

```
import numpy as np
import os
import matplotlib as mpl
import matplotlib.pyplot as plt
import pandas as pd
```

```
#merging all the csv files into 1 file
```

```
path = "/Users/muditkant/Downloads/Pandas-Data-Science-Tasks-master/SalesAnalysis/Sales_Data"
files = [file for file in os.listdir(path) if not file.startswith('.')] # ignore hidden files

all_months_data = pd.DataFrame()

for file in files:
    current_data = pd.read_csv(path+"/"+file)
    all_months_data = pd.concat([all_months_data, current_data])

all_months_data.to_csv("all_data.csv", index=False)
```

```
df = pd.read_csv("/Users/muditkant/Downloads/Pandas-Data-Science-Tasks-master/SalesAnalysis/Output/all_data.csv")
df.head()
```

2	176559	Bose SoundSport Headphones	1	99.99	04/07/19 22:30	682 Chestnut St, Boston, MA 02215
3	176560	Google Phone	1	600	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001
4	176560	Wired Headphones	1	11.99	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001

#checking for null values

```
#checking for null values
```

```
df.isnull().values.any()
```

```
True
```

```
df = df.dropna()
```

```
#checking for null values
```

```
df.isnull().values.any()
```

```
False
```

```
temp = df[df["Order Date"].str[:2] == "0r"]
temp.head()
```

Order ID	Product	Quantity Ordered	Price Each	Order Date
0	176558 USB-C Charging Cable	2	11.95	04/19/19 08:46
1	176559 Bose SoundSport Headphones	1	99.99	04/07/19 22:30
2	176560 Google Phone	1	600	04/12/19 14:38
4	176560 Wired Headphones	1	11.99	04/12/19 14:38

```
#order date consists of gibberish data
#removing gibberish values
```

```
df = df[df["Order Date"].str[:2] != "0r"]
df.head()
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	
0	176558	USB-C Charging Cable	2	11.95	04/19/19 08:46	917 1st St, Dallas, TX 75001
2	176559	Bose SoundSport Headphones	1	99.99	04/07/19 22:30	682 Chestnut St, Boston, MA 02215
3	176560	Google Phone	1	600	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001
4	176560	Wired Headphones	1	11.99	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001
5	176561	Wired Headphones	1	11.99	04/30/19 09:27	333 8th St, Los Angeles, CA 90001

#converting month values into int

```
df["Month"] = df["Order Date"].str[:2]
df.head()
```

2	176558	USB-C Charging Cable	2	11.95	04/19/19 08:46	917 1st St, Dallas, TX 75001	4
2	176559	Bose SoundSport Headphones	1	99.99	04/07/19 22:30	682 Chestnut St, Boston, MA 02215	4
3	176560	Google Phone	1	600	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001	4
4	176560	Wired Headphones	1	11.99	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001	4
5	176561	Wired Headphones	1	11.99	04/30/19 09:27	333 8th St, Los Angeles, CA 90001	4

Q1: Sales associated with each order

#converting Quantity Ordered and Price Each into numeric int

```
#converting month values into int
```

```
df["Month"] = df["Month"].astype('int32')
df.head()
```

2	176559	Bose SoundSport Headphones	1	99.99	04/07/19 22:30	682 Chestnut St, Boston, MA 02215	4
3	176560	Google Phone	1	600.00	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001	4
4	176560	Wired Headphones	1	11.99	04/12/19 14:38	669 Spruce St, Los Angeles, CA 90001	4
5	176561	Wired Headphones	1	11.99	04/30/19 09:27	333 8th St, Los Angeles, CA 90001	4

#Best month for sales and revenue generated

```
graph = df.groupby("Month").sum()
graph
```

```
# Q1: Sales associated with each order
```

```
#converting Quantity Ordered and Price Each into numeric int
```

```
df["Quantity Ordered"] = pd.to_numeric(df["Quantity Ordered"])
df["Price Each"] = pd.to_numeric(df["Price Each"])
```

