

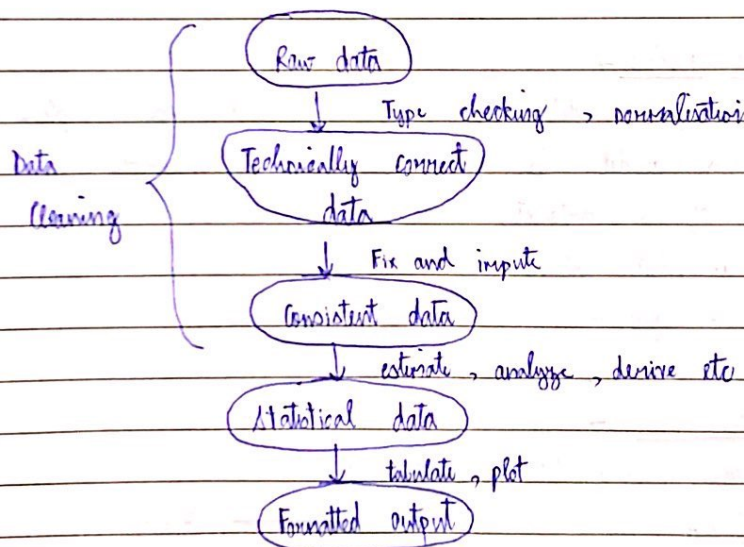
Assignment - 06

* Aim :- Perform different data cleaning operations using R/Python.

* Problem Statement → :- Perform the following operations using R/Python on the air quality and heart disease datasets.

1. Data cleaning
2. Data transformation - integration
3. Data integrative transformation
4. ETL connecting
5. Data model building

* Theory :-



• Data Cleaning →

Data cleaning is the process of preparing data for analysis by removing or identifying modifying data that is incorrect, incomplete, irrelevant, duplicate or improperly formatted.

Eg - Removing NAs.

```

airquality = read.csv("AirQuality.csv")
airquality $ Ozone = if else (is.na(airquality $ Ozone),
                             median(airquality $ Ozone, na.rm = TRUE), airquality)
  
```

33304

- Data Transformation →
It is the process of converting data from one format or structure.

Eg - $\text{airquality} \& \text{danger} = \text{airquality} \& \text{Temp} > 100$

- Data Integration →
Data integration is a method of combining data from different sources into a single, unified view. It begins with the ingestion process and includes steps such as cleansing ETL mapping and transformation.

// dataset 1 and dataset 2 is already available

$\text{dataset} = \text{Merge}(\text{dataset 1}, \text{dataset 2}, \text{by} = "FN")$

- Build data models →
Using simple models to better understand complex dataset. The goal of a model is to provide a simple low dimensional summary of a dataset.

Models:

1. list column
2. gapminder

- Error correction →
Eg - $\text{mean}(\text{airquality} \& \text{Temp})$

[1] NA

$\text{mean}(\text{airquality} \& \text{Temp}, \text{na.rm} = \text{TRUE})$

↑
Remove the error



33304

myCOMPANION

* Conclusion :- Thus, in this assignment we have performed data cleaning operation on air quality dataset using R