# Notes on Basics of Statistics:

It is the science of collecting, organizing, summarizing, and analyzing data to get meaningful insight from it to draw an optimal conclusion.

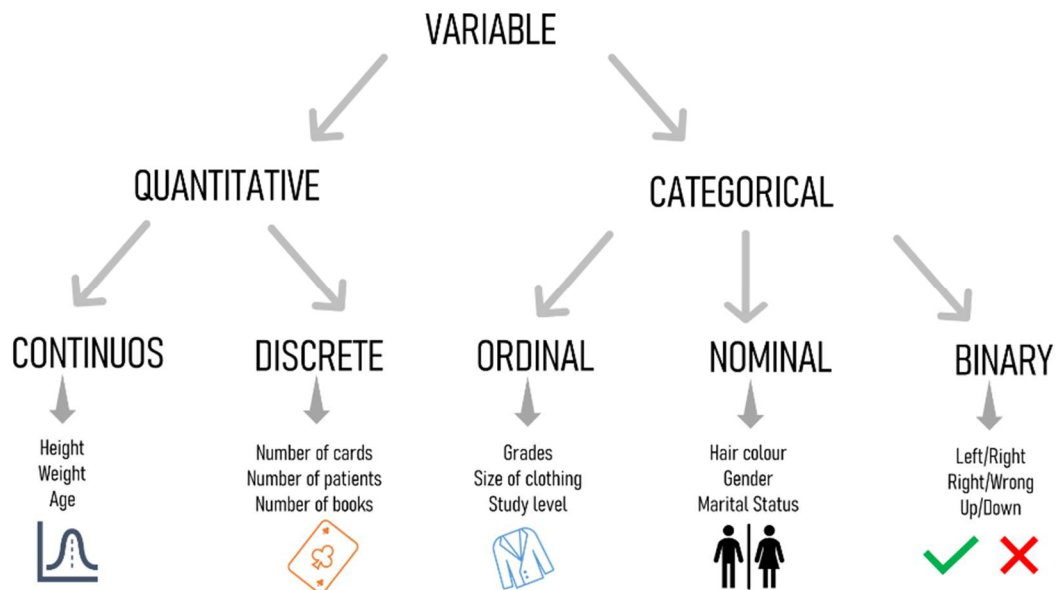## Type of Statistics:

1. **Descriptive Statistics:**
   This includes an analysis of the data at hand to understand the expected behavior using measures of central tendency (mean, median, mode, variance, standard deviation, skewness, and kurtosis).
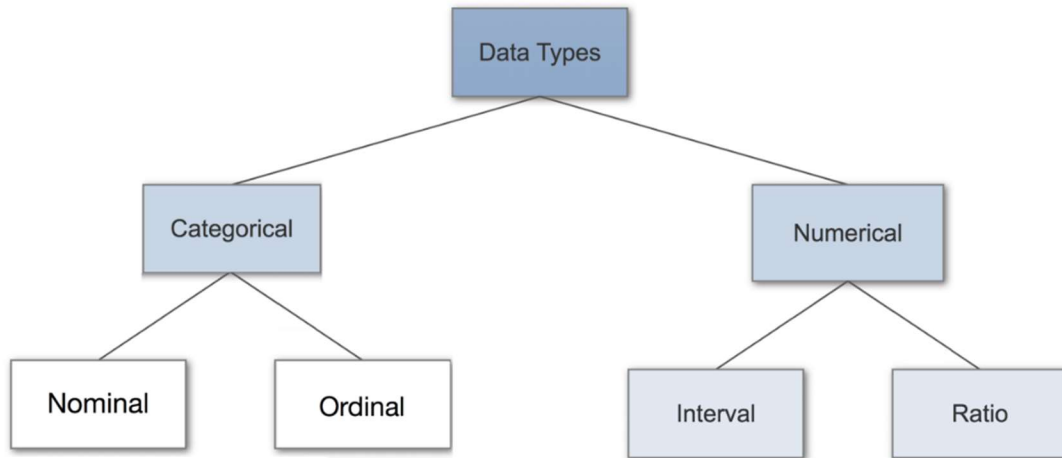
2. **Inferential Statistics:**

