**LONDON METROPOLITAN UNIVERSITY**

**islington college**
(इस्लिङ्टन कलेज)

**CC5067NI-Smart Data Discovery**


**60% of Individual Coursework**


**2022-23 Spring**


**Student Name: Bhuwani Bikram Nembang**

**London Met ID: 22015664**

**College ID: np01cp4s220164@islingtoncollege.edu.np**

**Assignment Due Date: Friday, March 3, 2023**

**Assignment Submission Date: Thursday, May 4, 2023**

**Word Count: 2140**

## Table of Contents

**Table Of Figures**

**Table of Tables**

# 1. Introduction to the project:

The coursework is about applying programming knowledge and skills to data analysis skills, demonstrating my skills for problem-solving and critical thinking, and evaluation. In this assignment, we analyse the data of ABC company for the year 2019. We are about to write Python programs and a technical report on data understanding, preparation, exploration, and initial analysis.

# Acknowledgement

I would like to extend my sincere appreciation and gratitude to ABC Company for providing me with the opportunity to complete the coursework on data analysis. The experience was invaluable and has enriched my knowledge in this field. I would like to express my gratitude to my instructors for their guidance, support, and valuable feedback throughout the coursework. Their expertise and insights were crucial in enhancing my understanding of data analysis techniques and tools. I would also like to acknowledge my fellow colleagues for their collaboration and teamwork during the project. Their input and participation were essential in achieving the project goals and objectives.

## Abstract

The coursework on data analysis conducted at ABC Company aimed to enhance the participants' knowledge and skills in the field of data analysis. The coursework included theoretical and practical training on various data analysis techniques and tools. The coursework provided a comprehensive overview of data analysis, including data cleaning, data transformation, data visualization, and statistical analysis. The participants gained hands-on experience with software tools such as Excel, Python, and R, and learned to use these tools for data analysis and visualization. The coursework also included a project where the participants applied their skills and knowledge to real-world data analysis problems. The project provided an opportunity for the participants to work collaboratively and to demonstrate their ability to apply data analysis techniques to real-world scenarios.

Overall, the coursework was an enriching experience that provided participants with practical skills and knowledge that can be applied in various industries and fields. The coursework was delivered by experienced instructors, and the participants received valuable feedback and support throughout the program. The coursework at ABC Company is a testament to the company's commitment to education and professional development, and it has equipped the participants with the necessary skills to excel in their future endeavours.

**Tools Used**

 **Anaconda:**

Anaconda software helps you create an environment for many different versions of Python and package versions. Anaconda is also used to install, remove, and upgrade packages in your project environments.



*Figure 1: figure of Anaconda.*

**Jupyter Notebook:**

The Jupyter Notebook is the original web application for creating and sharing computational documents. It offers a simple, streamlined, document-centric experience.



*Figure 2: figure of jupyter notebook.*

**Snipping tool:**

Snipping Tool is a Microsoft Windows screenshot utility included in Windows Vista and later. It can take still screenshots of an open window, rectangular areas, a free-form area, or the entire screen.



*Figure 3: figure of snipping tool.*

2201564 Bhuwani Bikram Nembang

## 2. Data Understanding:

### 2.1 Data set

A data set is a collection of related, discrete items of related data that may be accessed individually or in combination or managed as a whole entity (Contributer, 1999-2023). A dataset in smart data discovery refers to a collection of data that is used as the basis for analysis and insights generation. These datasets can come from a variety of sources, such as databases, spreadsheets, or other data storage systems.

In the context of smart data discovery, datasets are often pre-processed and transformed to make it easier for users to analyses the data and find patterns or insights. Smart data discovery platforms use advanced algorithms and machine learning techniques to automatically identify and extract relevant data from these datasets, and to generate insights and recommendations based on this data.

The characteristics of the resources used in data science can vary depending on the specific task or application. However, some common characteristics of resources used in data science include:

- **glob:** glob is a Python module that is used to find all the pathnames matching a specified pattern according to the rules used by the Unix shell, although results are returned in arbitrary order (Jain, 2023).


- **Matplotlib:** matplotlib.pyplot is a collection of functions that provide a convenient interface to create a variety of charts, plots, and visualizations in Python. It is a part of the Matplotlib library which is one of the most widely used visualization libraries in Python (Meghna, 2023).


- **.concat():** The .concat() method is a built-in method in JavaScript used to concatenate two or more arrays, creating a new array that contains all the elements of the original arrays in the order in which they were passed as arguments.

