

## Sales Data Project

```
In [126]: import pandas as pd
import os

In [127]: files = [file for file in os.listdir("C:\Users\946851\OneDrive - Cognizant\Desktop\Pandas-Data-Science-Tasks-master\SalesAnalysis\Sales_Data")]
all_months_data = pd.DataFrame()
for file in files:
    all_data.append(pd.read_csv("C:\Users\946851\OneDrive - Cognizant\Desktop\Pandas-Data-Science-Tasks-master\SalesAnalysis\Sales_Data/" + file))
all_months_data = pd.concat([all_months_data, df])

In [128]: all_data = pd.read_csv('salesfulldata.csv')

Out[128]:
```

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
...	...	...	...	...	...	...
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	700.00	09/02/19 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

186850 rows × 6 columns

## CLEANING

```
In [129]: #dropping NaN rows
all_data[all_data['Quantity Ordered'] == 0].dropna(inplace=True)
all_data[all_data['Price Each'] == 0].dropna(inplace=True)

Out[129]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
1	NaN	NaN	NaN	NaN	NaN
356	NaN	NaN	NaN	NaN	NaN
735	NaN	NaN	NaN	NaN	NaN
1433	NaN	NaN	NaN	NaN	NaN
1553	NaN	NaN	NaN	NaN	NaN

```
In [130]: #Removing OR character in month of orderdate column
all_data[all_data['Order Date'].str[8:2] == 'Or']
all_data = all_data[all_data['Order Date'].str[8:2] != 'Or']
all_data[all_data['Order Date'].str[8:2] == 'Or']

Out[130]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address
----------	---------	------------------	------------	------------	------------------

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

```
In [131]: all_data.head()

Out[131]:
```

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

```
In [132]: #Creating a month column and converting the data type from string to int
all_data['Month'] = all_data['Order Date'].str[8:2]
all_data['Month'] = all_data['Month'].astype('int32')

In [133]: #converting data types
all_data['Quantity Ordered'] = pd.to_numeric(all_data['Quantity Ordered'])
all_data['Price Each'] = pd.to_numeric(all_data['Price Each'])

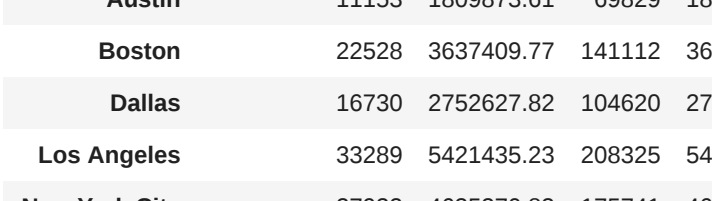
In [134]: #ADD SALES COLUMN
all_data['Sales'] = all_data['Price Each'] * all_data['Quantity Ordered']

In [135]: results = all_data.groupby('Month').sum()
results
```

```
Out[135]:
```

	Quantity Ordered	Price Each	Sales
Month			
1	10903	1911769.30	1822956.73
2	13449	218894.72	2200202.42
3	17005	2791207.63	2897103.36
4	20958	3367671.02	3366767.24
5	11867	313121.13	311536.75
6	15283	2662025.61	2577802.26
7	16072	2932539.56	2647775.76
8	13448	2320456.42	2344467.86
9	13109	2064992.09	2097560.12
10	22703	3715554.83	3736726.80
11	19798	3189600.68	3199020.20
12	28114	4598415.41	4633433.34

```
In [136]: #Seeing the results in chart
import matplotlib.pyplot as plt
Months = range(1,12)
plt.bar(Months, results['Sales'])
plt.xticks(Months)
plt.xlabel('Month')
plt.ylabel('Sales')
plt.show()
```



What city sold the most product?