   This includes the selection of samples from the population to the analysis of said samples to draw conclusions about the population.
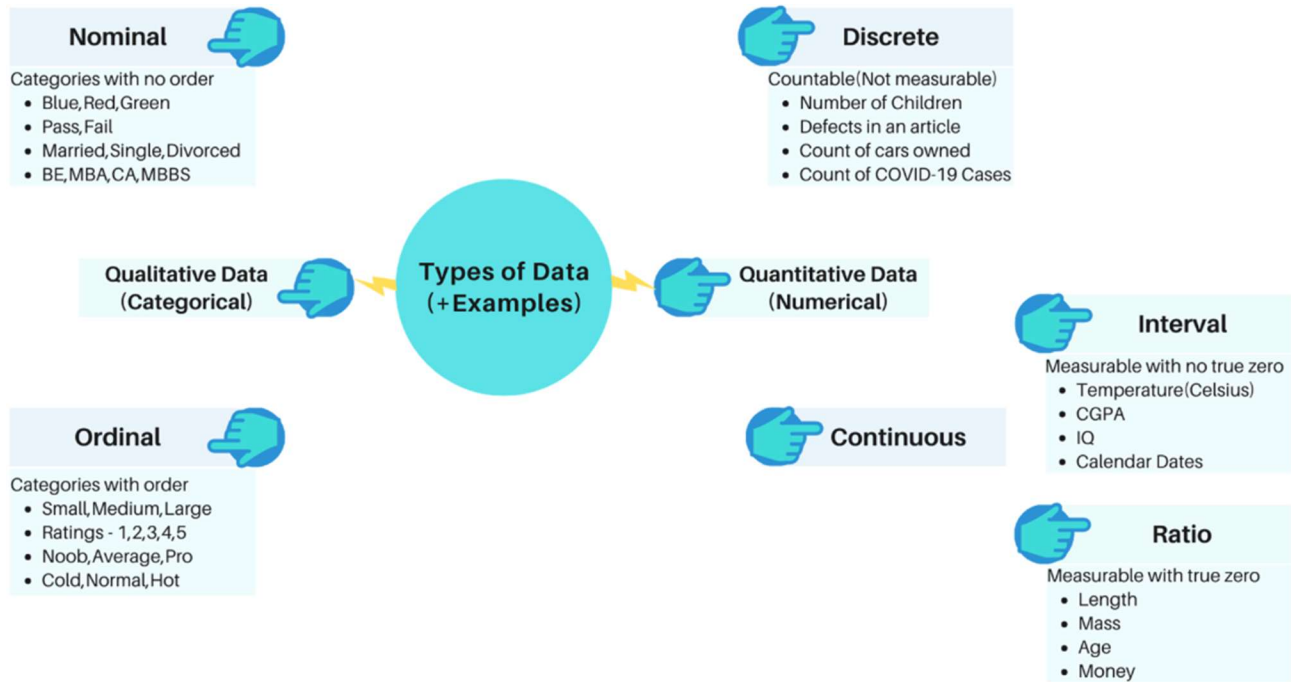
## Variables & Its type:

