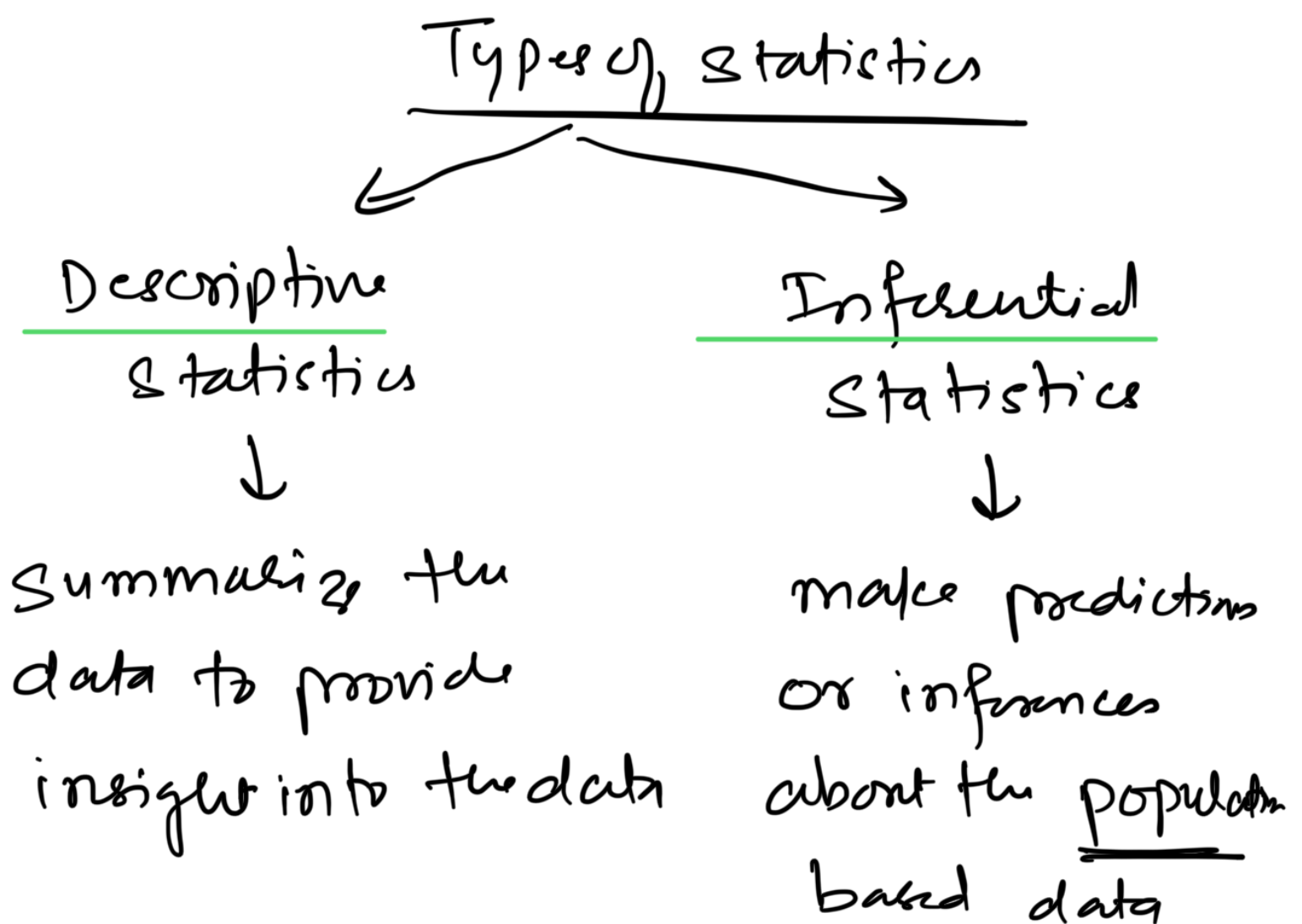# Statistics used in Data Science

Definition → Branch of mathematics for collection, analyzing, Interpretation of data & presentation of data.

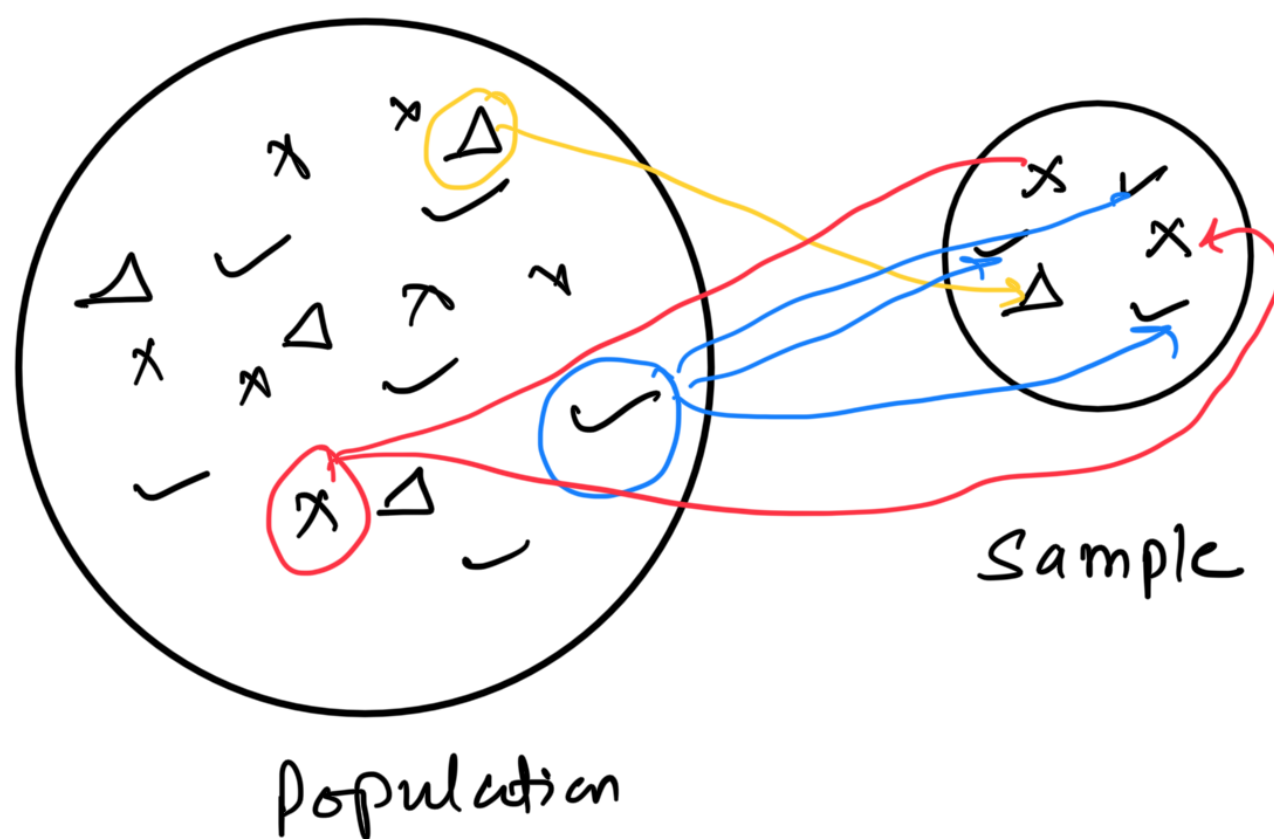## Importance -

Essential for extracting insights, making predictions and uncovering patterns in data science.

## Types of Statistics

| Descriptive Statistics | Inferential Statistics |
|---|---|
| ↓ | ↓ |
| Summarize the data to provide insight into the data | make predictions or inferences about the <u>population</u> based data |

Population (N) → is the entire data that you want to draw conclusions

Sample (n) → is the group

(part of population) from which you will collect data



Population

Sample

Sampling → Types of Sampling

## 1. Simple Random Sampling -

In this process of sampling where every member of the population has equal chance of being selected

## 2. Stratified Sampling -

In this method of sampling where population (N) is split into non-overlapping group.