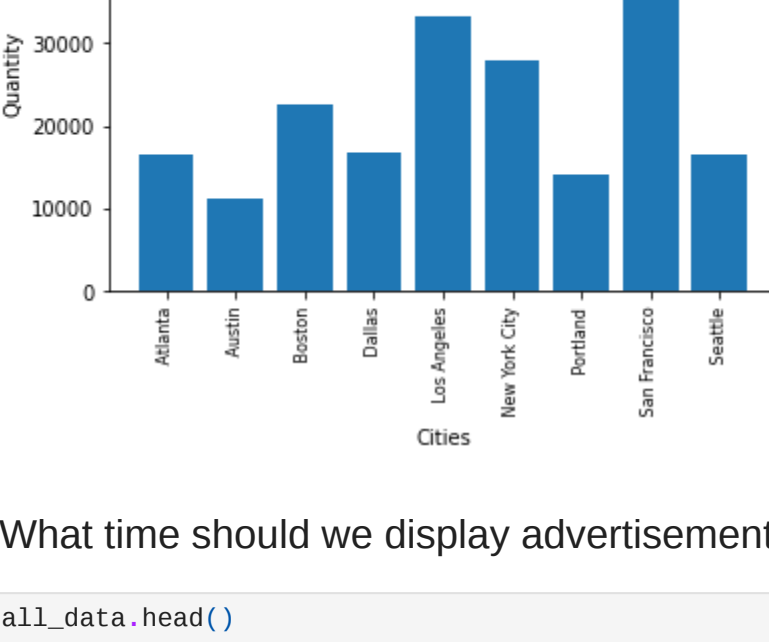
```
In [137]: #Taking city from address
def get_city(address):
    return address.split(',')[1]
all_data['City'] = all_data['Purchase Address'].apply(lambda x: get_city(x))

In [138]: results = all_data.groupby('City').sum()
results
#San Francisco
```

```
Out[138]:
```

	Quantity Ordered	Price Each	Month	Sales
City				
Atlanta	16602	2779908.20	104794	2795498.58
Austin	11153	1809873.61	69829	1819581.75
Boston	22528	3637409.77	141112	3661642.01
Dallas	16730	2752627.82	104620	2767975.40
Los Angeles	33289	5421435.23	208325	5452570.80
New York City	27932	4369370.83	175741	4664317.43
Portland	14063	2377747.47	87765	2330490.61
San Francisco	50239	8211461.74	315520	8326203.91
Seattle	16593	2732296.01	104941	2747755.48

```
In [139]: #Cities = all_data['City'].unique() # we did not get data in order
cities = [city for city, df in all_data.groupby('City')]
plt.bar(cities, results['Quantity Ordered'])
plt.xticks(rotation='vertical', size=8)
plt.xlabel('Cities')
plt.ylabel('Quantity')
plt.show()
```



What time should we display advertisements to maximize the likelihood of purchases?

```
In [140]: all_data.head()

Out[140]:
```

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

```
In [141]: #Take time from order date
all_data['Time'] = all_data['Order Date'].apply(lambda x: x.split(' ')[1].split(':')[0])
all_data

Out[141]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City	Time	
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	08
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99	Boston	22
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	14
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	14
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	09
...	...	...	...	...	...	...	...	...	...	
186845	259353	AAA Batteries (4-pack)	3	2.99	2019-09-17 20:56:00	840 Highland St, Los Angeles, CA 90001	9	8.97	Los Angeles	20
186846	259354	iPhone	1	700.00	2019-09-02 16:00:00	216 Dogwood St, San Francisco, CA 94016	9	700.00	San Francisco	16
186847	259355	iPhone	1	700.00	2019-09-23 07:39:00	220 12th St, San Francisco, CA 94016	9	700.00	San Francisco	07
186848	259356	34in Ultrawide Monitor	1	379.99	2019-09-19 17:30:00	511 Forest St, San Francisco, CA 94016	9	379.99	San Francisco	17
186849	259357	USB-C Charging Cable	1	11.95	2019-09-30 00:18:00	250 Meadow St, San Francisco, CA 94016	9	11.95	San Francisco	00

186850 rows × 10 columns

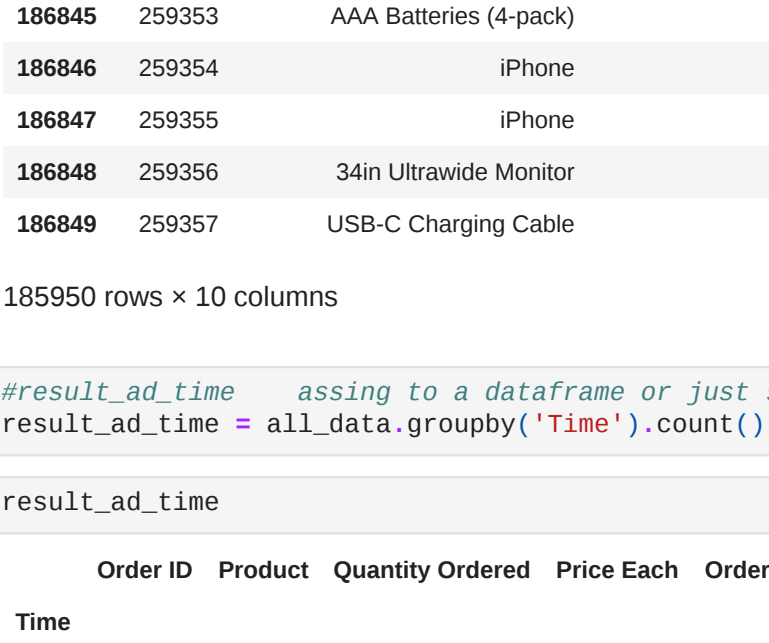
```
In [142]: #result_ad_time assing to a dataframe or just straightly can give in y axis
result_ad_time = all_data.groupby('Time').count()

In [143]: result_ad_time

Out[143]:
```

Time	Count
00	3910
01	2350
02	1243
03	831
04	854
05	1321
06	2482
07	4011
08	6256
09	8748
10	10944
11	12411
12	12587
13	12129
14	10984
15	10175
16	10384
17	10099
18	12280
19	12606
20	12278
21	10921
22	8822
23	6275

```
In [144]: Clock = [Time for Time, df in all_data.groupby('Time')]
plt.plot(Clock, result_ad_time['Order ID'])
plt.xticks(rotation='horizontal', size=8)
plt.xlabel('Time')
plt.ylabel('Count of Orders')
plt.show()
```



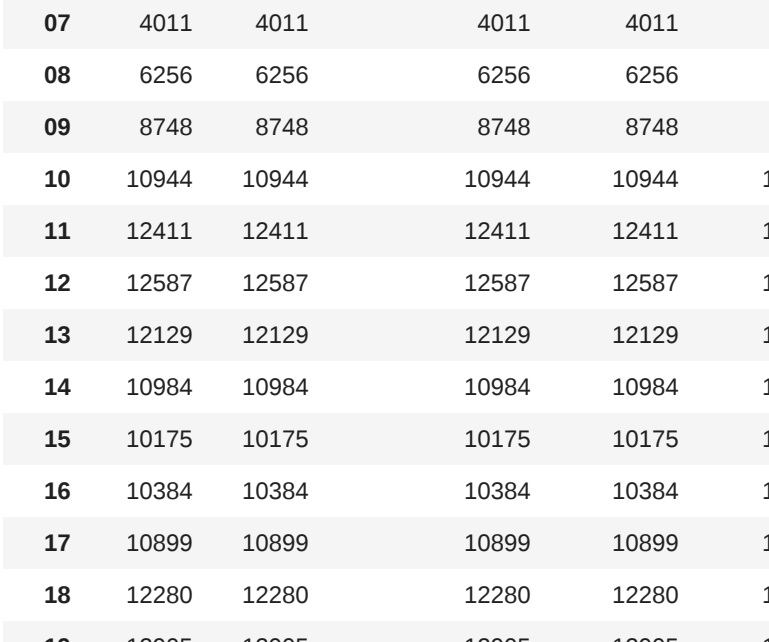
```
In [145]: #2nd method by converting data type of order date to date time
all_data['Order Date'] = pd.to_datetime(all_data['Order Date'])
all_data.head()

Out[145]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City	Time	
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	08
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99	Boston	22
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	14
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	14
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	09

```
In [147]: #Create a hours column
all_data['Hour'] = all_data['Order Date'].dt.hour
Result = all_data.groupby('Hour').count()

In [148]: Clock = [Time for Time, df in all_data.groupby('Hour')]
plt.bar(Clock, Result['Order ID'])
plt.xticks(Clock)
plt.show()
plt.plot(Clock, Result['Order ID'])
plt.grid()
plt.show()
#Earnings with maximum order 12,12,19
```



```
In [149]: #What's duplicated row also has based on order id to a new data frame
new_df = all_data[all_data['Order ID']. duplicated(keep=False)]
new_df

Out[149]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City	Time	Hour	
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	08	8
2	176559	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	682 Chestnut St, Boston, MA 02215	4	99.99	Boston	22	22
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	14	14
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	14	14
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	09	9

```
In [150]: #Taking duplicated row also has based on order id to a new data frame
new_df = all_data[all_data['Order ID']. duplicated(keep=False)]
new_df

Out[150]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City	Time	Hour	
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	14	14
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	14	14
18	176574	Google Phone	1	600.00	2019-04-03 19:42:00	20 Hill St, Los Angeles, CA 90001	4	600.00	Los Angeles	19	19
19	176574	USB-C Charging Cable	1	11.95	2019-04-03 19:42:00	20 Hill St, Los Angeles, CA 90001	4	11.95	Los Angeles	19	19
...	...	...	...	...	...	...	...	...	...	...	
186792	259303	AA Batteries (4-pack)	1	3.84	2019-09-20 20:18:00	106 7th St, Atlanta, GA 30301	9	3.84	Atlanta	20	20
186803	259314	Wired Headphones	1	11.99	2019-09-16 00:25:00	241 Highland St, Atlanta, GA 30301	9	11.99	Atlanta	00	0
186804	259314	AAA Batteries (4-pack)	2	2.99	2019-09-16 00:25:00	241 Highland St, Atlanta, GA 30301	9	5.98	Atlanta	00	0
186841	259350	Google Phone	1	600.00	2019-09-30 13:49:00	519 Maple St, San Francisco, CA 94016	9	600.00	San Francisco	13	13
186842	259350	USB-C Charging Cable	1	11.95	2019-09-30 13:49:00	519 Maple St, San Francisco, CA 94016	9	11.95	San Francisco	13	13

145449 rows × 11 columns

```
In [151]: #Group by based on similar orderid and concatenation of products bought
new_df['all Products'] = new_df.groupby('Order ID')['Product'].transform(lambda x: ', '.join(x))
new_df

In [151]: #C:\Users\946851\AppData\Local\Temp\7\ipykernel_11884\422882808.py:2: SettingWithCopyWarning:
A value is being 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
new_df['all Products'] = new_df.groupby('Order ID')['Product'].transform(lambda x: ', '.join(x))

Out[151]:
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	Sales	City	Time	Hour	all Products	
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	14	14	Google Phone,Wired Headphones
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	14	14	Google Phone,Wired Headphones
18	176574	Google Phone	1	600.00	2019-04-03 19:42:00	20 Hill St, Los Angeles, CA 90001	4	600.00	Los Angeles	19	19	Google Phone,USB-C Charging Cable
19	176574	USB-C Charging Cable	1	11.95	2019-04-03 19:42:00	20 Hill St, Los Angeles, CA 90001	4	11.95	Los Angeles	19	19	Google Phone,USB-C Charging Cable
30	176585	Bose SoundSport Headphones	1	99.99	2019-04-07 22:30:00	823 Highland St, Boston, MA 02215	4	99.99	Boston	11	11	Bose SoundSport Headphones,Bose SoundSport Head...
...	...	...	...	...	...	...	...	...	...	...	...	
186792	259303	AA Batteries (4-pack)	1	3.84	2019-09-20 20:18:00	106 7th St, Atlanta, GA 30301	9	3.84	Atlanta	20	20	34in Ultrawide Monitor,AA Batteries (4-pack)
186803	259314	Wired Headphones	1	11.99	2019-09-16 00:25:00	241 Highland St, Atlanta, GA 30301	9	11.99	Atlanta	00	0	Wired Headphones,AAA Batteries (4-pack)
186804	259314	AAA Batteries (4-pack)	2	2.99	2019-09-16 00:25:00	241 Highland St, Atlanta, GA 30301	9	5.98	Atlanta	00	0	Wired Headphones,AAA Batteries (4-pack)
186841	259350	Google Phone	1	600.00	2019-09-30 13:49:00	519 Maple St, San Francisco, CA 94016	9	600.00	San Francisco			