```
In [1]:  import pandas as pd
```

```
In [2]:  pd.__version__    #To Check version
```

```
Out[2]:  '2.2.2'
```

```
In [3]:  store = pd.read_csv(r'D:\Full Stack Data Scientist and AI\March 19 - Introductio
```

```
In [40]:  store    #store is the name of the object, created in Python. Called the variable
```

Out[40]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | O |
|---|---|---|---|---|---|---|---|
| 0 | Office Supplies | Houston | United States | Darren Powers | Message Book | 03-01-2020 | 2(103 |
| 1 | Office Supplies | Naperville | United States | Phillina Ober | GBC | 04-01-2020 | 2(112 |
| 2 | Office Supplies | Naperville | United States | Phillina Ober | Avery | 04-01-2020 | 2(112 |
| 3 | Office Supplies | Naperville | United States | Phillina Ober | SAFCO | 04-01-2020 | 2(112 |
| 4 | Office Supplies | Philadelphia | United States | Mick Brown | Avery | 05-01-2020 | 2(141 |
| ... | ... | ... | ... | ... | ... | ... | |
| 10189 | Office Supplies | New York City | United States | Patrick O'Donnell | Wilson Jones | 30-12-2023 | 2(143 |
| 10190 | Office Supplies | Fairfield | United States | Erica Bern | GBC | 30-12-2023 | 2(115 |
| 10191 | Office Supplies | Loveland | United States | Jill Matthias | Other | 30-12-2023 | 2(156 |
| 10192 | Technology | New York City | United States | Patrick O'Donnell | Other | 30-12-2023 | 2(143 |
| 10193 | Office Supplies | Charlottetown | Canada | Harry Olson | Wilson Jones | 30-12-2023 | 2(143 |

10194 rows × 19 columns

◄ ████████████_____ ►

In [5]: `id(store)` *#id gives address of memory allocation*

Out[5]:   2882727282448

In [6]:   `len(store)   #Numbers of Rows`

Out[6]:   10194

In [7]:   `store.shape   #shape gives -> Gives dimensions i.e Numbers of rows and columns in`

Out[7]:   (10194, 19)

In [8]:   `store.columns   #columns gives column names`

Out[8]:   Index(['Category', 'City', 'Country/Region', 'Customer Name', 'Manufacturer',
          'Order Date', 'Order ID', 'Postal Code', 'Product Name', 'Region',
          'Segment', 'Ship Date', 'Ship Mode', 'State/Province', 'Sub-Category',
          'Discount', 'Profit', 'Quantity', 'Sales'],
         dtype='object')

dtype='object' but actually datatype is int, float. But here system by default considered data type as object.

In [10]:   `len(store.columns)`

Out[10]:   19

# To check NULL Values

In [12]:   `store.isnull()   #Hey Python, is there any NULL value in the data set?`

Out[12]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Postal Code |
|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 10189 | False | False | False | False | False | False | False | False |
| 10190 | False | False | False | False | False | False | False | False |
| 10191 | False | False | False | False | False | False | False | False |
| 10192 | False | False | False | False | False | False | False | False |
| 10193 | False | False | False | False | False | False | False | False |

10194 rows × 19 columns

There is no null value, thus we get False. If there are any missing values then the answer is True

In [14]: `store.notnull()`  `#Hey Python, is there any Not NULL value in the data set?`

Out[14]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Postal Code |
|---|---|---|---|---|---|---|---|---|
| **0** | True | True | True | True | True | True | True | True |
| **1** | True | True | True | True | True | True | True | True |
| **2** | True | True | True | True | True | True | True | True |
| **3** | True | True | True | True | True | True | True | True |
| **4** | True | True | True | True | True | True | True | True |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| **10189** | True | True | True | True | True | True | True | True |
| **10190** | True | True | True | True | True | True | True | True |
| **10191** | True | True | True | True | True | True | True | True |
| **10192** | True | True | True | True | True | True | True | True |
| **10193** | True | True | True | True | True | True | True | True |

10194 rows × 19 columns

In [15]: `store.isnull().sum()`

Out[15]:
```
Category          0
City              0
Country/Region    0
Customer Name     0
Manufacturer      0
Order Date        0
Order ID          0
Postal Code       0
Product Name      0
Region            0
Segment           0
Ship Date         0
Ship Mode         0
State/Province    0
Sub-Category      0
Discount          0
Profit            0
Quantity          0
Sales             0
dtype: int64
```