```
df["Sales"] = df["Quantity Ordered"] * df["Price Each"]
df.head()
```

```
#visualizing using matplotlib

months = range(1,13)
plt.bar(months,graph["Sales"])
plt.xticks(months)
```

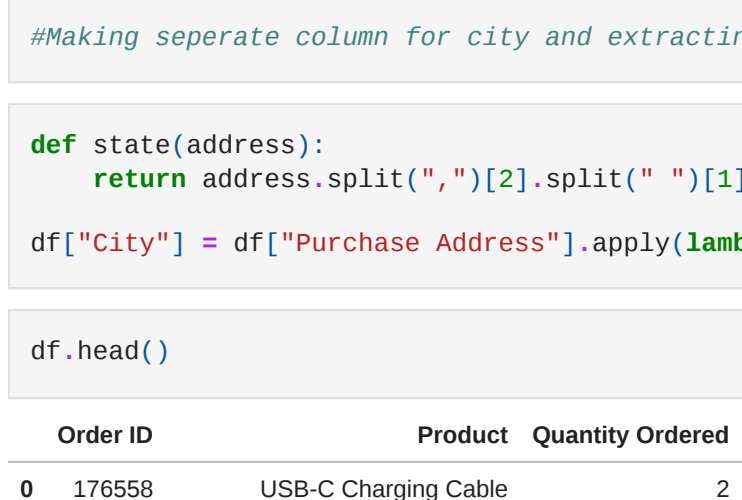
```
#Best month for sales and revenue generated
```

```
graph = df.groupby("Month").sum()
graph
```

Month	Quantity Ordered	Price Each	Sales
1	10903	1811768.38	182256.73
2	13449	2188884.72	220022.42
3	17005	2791207.83	2807100.38
4	20558	3367671.02	3390670.24
5	18667	3135125.13	3152606.75
6	15253	2562025.61	2577802.26
7	16072	2632539.56	264775.76
8	13448	2230345.42	2244467.88
9	13109	2084992.09	2097660.13
10	22703	3715554.83	3736726.88
11	19798	3180600.68	319603.20
12	28114	4588415.41	4613443.34

```
#visualizing using matplotlib
```

```
months = range(1,13)
plt.bar(months,graph["Sales"])
plt.xticks(months)
plt.ylabel('Sales in $ ->')
plt.xlabel('Months ->')
plt.show()
```



```
#Which city has best sales
```

```
df.head()
```

```
result = df.groupby("City").sum()
result
```

	Quantity Ordered	Price Each	Month	Sales
City				
Atlanta GA	16602	2779908.20	104794	2795498.58
Austin TX	11153	1609873.61	69829	1819561.75
Boston MA	23258	3637469.77	1411129	3661642.01
Dallas TX	16730	2752627.82	104620	2767975.40
Los Angeles CA	33289	5421435.23	208326	5452570.80

```
#Making seperate column for city and extracting value
```

```
def state(address):
    return address.split(",")[2].split(" ")[1]

df["City"] = df["Purchase Address"].apply(lambda x: x.split(',')[2] + " " + state(x))
```

```
df.head()
```

```
plt.xlabel("name ->")
plt.show()
```

Category	Sales in \$
Category 1	300
Category 2	400
Category 3	300
Category 4	550
Category 5	450
Category 6	850

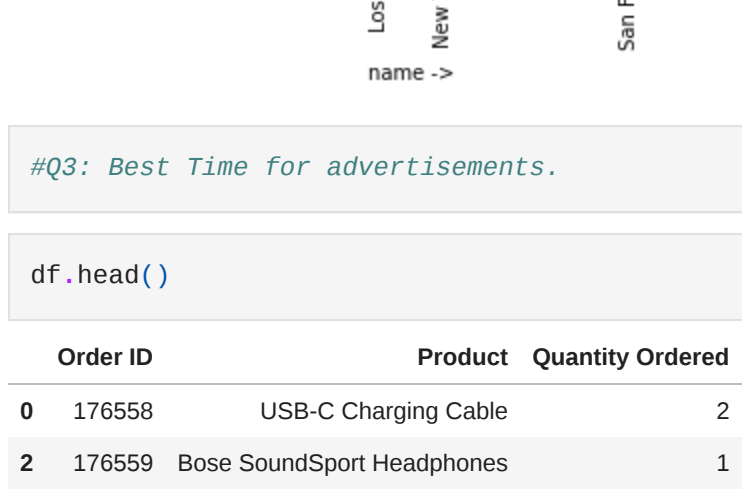
```
#Q2: Best sales in which city
```

