

Assignment No. 6

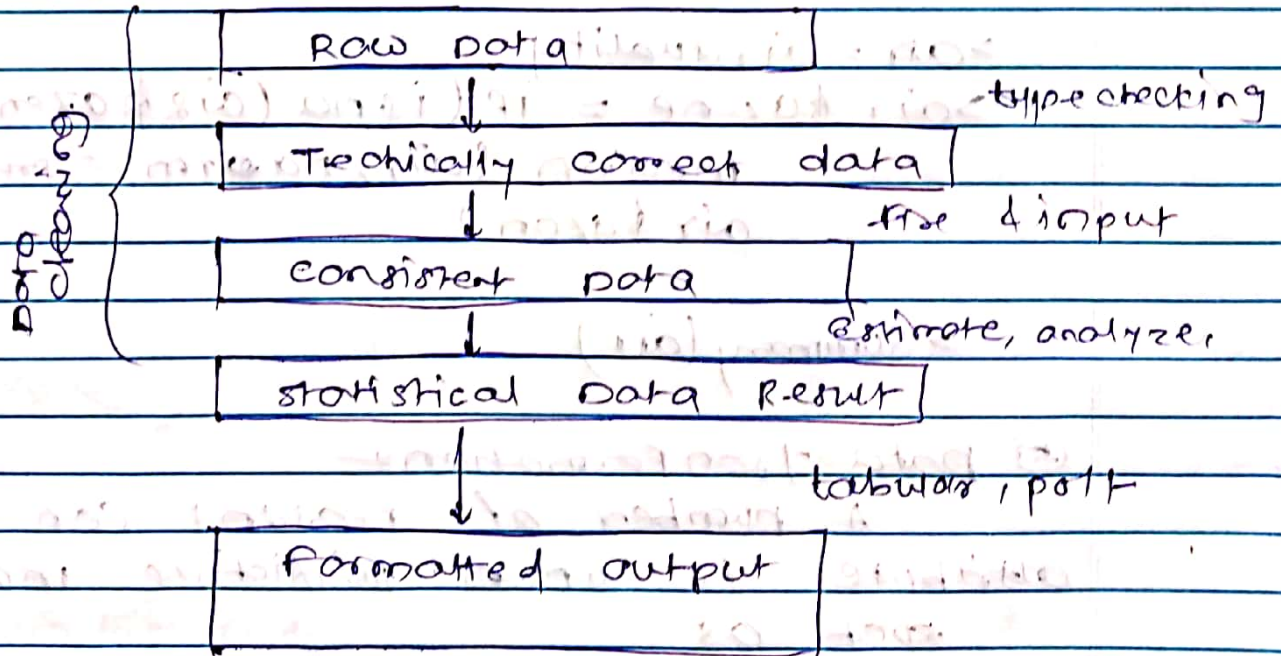
→ Problem Statement 2

Perform following operations using Python on Air quality & Heart disease dataset.

- data cleaning
- data integration
- data transformation
- error correcting
- data model building.

★ Theory :-

Diagram of statistical analysis value chain



① The dataset
* Airquality

② Errors & Corrections

```
> mean (airquality $ozone)
```

```
[1] NA
```

```
> mean (airquality $ozone, na.rm = TRUE)
```

```
[1] 42.12981
```

③ Check Summary

```
> summary (airquality)
```

④ Data cleaning (removing NAs)

```
> air = airquality
```

```
> air $ozone = if (is.na (air $ozone),  
median (air, na.rm = TRUE),  
air $ozone)
```

```
> summary (air)
```

⑤ Data Transformation

A number of reasons can be attribute to when predictive model such as

- Inadequate Data
- Inadequate model validation
- over-fitting

⑥ missing value treatment.

① Input value with median or mode

② Input value with kNN

eg.

```
> head
```

```
> air$solar_ranger = air$solar - R > 100
```

```
> head(air)
```

* Conclusion :-

Hence, I have studied & implemented the data preprocessing techniques.