0 means 0 missing values

In [17]: `store[:]`  `#Store slice -> Prints entire data set`

Out[17]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | O... |
|---|---|---|---|---|---|---|---|
| **0** | Office Supplies | Houston | United States | Darren Powers | Message Book | 03-01-2020 | 20 103 |
| **1** | Office Supplies | Naperville | United States | Phillina Ober | GBC | 04-01-2020 | 20 112 |
| **2** | Office Supplies | Naperville | United States | Phillina Ober | Avery | 04-01-2020 | 20 112 |
| **3** | Office Supplies | Naperville | United States | Phillina Ober | SAFCO | 04-01-2020 | 20 112 |
| **4** | Office Supplies | Philadelphia | United States | Mick Brown | Avery | 05-01-2020 | 20 141 |
| **...** | ... | ... | ... | ... | ... | ... | ... |
| **10189** | Office Supplies | New York City | United States | Patrick O'Donnell | Wilson Jones | 30-12-2023 | 20 143 |
| **10190** | Office Supplies | Fairfield | United States | Erica Bern | GBC | 30-12-2023 | 20 115 |
| **10191** | Office Supplies | Loveland | United States | Jill Matthias | Other | 30-12-2023 | 20 156 |
| **10192** | Technology | New York City | United States | Patrick O'Donnell | Other | 30-12-2023 | 20 143 |
| **10193** | Office Supplies | Charlottetown | Canada | Harry Olson | Wilson Jones | 30-12-2023 | 20 143 |

10194 rows × 19 columns

◄ ▬▬▬▬▬▬▬▬▬▬▬▬ ►

In [18]: 
```
store[0:10]  #This prints 0 to 10(n-1)th i.e 9th row (records). Row means record
```

Out[18]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Po C |
|---|---|---|---|---|---|---|---|---|
| 0 | Office Supplies | Houston | United States | Darren Powers | Message Book | 03-01-2020 | US-2020-103800 | 77 |
| 1 | Office Supplies | Naperville | United States | Phillina Ober | GBC | 04-01-2020 | US-2020-112326 | 60 |
| 2 | Office Supplies | Naperville | United States | Phillina Ober | Avery | 04-01-2020 | US-2020-112326 | 60 |
| 3 | Office Supplies | Naperville | United States | Phillina Ober | SAFCO | 04-01-2020 | US-2020-112326 | 60 |
| 4 | Office Supplies | Philadelphia | United States | Mick Brown | Avery | 05-01-2020 | US-2020-141817 | 19 |
| 5 | Furniture | Henderson | United States | Maria Etezadi | Global | 06-01-2020 | US-2020-167199 | 42 |
| 6 | Office Supplies | Henderson | United States | Maria Etezadi | Rogers | 06-01-2020 | US-2020-167199 | 42 |
| 7 | Office Supplies | Athens | United States | Jack O'Briant | Dixon | 06-01-2020 | US-2020-106054 | 30 |
| 8 | Office Supplies | Henderson | United States | Maria Etezadi | Ibico | 06-01-2020 | US-2020-167199 | 42 |

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Po C |
|---|---|---|---|---|---|---|---|---|
| **9** | Office Supplies | Henderson | United States | Maria Etezadi | Alliance | 06-01-2020 | US-2020-167199 | 42 |

**In Numpy we don't get 0, 1, 2, 3..There we count manually**

**But in Pandas dataFrame, we get index -> 0, 1, 2, 3..**

In [20]: `store[0:20:5]`

Out[20]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | |
|---|---|---|---|---|---|---|---|---|
| **0** | Office Supplies | Houston | United States | Darren Powers | Message Book | 03-01-2020 | US-2020-103800 | |
| **5** | Furniture | Henderson | United States | Maria Etezadi | Global | 06-01-2020 | US-2020-167199 | |
| **10** | Office Supplies | Henderson | United States | Maria Etezadi | Southworth | 06-01-2020 | US-2020-167199 | |
| **15** | Office Supplies | Huntsville | United States | Vivek Sundaresam | Acco | 07-01-2020 | US-2020-105417 | |