A variable is defined as the alphabetic character that expresses a numerical value or a number. In algebraic equations, a variable is used to represent an unknown quantity or quality of data.
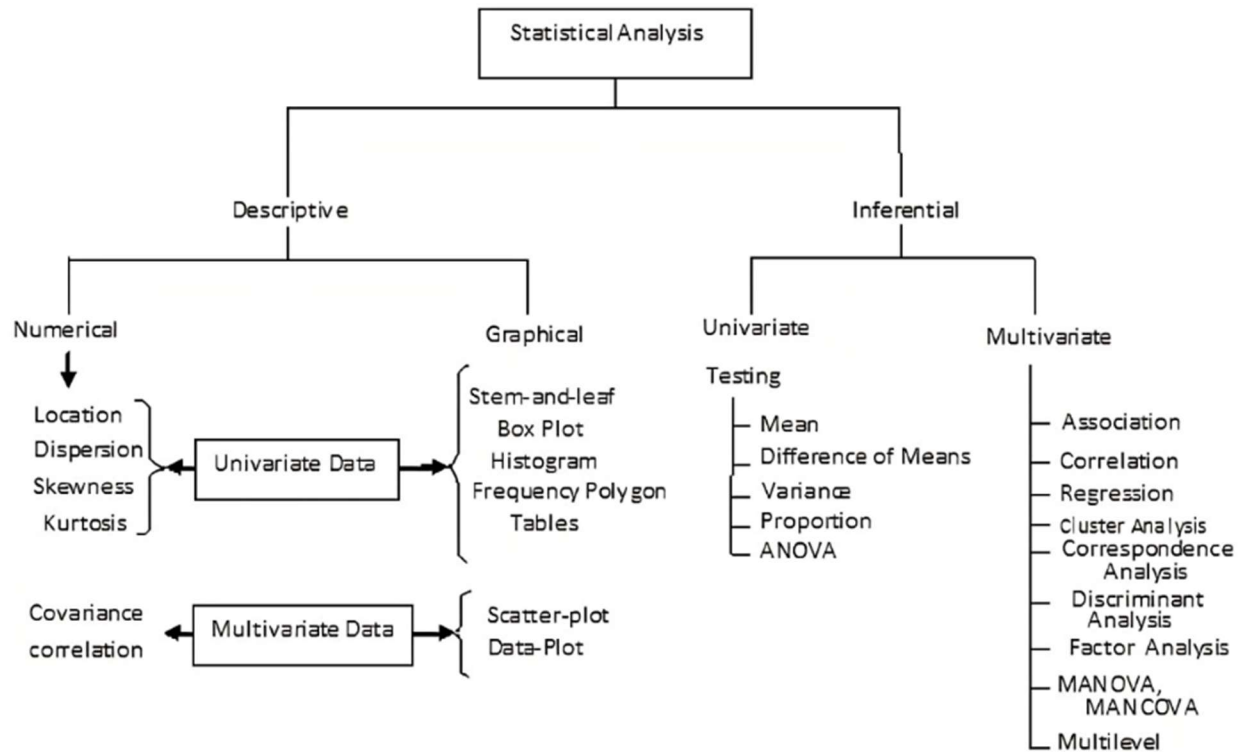
## Data Types:



## Understanding Data and Datatype:



**Nominal**
Categories with no order
- Blue, Red, Green
- Pass, Fail
- Married, Single, Divorced
- BE, MBA, CA, MBBS

**Discrete**
Countable (Not measurable)
- Number of Children
- Defects in an article
- Count of cars owned
- Count of COVID-19 Cases

**Qualitative Data (Categorical)**

**Types of Data (+Examples)**

**Quantitative Data (Numerical)**

**Interval**
Measurable with no true zero
- Temperature (Celsius)
- CGPA
- IQ
- Calendar Dates

**Ordinal**
Categories with order
- Small, Medium, Large
- Ratings - 1,2,3,4,5
- Noob, Average, Pro
- Cold, Normal, Hot

**Continuous**

**Ratio**
Measurable with true zero
- Length
- Mass
- Age
- Money

## Statistical Analysis



## <mark>The measure of central Tendency</mark>
- o **Mean**  (Mean affected by magnitude (outlier)
  - • Average of data point (sum of data point/ number of data points)
  - • <mark>Affected by outlier</mark>

    Mean of a sample:

$$\bar{x} = \frac{\sum\limits_{i=1}^{n} x_i}{n} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

- o **Median**
  - • <mark>This can be defined as Middle positional Value of given data set</mark> and can be derived in two steps,
    - (1) first by Arranging dataset (n) in ascending order then
    - (2) **for odd numbers of observations,** by selecting the observation at (n+1)/2 positions for an odd number of observation

**For an even number of observation,** calculating the average of two data points at nth/2 term and (n+1)/2 term,
[Median Formula | How To Calculate Median (Calculator, Excel Template) (educba.com)](#)

- <mark>Not affected by outlier</mark>

| Median (n=Odd) | Median (n=Even) |
|---|---|
| Median = $\dfrac{x_{(n+1)}}{2}$ | Median = $\dfrac{1}{2}\left(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}\right)$ |

- **Mode:** Mode is one of the values that indicate a central tendency of a set of data. <mark>Mode or modal value gives us an idea about which of the items in a data set is more likely to occur frequently</mark>. It is the measure of Central Tendency other than Mean and Median.
- Depending on the type of dataset given, you can find one<mark>, two three or even multiple modal values. Some data sets may have no mode value at all.</mark>
  - Most occurring value in the data set
  - Can be unimodal, bimodal, multimodal

$$\text{Mode} = l + \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \times h$$

Where,

l = lower limit of the modal class
h = size of the class interval
$f_1$ = frequency of the modal class
$f_0$ = frequency of the class preceding the modal class
$f_2$ = frequency of the class succeeding the modal class
([Mode Formula: Learn to Calculate Mode For Set of Data - Embibe](#))

