# Data Analysis Process

- **Prudhvi Vardhan Notes**



## 1. Asking questions

- What features will Contribute to my analysis?
- What Features are not important for my analysis?
- Which of the features have a strong Correaltion?
- Do i need Data Preprocessing?
- What Kind of Feature engineering / Maupulation is required?

**how to Ask Better Questions?**

- subject Matter Expertise
- Experience

## 2. Data Wrangling / Munging

Data Wrangling , sometimes referred to as data munging , it is the process of **Transforming and Mapping data** from on "raw" data form into another format with intnet of making it more appropriate and valuable for a variety of downstream purposes such as analytics

- Gathering data
- Assesing data
- Cleaning data

## 2a : Gathering Data



CSV FILES          API          WEB SCRAPING          DATABASES

## 2b : Assessing Data

1. Finding the number of rows/columns( shape)
2. Data types of various columns (info())
3. Checking for missing values (info())
4. Check for duplicate data (is_unique)
5. Memory occupied by the dataset (info)
6. High level mathematical overview of the data (describe)

## 2c : Cleaning Data

1. Missing Data (e.g mean)
2. Remove duplicate data (drop_duplicates)
3. Incorrect data type (astype)

## 3. Exploratory Data Analysis

To analyze and investigate data sets and summarize their main characteristics, often employing data visualization methods.

- Exploring Data
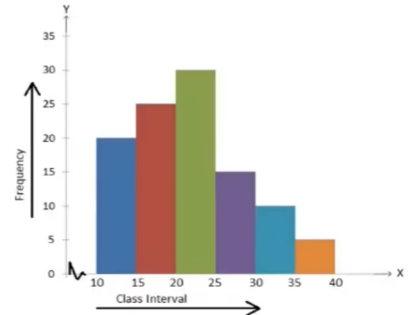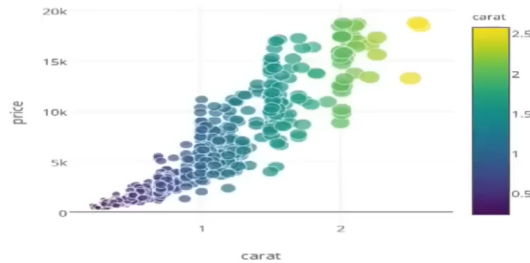- Augmenting Data

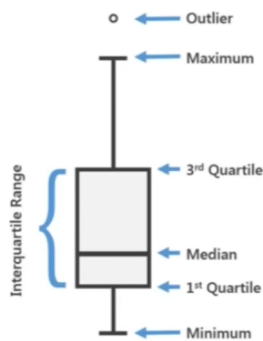## Step 3 : Exploratory Data Analysis



Explore



Augment

## 3a : Exploring Data

1. Finding Correlation and Covariance
2. Doing univariate and multivariate analysis
3. Plotting graphs( data visualization)



## 3b : Augmenting Data



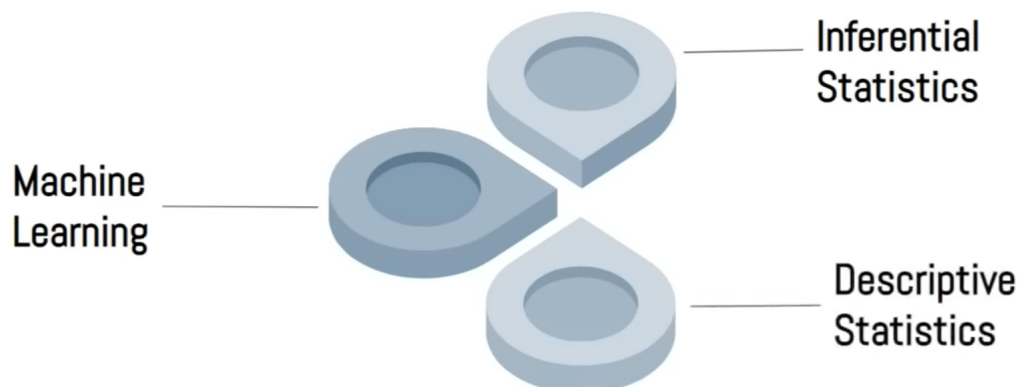Removing Outliers          Merging Dataframes          Adding new Column

These operations are collectively called as **Feature Engineering**

## 4. Drawing Conclusions

- FROM Machine learning Algorithms
- Descriptive statistics
- Inferential Statistics

## Step 4 : Drawing Conclusions



### 5. Communicating Result / story Telling

- Using PowerBI
- Tableau

## Step 5 : Communicating Results/ Data Storytelling



.

.

.