```
result = df.groupby("City").sum()
result
```

City	Quantity Ordered	Price Each	Month	Sales
Atlanta GA	16602	2779908.20	104794	2795498.58
Austin TX	11153	1809073.61	69829	1819581.75
Boston MA	22528	3637409.77	141112	3661842.01
Dallas TX	16730	2752627.82	104620	2767975.40
Los Angeles CA	33289	5421435.23	208325	5452570.80
New York City NY	27932	4635370.83	175741	4664317.43
Portland ME	2750	447189.25	17144	449758.27
Portland OR	11303	1860568.22	70621	1870732.34
San Francisco CA	50239	8211461.74	315520	8262203.91
Seattle WA	16553	2733296.01	104941	2747755.48

```
#visualizing using matplotlib
```

```
cities = df["City"].unique()
plt.barcities,result["Sales"])
plt.xticks(cities,rotation = "vertical")
plt.ylabel('Sales in $ ->')
plt.xlabel("name ->")
plt.show()
```

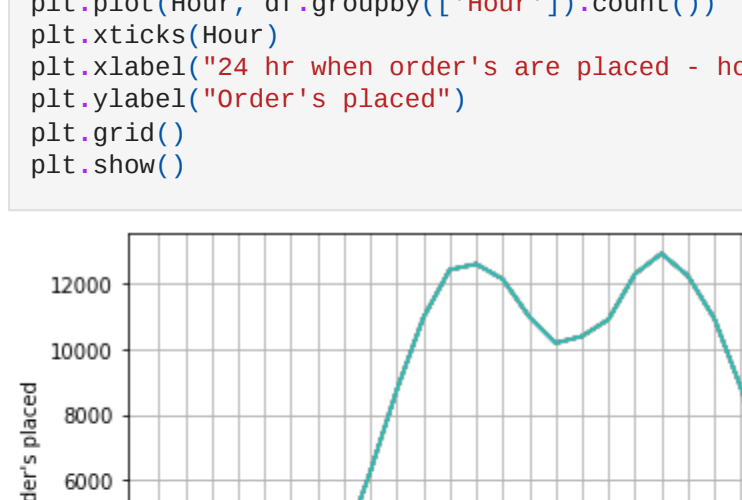


```
# In city colums it shows San Francisco CA bes sales
# While in visulization it shows, austin TX
```

```
## Need to search why this happened.
```

```
## X data and Y data needs to be in same order. That's why it's causing
```

```
cities = [city for city, df in df.groupby(['City'])]
plt.bar(cities,result["Sales"])
plt.xticks(cities,rotation = "vertical")
plt.ylabel('Sales in $ ->')
plt.xlabel("name ->")
plt.show()
```



```
#Q3: Best Time for advertisements.
```

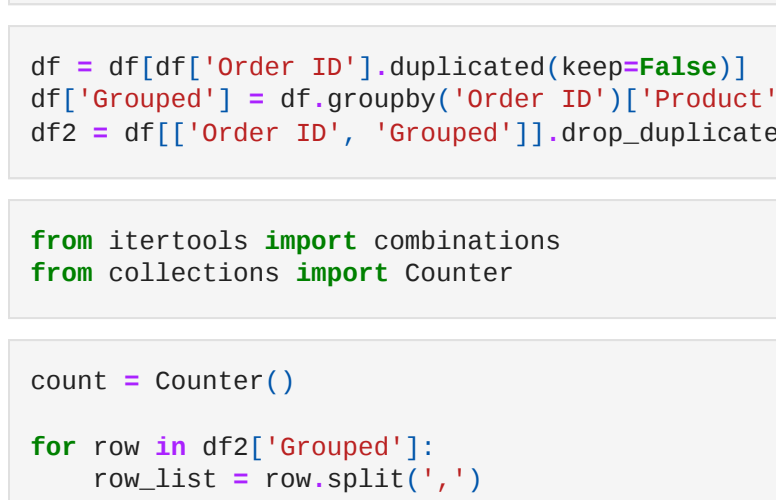
```
df.head()
```

df.head(20)

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	
0	176558	USB-C Charging Cable	2	11.95	2019-04-19 08:46:00	917 1st St, Dallas, TX 75001	4	23.90
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99
3	176560	Google Phone	1	600.00	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	600.00
4	176560	Wired Headphones	1	11.99	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	11.99
5	176561	Wired Headphones	1	11.99	2019-04-30 09:27:00	333 8th St, Los Angeles, CA 90001	4	11.99