- Shows Expected Behavior of Data Set at the center

**Spread (Overall spread of Data)**
- **Range (Max-Min)**
  - The range for any distribution is given by = highest value – lowest value.
    <mark>Range = Highest Value – Lowest Value</mark>
  - Doesn't show variation in between
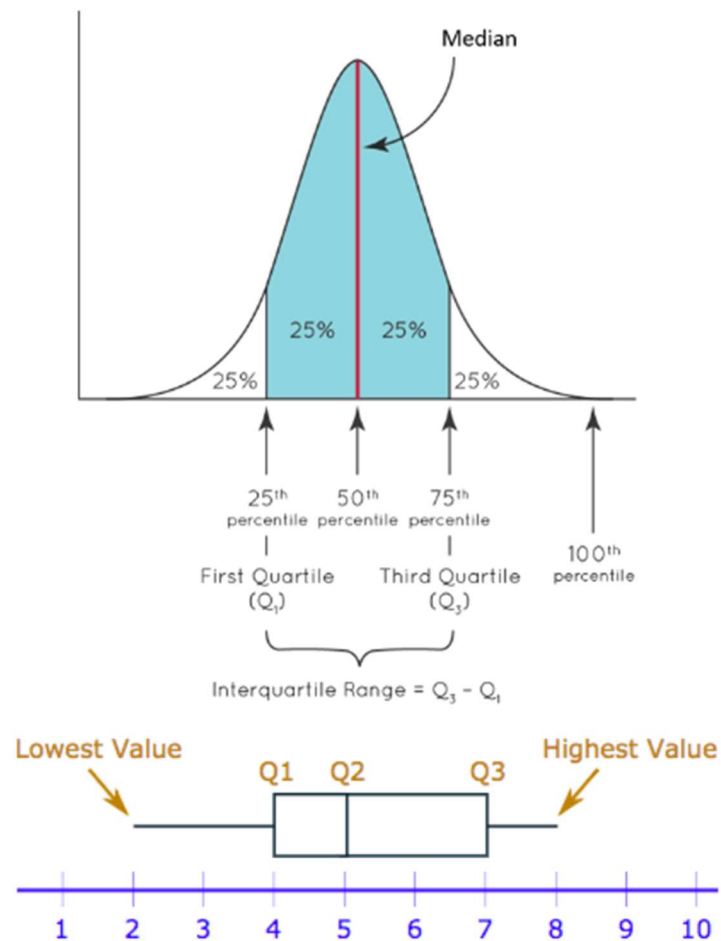
## IQR (Inter Quartile Range)

When the set of observations are arranged in ascending order the quartiles are represented as,

| | |
|---|---|
| $Q_1 = \left(\dfrac{n+1}{4}\right)^{th}$ term | Excel Formula : **QUARTILE(( ), 1)** |
| $Q_2 = \left(\dfrac{2(n+1)}{4}\right)^{th}$ term | Excel Formula : **QUARTILE(( ), 2)** |
| $Q_3 = \left(\dfrac{3*(n+1)}{4}\right)^{th}$ term | Excel Formula : **QUARTILE(( ), 3)** |

The interquartile range (IQR) = Upper Quartile (Q3) – Lower Quartile (Q1)
Lower Bound Limit = Q1 - 1.5 x IQR
Upper Bound Limit = Q3 + 1.5 x IQR

## Measures of Dispersion:

- **Variance :**
  Average squared deviation of the data points from their mean.
  - ○ Affected by extreme

| Population variance: | Sample variance: |
|---|---|
| $$\sigma^2 = \frac{\sum_{i=1}^{N}(x_i - \mu)^2}{N}$$ where $\mu$ is the population mean. | $$s^2 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n - 1}$$ |

