

```
import pandas as pd
import numpy as np
```

```
df=pd.read_csv('Sales_April_2017.csv')
df.head()
```

```
['Sales_April_2019.csv',
 'Sales_August_2019.csv',
 'Sales_December_2019.csv',
```

```
'Sales_June_2019.csv',
'Sales_March_2019.csv',
'Sales_May_2019.csv',
'Sales_February_2019.csv',
```

```
df.append.dtypes
```

Order ID	object
Product	object
Quantity Ordered	object
Price Each	object
Order Date	object
Purchase Address	object
dtype:	object

## Data clean-up: deleting duplicate rows

```
df_temp=df.append(df.append('Order Date'].str[0:2]=='Or')
df_temp.head()
```

1149	Order ID	Product	Quantity O
1155	Order ID	Product	Quantity O
2878	Order ID	Product	Quantity O

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)

```
df_new['Price Each'] = df_new['Price Each']
```

	Product	Quantity Ordered	Price Each	Purchase Address
Order ID				
1	USB-C Charge Cable	2	11.95	041819 08 46 417 1st St, Dallas, TX 75001
2	Base SoundPort Headphones	1	99.99	040710 22 30 682 Chestnut St, Boston, MA 02215
3	Google Phone	1	600.00	041219 14 38 669 Spruce St, Los Angeles, CA 90001
4	Wired Headphones	1	11.99	041219 14 38 669 Spruce St, Los Angeles, CA 90001
5	Wired Headphones	1	11.99	043019 09 27 333 Rth St, Los Angeles, CA 90001

Adding additional columns

```
df_new['Month']=df_new['Order Date'].str[0:2]
```

## Adding city column

```

#using apply function
def get_city(address):
    return address.split(',')[1]

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

#new['city'] = df.new['Purchase Address'].apply(lambda x: get_city(x))
#df.new.head()

#User's data AppData\LocalTemp\ipykernel_5168\3727653401.py:8: SettingWithCopyWarning:
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using loc[row_index,col_indexer] = value instead

See the warnings in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df.new['city'] = df.new['Purchase Address'].apply(lambda x: get_city(x))

```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month sales	city
1	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
2	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
3	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
4	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
5	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
6	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
7	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
8	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
9	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco
10	Chips and Beer	1	10.00	2016-01-01	1234 Main St, San Francisco, CA 94102	1	San Francisco

3	176560	Google Phone
4	176560	Wired Headphones
5	176560	Wired Headphones

```
months=range(1,13)
plt.figure(figsize=(9,4))
plt.bar(months,results['sales'])
plt.xticks(months)
plt.ylabel('Total sales in a month in USD')
plt.grid(axis='y')
plt.show()
```

Month	Total sales in a month in USD (x 1e6)
1	~1.5
2	~2.5
3	~2.5
4	~3.5
5	~2.5
6	~3.5
7	~3.5
8	~3.5
9	~3.5
10	~3.5
11	~3.5
12	~4.5

Category	Sales (millions of dollars)
1	1.5
2	2.0
3	2.0
4	2.0

```
cities=results.index
plt.figure(figsize=(9,4))
plt.bar(cities,results['sales'])
plt.xticks(cities,rotation='vertical',size=8)
plt.ylabel('Total sales in a city in USD')
plt.xlabel('city name')
plt.grid(axis='y')
plt.show()
```

Year	Total sales in millions of dollars
2006	2.8
2007	3.7
2008	3.7

	source	device	lat	lon	time	city	state	zip
4	176560	Google Phone	600.00	-121.840	2019-04-12 14:30:00	669 S Normandie St	Los Angeles, CA	90001
4	176560	Wired Headphones	1	11.99	2019-04-12 09:00:00	669 S Normandie St	Los Angeles, CA	90001
5	176561	Wired Headphones	1	11.99	2019-04-30 09:27:00	333 Bth St	Los Angeles, CA	90001

```
df_new['hour'] = df_new['Order Date'].dt.hour
df_new['minute'] = df_new['Order Date'].dt.minute
df_new.head()
```

C:\Users\avekita\AppData\Local\Temp\ipykernel\_5168\2136656433.py:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using loc[row\_index,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)  
df\_new['hour'] = df\_new['Order Date'].dt.hour  
C:\Users\avekita\AppData\Local\Temp\ipykernel\_5168\2136656433.py:2: SettingWithCopyWarning:

```
df_new['Minute']=df_new['Order Date'].dt.minute
```

[illegible]

```
plt.plot(hours, results['Order ID'])
plt.xticks(hours)
plt.xlabel('Hour of the day')
plt.ylabel('Total number of sales in an hour')
```

Order ID	Product	Quantity Ordered	Price Each	Order Date	Purchase Address	Month	sales	City	Hour	Minute
0 176558	USB-C Charging Cable	2	11.95	2019-04-10 09:46:00	517 1st St, Dallas, TX 75001	4	23.90	Dallas(TX)	8	46
2 176559	Base Soundstip/Split/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 Susan 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 Susan St, Los Angeles, CA 90001	4	11.99	Los Angeles(CA)	14	38
5 176561	Wired Headphones	1	11.99	2019-04-10 09:47:00	333 8th St, Los Angeles, CA 90001	4	11.99	Los Angeles(CA)	9	27
6 176562	USB-C Charging Cable	1	11.95	2019-04-29 13:03:00	381 Wilton St, San Francisco, CA 94104	4	11.95	San Francisco(CA)	13	3
7 176563	Base Soundstip/Split/Headphones	1	99.99	2019-04-02 09:50:00	666 Court St, Seattle, WA 98101	4	99.99	Seattle(WA)	7	46
8 176564	JURIS-C/Charging Cable	1	11.95	2019-04-17 09:27:00	270 Bldg St, Atlanta, GA 30301	4	11.95	Atlanta(GA)	10	58

Here we carefully observing the data , we can say that Order ID of two or more rows are same ,these corre

```
2890 179328 Wired Headphones, AA Batteries (4 pack)
100 rows x 2 columns

Now we can count number of occurrences of the combinations

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)
```

```

('Google Phone', 'Wired Headphones') 413
('Vareebadd Phone', 'USB-C Charging Cable') 361
('iPhone', 'Apple AirPods Headphones') 360
('Google Phone', 'Bose SoundSport Headphones') 220

```

```

USB-C Charging Cable      23901
Varadero Phone            2050
Wired Headphones          20524
iPhone                    6847
Name: Quantity Ordered, dtype: int64

products=quantity_ordered.index
plt.bar(products, quantity_ordered)
plt.xticks(products, rotation='vertical', size=8)
plt.ylabel('Number of units sold')
plt.xlabel('Products')
plt.show()

```

Product	Quantity Ordered
USB-C Charging Cable	23901
Varadero Phone	2050
Wired Headphones	20524
iPhone	6847

Product	Units Sold
Product A	22,000
Product B	20,000
Product C	18,000
Product D	15,000

AA Batteries (4-pack)	2.84
AAA Batteries (4-pack)	2.25
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	900.00
USB-C Charging Cable	11.95
Verablast Phone	400.00
Wired Headphones	11.99
iPhone	700.00
Name: Price Each, dtype: float64	

```
ax1.set_ylabel('Total quantity ordered',color='y')
ax1.set_xlabel('Products')

ax2.plot(products,prices,color='r')
```

we can see from the plot that there is an inverse correlation in quantity ordered and mean prices.