```
#visualizing using matplotlib
```

```
Hour = [Hour for Hour, df in df.groupby(['Hour'])]
plt.plot(Hour, df.groupby(['Hour']).count())
plt.xticks(Hour)
plt.xlabel("24 hr when order's are placed - hours")
plt.ylabel("Order's placed")
plt.grid()
plt.show()
```



```
#delivering advertisements just before Peak hours of sale: 11 - 12 & 18 - 19 hrs
```

```
#Q4: What are the most often products sold in a groupof 2 or 3?
```

```
df.head(20)
```

```
Q5: What product was sold the most and why
```

```
product_group = df.groupby('Product')
quantity_ordered = product_group.sum()['Quantity Ordered']

keys = [pair for pair, df in product_group]
plt.bar(keys, quantity_ordered)
plt.xticks(keys, rotation='vertical', size=8)
plt.show()
```

Product Group	Quantity Ordered
USB-C Charging Cable	1
Bose SoundSport Headphones	1
Google Phone	1
Wired Headphones	1
Macbook Pro Laptop	1
AA Batteries (4-pack)	1
Lightning Charging Cable	1
Apple AirPods Headphones	1
AAA Batteries (4-pack)	1
USB-C Charging Cable	2

```
#order's have duplicate order ID:
```

```
df = df[df['Order ID'] duplicated(keep=False)]
df[["grouped"] = df.groupby("Order ID")["Product"].transform(lambda x: ','.join(x))
df2 = df[["Order ID", "grouped"]].drop_duplicates()
```

```
from itertools import combinations
from collections import Counter

count = Counter()

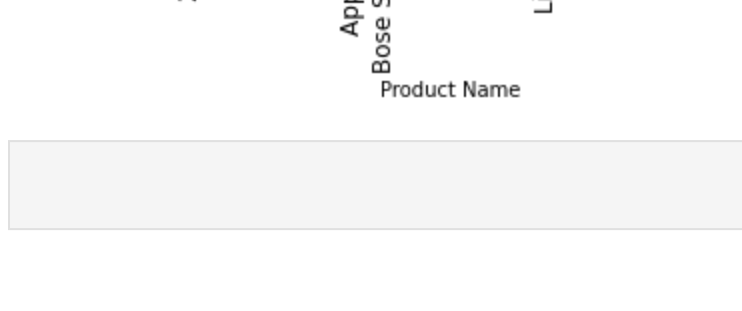
for row in df2[["grouped"]]:
    row_list = row.split(',')
    count.update(Counter(combinations(row_list, 1)))

for key, value in count.most_common(10):
    print(key, value)
```

```
('USB-C Charging Cable',) 2111
('iPhone',) 1867
('Lightning Charging Cable',) 1827
('Wired Headphones',) 1674
('Google Phone',) 1639
('Apple AirPods Headphones',) 929
('Bose SoundSport Headphones',) 820
('AAA Batteries (4-pack)',) 815
('AA Batteries (4-pack)',) 768
('Vareebaddi Phone',) 691
```

```
# Q5: What product was sold the most and why
```

```
product_group = df.groupby('Product')
quantity_ordered = product_group.sum()["Quantity Ordered"]
keys = pair for pair, df in product_group)
plt.bar(keys, quantity_ordered)
plt.xticks(keys, rotation='vertical', size=8)
plt.show()
```



```
prices = df.groupby('Product').mean()["Price Each"]
fig, ax1 = plt.subplots()

ax2 = ax1.twinx()
ax1.bar(keys, quantity_ordered, color='g')
ax2.plot(keys, prices, color='b')

ax1.set_xlabel('Product Name')
ax1.set_ylabel('Quantity Ordered
```