

# Statistics

Defn : Statistics is the science of collecting, organizing and analyzing data.



Decision Making

Data : "facts or pieces of information"  $\Rightarrow$  Measured, Collected, Analyzed

Eg: Weights of students in the class

{ 60kg, 50kg, 45kg, 30kg, ... }.

IQ of the students of incly

{ 100, 90, 95, 99, ... }.

## House Price Dataset

City	Area	No. of Rooms	Price	Analysing this Data
Bangalore	1000	2	45 Lakhs	
New York	1250	2.5	50 Lakhs	Data Scientist $\rightarrow$ Model
Mumbai	-	-	-	$\downarrow \uparrow$ Price

Data Analyst  $\rightarrow$  Report  $\rightarrow$  Visualization  $\rightarrow$  Meaning Decisions.

$\hookrightarrow$  Project  $\rightarrow$

## Application

- ① Data Exploration And Summarize
- ② Model Building And Validation
- ③ Statistical Analysis  $\rightarrow$  Sample Data  $\rightarrow$  Population Data.
- ④ Hypothesis Testing
- ⑤ Optimization And Efficiency.
- ⑥ Reporting

# Types Of Statistics

## (1) Descriptive Statistics .

Descriptive statistics involves methods for summarizing and organizing data to make it understandable. This type of statistics helps to describe the basic features of the data in a study.

### 1) Measure Of Central Tendency

[Mean, Median, Mode]

### 2) Measure of Dispersion [Variance, Standard deviation]

### 3) Data Distribution

i) Histograms

ii) Box plot

iii) Pie Chart

iv) PDF , PMF

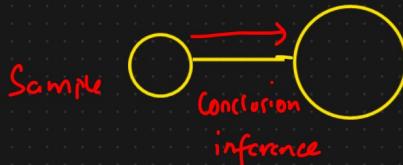
### 4) Summary Statistics

#### i) Five Number Summary

Q1, Q2, Q3, Maximum

## (2) Inferential Statistics

Inferential statistics involves methods for making predictions or inferences about a population based on a sample of data. It allows for hypothesis testing, estimation, and drawing conclusions.



### ① Hypothesis Testing

#### ② P value

#### ③ Confidence Interval

#### ④ Statistical Analysis Test

##### ① Z test

##### ② t test

##### ③ ANOVA → F test

##### ④ Chi Square.

## Summary

Type of Statistics	Key Concepts	Examples
Descriptive Statistics	Measures of Central Tendency (Mean, Median, Mode), Measures of Dispersion (Range, Variance, Standard Deviation), Data Distribution (Histograms, Box Plots), Summary Statistics (Five-number Summary)	Mean score of students, Range of temperatures, Histogram of ages
Inferential Statistics	Hypothesis Testing (Null and Alternative Hypotheses, P-value), Confidence Intervals, Regression Analysis (Simple and Multiple Linear Regression), ANOVA, Chi-Square Test	P-value in test scores comparison, 95% confidence interval for average height, Predicting house prices, Comparing test scores of different schools, Association between gender and product preference

Eg: Let say there are 20 statistics class in your College, and you have collected the height of students in the class.

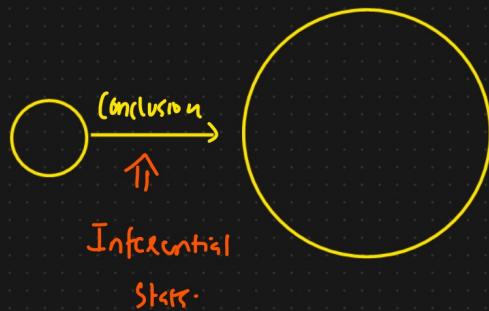
Heights are recorded [175cm, 180cm, 140cm, 135, 160cm, 120cm]

## Descriptive Question

"What is the average height of the entire classroom" ← Measure of Central Tendency

## Inferential Question

"Are the height of the sample students in classroom similar to what you expect in the entire College"?



## ③ Population And Sample Data

### Population

**Definition:** A population is the entire set of individuals or objects of interest in a particular study. It includes all members of a defined group that we are studying or collecting information on.

### Sample Data

**Definition:** A sample is a subset of the population that is used to represent the entire group. Sampling involves selecting a group of individuals or observations from the population to draw conclusions about the whole population.

### Characteristics

1) Complete Set : Contains all the observation of interest

2) Parameter : A numerical value summarizing the entire population

i) Population mean ( $\mu$ )

2) Population Variance ( $\sigma^2$ )

### Example

1) Population in a School Study

### Characteristics

1) Subset : Represents a portion of the population.

2) Statistic : A numerical value summarizing the sample data [Sample mean, Sample Variance].