- **Standard Deviation (SD) : How data variating inside the range**
  **Sample Standard Deviation:**
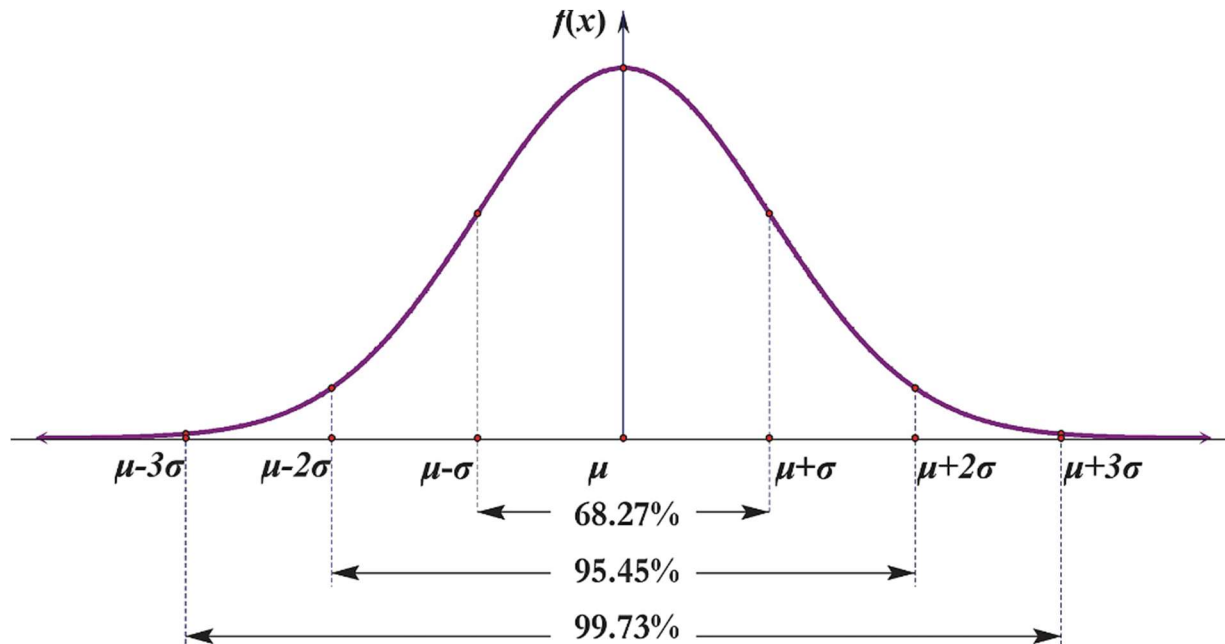
$$\sigma = \sqrt{\frac{\sum (x - u)^2}{N}} \qquad s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n - 1}}$$

(Note: for sample data n-1 is considered for **Basel correction**)

$1\sigma = 68.26\,\%$, Data spread both lef and right from center is 68.26 %
$2\sigma = 95.44\,\%$, Data spread both lef and right from center is 95.44 %
$3\sigma = 99.74\%$, Data spread both lef and right from center is 97.74 %
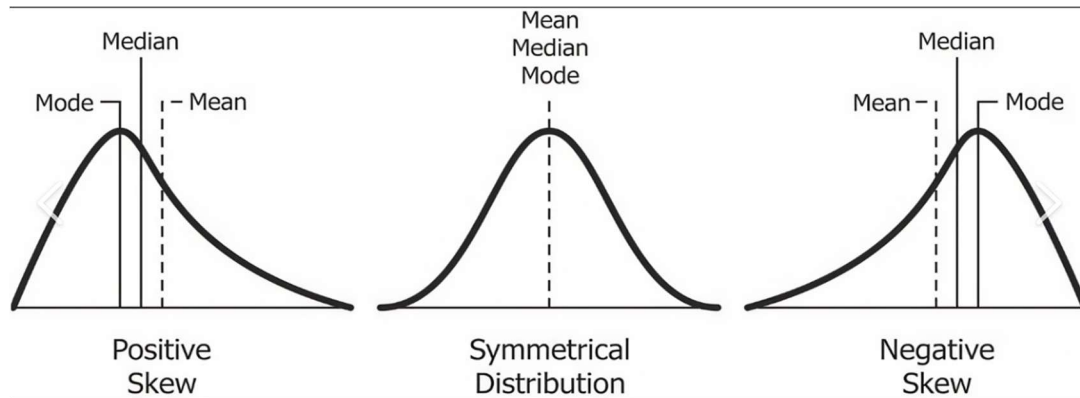


**The larger the value of *n* is, the closer that the population and sample standard deviations will be.**

**Shape of Distribution**

- **Symmetricity / Normal distribution**

  Mean = Median = Mode



| Positive Skew | Symmetrical Distribution | Negative Skew |

$$\text{Skewness} = \frac{\sum_i^N (X_i - \overline{X})^3}{(N-1) * o^3}$$

**Skewness (skew: where it has been flattened)**

    o  Positive skewed **(mean>median)/mean-median= + ve** (i.e. Salary in an Organization) - flattened at right

    o  **Negative Skewed (Median > Mean)/mean-median= - ve** (i.e. easy Exam)

