

## Data Analytics Problems on Sales\_Data

### 1. Import necessary libraries for sales Data Analysis

```
In [3]: import pandas as pd
import os
```

### 2.Concatenate each month sale Data ([https://github.com/svkarthik86/Assignment/tree/main/Sales\\_Data](https://github.com/svkarthik86/Assignment/tree/main/Sales_Data)) into one dataframe and save the dataframe to annual\_sale.csvs

```
In [4]: os.listdir('C:\\Users\\user\\Sales_Data')

Out[4]: ['annual_sale.csv',
'Sales_April_2019.csv',
'Sales_August_2019.csv',
'Sales_December_2019.csv',
'Sales_February_2019.csv',
'Sales_January_2019.csv',
'Sales_July_2019.csv',
'Sales_June_2019.csv',
'Sales_March_2019.csv',
'Sales_May_2019.csv',
'Sales_November_2019.csv',
'Sales_October_2019.csv',
'Sales_September_2019.csv']

In [24]: path=r"C:\Users\user\Sales_Data\\"
files=[s for s in os.listdir('C:\\Users\\user\\Sales_Data')]
sales_data=pd.DataFrame()
for s in files:
    annual_sale=pd.read_csv(path+s)
    annual_sale=pd.concat([sales_data,annual_sale])
    annual_sale.to_csv("annual_sale.csv",index=False)
    annual_sale
```

	1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
	2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
	3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
	4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301
	...		...	...		...	...
	11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
	11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
	11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
	11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
	11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016
11686 rows × 6 columns							

11686 rows × 6 columns

### 3.Read the annual\_sale.csv fom current working directory and store it as annual\_sale dataframe , and display the first five row of the dataframe

```
In [25]: annual_sale.head(5)

Out[25]:
```

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301

### 4.show the metadata information of the annual\_sale data frame and check if data is missing or not, if yes How many data are missing.

```
In [26]: print(annual_sale.info())
annual_sale.isna().sum().sum()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11686 entries, 0 to 11685
Data columns (total 6 columns):
# Column Non-Null Count Dtype
---
0 Order ID 11646 non-null object
1 Product 11646 non-null object
2 Quantity Ordered 11646 non-null object
3 Price Each 11646 non-null object
4 Order Date 11646 non-null object
5 Purchase Address 11646 non-null object
dtypes: object(6)
memory usage: 547.9+ KB
None
248

Out[26]:
```

### 5.Clean up the data!

#### 5.1 Verify all the column names are in a valid format, if any space between the column name then rename the column names

example: Product ID as Product\_ID

```
In [28]: annual_sale.columns=[i.replace(" ", "_") for i in annual_sale.columns]
annual_sale
```

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301
...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016

11686 rows × 6 columns

#### 5.2 check the isnan is present in dataframe, if there is nan is present remove the nan using dropna() function

```
In [29]: annual_sale.dropna(inplace=True)

Out[29]:
```

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301
...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016

11646 rows × 6 columns

#### 5.3 Find The duplicated data present in the data frame, and remove the duplicated data from the dataframe

```
In [33]: annual_sale=annual_sale[~(annual_sale.Price_Each=="Price_Each")]
annual_sale
```

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301
...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016

11629 rows × 6 columns

```
In [34]: annual_sale.drop_duplicates(inplace=True)
annual_sale
```

C:\Users\user\AppData\Local\Temp\ipykernel\_1580\572220841.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

annual\_sale.drop\_duplicates(inplace=True)

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301
...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016

11611 rows × 6 columns

### 6.memory\_usage

check memory\_usage of Product column , type cast the Product column as "category" type and then check memmory\_usage compare the memory utlization and how much percentage effectively reduce the storage space?

```
In [35]: annual_sale.Product.memory_usage()
```

Out[35]: 185776

```
In [36]: annual_sale.Product=annual_sale.Product.astype("category")
annual_sale.Product.memory_usage()
```

C:\Users\user\AppData\Local\Temp\ipykernel\_1580\1369973099.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

annual\_sale.Product=annual\_sale.Product.astype("category")

Out[36]: 105297

### 7.Create and add a new column

Add month column to annual\_sale DataFrame object from 'Order Date' column using Series.str method

annual\_sale["month"]

```
In [44]: annual_sale['month']=annual_sale.Order_Date.str[:2]
annual_sale
```

C:\Users\user\AppData\Local\Temp\ipykernel\_1580\2994047926.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

annual\_sale['month']=annual\_sale.Order\_Date.str[:2]

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address	month	(month, date)
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001	09	09
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301	09	09
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001	09	09
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101	09	09
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301	09	09
...	...	...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001	09	NaN
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016	09	NaN
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016	09	NaN
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016	09	NaN
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016	09	NaN

11611 rows × 8 columns

#### 7.1 converts the datatype of month column as int using astype('int32')

```
In [45]: annual_sale.month=annual_sale.month.astype('int32')
annual_sale
```

C:\Users\user\AppData\Local\Temp\ipykernel\_1580\2799170811.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

annual\_sale.month=annual\_sale.month.astype('int32')

	Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address	month	(month, date)
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001	9	09
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301	9	09
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001	9	09
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101	9	09
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301	9	09
...	...	...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001	9	NaN
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016	9	NaN
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016	9	NaN
11684	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016	9	NaN
11685	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016	9	NaN

11611 rows × 8 columns

#### 7.2 Add sale column to annual\_sale DataFrame object using following calculation

sales = Quantity\_Ordered \* Price\_Each

```
In [46]: annual_sale['sale']=annual_sale.Quantity_Ordered.astype(float)*annual_sale.Price_Each.astype(float)
annual_sale
```

C:\Users\user\AppData\Local\Temp\ipykernel\_1580\36499044914.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

annual\_sale['sale']=annual\_sale.Quantity\_Ordered.astype(float)\*annual\_sale.Price\_Each.astype(float)

annual_sale['sale']=annual_sale.Quantity_Ordered.astype(float)*annual_sale.Price_Each.astype(float)									
Order_ID	Product	Quantity_Ordered	Price_Each	Order_Date	Purchase_Address	month	(month, date)	sale	
0	248151	AA Batteries (4-pack)	4	3.84	09/17/19 14:44	380 North St, Los Angeles, CA 90001	9	09	15.36
1	248152	USB-C Charging Cable	2	11.95	09/29/19 10:19	511 8th St, Austin, TX 73301	9	09	23.90
2	248153	USB-C Charging Cable	1	11.95	09/16/19 17:48	151 Johnson St, Los Angeles, CA 90001	9	09	11.95
3	248154	27in FHD Monitor	1	149.99	09/27/19 07:52	355 Hickory St, Seattle, WA 98101	9	09	149.99
4	248155	USB-C Charging Cable	1	11.95	09/01/19 19:03	125 5th St, Atlanta, GA 30301	9	09	11.95
...	...	...	...	...	...	...	...	...	...
11681	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001	9	NaN	8.97
11682	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016	9	NaN	700.00
11683	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016	9	NaN	700.00
11684	260266	24in U-Brasside Monitor	1	230.00	09/01/19 13:30	511 Export St, San Francisco, CA 94016	9	NaN	230.00