3) Random Sampling : Samples should be randomly selected to avoid bias

### Example

1) Sample In a School Study

- \* All Students enrolled in a School
- \* Determine the avg height of Student, **Population mean**
- i) A group of 50 Students from School
- ii) Use case: Estimate the average height of students in a School

## 2) Population in Market Research

- \* All consumers in a city.
- \* To understand the purchasing behaviour of all consumers.

## 2) Sample in Market Research

- \* 500 consumers from the city
- ↓
- \* Behaviour → Population.

## 3) Population in a Medical Study

## 3) Sample in a Medical Study

- \* All the patients with a specific disease
- \* To study the effectiveness of a drug.

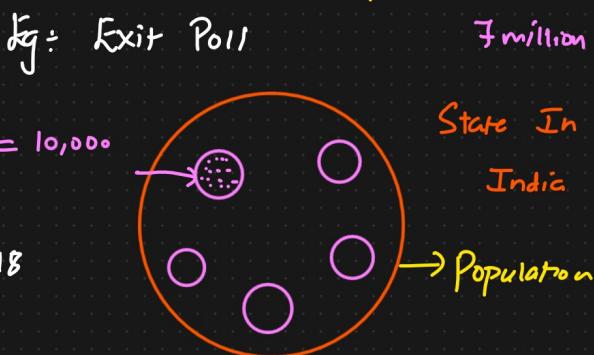
- \* 200 Patients
- \* Test the effectiveness of the drug

## ④ Types of Sampling Techniques

1. Probability Sampling
2. Non Probability Sampling

Sample = 10,000

Age  $\geq 18$



### 1. Probability Sampling:

#### a) Simple Random Sampling :

Every member of the population has an equal chance of being selected.

Eg: Selecting people randomly

Draw names random from a class of students.

#### b) Systematic Sampling

Select every  $n^{\text{th}}$  member of the population after a random starting point.

Eg: Airport → Credit Card → 5<sup>th</sup> person, 10<sup>th</sup> person, 15<sup>th</sup> person

Feedback Survey  $\rightarrow$  Selected every 11<sup>th</sup> member  $\rightarrow$  Feedback Survey.

### c) Stratified Sampling:

Divide the population into Strata (groups) based on specific characteristics and then randomly sampling from each Strata.

Eg: Divide employees by department and then randomly select a proportional number from each department to form a survey sample.

$$\text{Eg: Age} \rightarrow \frac{<12}{\uparrow\uparrow} \quad \frac{12-18}{\uparrow\uparrow} \quad \frac{>18}{\uparrow\uparrow} \quad \{ \text{Politics} \}$$

### d) Cluster Sampling

Divide the population into clusters, randomly selecting clusters, then sampling all the members from the selected clusters

Eg: Randomly selecting several schools from a district, and surveying all teachers within those schools.

### e) Multi Stage Sampling

Combining several sampling methods. Usually involves selecting clusters, then randomly sampling within those clusters

Eg: Randomly selecting cities, each selected city randomly selecting households to survey.

## ② Non Probability Sampling

Select individuals who are easiest to reach.

Eg: Surveying people at mall.

### ③ Convenience Sampling

Selecting individual who are easiest to reach

### ④ Judgmental (Purposive) Sampling

Select individual based on the researcher's judgement → Useful or Representative

Eg: Choose experts in a field to participate {Data Science}

### ⑤ Snowball Sampling

Existing Study Subjects recruit future Subjects from among their Acquaintances.

Eg: Survey members of a rare disease.