- **.to_csv():** .to_csv() is a method used to write a Data Frame to a CSV (Comma-Separated Values) file. A Data Frame is a 2-dimensional labelled data structure with columns of potentially different types.

- **.dropna():** .dropna() method is used to remove missing or null values from a Pandas Data Frame. The method drops rows or columns that contain missing values, depending on the value of the axis parameter.

- **.astype (int):** .astype(int) is a method used to convert a Pandas Data Frame or Series to an integer data type. The method converts each element of the Data Frame or Series to an integer, rounding down any decimal values.

- **.str.split():** .str.split() method is used to split a string column of a Pandas Data Frame into multiple columns based on a specified delimiter. The method returns a new Data Frame with the original string column split into multiple columns.

- **.corr():** .corr() method is used to compute the correlation matrix of a Pandas Data Frame. The correlation matrix shows the pairwise correlation coefficients between all pairs of columns in the Data Frame.

Overall, datasets play a critical role in smart data discovery, as they provide the foundation for data-driven insights and decision-making. By using advanced analytics tools and techniques to analyses these datasets, businesses can gain a deeper understanding of their operations, customers, and markets, and use this knowledge to drive growth and improve their bottom line.

| S.N | Column Name | Data Type | Description |
|---|---|---|---|
| 1 | Order Id | Float64 | All the Order Id are stored in this column |
| 2 | Product | object | All the product details are stored here, object as a data type. |
| 3 | Quantity Ordered | Int32 | All the quantities ordered are stored here as Integer Numbers. |
| 4 | Price Each | Int32 | All the prices of the product are stored datatype as Integer Numbers. |
| 5 | Ordered Date | object | All the ordered date are stored here, object as a data type. |
| 6 | Purchase Address | object | All the Purchase addresses are stored here, object as a data type. |
| 7 | Month | Int32 | Only months are stored here from Ordered Date as Integer Numbers. |
| 8 | City | object | All the cities from the purchase address are stored here, object as a data type. |
| 9 | Sales | Int32 | All the total sales of the product are stored datatype as Integer Numbers. |

*Table 1: Description of data frame data.*

**Description:** In the above table, description of the data frame and updated CSV file is elaborated in brief way.

## 3. Data Preparation

### 3.1. Write a Python program to merge data from each month into one CSV and read in the updated data frame.

```
In [1]: import pandas as pd
        import numpy as np
        import glob
```

```
In [2]: #defining the path to the CSV file.
        path = 'C:/Users/Acer/Desktop/sem-2/smart data recovery/2201564_BhuwaniBikramNembang/*.csv'
```

```
In [3]: #get a list of all csv file, read it into a dataframe
        files = glob.glob(path)
```

```
In [4]: #Initialize an empty list to store the dataframes.
        dataframe = []
```

```
In [5]: #Loop through each csv file, read it into a dataframes
        for file in files:
            df = pd.read_csv(file)
            dataframe.append(df)
```

```
In [6]: # Concatenate all dataframes into a single dataframe
```

*Figure 4: importing file data and merged in to one csv file.*

```
In [6]: # Concatenate all dataframes into a single dataframe
        merged_df = pd.concat(dataframe)
```

```
In [7]: #merging data into one csv file name as newrecord.
        merged_df.to_csv('C:/Users/Acer/Desktop/sem-2/smart data recovery/2201564_BhuwaniBikramNembang/newrecord.csv')
```

```
In [8]: merged_df
```

Out[8]:

|  | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address |
|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2.0 | 11.95 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 |
| 1 | NaN | NaN | NaN | NaN | NaN | NaN |
| 2 | 176559.0 | Bose SoundSport Headphones | 1.0 | 99.99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 |
| 3 | 176560.0 | Google Phone | 1.0 | 600.00 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| 4 | 176560.0 | Wired Headphones | 1.0 | 11.99 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3.0 | 2.99 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 |
| 11682 | 259354.0 | iPhone | 1.0 | 700.00 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 |
| 11683 | 259355.0 | iPhone | 1.0 | 700.00 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1.0 | 379.99 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 |
| 11685 | 259357.0 | USB-C Charging Cable | 1.0 | 11.95 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 |

186850 rows × 6 columns

*Figure 5: printing the data from the updated .csv file.*

m-2 > smart data recovery > 2201564_BhuwaniBikramNembang



*Figure 6: proved as new csv file is generated from the overall csv data.*

- ▪ **Description:** In these above figures, I merge data from each month into one CSV and read in the updated data frame. At first, I combined data from all the csv files into one data frame.

## 3.2. Write a Python program to remove NAN missing values from the updated data frame.

```
In [9]:   # dropping all NaN values from dataframe.
          merged_df = merged_df.dropna()

In [10]:  merged_df

Out[10]:
```

|  | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address |
|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2.0 | 11.95 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 |
| 2 | 176559.0 | Bose SoundSport Headphones | 1.0 | 99.99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 |
| 3 | 176560.0 | Google Phone | 1.0 | 600.00 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| 4 | 176560.0 | Wired Headphones | 1.0 | 11.99 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| 5 | 176561.0 | Wired Headphones | 1.0 | 11.99 | 4/30/2019 9:27 | 333 8th St, Los Angeles, CA 90001 |
| ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3.0 | 2.99 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 |
| 11682 | 259354.0 | iPhone | 1.0 | 700.00 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 |
| 11683 | 259355.0 | iPhone | 1.0 | 700.00 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1.0 | 379.99 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 |
| 11685 | 259357.0 | USB-C Charging Cable | 1.0 | 11.95 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 |

185950 rows × 6 columns

*Figure 7: removing NAN missing values from the updated data frame.*

- **Description:** In the above figure, all the NaN values are dropped and presented in a tabular form.

## 3.3. Write a Python program to convert Quantity Ordered and Price Each to numeric.

```
In [11]:  #changing datatype as int of quantity ordered and price each
          merged_df["Quantity Ordered"]=merged_df["Quantity Ordered"].astype(int)
          merged_df["Price Each"]=merged_df["Price Each"].astype(int)

In [12]:  print(merged_df.dtypes)

          Order ID            float64
          Product              object
          Quantity Ordered      int32
          Price Each            int32
          Order Date           object
          Purchase Address     object
          dtype: object

In [13]:  merged_df

Out[13]:
```

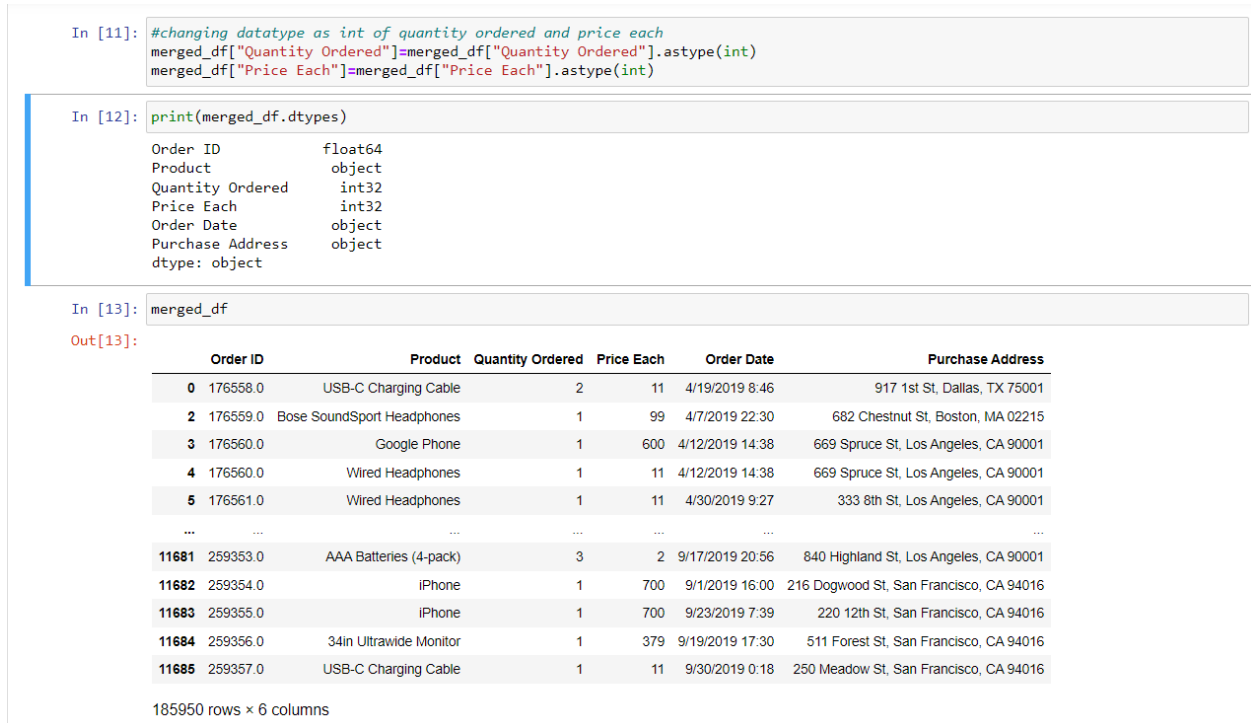|  | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address |
|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2 | 11 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 |
| 2 | 176559.0 | Bose SoundSport Headphones | 1 | 99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 |
| 3 | 176560.0 | Google Phone | 1 | 600 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| 4 | 176560.0 | Wired Headphones | 1 | 11 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 |
| 5 | 176561.0 | Wired Headphones | 1 | 11 | 4/30/2019 9:27 | 333 8th St, Los Angeles, CA 90001 |
| ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3 | 2 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 |
| 11682 | 259354.0 | iPhone | 1 | 700 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 |
| 11683 | 259355.0 | iPhone | 1 | 700 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1 | 379 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 |
| 11685 | 259357.0 | USB-C Charging Cable | 1 | 11 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 |

185950 rows × 6 columns

*Figure 8: converting quantity ordered and price each to numeric.*

▪ **Description:** In the above figure, quantity ordered and price each datatype is converted to int and displayed.

2201564 Bhuwani Bikram Nembang

## 3.4. Create a new column named Month from Ordered Date of the updated data frame and convert it to integer as data type.

```
In [14]: # adding new column month from Orederd date
         merged_df['Month'] = merged_df['Order Date'].str.split('/').str[0]
         merged_df['Month']= merged_df['Month'].astype(int)
```

```
In [15]: print(merged_df.dtypes)

         Order ID          float64
         Product            object
         Quantity Ordered    int32
         Price Each          int32
         Order Date         object
         Purchase Address   object
         Month               int32
         dtype: object
```

*Figure 9: creating a new column named month from Ordered date of the update data frame.*

```
In [16]: merged_df

Out[16]:
```

| | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address | Month |
|---|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2 | 11 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 | 4 |
| 2 | 176559.0 | Bose SoundSport Headphones | 1 | 99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 | 4 |
| 3 | 176560.0 | Google Phone | 1 | 600 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 |
| 4 | 176560.0 | Wired Headphones | 1 | 11 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 |
| 5 | 176561.0 | Wired Headphones | 1 | 11 | 4/30/2019 9:27 | 333 8th St, Los Angeles, CA 90001 | 4 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3 | 2 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 | 9 |
| 11682 | 259354.0 | iPhone | 1 | 700 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 | 9 |
| 11683 | 259355.0 | iPhone | 1 | 700 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 | 9 |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1 | 379 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 | 9 |
| 11685 | 259357.0 | USB-C Charging Cable | 1 | 11 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 | 9 |

185950 rows × 7 columns

*Figure 10: column created and added in updated data frame.*

- **Description:** In the above figure, a new column month from ordered date of the updated data frame and converted into integer as data type.

## 3.5. Create a new column named City from Purchase Address based on the value in the updated data frame.



```
In [18]: # creating a new column from purchase address.
         merged_df['City'] = merged_df['Purchase Address'].str.split(',').str[1]

In [19]: merged_df
```

Out[19]:

| | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address | Month | City |
|---|---|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2 | 11 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 | 4 | Dallas |
| 2 | 176559.0 | Bose SoundSport Headphones | 1 | 99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 | 4 | Boston |
| 3 | 176560.0 | Google Phone | 1 | 600 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 | Los Angeles |
| 4 | 176560.0 | Wired Headphones | 1 | 11 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 | Los Angeles |
| 5 | 176561.0 | Wired Headphones | 1 | 11 | 4/30/2019 9:27 | 333 8th St, Los Angeles, CA 90001 | 4 | Los Angeles |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3 | 2 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 | 9 | Los Angeles |
| 11682 | 259354.0 | iPhone | 1 | 700 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 | 9 | San Francisco |
| 11683 | 259355.0 | iPhone | 1 | 700 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 | 9 | San Francisco |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1 | 379 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 | 9 | San Francisco |
| 11685 | 259357.0 | USB-C Charging Cable | 1 | 11 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 | 9 | San Francisco |

185950 rows × 8 columns

*Figure 11: creating a new column named City from Purchase Address based on the value in the updated data frame.*

- **Description:** In the above figure, a new column named City from purchase address Is created based on the values in the updated data frame using split function.

## 4. Data analysis

## 4.1. Write a Python program to show summary statistics of sum, mean, standard deviation, skewness, and kurtosis of any chosen variable.

```
In [20]:  # creating a new local variables for summary statistics.
          Quantity_Ordered_sum =merged_df['Quantity Ordered'].sum()
          Quantity_Ordered_mean =merged_df['Quantity Ordered'].mean()
          Quantity_Ordered_std =merged_df['Quantity Ordered'].std()
          Quantity_Ordered_skew =merged_df['Quantity Ordered'].skew()
          Quantity_Ordered_kurt =merged_df['Quantity Ordered'].kurt()
```

```
In [21]:  # printing all summary statistics of sum, mean, standard deviation, skewness, and kurtosis of Quantity Ordered.
          print("summary statistics of sum, mean, standard deviation, skewness, and kurtosis of Quantity Ordered")
          print('')
          print('Sum:',Quantity_Ordered_sum)
          print("Kurtosis",Quantity_Ordered_skew)
          print("SKewness:",Quantity_Ordered_kurt)
          print("Mean:",Quantity_Ordered_mean,)
          print("Standard:",Quantity_Ordered_std)

          summary statistics of sum, mean, standard deviation, skewness, and kurtosis of Quantity Ordered

          Sum: 209079
          Kurtosis 4.833164172577953
          SKewness: 31.82048892027536
          Mean: 1.1243828986286637
          Standard: 0.44279262402849096
```

*Figure 12: summary statistics of sum, mean, std, skew, and kurtosis.*

- • **Description:** In the above figure, summary statistics of sum, mean, standard deviation, skewness, and kurtosis is displayed.

**.mean():** .mean() method is used to compute the arithmetic mean (average) of a set of numbers. This method is commonly used with Pandas Data Frames and Series objects to compute the mean of a column or row of data.

**.sum():** .sum() method is used to compute the sum of a set of numbers. This method is commonly used with Pandas Data Frames and Series objects to compute the sum of a column or row of data.

**.std():** .std() method is used to compute the standard deviation of a set of numbers. This method is commonly used with Pandas DataFrames and Series objects to compute the standard deviation of a column or row of data.

**.skew():** .skew() method is used to compute the skewness of a set of numbers. This method is commonly used with Pandas DataFrames and Series objects to compute the skewness of a column or row of data.

**.kurt():** .kurt() method is used to compute the kurtosis of a set of numbers. This method is commonly used with Pandas DataFrames and Series objects to compute the kurtosis of a column or row of data.

## 4.2. Write a Python program to calculate and show correlation of all variables.

```
In [22]: corr_matrix=merged_df.corr()
```

```
In [23]: print("Corredlation matrix:")
         print(corr_matrix)
```

```
Corredlation matrix:
                   Order ID  Quantity Ordered  Price Each
Order ID           1.000000          0.000702   -0.002861
Quantity Ordered   0.000702          1.000000   -0.148335
Price Each        -0.002861         -0.148335    1.000000
```

*Figure 13: calculating and showing correlation of all variables.*

```
In [24]: import seaborn as sns
         # Create a correlation matrix of Quantity Ordered and Price Each
         correlation = merged_df[['Quantity Ordered', 'Price Each']].corr()

         # Plot a heatmap of the correlation matrix
         sns.set_theme(style="whitegrid")
         heatmap = sns.heatmap(data=correlation, annot=True, cmap='RdYlBu_r', fmt='.4g')
```
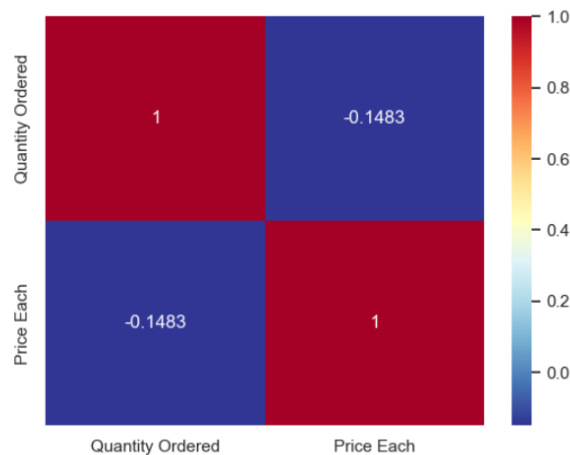


*Figure 14: heatmap of correlation matrix.*

- ▪ **Description:** In the above figures, correlation matrix is displayed using heat map.

**.corr():** .corr() method is used to compute the correlation between columns in a DataFrame.

2201564 Bhuwani Bikram Nembang

## 5. Data exploration

### 5.1. Which Month has the best sales? and how much was the earning in that month? Make a bar graph of sales as well.

```
In [25]: # adding  new column sales.
         merged_df['Sales'] = merged_df['Quantity Ordered']* merged_df['Price Each']
```

```
In [26]: merged_df
```

Out[26]:

| | Order ID | Product | Quantity Ordered | Price Each | Order Date | Purchase Address | Month | City | Sales |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 176558.0 | USB-C Charging Cable | 2 | 11 | 4/19/2019 8:46 | 917 1st St, Dallas, TX 75001 | 4 | Dallas | 22 |
| 2 | 176559.0 | Bose SoundSport Headphones | 1 | 99 | 4/7/2019 22:30 | 682 Chestnut St, Boston, MA 02215 | 4 | Boston | 99 |
| 3 | 176560.0 | Google Phone | 1 | 600 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 | Los Angeles | 600 |
| 4 | 176560.0 | Wired Headphones | 1 | 11 | 4/12/2019 14:38 | 669 Spruce St, Los Angeles, CA 90001 | 4 | Los Angeles | 11 |
| 5 | 176561.0 | Wired Headphones | 1 | 11 | 4/30/2019 9:27 | 333 8th St, Los Angeles, CA 90001 | 4 | Los Angeles | 11 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 11681 | 259353.0 | AAA Batteries (4-pack) | 3 | 2 | 9/17/2019 20:56 | 840 Highland St, Los Angeles, CA 90001 | 9 | Los Angeles | 6 |
| 11682 | 259354.0 | iPhone | 1 | 700 | 9/1/2019 16:00 | 216 Dogwood St, San Francisco, CA 94016 | 9 | San Francisco | 700 |
| 11683 | 259355.0 | iPhone | 1 | 700 | 9/23/2019 7:39 | 220 12th St, San Francisco, CA 94016 | 9 | San Francisco | 700 |
| 11684 | 259356.0 | 34in Ultrawide Monitor | 1 | 379 | 9/19/2019 17:30 | 511 Forest St, San Francisco, CA 94016 | 9 | San Francisco | 379 |
| 11685 | 259357.0 | USB-C Charging Cable | 1 | 11 | 9/30/2019 0:18 | 250 Meadow St, San Francisco, CA 94016 | 9 | San Francisco | 11 |

185950 rows × 9 columns

*Figure 15: the best sales month with a price of each.*

- **Description:** In the above figure, Sales column is created with the value from Quantity Ordered and Price Each which is multiply and set values to Sales.

```
In [27]: monthly = merged_df.groupby('Month')['Sales'].sum()
         print("the monthly sales is:")
         monthly

         the monthly sales is:
```

```
Out[27]: Month
         1     1813956
         10    3719205
         11    3184394
         12    4591824
         2     2191696
         3     2794068
         4     3374951
         5     3138287
         6     2566187
         7     2635443
         8     2234194
         9     2087435
         Name: Sales, dtype: int32
```

*Figure 16: Highest sale month and max ordered quantity.*

- **Description:** In the above figure, all the months and their sales are displayed.

  The best month from by the sales is 12 and the earning of the month is 4591824.

```
In [28]: from matplotlib import pyplot as plt
         # calculate the total sales for each month
         monthly = merged_df.groupby('Month')['Sales'].sum()

         # sort the monthly sales by month in ascending order
         monthly = monthly.sort_index()

         # create a bar graph of the monthly sales
         plt.figure(figsize = (10,5))
         plt.bar(monthly.index, monthly.values)

         # add a title to the plot
         plt.title("Total Sales by Month")

         # show the plot
         plt.show()
```
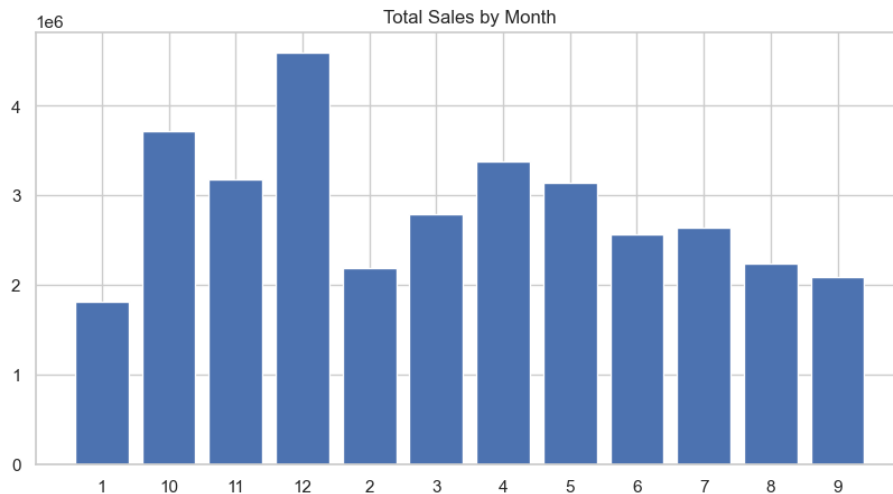


*Figure 17: best sales month in bar graph.*

- ▪ **Description:** In the above figure, total sales of the month is displayed in bar Graph.

