

Introduction part:

• check the quality of the data

→ Descriptive Statistics.

• To make a statement or a conclusion

→ Inferential Statistics

\* What is Statistics?

It is the science of collecting, organising and analysing data (is better decision making)

\* What is data?

Data means facts and pieces of information that can be measured.

eg: the IQ of a class students.

{ 98, 97, 68, 57, 110 }

avg IQ min, max

\* Descriptive Statistics

It consists of organising and summarising data.

Inferential Statistics

Techniques where we used the data that we have measured to form conclusion

&

To make a statement / conclusion on a descriptive statistics we use inferential statistics

eg question

Are the avg marks of Java class students is same as python class students or not?

→ Inferential Statistics

What is the avg marks of students

→ Descriptive Statistics

population (N) and Sample (n)

The entire group of the data we call it as population.

eg All people in India

A subset of a population we call it as a sample

eg: 1 lakh people from different region of India.

Key points:

- populations are larger than samples
- Samples should be representative of the population
- Sample allow for easier faster and less costly collection.

Types of Sampling techniques

1. Simple random sampling

every member of a population has an equal

chance of being selected for our sample

eg: Avg mileage of a bike.

Avg ratio of married people in bnglr

## 2. Stratify Sampling

where the population is split into non overlapping group

eg: people is male or female

## 3. Systematic Sampling

From the population every  $n$ th sample we are going to select.

eg: ~~the person~~ while doing survey in the mall on the topic

of modernisation, collecting information of every fifth person who is coming out from mall

## Convenience Sampling

The sample is collected based on our convenience from the particular domain experts.

ie: Sampling technique selection always depends on problem statement.

Variable: It is a property that can take any value.

## Two kinds of Variable

1) Quantitative (Numerical) Variables.

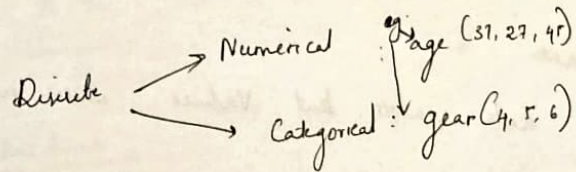
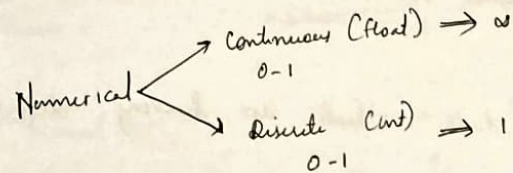
2) Qualitative (Categorical) Variable

### ① Quantitative Variable

A value can be measured and we can perform mathematical operation like  $(A, S, M, D)$

eg: mpg, weight, height.

gear  $\rightarrow$  discrete Categorical data.





## 2) Qualitative / Categorical

Non-measurable data and based on some characteristic we can derive categorical variable.

eg: gender

- male
- female
- other

working type

- IT
- non-IT

Blood group

- A<sup>+</sup>
- B<sup>+</sup>
- O<sup>+</sup>
- AB<sup>+</sup>

## Variable measurement scales

types of measured variables

### 1) Nominal data.

The categorical data which are having different classes.

### Ordinal data.

Order of the data matters but values does not.

80	46	32	98	57
↓	↓	↓	↓	↓
2nd	4th	5th	1st	3rd
68	92	88	49	70
↓	↓	↓	↓	↓
4th	1st	2nd	5th	3rd

### 3) Interval data

Order matters and value also matters but natural zero is not present.

→ (that value never reaches zero)

### 4) Ratio data

The ratio data can be measured. order equal and have meaningful zero.

## Descriptive Statistics

1. measure of central tendency

mean                      median                      mode

population mean ( $\mu$ )

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

sample mean ( $\bar{x}$ )

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

### Median

- Sort the value either asc or desc
- Choose the mid value
- if you get mid 2 value take avg of those

[1, 2, 2, 3, 4, 5]

[1, 2, 2, 3, 4, 5, 100]

mean: 2.8

mean: 16.7

median: 2.5

median: 3

mean will be affected by outlier where as

median won't affect by outlier

for null value imputation (mean and median is used)

Mode

most repetitive Value

Measure of Dispersion

eg:

	Person 1	person 2
Mon	7:30 am	8 am
Tues	7:45 am	11 am
wed	8 am	9 am
Thurs	7:15 am	7 am
	7 am	10 am
	?	1
	7-8	9-10

Variance high  
prediction low

Variance

Standard deviation

range

population Variance ( $\sigma^2$ )

Sample Variance ( $S^2$ )

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

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

$n-1$  = degree of freedom

Calculate the Variance

{1, 2, 2, 3, 4, 5}

8.02

$$\sigma^2 = \frac{(1-2.8)^2 + (2-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{6}$$

$$\frac{3.36 + 0.69 + 0.69 + 0.02 + 1.36 + 4.69}{6} = \frac{10.83}{6} = 1.805$$

$$\sqrt{\sigma^2} = \sqrt{1.805} = 1.34$$

$\mu = 2.8$

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

population = SD

$$\sigma = \sqrt{\sigma^2}$$

$$\sigma = \sqrt{1.8}$$

$$\sigma = 1.34$$

$$\sigma = 1.34 \text{ km} \downarrow \text{8km}$$

Sample

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

Sample = SD

$$S = \sqrt{s^2}$$

$$S = \sqrt{2.16}$$

$$S = 1.46$$

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

$$= 5 - 1$$

$$= 4$$

Percentile and Quartile.

Percentile is a value below which a certain percentage of observations will come under.

{ 2, 3, 4, 5, 5, 6, 7, 7, 8 }

How much % of data will come below value 6?

$$\text{ile rank of } x = \frac{\text{No of value below } x}{N} \times 100$$

$$= \frac{7}{11} \times 100$$

$$= 63\% \quad \text{Observation data}$$

$$is < 6.$$



• Quantile helps to find the value which is present at the given percentile rank.

$\Sigma 1, 1, 2, 3, 4, 5, 5, 6, 7, 7, 8$ .

which Value is percent at 25%?

$$\text{Value} = \frac{\text{percentile}}{100} \times \frac{n+1}{100}$$

$$\frac{25}{100} \times \frac{12}{100}$$

$$= 3 \rightarrow \text{index}$$

$$= \text{Value} = 2$$

$$\begin{array}{l} \text{if it is } 90\% \\ \frac{90}{100} \times 12 \end{array}$$

$$10.8 = \text{index}$$

$$10 \rightarrow \text{index}$$

$$7 = \text{Value.}$$