

What is statistics  $\div$  It is science, of collecting, organising and analyzing the data.

Data  $\div$  Facts (or) pieces of information that can be measured.

Types of Statistics  $\div$

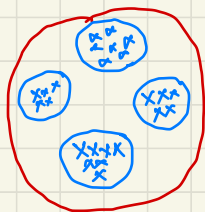
(i) Descriptive stats  $\div$  it consists of organizing and summarizing data.

(ii) Inferential stats  $\div$  it is a technique where we use the data that we have measured to form conclusions.

Population and sample  $\div$  Population denoted by  $N$   
Sample denoted by  $n$

example  $\div$  Elections  $\div$  Goa.

For exit polls the news reporters cannot ask each and everyone whom you have voted so what they do is take random people from random regions & whichever party gets more votes they conclude.



The red circle is population of goa & the blue circle are the samples

Sampling Techniques  $\div$

① Simple Random Technique  $\div$  just pick random data.

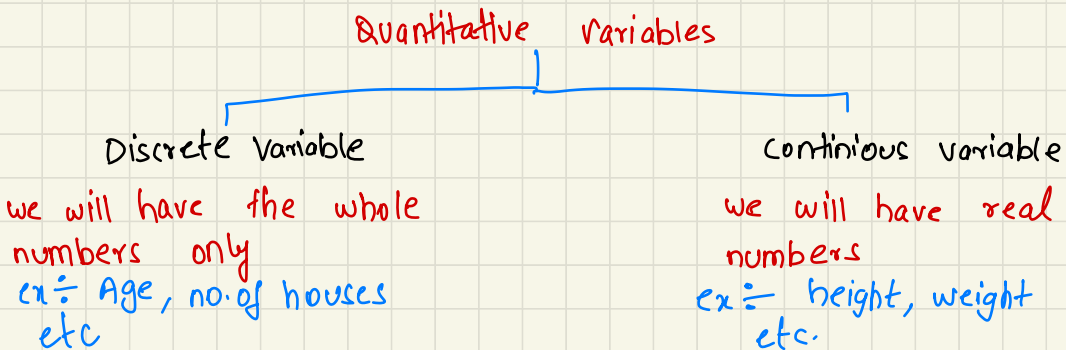
In simple random sampling every member of the population equal probability of getting selected

- ③ Stratified Sampling  $\div$  where the population ( $N$ ) is split into non-overlapping groups. (strata)
- ③ Systematic sampling  $\div$  From the population ( $N$ ) we pick every  $n$ th individual.
- ④ convenience sampling  $\div$  we choose only those people who are experts (or) interested in doing this.

**Variables  $\div$**  it is a property that can take values  
ex: height, weight.

**Two kinds of variable  $\div$**

- (1) Quantitative variable  $\div$  measured numerically  
ex: age, marks etc
- (2) Qualitative variable / categorical variable  $\div$  Here we will have categories for the variable & we cannot perform mathematical operations  
Ex: male, female.



**variable measurement Scales**  $\div$  4 types of measured variable

(i) **Nominal**  $\div$  The data is categorical we cannot perform any mathematical operations ex: **Colours, gender, Subjects.**

(ii) **ordinal**  $\div$  The order of the data matters but value doesnot. ex: **5 Student marks Rank**

**100**  
**96**  
**57**  
**85**  
**44**

**1**  
**2**  
**4**  
**3**  
**5**

} Here we are more worried about the Rank rather than the number of marks.

(iii) **interval**  $\div$  The value & order of the data matters, natural zero is not present.

ex: **temperature**

**70-80    80-90    90-100**

(iv) **Ratio data**  $\div$

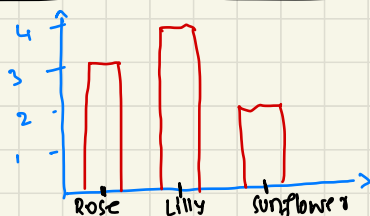
**\* Frequency distribution**  $\div$

Sample dataset  $\div$  **Rose, Lilly, Sunflower, Rose, Lilly, Sunflower, Rose, Lilly, Lilly.**

