# *Analyzing Amazon Sales Data*

## *Wireframe Documentation*
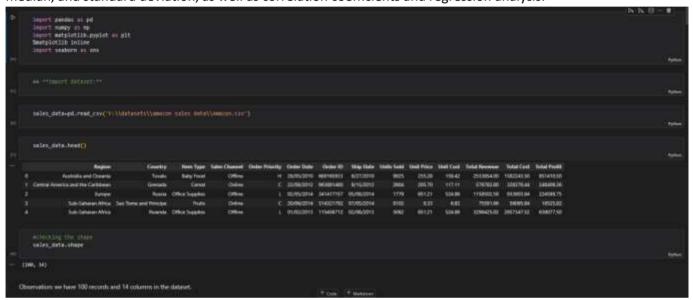
**Revision Number - 1.0**

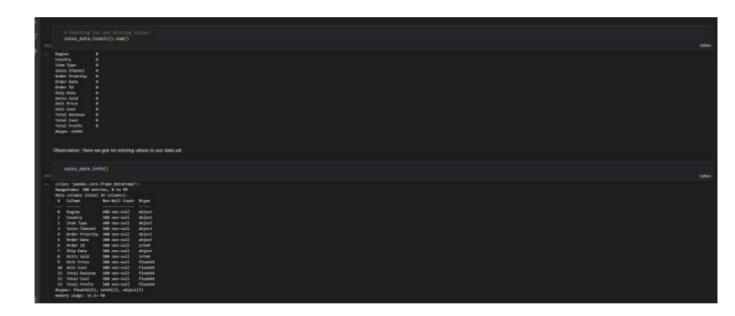**Last Date of Revision - 14/04/2023**

**Mohammad Aadil**

# Performed Exploratory Data Analysis :-

Exploratory Data Analysis (EDA) is an approach to analyzing data that emphasizes gaining an understanding of the data before formal modeling or hypothesis testing is carried out. EDA is a crucial step in the data analysis process, as it can help identify patterns, anomalies, and relationships between variables, and provide insights into the underlying structure of the data

EDA typically involves a combination of visual and numerical methods to explore the data. Visual methods include techniques such as scatterplots, boxplots, histograms, and heatmaps, which allow us to visualize the distribution of variables and the relationships between them. Numerical methods include summary statistics such as mean, median, and standard deviation, as well as correlation coefficients and regression analysis.

```python
# Let's see the columns name
sales_data.columns
```

```
Index(['Region', 'Country', 'Item Type', 'Sales Channel', 'Order Priority',
       'Order Date', 'Order ID', 'Ship Date', 'Units Sold', 'Unit Price',
       'Unit Cost', 'Total Revenue', 'Total Cost', 'Total Profit'],
      dtype='object')
```

```python
sales_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   Region          100 non-null    object
 1   Country         100 non-null    object
 2   Item Type       100 non-null    object
 3   Sales Channel   100 non-null    object
 4   Order Priority  100 non-null    object
 5   Order Date      100 non-null    object
 6   Order ID        100 non-null    int64
 7   Ship Date       100 non-null    object
 8   Units Sold      100 non-null    int64
 9   Unit Price      100 non-null    float64
 10  Unit Cost       100 non-null    float64
 11  Total Revenue   100 non-null    float64
 12  Total Cost      100 non-null    float64
 13  Total Profit    100 non-null    float64
dtypes: float64(5), int64(2), object(7)
memory usage: 11.1+ KB
```

Observation In the information of the dataset we got to know the number of columns, columns name and the data type of the columns.

```python
# Checking for any missing values
sales_data.isnull().sum()
```

```
Region            0
Country           0
Item Type         0
Sales Channel     0
Order Priority    0
Order Date        0
Order ID          0
Ship Date         0
Units Sold        0
Unit Price        0
Unit Cost         0
Total Revenue     0
Total Cost        0
Total Profit      0
dtype: int64
```

Observation : Here we got no missing values in our data set

```python
sales_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   Region          100 non-null    object
 1   Country         100 non-null    object
 2   Item Type       100 non-null    object
 3   Sales Channel   100 non-null    object
 4   Order Priority  100 non-null    object
 5   Order Date      100 non-null    object
 6   Order ID        100 non-null    int64
 7   Ship Date       100 non-null    object
 8   Units Sold      100 non-null    int64
 9   Unit Price      100 non-null    float64
 10  Unit Cost       100 non-null    float64
 11  Total Revenue   100 non-null    float64
 12  Total Cost      100 non-null    float64
 13  Total Profit    100 non-null    float64
dtypes: float64(5), int64(2), object(7)
memory usage: 11.1+ KB
```

## Visualizing Sales Trend:-

Visualizing sales trends year and month wise can provide a more granular view of sales data and help identify patterns and seasonal trends. This type of analysis can be particularly useful for businesses that experience fluctuations in sales due to seasonal or other factors.

One way to visualize sales trends year and month wise is to use a line plot or a heatmap. A line plot can show how sales vary over time, while a heatmap can show which months are typically high or low in sales.