# Data Preprocessing

★ Process of transforming raw data into an understandable format.

★

```
                    ┌──────────┐
                    │   Data   │
                    │ Cleaning │
                    └────┬─────┘
                         │
┌──────────────┐  ┌──────┴────────┐  ┌───────────┐
│     Data     │  │     Data      │  │   Data    │
│Transformation│  │ Pre-processing│  │ Reduction │
└──────────────┘  └──────┬────────┘  └───────────┘
                         │
                  ┌──────┴────────┐
                  │     Data      │
                  │ Discretization│
                  └───────────────┘
```

1. Data Cleaning

① Handling Missing Values

ⓐ Replace missing values manually

ⓑ Replace missing values by zero

ⓒ Dropping or ignoring rows with missing values.

ⓓ Use Global constants to fill missing values.

⑪ Noisy Data

ⓐ Duplicate entry

ⓑ Multiple entries for a Single Entity

ⓒ Nulls

ⓓ Huge Outliers.

## 2] Data Transformation

=> Process of converting data from one for
to another.

Techniques:-

① Rescaling Data [0-1 or 0-100].
② Normalizing Data. [0.0-1.0 or -1.0 to
③ Binarizing Data [0 or 1]
④ Standardizing Data [mean=0 & S.D=1
⑤ Label Encoding [text → numeric]
⑥ One Hot Encoding [spliting coloums in
many columns]

|        | Age |
|--------|-----|
| India  | 34  |
| Japan  | 15  |

|   | Age |
|---|-----|
| 0 | 34  |
| 1 | 15  |

Label Encoding

| 0 | 1 | Age |
|---|---|-----|
| 1 | 0 | 34  |
| 0 | 1 | 15  |

One hot Encoding

3] Data Discretization

⇒ method of attribute of continuous data into finite set of intervals

| Age | 10 | 20 | 21 | 50 | 51 | 60 | 70 |

⇓

Data Discretization

| Age | Age | Age |
|---|---|---|
| 10,20,21 | 50,51 | 60,70 |
| Young | Mature | old |

4] Data Reduction

⇒ reducing a large capacity of the data into small datasets.

Techniques
① Data Cube Aggregation.
⑪ Numerosity reduction
⑪⑪ Data compression

* Types of Attributes
  (i) Nominal [names or symbols]
  (ii) Binary [0 or 1]
  (iii) Ordinal [ranking among them]
  (iv) Numeric [measurable quantity, represented in integer or real values]