3. **Systematic Sampling -**

In this method of Sampling a probability sampling method where researchers select members from population at $n^{th}$ interval

4. **Convenience Sampling -**

In this method of Sampling a process of taking Sample data from those who has knowledge / expertise on the domain area.
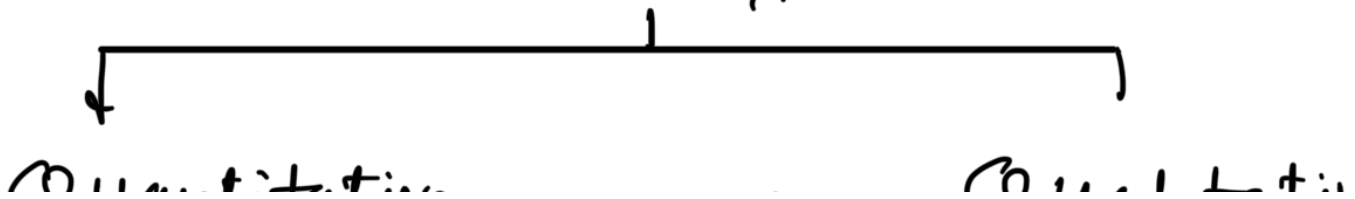
## Data Used in Data Science

**Data:** Data is a collection of facts, such as numbers, words, measurements, observation or even descriptions of things.

- In the context of data science, data is used to make decision, prediction and inference

### Data Types

| Quantitative (number) | Qualitative (categorical) (textual data) |
|---|---|
| — Take on numeric values | — Take names or labels |
| _Eg_ No. of students in a class | _Eg_ Eye color |
| — Square feet in a house | — Gender |
| — Population of a city | — Breed of dog |
| — Height of Students | — Level of Educat" |
| | — Marital status |
| (Measure it numerically) | (Based on some Categorical value) |
| — Add, Multiply, Sub, div | |

Data
```
         Data
        /    \
Quantitative   Qualitative
   /    \
Discrete  Continuous
  ↓          ↓
whole numbers   infinite numbers
```

<u>Discrete Data</u> — Countable and can be taken on some specific value.

eg. No of students in class

No of cars park in parkings

## Continuous Data - Can take any value within a range and its often measured rather than count.

eg. Height of students

Temperature Readings

## Qualitative Data → Types of data are

1. Nominal ⟶ Categorical data

2. Ordinal ⟶ Order of data matters

3. Interval ⟶ Order matters, {o, max)

4. Ratio ⟶ Something measured on ratio scale.

eg. Students Marks and Rank

| | St Marks | Ranks | |
|---|---|---|---|
| Pass | 100 | (first) 1 | Ordinal data |
| Pass | 96 | (second) 2 | |

Pass     57     4

Pass     85     3

Nominal

    Pass     44     5

Fail     ← 30

eg. <u>Temperature</u>

<u>Fahrenheit</u>     70-80    } Interval
               80-90
               90-100

eg. <u>Ratio</u> -

      <u>Percentile</u>

---

## <u>Data Classifications</u>

Structured            Unstructured
Data                 Data

<u>Structured Data:</u>

- is organized in a fixed format, often in rows and columns, making it easy for searching & analyzing.

- Eg. Database, Relational Databases

Excel sheets (Spreadsheets)

eg. data = {

'Name': ['Amit', 'Sai', 'Aditya'],

'salary': ['20k', '36k', '50k']



Rows & col
Structured
data

## Unstructured Data

- data lacks a predefined format, and is not easily searchable.

- It includes huge variety of data formats as text, image, videos, audios, and social meadia post.

eg - Text documents

- Audio files
- Blog post

eg. text data = [" I am Kiran."

" Data science is [ scinating".

" Python is most popular

language."]

## Semi-structured data

- It is not organized in a rigid structure like structured data, but it contains tags or markers to separate elements,

- XML language
  JSON files

## Summary

| Data Type | Description | Example |
|---|---|---|
| Discrete | Countable specific values | No of students |
| Continuous | Any value within range | Temperature |

| | | |
|---|---|---|
| Nominal | Categories | Eye color |
| Ordinal | Ordered categories | Age =f teenage, middleage, old age ¿ |
| Interval | Range, | [0, 100] |
| Ratio | selceted values ont of total value | Procentie |
| Structured data | Organiz in row & Cols | Excel SQL db |
| Unstructured data | No predefined format | Text document |
| Semi-Structured data | Tags or markers | JSON files XML |

| Files