Frequency distribution table  $\div$

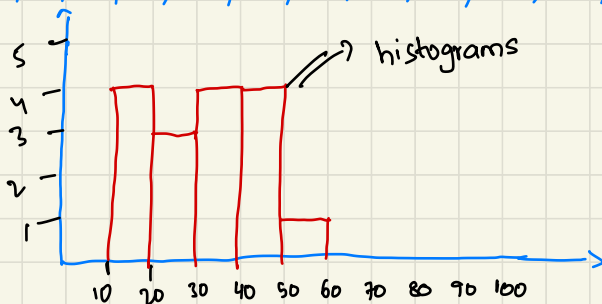
Flower	frequency	cumulative-freq
<b>Rose</b>	3	3
<b>Lilly</b>	4	7
<b>sunflower</b>	2	9

**Note:** \* if the variable is discrete then we use bar graph  
if the variable is continuous the we will use histogram



② Histograms  $\div$  we will divide in bins, the size is 10

Age:  $\{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51\}$



\* Measure of central tendency  $\div$  It is used to measure (or) determine the center of the data.

(i) Arithmetic mean for population and sample  $\div$

Population  $\div$  (N)  
 $x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$
$$= \frac{32}{10} = 3.2$$

Sample (n)

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$
$$= 3.2$$

\* Median  $\div$

$x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$

mean = 3.2  $\Rightarrow$  This was before we added 100

$$\text{mean} = \frac{32 + 100}{11} = 12$$

$\Rightarrow$  This is an outlier & it is also affecting the data adversely.

Now we can see that there is huge difference in mean from before adding 100 (3.2) to after adding 100 (12)

Median  $\div$   $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$

1) While doing median we have to sort the numbers

2) For odd numbers  $\div$  The middle element is called median  
 $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$

For the above data 3 is the median.

3) For even numbers  $\div$  The middle 2 elements Average is the median  $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112\}$

For the above data  $\frac{3+4}{2}$  is the median.

\* mode  $\div$   $\{1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 200\}$

most frequent element in the data is called the mode

$\therefore$  In the above data my mode is '6'

mode is also used to fill the missing values in the data, But works good for categorical variables.

\* Measure of Dispersion  $\div$  How well spread your data is.

(i) variance  $\div$

Population variance

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

Sample variance

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

example  $\div$

$x$	$\mu$	$x - \mu$	$(x - \mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	-0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
			<u>10.84</u>

$$\text{population variance} = \frac{10.84}{6} = 1.81 //$$

Standard deviation  $\div$  it is the square root of variance.

$$\sqrt{1.81} = 1.34$$

\* with the help of variance we can say how the data is spread.

\* with the help of S.D we can give range in which a particular value is present.

### Percentiles and Quartiles :-

$$\text{percentile Rank of } x = \frac{\text{\# of values below } x}{n} \times 100$$

ex:- { 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12 }

what is percentile Ranking of 10?

$$= \frac{\frac{16}{20} \times 100}{100} = 80\% \quad \left[ \text{Here we can say that 80 percentage of the values are less than 10} \right]$$

what is percentile Ranking of 11?

$$= \frac{17}{20} \times \frac{100}{100} = 85\% \quad \left[ \text{Hence we can say that 85\% of the values are less than 11} \right]$$

### Five number Summary :-

- 1) minimum
- 2) First Quartile ( $Q_1$ )
- 3) median
- 4) Third Quartile ( $Q_3$ )
- 5) maximum

## Removing the outliers :-

whenever we want to remove an outlier we need to have a lower fence & higher fence. that means.

- (i) All the numbers above the higher fence are outliers
- (ii) All the numbers below the lower fence are outliers.

$$\text{Lower fence} = Q_1 - 1.5 * (IQR)$$

$$\text{upper fence} = Q_3 + 1.5 * (IQR)$$

$$IQR (\text{inter quartile range}) = Q_3 - Q_1 \\ (75^{\text{th}}) - (25^{\text{th}})$$

ex:  $\{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 \}$   $n = 19$

Calculate the higher and the lower fence?

