

```
import pandas as pd
import os
import matplotlib.pyplot as plt
```

```
[72]: df = pd.read_csv("../Sales/Data/Sales_April_2019.csv")
files = [file for file in os.listdir('../Sales/Data')] #create a list for all files
all_months_data = pd.DataFrame() #create empty data frame

for file in files:
    df = pd.read_csv("../Sales/Data/"+file) #start a loop, iterates each file
    all_months_data = pd.concat([all_months_data, df]) #read each csv file iterated in loop and store content in data frame df
    #combine the data from all csv file into a single data frame

all_months_data.to_csv("all_data.csv", index=False) #exports the combine data frame to a csv file
```

Read in updated data frame

```
In [73]: all_data = pd.read_csv("all_data.csv")
all_data.head(6)
```

```
Out[73]:   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         600.00  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 [74]: all_data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 186850 entries, 0 to 186849
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype  
---  --
0   Order ID               186395 non-null  object 
1   Product                186395 non-null  object 
2   Quantity Ordered       186395 non-null  object 
3   Price Each             186395 non-null  object 
4   Order Date             186395 non-null  object 
5   Purchase Address       186395 non-null  object 
dtypes: object(6)
memory usage: 8.4+ MB
```

In [75]: nan_df = all_data[all_data.isna().any(axis=1)] #To find Nan values in Data
nan_df.head()

```
Out[75]:   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
1423     NaN        NaN              NaN          NaN          NaN              NaN
1553     NaN        NaN              NaN          NaN          NaN              NaN
```

In [76]: all_data.dropna(inplace=True) #drop Nan values

In [77]: all_data.head()

```
Out[77]:   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         600.00  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
```

Find 'Or' and delete it

```
In [78]: temp_df = all_data[all_data['Order Date'].str[0:2] != 'or'] #to check if first two characters in order date column have 'or'
temp_df.head()
```

```
Out[78]:   Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
519      519      Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
1149      1149      Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
1155      1155      Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
2878      2878      Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
2893      2893      Order ID      Product  Quantity Ordered  Price Each  Order Date      Purchase Address
```

In [79]: all_data = all_data[all_data['Order Date'].str[0:2] != 'or'] #filtering out rows for 'or'

Convert columns to the correct data type

```
In [80]: all_data['Quantity Ordered'] = pd.to_numeric(all_data['Quantity Ordered']) #Make int
all_data['Price Each'] = pd.to_numeric(all_data['Price Each']) #Make float
```

Augment data with additional columns

Add Month Column

```
In [81]: all_data['Month'] = all_data['Order Date'].str[0:2] #create new column Month
all_data['Month'] = all_data['Month'].astype('int32') #change data type
all_data.head()
```

```
Out[81]:   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
```

Add a Sales column

```
In [82]: all_data['Sales'] = all_data['Quantity Ordered'] * all_data['Price Each'] #calculate total sales for each row
all_data.head()
```

```
Out[82]:   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
```

Add a city column

```
In [88]: #use apply
def get_city(address):
    return address.split(',')[1] #split address to get city from

def get_state(address):
    return address.split(',')[2].split('.')[1] #split address to get state

all_data['City'] = all_data['Purchase Address'].apply(lambda x: f'{get_city(x)} ({get_state(x)})') #create City column and concatenated city and state using
all_data.head()
```

```
Out[88]:   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 (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)
```

Question1: What was the best month for sales? How much was earned that month?

```
In [84]: results = all_data.groupby('Month').sum()
results.head(12)
```

```
Out[84]:   Quantity Ordered  Price Each  Sales
Month
1          10903    1811768.38    1822256.73
2          13449    2188884.72    2202022.42
3          17005    2791207.83    2807100.38
4          20558    3367671.02    3380670.24
5          18667    3135125.13    3152606.75
6          15253    2562025.61    2577802.26
7          16072    2632039.56    2647775.76
8          13448    2230345.42    2244467.88
9          13109    2084902.09    2087560.13
10         22703    3715554.83    3736726.88
11          19798    2180600.68    2198603.20
12         28114    4588415.41    4613443.34
```

```
In [85]: Months = range(1,13)
plt.bar(Months, results['Sales'])
plt.xticks(Months)
plt.ylabel('Sales in USD ($)')
plt.xlabel('Month Number')
plt.ticklabel_format(style='plain') #disable the scientific notation on the y-axis

plt.show()
```

December was the best month for sales

Question2: What city had the highest number of sales?

```
In [89]: results = all_data.groupby('City').sum()
results
```

```
Out[89]:   Quantity Ordered  Price Each  Month  Sales
City
Austin (GA)          16602    2779908.20    104794    2795498.58
Atlanta (TX)         11153    1803873.61     69829    1819581.75
Boston (MA)          22528    3637409.77    141112    3661642.01
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    1800558.22     70621    1870732.34
San Francisco (CA)   50239    8211461.74    315520    8262203.91
Seattle (WA)         16553    2723926.01    104841    2747755.48
```

```
In [96]: cities = [city for city,
```