# Data Science

- Data science : multidisciplinary → uses - scientific methods, Algorithms, statistics, Visualisation, extract info from structured, Unstructured data.

- Big data : collection of data i.e., volume, Variety, Veracity, Velocity, Value (5Vs)
  ↓
  Uncertain & Inconsistency.

→ AI > ML > NN > DL .

→ <u>Data types and scales</u> :-

Data types :- Structured, Unstructured, semi structured, Quasi structured
  (Textual data + inconsistent)

  :- Cross-sectional data  — Time series data  — Panel data
       ↓                        ↓                    ↓
    Many variables          Single variable      Several variables
    collected same time     collected several    collected several
    or durration            times (week, month)  times. (Ex: Unemployee
                                                  rate of countries)

→ Object ; Attributes :-   Qualitative

Attributes :- Nominal  — Binary — ordinal — Continuous and discontinuous
  — Numeric : (interscaled, Ratio scales). ——→ Quantitative.

Nominal :- Name of things - Name of symbols: Ex : Lecturer, White, Professor.

Ordinal :- Values with meaningful sequence (Ranking) Ex: Medium, High, Low.

Binary :- ⎡ Symmetric : Both values are equal Ex: Gender
(yes/No)  ⎣ Asymmetric: Both values not equal Ex: Result.

Numeric :- Real/integer value, measure quantity.
  ⎣ ⎡ Interval scaled :- Difference of adjacent values consistent : Ex: Temparature,
        - No predefined starting point (True zero value) Dates.
    ⎣→ Ratio - scaled :
        - Difference of adjacent values    Ex: Height, weight, Length
        - Predefined. True zero value    No zero in difference of Height.

Continuous :- Infinite values, measured, subdivisible, float Ex: Height, weight, time

Discontinuous : Specific values, counted, non subdivisible Ex: No of students, No of cars.

§→ Population :- set of possible observations
Sample :- Logical subset of the population.

§ Measures of centraltendency :-
  1. Mean :
     $\mu = \dfrac{\Sigma x}{N} = \dfrac{x_1 + x_2 \cdots}{N}$    2, Weighted mean   3, Median   4, Mode :
                                                                    $= \dfrac{n+1}{2}$   Most often occures.

# Types of Data Analytics :-

## 1. Descriptive

- summarize historical data to identify patterns
- Used Aggregate functions on database.

## 2. Predictive

- Predicts future
- Predict probability of future occurrences
- Ex: Regression classification

## 3. Prescriptive

- Choose optimal actions to perform on insights from Descriptive and predictive analysis

Ex: Linear programming meta-heuristics Alg.

---

## Steps in Data Science :

**Step 1: Setting Research goal :-** ⎡ Define research goal ⇒ Well defined, Deliverable
⎣ Create project charter ⇒ objectives, Resources, Timeline.

**Step 2 Retrieving Data :** ⎡ Internal Data Ex: Servers ⎤ → ⎡ Data retrieve
⎣ External data Ex: Facebook          ⎣ Data owner - ship

**Step 3: Data Preparation: (Preprocessing).**

**Step 4: Data Exploration:** — ⎡ Simple graphs
⎢ Combined graphs
⎢ Link and brush
⎣ Non-graphic Techniques

1. Histogram
2. scatter plot
3. Box plot
4. pie charts
5. Bar.

**Step 5 Build the Model:** — ⎡ Model selection
⎢ Model execution
⎣ Diagnosis and model comparison.

$$MSE = \frac{1}{n} \Sigma (y - \hat{y})^2$$

**Step 6 Presentation and Automation :-** ⎡ Presenting Data
⎣ Automating Data Analysis