$$Q_1 = \frac{25}{100} \times (27) = 5 = 3$$

$$Q_3 = \frac{75}{100} \times (27) = 15 = 7$$

$$\text{Lower fence} = 3 - (1.5) (4)$$

$$= 3 - 6$$

$$= -3$$

$$\text{upper fence} = 7 + (1.5) * 4$$

$$= 7 + 6$$

$$= 13$$

## Data after removing outliers :-

$\{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9 \}$

$$\text{min} = 1$$

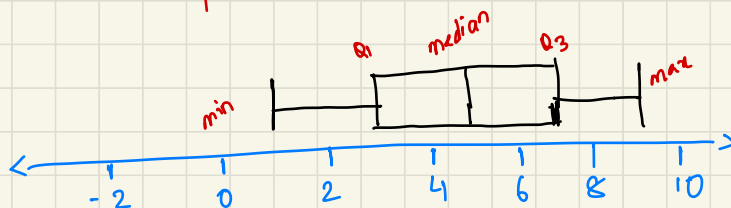
$$Q_1 = 3$$

$$\text{median} = 5$$

$$Q_3 = 7$$

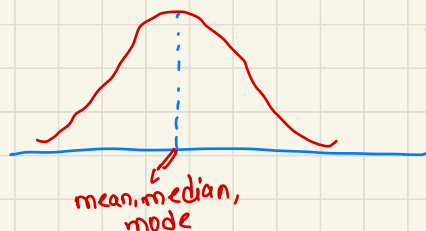
$$\text{max} = 9$$

## Box plot :-

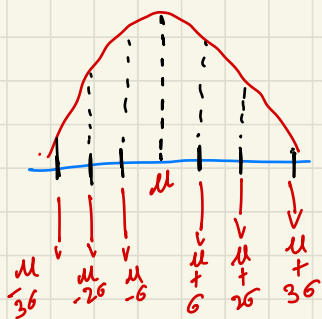


① Distribution  $\div$  it is used to get a brief idea about the dataset.

(i) Gaussian/normal distribution  $\div$



The 2 halves of the bell curve are symmetrical.



empirical formula  $\div$

68-95-99.7% Rule

Between the  $\mu \pm \sigma$  (1<sup>st</sup> SD) 68% of the data is present.  
 $\mu \pm 2\sigma$  (2<sup>nd</sup> SD) 95% of the data is present.  
 $\mu \pm 3\sigma$  (3<sup>rd</sup> SD) 99.5% of the data is present.

Z-Score  $\div$  it will help us in telling how much standard deviation it is away from the mean.

$$Z \text{ score} = \frac{x_i - \mu}{\sigma}$$

Note  $\div$  when you apply zscore for a particular distribution it will be converted to standard normal distribution.

In standard normal distribution mean = 1 & S.D = 0.



## practical Application $\div$

Dataset  $\div$

$\Rightarrow$ years	$\Rightarrow$ Rs	$\Rightarrow$ kg
Age	Salary	weight
24	40K	70
24	80K	80
26	60K	55
27	70K	45

For this dataset we can see that the features have different units hence to make them similar we do standardization, such that  $\text{mean} = 0$  &  $\text{std} = 1$

\* The process of converting your data in a way that  $\text{mean} = 0$  and  $\text{SD} = 1$  is called as standardization

Normalization  $\div$  Here you will have an option to convert all the values b/w zero and one (0, 1)

$$\text{min, max scalar} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

## practical Applications $\div$ { Ind vs SA }

(2021)  
① ODI Series  $\div$

Scores Avg = 250  
SD = 10  
R. pant score = 240

Q) In which year R. pant has better score?

(2020)  
② ODI Series  $\div$

Avg = 260  
SD = 12  
R. pant = 245

2021 (Z score)

$$= \frac{-10}{10} = -1$$

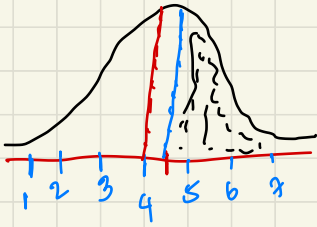
2020 (Z score)

$$= \frac{-15}{12} = -1.25$$

interview problem :-

Q) what percentage of scores fall above 4.25?

$\mu$



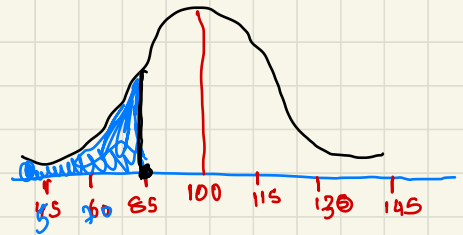
1) do the z score =  $\frac{4.25 - 4}{1} = 0.25$  [ look this at z-table ]

2) Since it is bell curve total consider as 1

So percentage above 4.25 is =  $1 - (0.25 \text{ value in z-table})$   
=  $1 - 0.5987$   
= 40.1%

Q) in india the average IQ is 100, with a standard deviation of 15. what percentage of the population would you expect to have an IQ lower than 85?

$$\begin{aligned} z_{\text{score}} &= \frac{x - \bar{x}}{s} ; \\ &= \frac{85 - 100}{15} \\ &= -1 \end{aligned}$$



$\therefore$  Lower IQ people are 24.20%