

# DATA SCIENCE

0287

Data: Data is information such as facts and numbers used to analyse something or make decisions.

Data Science: Data science is the field of study, that combines domain expertise, programming skills and knowledge of mathematics and statistics to extract meaningful insights from data.

## Types of Data:

- ① Qualitative
  - Nominal Data
  - Ordinal Data
- ② Quantitative
  - Discrete data
  - Continuous data

Nominal data: Nominal data is used just for labelling variables, without any type of quantitative value.

Ordinal data: Ordinal data is data which is placed into some kind of order by their position on a scale.

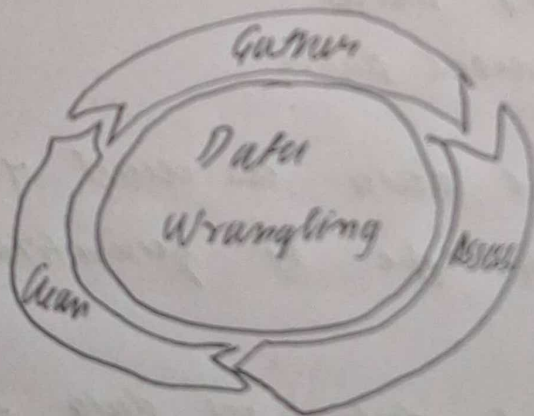
Data source : the location where data is being used to originate from.

Types of Data sources:

- { Relational
- { Multi dimensional
- { Dimensional modelled relational

### Data Wrangling

Data wrangling is the process of removing errors and combining complex data sets to make them more accessible and easier to analyse.



Main concept of Data source is collection of data, design and analysis.



## Big Data

Big data refers to significant volumes of data that cannot be processed effectively with the traditional applications that are used.

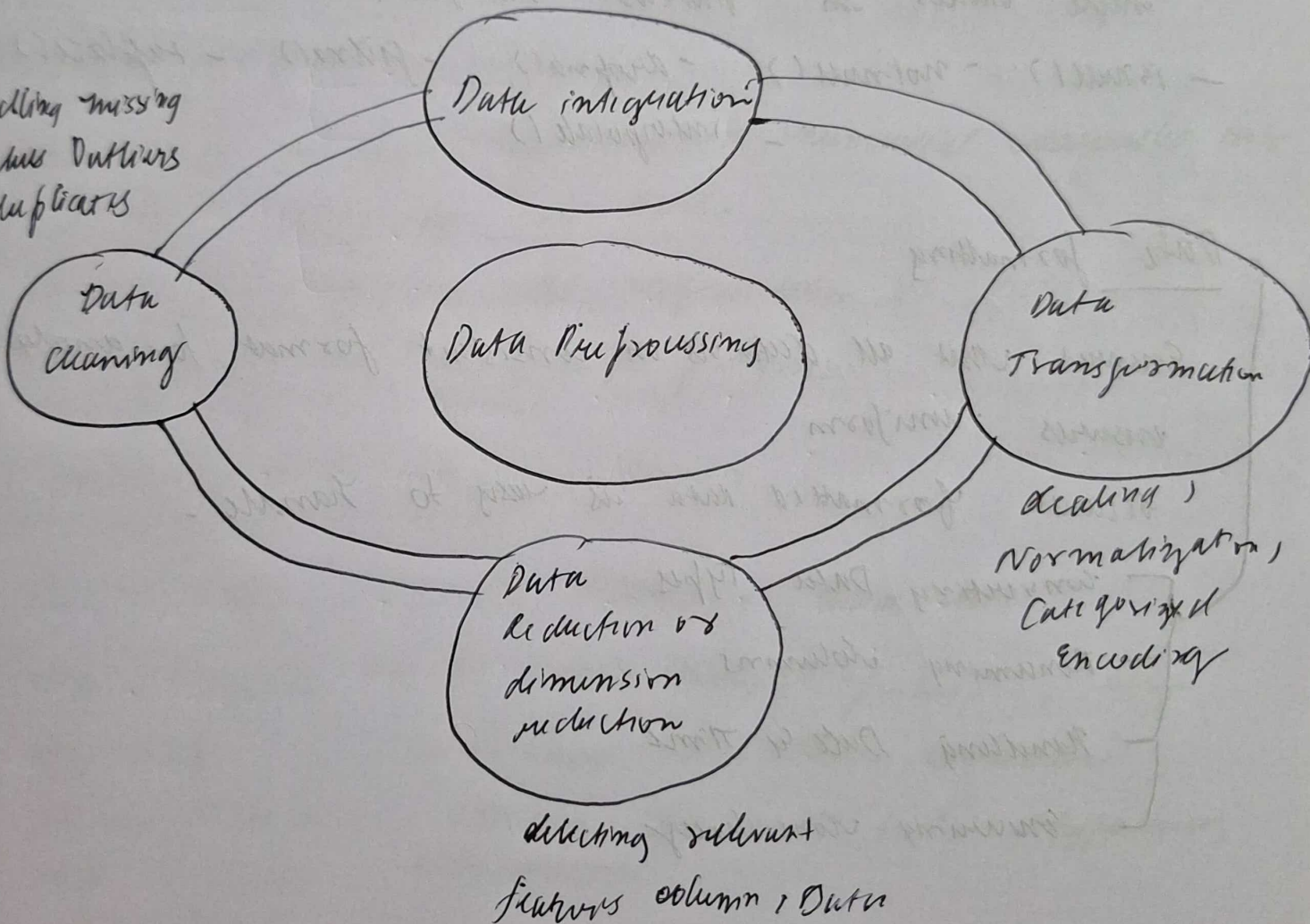
Velocity → speed generating data

Volume → size - KB - PB

Variety →  
Structured  
unstructured  
semi-structured

Complex & large datasets.

Handling missing  
values Outliers  
duplicates



## DATA CLEANING

involves the systematic identification and correction of errors, inconsistencies and inaccuracies within a dataset.

- handling missing values
- removing duplicates
- addressing outliers.

essential because raw data is messy.

### Handling missing values

several useful functions for detecting, removing or replacing null values in Pandas Dataframe.

- isnull()
- notnull()
- dropna()
- fillna()
- replace()
- interpolate()

### Data formatting

Ensures that all data is in consistent format for analysis.  
ensures uniform

because formatted data is easy to handle.

- converting Data types
- Renaming columns
- Handling Date & Time
- Ensuring consistency



## Data Normalization

process used to scale numerical features in a dataset to a standard range.

Ensures no single feature to dominate due to weight.

Why?

Improves Model Performance

Speeds up Training

Prevents Numerical Instability.

Methods

Min-Max Normalization

⇒ Z-score Normalization (Standardization)

## Binning

process of converting continuous numerical variables into discrete intervals or bins.

Useful in simplifying the representation of data.

Also helps in identifying patterns.

## Exploratory Data Analysis (EDA)

Is an approach to analyzing data sets.

crucial step in the data analysis process, help analysts to understand the data's structure, detect patterns, identify anomalies, test hypotheses and check assumptions.

helps making informed decisions.

## Components of EDA

- Descriptive Statistics - summary statistics
- Data Visualization
- Data Cleaning
- Data Transformation
- Correlation Analysis

## Correlation

Is a statistical technique that measures the strength and direction of the linear relationship between two quantitative variables

- Pearson Correlation
- Spearman Correlation
- Kendall Correlation