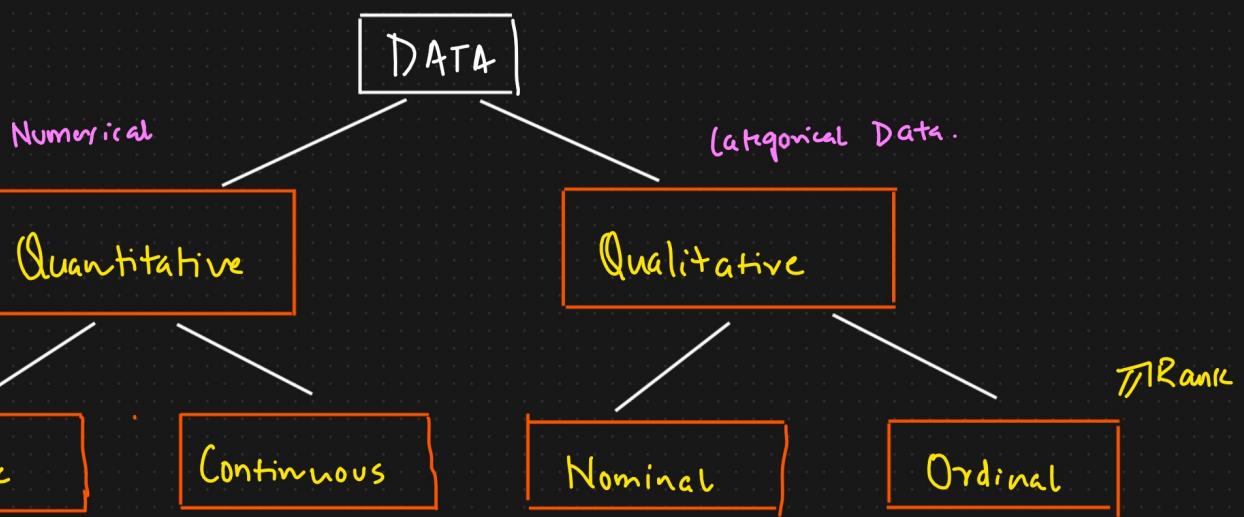
### ⑥ Quota Sampling

Age, group, gender, caste.

Selecting the Sampling technique depends on Usability

### ⑤ Types of Data

DMC	DC	ISI	BUI	FWI	Classes	Region
3.4	7.6	1.3	3.4	0.5 not fire	0	
4.1	7.6	1	3.9	0.4 not fire	0	
2.5	7.1	0.3	2.7	0.1 not fire	0	
1.3	6.9	0	1.7	0 not fire	0	
3	14.2	1.2	3.9	0.5 not fire	0	
5.8	22.2	3.1	7	2.5 fire	0	
9.9	30.5	6.4	10.9	7.2 fire	0	
12.1	38.3	5.6	13.5	7.1 fire	0	



Whole numbers	Any value	Eg: gender	Eg: Customer Feedback
Eg: No. of bank accounts	Eg: Weights	M, F	Good, Bad, Better
No. of children in a family	Height	Blood Group	1 $\frac{1}{3}$ 2
	Temperature	Pinode	
	Speed		

## ⑥ Scales of Measurement of Data

The scales of measurement describe the nature of information within the values assigned to variables.

### 4 Primary Scales of Measurement

- ① Nominal Scale
- ② Ordinal Scale
- ③ Interval
- ④ Ratio

### ① Nominal Scale

Dcfn: This scale classifies data into distinct categories that do not have an intrinsic order

Qualitative/Categorical data.

#### Characteristics

- i) Data is categorized based on labels, names or qualities
- ii) Categories are mutually exclusive
- iii) No logical order among categories [No Rank]

#### Example

Gender	Colors
→ M	Red → 5 50%
→ F	Blue → 4 40%
	Pink → 1 10%

#### Types of Cuisines

{ Italian }  
Chinense  
Mexican

## (2) Ordinal Scale

Defn: This scale classifies the data into categories that can be ranked or ordered.

### Characteristics

- i) Data is categorized and ranked in a specific order.
- ii) The intervals between ranks are not necessarily equal.

### Example

<u>Education level</u>	$\Rightarrow$ F.E.	<u>Ranks</u>	<u>Customer Feedback</u>	$\Rightarrow$
High School		1	Satisfied	1
Bachelor		2	Very Satisfied	2
Masters		3		
Doctorate		4	Not Satisfied	0

### Socio Economic Status

Low	2
Middle	1
High	0

## (3) Interval Scale

Defn: The interval scale not only categorizes and orders but also specifies the exact difference between intervals. It lacks a true zero point.

### Characteristics

- 1) Data is ordered with consistent interval between values.
- 2) Allows for meaningful comparison of differences [Ratio cannot be measured].
- 3) No true zero point.

## Example

Temperature in Fahrenheit:

$\rightarrow 10^{\circ}\text{F}, 20^{\circ}\text{F}, 30^{\circ}\text{F}$

$$20 - 10 = 10$$

$\Rightarrow 0^{\circ}\text{F}$  "No Temperature"

IQ Scores

$90, 100, 110$

$$\hookrightarrow \text{Difference } 100 - 90 = 10$$

$$30 - 20 = 10$$

$$\frac{30}{10} = \boxed{3\text{ : }1} \leftarrow$$

Calendar Years

$2024, 2020, 2016 \Rightarrow \underline{\underline{0 \text{ years}}}$

## 4) Ratio Scale

Eg: Student marks in a class

Ⓐ The order matters

$0, 90, 60, 30, 75, 45$

Ⓑ Differences Are measurable [Ratio]  $\text{ASC} = 0, 30, 45, 60, 75, 90 \leftarrow$   
(can be meaningful)

⓫ Contains a 0 starting point.

$$45 - 30 = 15$$

$$60 - 30 = 30$$

$$\text{Ratio} = \frac{90}{30} = \frac{3}{1}$$

$$\boxed{3:1}$$

Assignment

Eg: Weights

① length of Different Rivers In the World?

$\rightarrow 10, 20, 30, 40$

② Favorite food based on Gender?

Income:  $10,000, \leftarrow$

③ Marital Status?

$20,000, \leftarrow$   
 $30,000, \leftarrow$   
 $40,000$

④ IQ Measurement?