◄ ████████████████ ►

In [38]: `store.head()`  *#head() function gives top 5 rows/ displays top 5 records*

Out[38]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Po C |
|---|---|---|---|---|---|---|---|---|
| 0 | Office Supplies | Houston | United States | Darren Powers | Message Book | 03-01-2020 | US-2020-103800 | 77 |
| 1 | Office Supplies | Naperville | United States | Phillina Ober | GBC | 04-01-2020 | US-2020-112326 | 60 |
| 2 | Office Supplies | Naperville | United States | Phillina Ober | Avery | 04-01-2020 | US-2020-112326 | 60 |
| 3 | Office Supplies | Naperville | United States | Phillina Ober | SAFCO | 04-01-2020 | US-2020-112326 | 60 |
| 4 | Office Supplies | Philadelphia | United States | Mick Brown | Avery | 05-01-2020 | US-2020-141817 | 19 |

◄ ━━━━━━━━━━━━━━━━━━━━━━━ ►

In [42]:
```
store.tail()  #tail() function gives bottom 5 rows/ displays bottom 5 records
```

Out[42]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | O |
|---|---|---|---|---|---|---|---|
| **10189** | Office Supplies | New York City | United States | Patrick O'Donnell | Wilson Jones | 30-12-2023 | 20 143 |
| **10190** | Office Supplies | Fairfield | United States | Erica Bern | GBC | 30-12-2023 | 20 115 |
| **10191** | Office Supplies | Loveland | United States | Jill Matthias | Other | 30-12-2023 | 20 156 |
| **10192** | Technology | New York City | United States | Patrick O'Donnell | Other | 30-12-2023 | 20 143 |
| **10193** | Office Supplies | Charlottetown | Canada | Harry Olson | Wilson Jones | 30-12-2023 | 20 143 |

◀ ▶

store. #store . tab -> Displays all functionalities of pandas

In [54]: 
```
store.isna()  #isna() and isnull() both are same
```

Out[54]:

| | Category | City | Country/Region | Customer Name | Manufacturer | Order Date | Order ID | Postal Code |
|---|---|---|---|---|---|---|---|---|
| **0** | False | False | False | False | False | False | False | False |
| **1** | False | False | False | False | False | False | False | False |
| **2** | False | False | False | False | False | False | False | False |
| **3** | False | False | False | False | False | False | False | False |
| **4** | False | False | False | False | False | False | False | False |
| **...** | ... | ... | ... | ... | ... | ... | ... | .. |
| **10189** | False | False | False | False | False | False | False | False |
| **10190** | False | False | False | False | False | False | False | False |
| **10191** | False | False | False | False | False | False | False | False |
| **10192** | False | False | False | False | False | False | False | False |
| **10193** | False | False | False | False | False | False | False | False |

10194 rows × 19 columns

◀ ▶

# Introduction to Statistical Concept in Pandas

Pandas is a library which handles rows, columns and series.

- Excel sheet means either number or text

- Number -> Numerical Data

- Text data -> Categorical Data

- Dataset is formed with a combination of numerical data and categorical data.

- Numerical data and Categorical data is called Statistical World.

- If dataset is number -> We call it numerical

- If dataset is text -> We call it categorical

In [62]: 
```
store.describe()
```

Out[62]:

|       | Discount | Profit | Quantity | Sales |
|-------|----------|--------|----------|-------|
| count | 10194.000000 | 10194.000000 | 10194.000000 | 10194.000000 |
| mean | 0.155385 | 28.673417 | 3.791838 | 228.225854 |
| std | 0.206249 | 232.465115 | 2.228317 | 619.906839 |
| min | 0.000000 | -6599.978000 | 1.000000 | 0.444000 |
| 25% | 0.000000 | 1.760800 | 2.000000 | 17.220000 |
| 50% | 0.200000 | 8.690000 | 3.000000 | 53.910000 |
| 75% | 0.200000 | 29.297925 | 5.000000 | 209.500000 |
| max | 0.800000 | 8399.976000 | 14.000000 | 22638.480000 |

- describe() refers to descriptive statistics

- Only Discount, Profit, Quantity, Sales have numbers.

- Thus, describe() displays these attributes as describe() by default displays only numerical data

In [ ]: