

Sales Analysis

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

Merging all the 12 months data into a single dataframe

```
In [2]: files = [file for file in os.listdir('Pandas-Data-Science-Tasks-master/Pandas-Data-Science-Tasks-master/SalesAnalysis/Sales_Data')]

all_months_data = pd.DataFrame()

for file in files:
    df = pd.read_csv('Pandas-Data-Science-Tasks-master/Pandas-Data-Science-Tasks-master/SalesAnalysis/Sales_Data/' + file)
    all_months_data = pd.concat([all_months_data,df])

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

Reading in updated dataframe

```
In [3]: all_data = pd.read_csv('all_data.csv')
all_data.head()
```

Out[3]:

	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
1	NaN	NaN	NaN	NaN	NaN	NaN
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

```
In [4]: nan_df = all_data[all_data.isna().any(axis=1)]
nan_df
```

Out[4]:

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
1	NaN	NaN	NaN	NaN	NaN	NaN
356	NaN	NaN	NaN	NaN	NaN	NaN
735	NaN	NaN	NaN	NaN	NaN	NaN
1433	NaN	NaN	NaN	NaN	NaN	NaN
1553	NaN	NaN	NaN	NaN	NaN	NaN
...
185176	NaN	NaN	NaN	NaN	NaN	NaN
185438	NaN	NaN	NaN	NaN	NaN	NaN
186042	NaN	NaN	NaN	NaN	NaN	NaN
186548	NaN	NaN	NaN	NaN	NaN	NaN
186826	NaN	NaN	NaN	NaN	NaN	NaN

545 rows × 6 columns

Cleaning Up the data

Dropping Nan

```
In [5]: all_data = all_data.dropna(how='all')
all_data.head()
```

Out[5]:

	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

Finding errors in dataframe 'Or'

```
In [6]: temp_df = all_data[all_data['Order Date'].str[0:2] == 'Or']
temp_df
```

Out[6]:

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	
	519	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	1149	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	1155	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	2878	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	2893	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address

	185164	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	185551	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	186563	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	186632	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
	186738	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address

355 rows × 6 columns

Removing of 'Or' error

In [7]:

```
all_data = all_data[all_data['Order Date'].str[0:2] != 'Or']
all_data
```

Out[7]:

	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

	186845	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001
	186846	259354	iPhone	1	700	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016
	186847	259355	iPhone	1	700	09/23/19 07:39	220 12th St, San Francisco, CA 94016
	186848	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016
	186849	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016

185950 rows × 6 columns

Converting columns to correct types

In [8]:

```
all_data['Quantity Ordered'] = pd.to_numeric(all_data['Quantity Ordered'])

all_data['Price Each'] = pd.to_numeric(all_data['Price Each'])
```

Adding additional colums

Task 2: Adding month column

In [9]:

```
all_data['Month']= all_data['Order Date'].str[0:2]
all_data['Month'] = all_data['Month'].astype('int32')
all_data
```

Out[9]:

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	
	0	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	600.00	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

	186845	259353	AAA Batteries (4-pack)	3	2.99	09/17/19 20:56	840 Highland St, Los Angeles, CA 90001	9
	186846	259354	iPhone	1	700.00	09/01/19 16:00	216 Dogwood St, San Francisco, CA 94016	9
	186847	259355	iPhone	1	700.00	09/23/19 07:39	220 12th St, San Francisco, CA 94016	9
	186848	259356	34in Ultrawide Monitor	1	379.99	09/19/19 17:30	511 Forest St, San Francisco, CA 94016	9
	186849	259357	USB-C Charging Cable	1	11.95	09/30/19 00:18	250 Meadow St, San Francisco, CA 94016	9

185950 rows × 7 columns

Task 3: Adding a Sales Column

```
In [10]: all_data['Sales'] = all_data['Quantity Ordered'] * all_data['Price Each']
all_data.head()
```

Out[10]:

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

Task 4: Add a City Column

```
In [11]: def city(address):
        return address.split(',')[1]

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

all_data['City Name'] = all_data['Purchase Address'].apply(lambda x: f"{city(x)} ({state(x)})")
all_data.head()
```

Out[11]:

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

Data Exploration

Question 1: What was the best month for sales? How much was earned that month?

```
In [12]: results = all_data.groupby('Month').sum()['Sales']
results
```

Out[12]:

Month

1	1822256.73
2	2202022.42
3	2807100.38
4	3390670.24
5	3152606.75
6	2577802.26
7	2647775.76
8	2244467.88
9	2097560.13
10	3736726.88
11	3199603.20
12	4613443.34

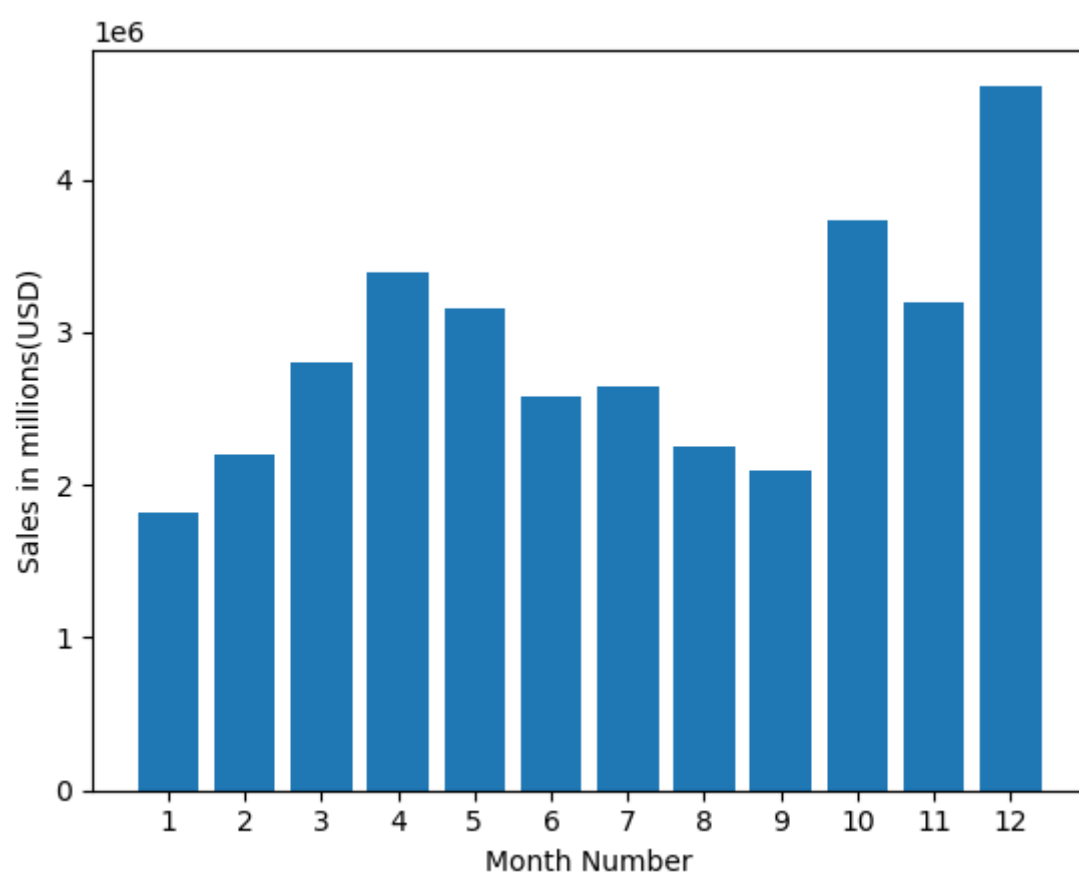
Name: Sales, dtype: float64

```
In [13]: import matplotlib.pyplot as plt
%matplotlib inline

months = range(1,13)
plt.bar(months,results)
plt.xticks(months)
plt.xlabel('Month Number')
plt.ylabel('Sales in millions(USD)')
```

Out[13]:

Text(0, 0.5, 'Sales in millions(USD)')



Question 2:What city had the highest number of sales?

```
In [14]: results = all_data.groupby('City Name').sum()['Sales']
results
```

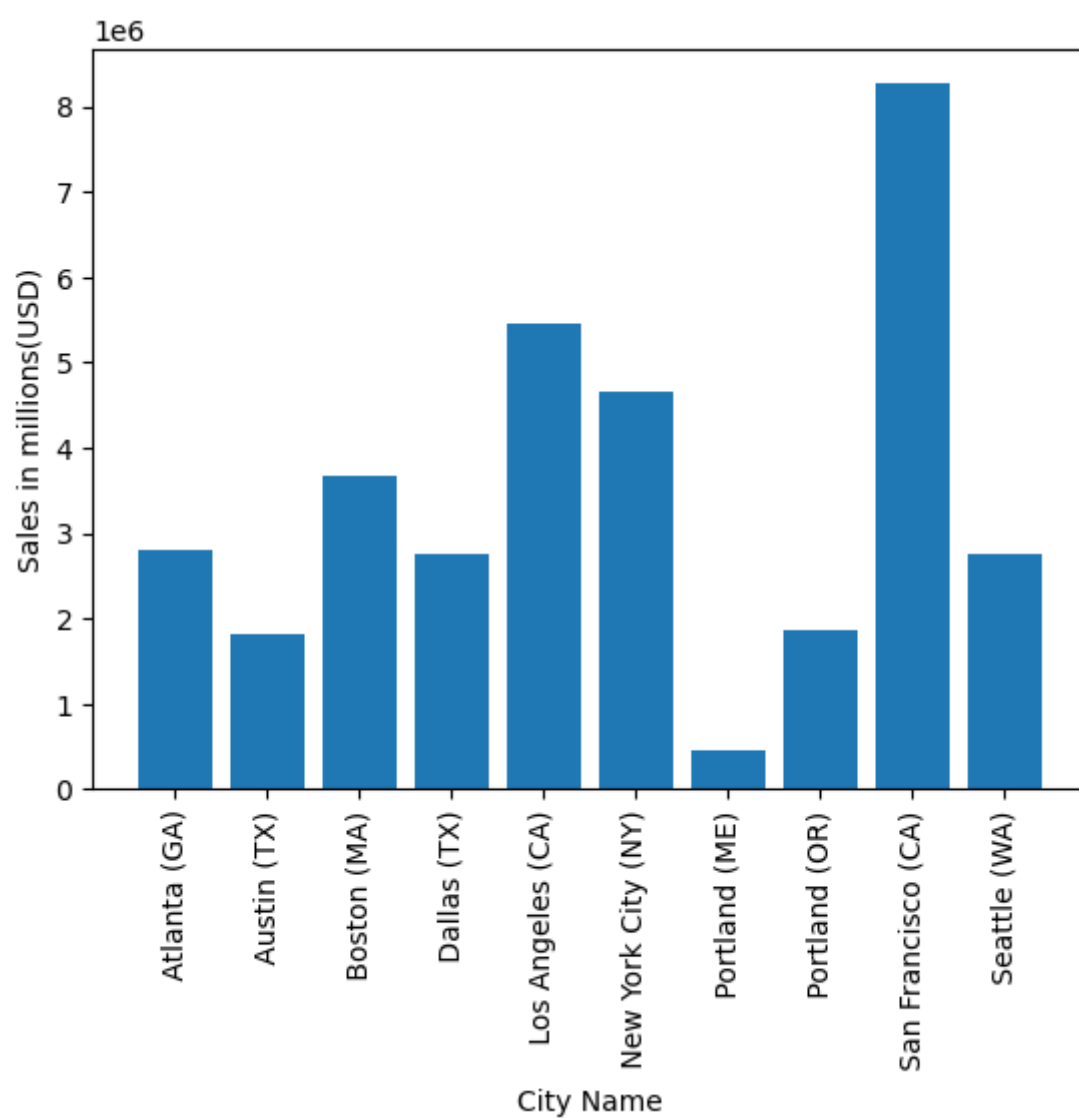
```
Out[14]: City Name
Atlanta (GA)      2795498.58
Austin (TX)       1819581.75
Boston (MA)       3661642.01
Dallas (TX)       2767975.40
Los Angeles (CA)  5452570.80
New York City (NY) 4664317.43
Portland (ME)      449758.27
Portland (OR)     1870732.34
San Francisco (CA) 8262203.91
Seattle (WA)      2747755.48
Name: Sales, dtype: float64
```

```
In [15]: import matplotlib.pyplot as plt
%matplotlib inline

# To match the city name correctly with the values of sales
cities = [city for city, df in all_data.groupby('City Name')]

plt.bar(cities,results)
plt.xticks(cities, rotation = 'vertical')
plt.xlabel('City Name')
plt.ylabel('Sales in millions(USD)')
```

```
Out[15]: Text(0, 0.5, 'Sales in millions(USD)')
```



Question 3: What time should we display advertisements to maximize likelihood of customers buying products?

In [16]:

```
all_data['Order Date'] = pd.to_datetime(all_data['Order Date'])

all_data['Hour'] = all_data['Order Date'].dt.hour
all_data['Minute'] = all_data['Order Date'].dt.minute

all_data.head()
```

Out[16]:

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City Name	Hour	Minute
0	176558	USB-C Charging Cable	2	11.95	2019-04-19 08:46:00	917 1st St, Dallas, TX 75001	4	23.90	Dallas (TX)	8	46
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99	Boston (MA)	22	30
3	176560	Google Phone	1	600.00	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	600.00	Los Angeles (CA)	14	38
4	176560	Wired Headphones	1	11.99	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	11.99	Los Angeles (CA)	14	38
5	176561	Wired Headphones	1	11.99	2019-04-30 09:27:00	333 8th St, Los Angeles, CA 90001	4	11.99	Los Angeles (CA)	9	27

In [17]:

```
results = all_data.groupby('Hour').count()
results
```

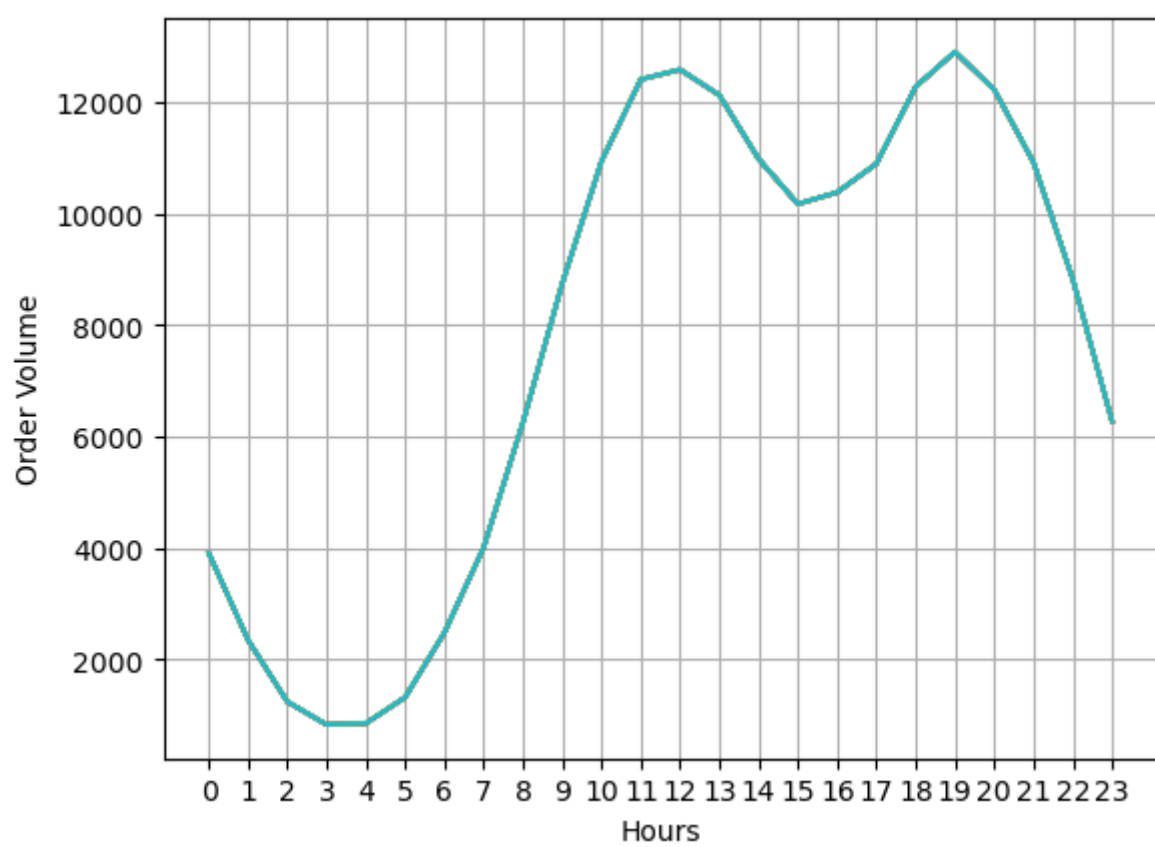
Out[17]:

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City Name	Minute
Hour										
0	3910	3910	3910	3910	3910	3910	3910	3910	3910	3910
1	2350	2350	2350	2350	2350	2350	2350	2350	2350	2350
2	1243	1243	1243	1243	1243	1243	1243	1243	1243	1243
3	831	831	831	831	831	831	831	831	831	831
4	854	854	854	854	854	854	854	854	854	854
5	1321	1321	1321	1321	1321	1321	1321	1321	1321	1321
6	2482	2482	2482	2482	2482	2482	2482	2482	2482	2482
7	4011	4011	4011	4011	4011	4011	4011	4011	4011	4011
8	6256	6256	6256	6256	6256	6256	6256	6256	6256	6256
9	8748	8748	8748	8748	8748	8748	8748	8748	8748	8748
10	10944	10944	10944	10944	10944	10944	10944	10944	10944	10944
11	12411	12411	12411	12411	12411	12411	12411	12411	12411	12411
12	12587	12587	12587	12587	12587	12587	12587	12587	12587	12587
13	12129	12129	12129	12129	12129	12129	12129	12129	12129	12129
14	10984	10984	10984	10984	10984	10984	10984	10984	10984	10984
15	10175	10175	10175	10175	10175	10175	10175	10175	10175	10175
16	10384	10384	10384	10384	10384	10384	10384	10384	10384	10384
17	10899	10899	10899	10899	10899	10899	10899	10899	10899	10899
18	12280	12280	12280	12280	12280	12280	12280	12280	12280	12280
19	12905	12905	12905	12905	12905	12905	12905	12905	12905	12905
20	12228	12228	12228	12228	12228	12228	12228	12228	12228	12228
21	10921	10921	10921	10921	10921	10921	10921	10921	10921	10921
22	8822	8822	8822	8822	8822	8822	8822	8822	8822	8822
23	6275	6275	6275	6275	6275	6275	6275	6275	6275	6275

In [18]:

```
Hours = [hour for hour, df in all_data.groupby('Hour')]

plt.plot(Hours,results)
plt.xticks(Hours)
plt.xlabel('Hours')
plt.ylabel('Order Volume')
plt.grid()
```



Question 4: What products are most often sold together?

```
In [19]: df = all_data[all_data['Order ID'].duplicated(keep = False)]

df['Grouped'] = df.groupby('Order ID')['Product'].transform(lambda x: ','.join(x))

df = df[['Order ID', 'Grouped']]. drop_duplicates()

df.head(25)
```

C:\Users\SOUPTIK DAS\AppData\Local\Temp\ipykernel_13928\1318596486.py:3: 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
df['Grouped'] = df.groupby('Order ID')['Product'].transform(lambda x: ','.join(x))

Out[19]:

	Order ID	Grouped
3	176560	Google Phone,Wired Headphones
18	176574	Google Phone,USB-C Charging Cable
30	176585	Bose SoundSport Headphones,Bose SoundSport Hea...
32	176586	AAA Batteries (4-pack),Google Phone
119	176672	Lightning Charging Cable,USB-C Charging Cable
129	176681	Apple Airpods Headphones,ThinkPad Laptop
138	176689	Bose SoundSport Headphones,AAA Batteries (4-pack)
189	176739	34in Ultrawide Monitor,Google Phone
225	176774	Lightning Charging Cable,USB-C Charging Cable
233	176781	iPhone,Lightning Charging Cable
250	176797	Google Phone,Bose SoundSport Headphones,Wired ...
260	176805	Google Phone,USB-C Charging Cable
264	176808	Google Phone,Wired Headphones
270	176813	Google Phone,Wired Headphones
394	176935	AAA Batteries (4-pack),27in FHD Monitor
435	176975	USB-C Charging Cable,AAA Batteries (4-pack)
450	176989	Google Phone,USB-C Charging Cable
455	176993	iPhone,Wired Headphones
485	177022	iPhone,Wired Headphones
567	177102	iPhone,27in 4K Gaming Monitor
581	177115	iPhone,Lightning Charging Cable
584	177117	ThinkPad Laptop,AAA Batteries (4-pack)
635	177167	iPhone,Apple Airpods Headphones,AAA Batteries ...
648	177178	iPhone,Lightning Charging Cable
652	177181	Wired Headphones,Apple Airpods Headphones

```
In [30]: from itertools import combinations
from collections import Counter

count = Counter()
```

```
for row in df['Grouped'].  
    row_list = row.split(',')  
    count.update(Counter(combinations(row_list,2)))  
  
for key, value in count.most_common(10):  
    print(key,value)
```

```
('iPhone', 'Lightning Charging Cable') 1005  
( 'Google Phone', 'USB-C Charging Cable') 987  
( 'iPhone', 'Wired Headphones') 447  
( 'Google Phone', 'Wired Headphones') 414  
( 'Vareebadd Phone', 'USB-C Charging Cable') 361  
( 'iPhone', 'Apple Airpods Headphones') 360  
( 'Google Phone', 'Bose SoundSport Headphones') 220  
( 'USB-C Charging Cable', 'Wired Headphones') 160  
( 'Vareebadd Phone', 'Wired Headphones') 143  
( 'Lightning Charging Cable', 'Wired Headphones') 92
```

Question 5: What product sold the most? Why do you think the product was sold the most?

In [20]: all_data.head()

	Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City Name	Hour	Minute
0	176558	USB-C Charging Cable	2	11.95	2019-04-19 08:46:00	917 1st St, Dallas, TX 75001	4	23.90	Dallas (TX)	8	46
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99	Boston (MA)	22	30
3	176560	Google Phone	1	600.00	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	600.00	Los Angeles (CA)	14	38
4	176560	Wired Headphones	1	11.99	2019-04-12 14:38:00	669 Spruce St, Los Angeles, CA 90001	4	11.99	Los Angeles (CA)	14	38
5	176561	Wired Headphones	1	11.99	2019-04-30 09:27:00	333 8th St, Los Angeles, CA 90001	4	11.99	Los Angeles (CA)	9	27

In [21]: results = all_data.groupby('Product').sum()['Quantity Ordered']
results

Out[21]:

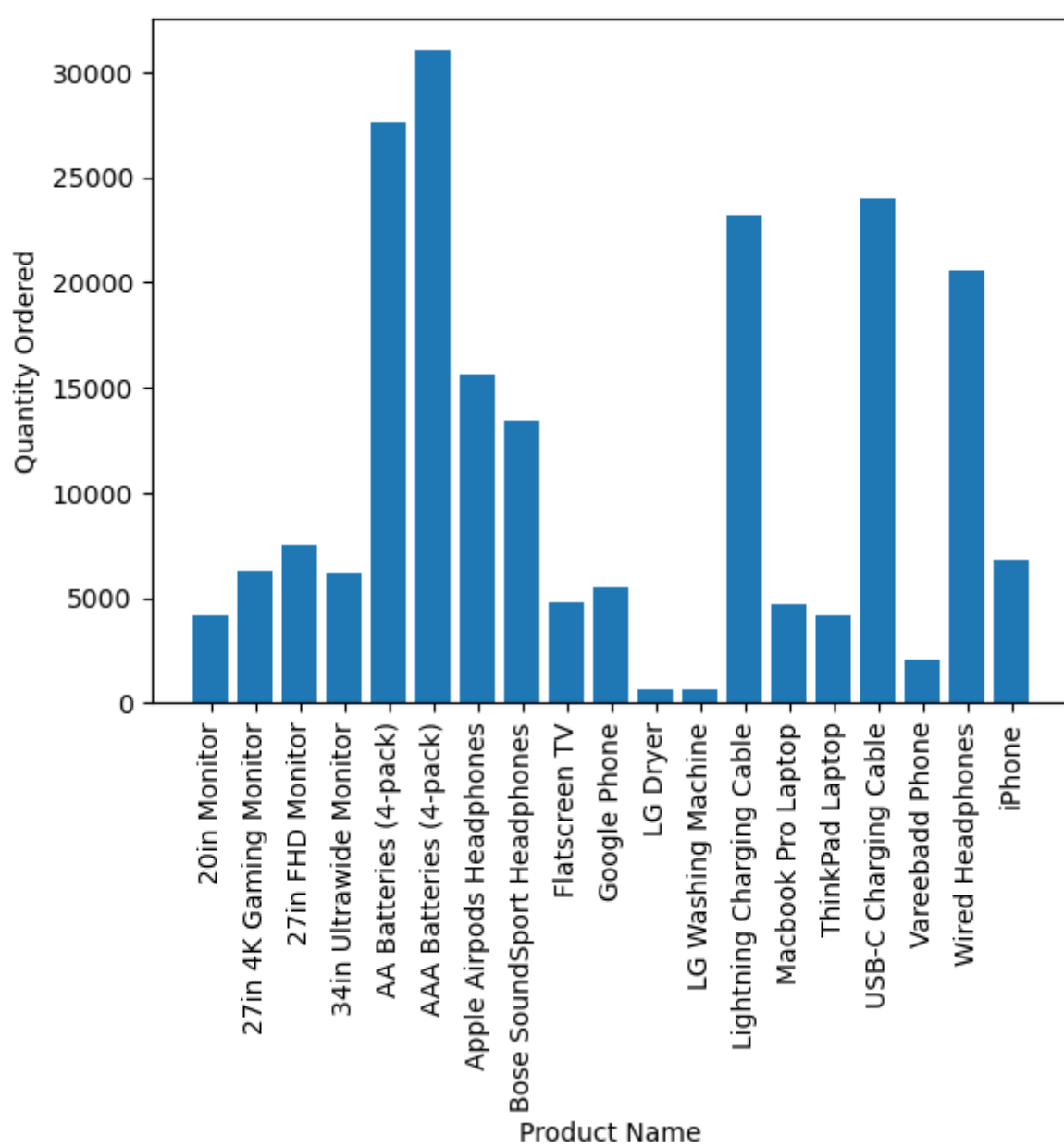
Product	
20in Monitor	4129
27in 4K Gaming Monitor	6244
27in FHD Monitor	7550
34in Ultrawide Monitor	6199
AA Batteries (4-pack)	27635
AAA Batteries (4-pack)	31017
Apple Airpods Headphones	15661
Bose SoundSport Headphones	13457
Flatscreen TV	4819
Google Phone	5532
LG Dryer	646
LG Washing Machine	666
Lightning Charging Cable	23217
Macbook Pro Laptop	4728
ThinkPad Laptop	4130
USB-C Charging Cable	23975
Vareebadd Phone	2068
Wired Headphones	20557
iPhone	6849

Name: Quantity Ordered, dtype: int64

In [22]: Products = [product for product, df in all_data.groupby('Product')]

```
plt.bar(Products,results)  
plt.xticks(Products, rotation = 'vertical')  
plt.xlabel('Product Name')  
plt.ylabel('Quantity Ordered')
```

Out[22]: Text(0, 0.5, 'Quantity Ordered')



Plotting line chart against bar graph to determine quantity sold and price of product w.r.t product name

```
In [23]: price = all_data.groupby('Product').mean()['Price Each']
price
```

```
Out[23]: Product
20in Monitor          109.99
27in 4K Gaming Monitor  389.99
27in FHD Monitor       149.99
34in Ultrawide Monitor  379.99
AA Batteries (4-pack)    3.84
AAA Batteries (4-pack)   2.99
Apple AirPods Headphones 150.00
Bose SoundSport Headphones 99.99
Flatscreen TV          300.00
Google Phone           600.00
LG Dryer               600.00
LG Washing Machine     600.00
Lightning Charging Cable 14.95
Macbook Pro Laptop     1700.00
ThinkPad Laptop        999.99
USB-C Charging Cable    11.95
Vareebadd Phone        400.00
Wired Headphones       11.99
iPhone                 700.00
Name: Price Each, dtype: float64
```

Plotting the two graphs together

```
In [25]: fig, ax1 = plt.subplots()

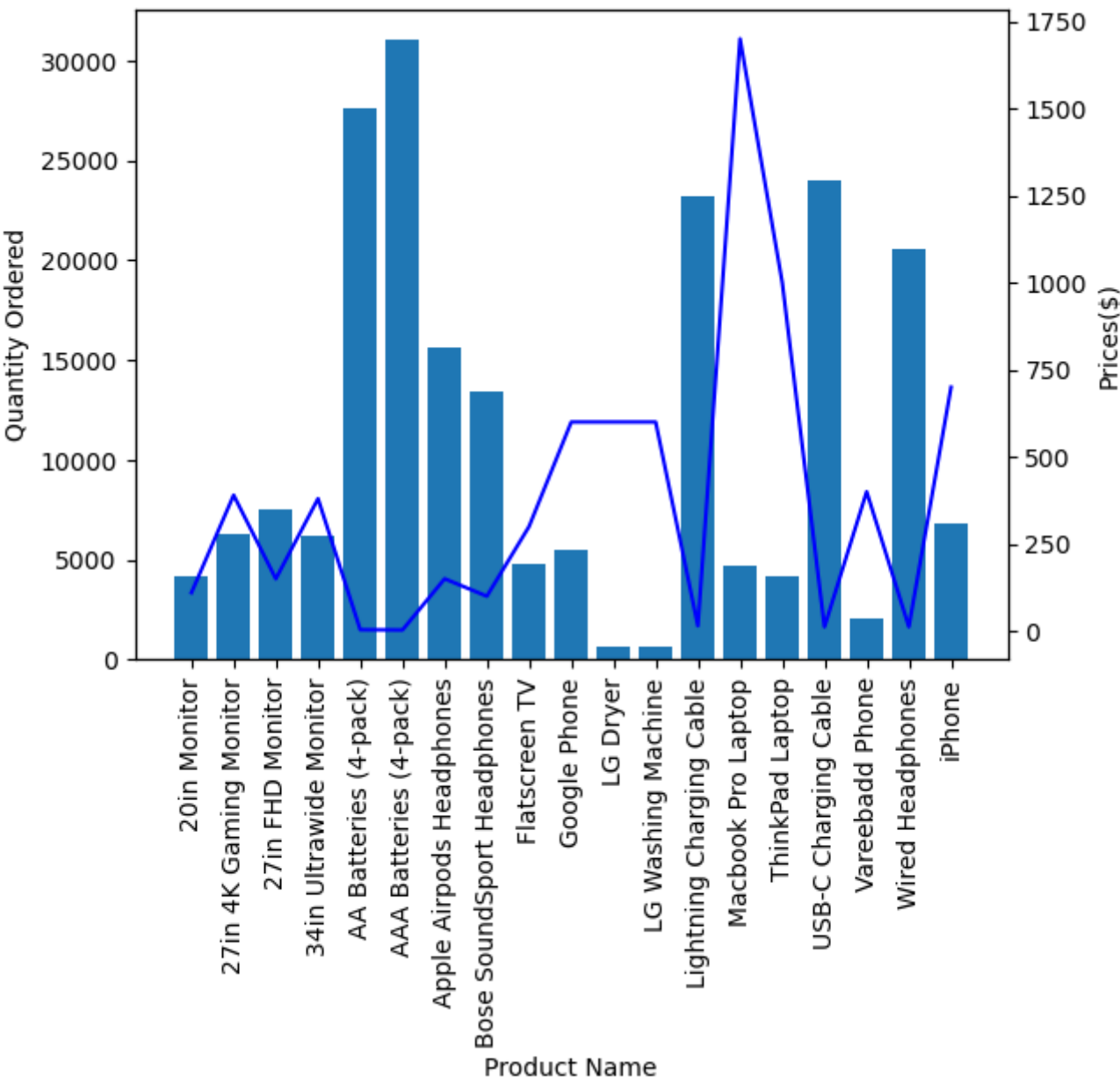
ax2 = ax1.twinx()

ax1.bar(Products,results)
ax2.plot(Products,price, 'b-')

ax1.set_xlabel('Product Name')
ax1.set_ylabel('Quantity Ordered')
ax2.set_ylabel('Prices($)')
ax1.set_xticklabels(Products, rotation='vertical')
```

C:\Users\SOUPTIK DAS\AppData\Local\Temp\ipykernel_13928\1795160273.py:11: UserWarning: FixedFormatter should only be used together with FixedLocator
ax1.set_xticklabels(Products, rotation='vertical')


```
Out[25]: [Text(0, 0, '20in Monitor'),
Text(1, 0, '27in 4K Gaming Monitor'),
Text(2, 0, '27in FHD Monitor'),
Text(3, 0, '34in Ultrawide Monitor'),
Text(4, 0, 'AA Batteries (4-pack)'),
Text(5, 0, 'AAA Batteries (4-pack)'),
Text(6, 0, 'Apple Airpods Headphones'),
Text(7, 0, 'Bose SoundSport Headphones'),
Text(8, 0, 'Flatscreen TV'),
Text(9, 0, 'Google Phone'),
Text(10, 0, 'LG Dryer'),
Text(11, 0, 'LG Washing Machine'),
Text(12, 0, 'Lightning Charging Cable'),
Text(13, 0, 'Macbook Pro Laptop'),
Text(14, 0, 'ThinkPad Laptop'),
Text(15, 0, 'USB-C Charging Cable'),
Text(16, 0, 'Vareebadd Phone'),
Text(17, 0, 'Wired Headphones'),
Text(18, 0, 'iPhone')]
```



```
In [ ]:
```