**plt.figure():** plt.figure() function is used to create a new figure for plotting. A figure is a container that holds all the subplots, axes, images, and other elements of a plot.

**plt.bar():** plt.bar() function is used to create a bar chart. A bar chart is a chart that represents categorical data with rectangular bars, where the height or length of the bar represents the value of the category.

**plt.title():** plt.title() function is used to add a title to a plot. The title is used to describe the contents of the plot or to provide a brief summary of the data being presented.

**plt.show():** plt.show() function is used to display a plot that has been created using the other Matplotlib functions.

## 5.2. Which city has sold the highest product?

```
In [29]: # which city has highest product
         merged_df['City'].value_counts()

Out[29]:  San Francisco    44732
          Los Angeles      29605
          New York City    24876
          Boston           19934
          Atlanta          14881
          Dallas           14820
          Seattle          14732
          Portland         12465
          Austin            9905
          Name: City, dtype: int64
```
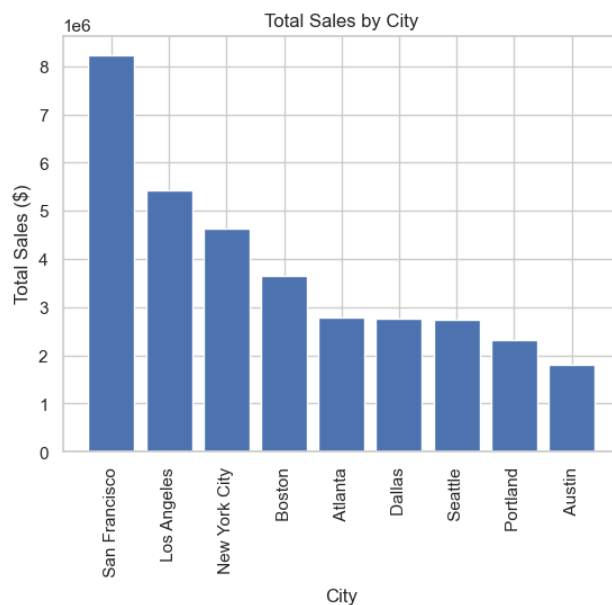
*Figure 18: city with the highest product sold.*

▪ **Description***:* In the above figure, all the cities with their sales is displayed.

**.valuecounts():** .value_counts() method is a Pandas function that is used to count the unique values in a Series object. It returns a new Pandas Series object containing the counts of each unique value in the original Series.

```
In [30]: city_sales = merged_df.groupby('City')['Sales'].sum().sort_values(ascending=False)

         # Create a bar plot of the total sales by city
         plt.bar(city_sales.index, city_sales.values)
         plt.title('Total Sales by City')
         plt.xlabel('City')
         plt.ylabel('Total Sales ($)')
         plt.xticks(rotation=90)
         plt.figure(figsize = (10,5))
         plt.show()
```



```
<Figure size 1000x500 with 0 Axes>
```

*Figure 19: sales with best city in bar graph.*

- **Description:** In the above figure, the city with the highest product sold is San Francisco which is also displayed in a bar graph.