**Kurtosis: (Strength of Relationship with Standard Deviation)**

It is statistical measures which define shape of each tails of distribution where it is heavy tailed (presence of outliers) or light tailed (paucity of outliers) compared to normal distribution.
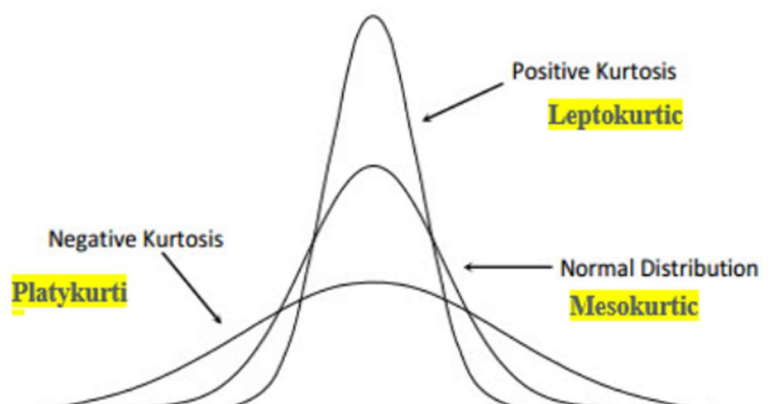
If **Kurtosis > ±3 SD then it is to be called Leptokurtic** (short tailed) with Low Standard Deviation

if **Kurtosis < ±3 SD then it is called Platykurtic** (long Tailed) with High Standard Deviation

If **Kurtosis = ± 3 SD then it is called Mesokurtic**

$$\text{Kurtosis} = n * \frac{\sum_i^n (Y_i - \overline{Y})^4}{\sum_i^n (Y_i - \overline{Y}^2)^2}$$

## Getting Major Statistical Details in Excel / Python / R:

**R-Program / Library (Pastecs)**
**stat.desc(filename$column_name) to get all the statistical measures:**

```
> library(pastecs)
> stat.desc(Meetings$MeetingTime)
     nbr.val      nbr.null       nbr.na           min           max        range          sum       median         mean
   70.000000      0.000000      0.000000     30.000000    150.000000   120.000000  4520.000000    60.000000    64.571429
     SE.mean  CI.mean.0.95          var       std.dev      coef.var
    2.757904      5.501868    532.422360     23.074279      0.357345
```
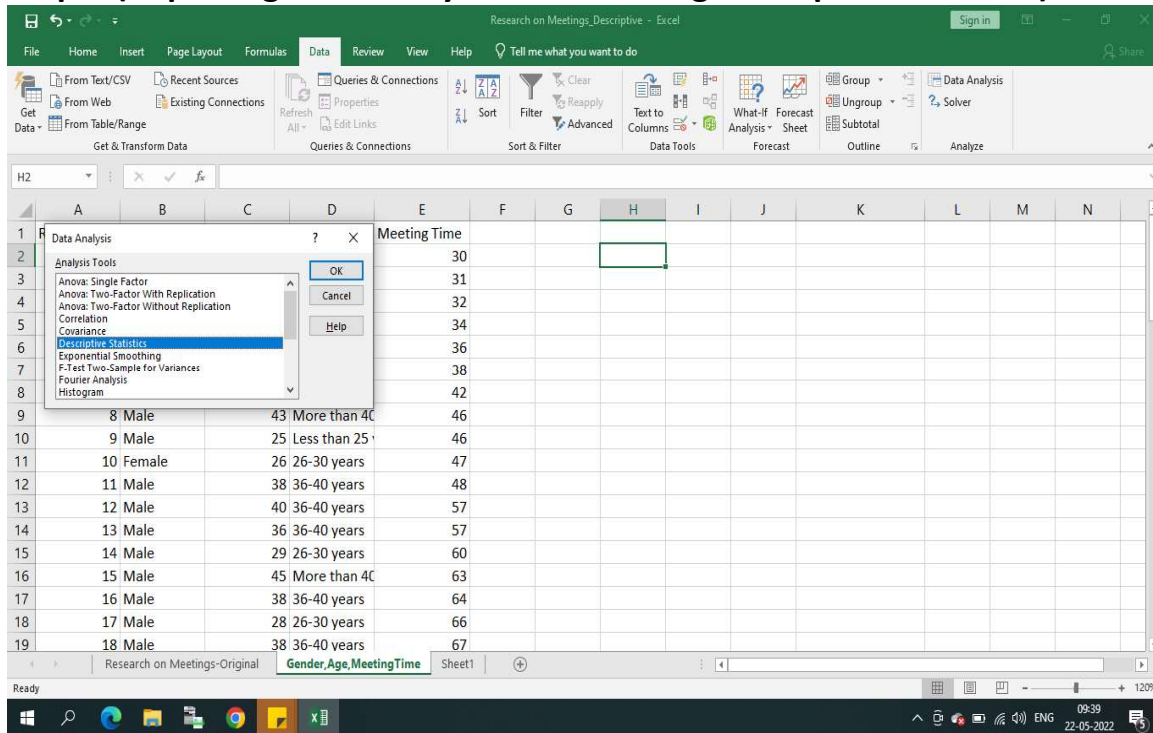
**Python /library (pandas)**
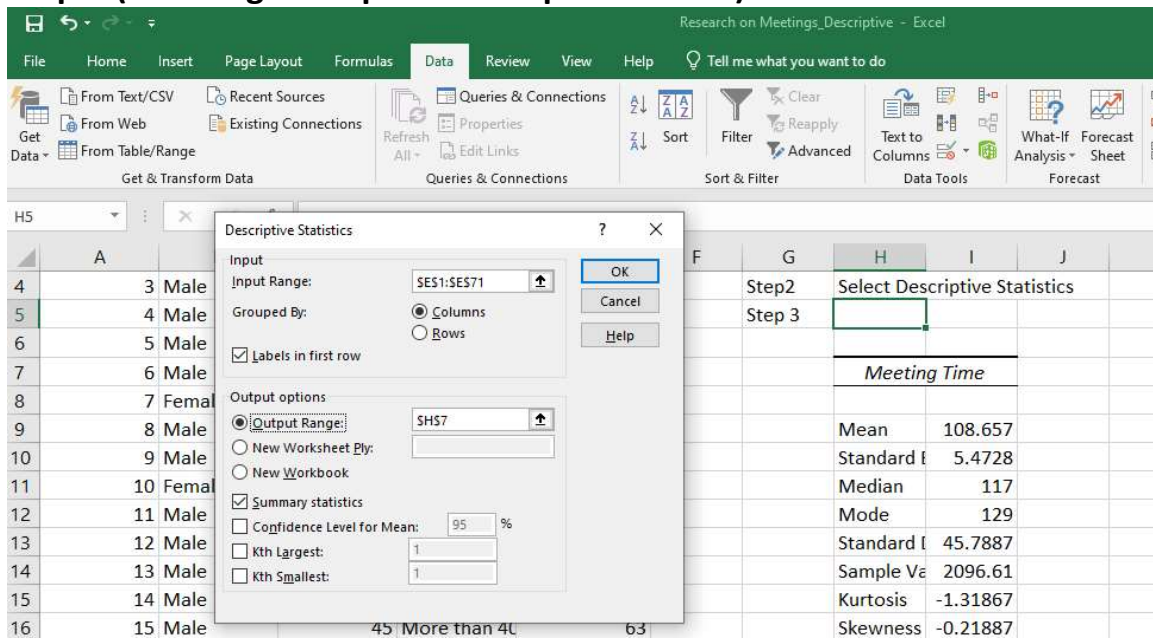**df. describe() to get all the statistical measures:**

```
dfm.describe()
✓  0.7s
```

|       | Respondents | Age(inyears) | MeetingTime |
|-------|-------------|--------------|-------------|
| count | 70.000000   | 70.000000    | 70.000000   |
| mean  | 35.500000   | 32.171429    | 64.571429   |
| std   | 20.351085   | 5.036019     | 23.074279   |
| min   | 1.000000    | 24.000000    | 30.000000   |
| 25%   | 18.250000   | 28.000000    | 46.250000   |
| 50%   | 35.500000   | 31.500000    | 60.000000   |
| 75%   | 52.750000   | 36.000000    | 75.000000   |
| max   | 70.000000   | 45.000000    | 150.000000  |

## Step-1 (Importing Data Analysis and selecting Descriptive Statistics)



## Step-2 (selecting the input and output columns)

# Type of Plots