## 5.3. Which product was sold the most in overall? Illustrate it through bar graph.

```
In [31]: # Which product was sold the most in overall? Illustrate it through bar graph.
         product = merged_df.groupby('Product')['Sales'].sum().sort_values(ascending=False)
         plt.bar(product.index, product.values)
         plt.title('Graph')
         plt.xlabel('Name of product')
         plt.ylabel('Total Sales ($)')
         plt.xticks(rotation=90)
         plt.figure(figsize = (10,8))
         plt.show()
```
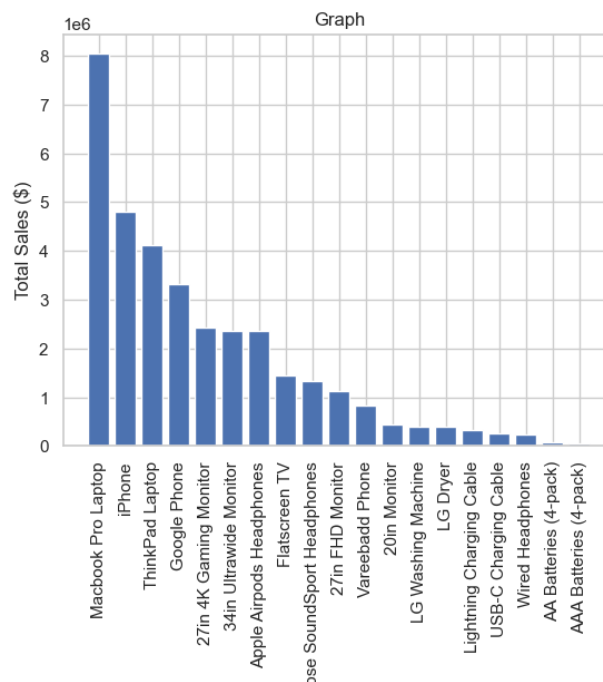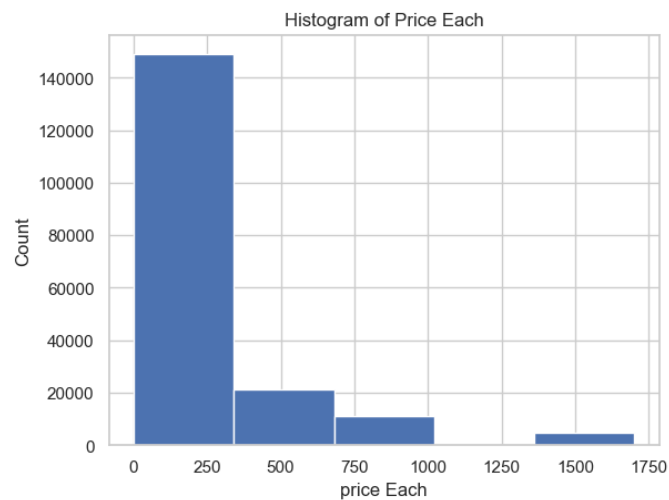


*Figure 20: best sale product in overall.*

▪ **Description:** In the above figure, it is displayed that MacBook Pro Laptop is the most sold product overall from the data frame which is also displayed in Bar graph.

2201564 Bhuwani Bikram Nembang

**5.4. Write a Python program to show histogram plot of any chosen variables. Use proper labels in the graph.**

```
In [32]: #Python program to show histogram plot of any chosen variables. Use proper labels in the grap
         plt.hist(merged_df['Price Each'], bins=5)
         plt.title('Histogram of Price Each')
         plt.xlabel("price Each")
         plt.ylabel('Count')
         plt.figure(figsize = (10,5))
         plt.show()
```



```
<Figure size 1000x500 with 0 Axes>
```

*Figure 21: histogram plot of proper labels in graph.*

▪ **Description:** In the above figure, histogram plot of price each is variable is displayed which is also in a bar graph.

2201564 Bhuwani Bikram Nembang

## Conclusion

In conclusion, the sales data analysis of ABC company using Python Pandas has provided valuable insights into the business's performance. By cleaning, transforming, and visualizing the sales data using Pandas, we were able to identify patterns and trends that can help ABC company make data-driven decisions.

We started by importing the sales data into a Pandas Data Frame using the read_csv() function, then cleaned and transformed the data by dropping missing values, merging datasets, and creating new columns.

Next, we performed exploratory data analysis by visualizing the sales data using Pandas plotting functions such as plot (), bar (), and plot.show (). These visualizations helped us identify which products were selling the most, which sales channels were the most effective, and which regions were performing well.

Finally, we used Pandas functions such as groupby() to aggregate the sales data and calculate various performance metrics such as total revenue, average order value, and customer retention rate.

Through our analysis, we identified opportunities for ABC company to increase revenue by focusing on high-performing products, optimizing sales channels, and targeting underperforming regions. We also identified areas for improvement, such as reducing customer churn and increasing the average order value.

Overall, the skills and techniques learned in this coursework can be applied to a wide range of sales data analysis tasks and can help ABC company and other businesses gain a competitive advantage by making data-driven decisions.

## References

Contributer, T. (1999-2023). *TechTarget*. Retrieved from TechTarget:
https://www.techtarget.com/whatis/definition/data-set

Jain, K. (2023, feb 27). *favtutor*. Retrieved from favtutor: https://favtutor.com/blogs/glob-
python

Meghna, K. (2023, mar 22). *GeekforGeeks*. Retrieved from GeekforGeeks:
https://www.geeksforgeeks.org/python-introduction-matplotlib/

2201564 Bhuwani Bikram Nembang