

OBJECTIVE

1.What is the overall sales trend?

2.Which are the Top 10 products by sales?

3.Which are the Most Selling Products?

4.Which is the most preferred Ship Mode?

5.Which are the Most Profitable Category and Sub-Category?

IMPORTING REQUIRED LIBRARIES

```
In [8]: # Data Manipulation
import pandas as pd

# Data Visualisation
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

In [9]: # Importing dataset
df = pd.read_excel('superstore_sales.xlsx')

In [9]: # To display First five rows of the dataset
df.head()
```

	order_id	order_date	ship_date	ship_mode	customer_name	segment	state	country	market	region	...	category	sub_category	product_name	sales	quantity	discount	profit
0	AC-2011-47883	2011-01-01	2011-01-06	Standard Class	Roby Braunhardt	Consumer	Constantine	Algeria	Africa	Africa	Office Supplies	Storage	Tenex Lockers, Blue	408.300	2	0.0	106.14
1	IN-2011-37913	2011-01-01	2011-01-08	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC	Oceania	Office Supplies	Supplies	Acme Trimmer, High Speed	120.366	3	0.1	36.06
2	HU-2011-1220	2011-01-01	2011-01-05	Second Class	Annie Thurman	Consumer	Budapest	Hungary	EMEA	EMEA	Office Supplies	Storage	Tenex Box, Single Width	66.120	4	0.0	29.64
3	IT-2011-3647632	2011-01-01	2011-01-05	Second Class	Eugene Moren	Home Office	Stockholm	Sweden	EU	North	Office Supplies	Paper	Ennemax Note Cards, Premium	44.865	3	0.5	-26.05
4	IN-2011-47883	2011-01-01	2011-01-08	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC	Oceania	Furnishings	Furnishings	Eldon Light Bulb, Duo Pack	113.670	5	0.1	37.70

5 rows × 21 columns

```
In [10]: # To display Last five rows of the dataset
df.tail()
```

	order_id	order_date	ship_date	ship_mode	customer_name	segment	state	country	market	region	...	category	sub_category	product_name	sales	quantity	discount	profit
51285	CA-2014-12-31	2014-12-31	2015-01-04	Standard Class	Erica Bern	Corporate	United States	US	West	Office Supplies	Binders	Cardinal Slant-D Ring Binder, Heavy Gauge Vinyl	13.904	2	0.2	0.2
51286	MO-2014-12-31	2014-12-31	2015-01-05	Standard Class	Liz Price	Consumer	Illinois	Moscow	Moscow	Africa	...	Office Supplies	Binders	Wilson Jones Hole Reinforcements, 3.990	1	0.0	0.0	
51287	MX-2014-12-31	2014-12-31	2015-01-02	Second Class	Charlotte Nelson	Consumer	Managua	Nicaragua	LATAM	Central	...	Office Supplies	Labels	Hon Color Coded Labels, 5000 Label Set	26.400	3	0.0	1.0
51288	MX-2014-12-31	2014-12-31	2015-01-06	Standard Class	Tamara Dahlen	Consumer	Chihuahua	Mexico	LATAM	North	...	Office Supplies	Labels	Hon Legal Exhibit Labels, Alphabetical	1.120	1	0.0	0.0
51289	CA-2014-12-31	2014-12-31	2015-01-04	Standard Class	Jill Matthias	Consumer	Colorado	United States	US	West	...	Office Supplies	Fasteners	Bagged Rubber Bands	3.024	3	0.2	0.2

5 rows × 21 columns

```
In [11]: # To display Shape of the dataset(No. of rows and columns)
df.shape
```

Out[11]: (51290, 21)

```
In [12]: # Columns present in the dataset
df.columns
```

```
In [12]: Index(['order_id', 'order_date', 'ship_date', 'ship_mode', 'customer_name',
        'segment', 'state', 'country', 'market', 'region', 'category', 'sub_category', 'product_name', 'sales', 'quantity', 'discount', 'profit'],
        dtype='object')
```

```
In [13]: # A concise summary of the dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51290 entries, 0 to 51289
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
--  --
0   order_id              51290 non-null  object
1   order_date            51290 non-null  datetime64[ns]
2   ship_date             51290 non-null  datetime64[ns]
3   ship_mode             51290 non-null  object
4   customer_name         51290 non-null  object
5   segment               51290 non-null  object
6   state                51290 non-null  object
7   country               51290 non-null  object
8   market               51290 non-null  object
9   region               51290 non-null  object
10  product_id            51290 non-null  object
11  category              51290 non-null  object
12  sub_category          51290 non-null  object
13  product_name          51290 non-null  object
14  sales                 51290 non-null  float64
15  quantity              51290 non-null  int64
16  discount              51290 non-null  float64
17  profit                51290 non-null  float64
18  shipping_cost          51290 non-null  float64
19  order_priority         51290 non-null  object
20  year                  51290 non-null  int64
dtypes: datetime64[ns](2), float64(4), int64(2), object(13)
memory usage: 8.2+ MB
```

```
In [16]: #To check if there are any missing values in the entire data set we use the isnull function, then see if there are any values.
df.isna().sum()
```

```
Out[16]:
order_id      0
order_date    0
ship_date     0
ship_mode     0
customer_name  0
segment       0
state         0
country       0
market        0
region        0
product_id    0
category      0
sub_category  0
product_name  0
sales         0
quantity      0
discount      0
profit        0
shipping_cost  0
order_priority 0
year          0
dtype: int64
```

```
In [17]: # Luckily there are no null values
```

Some descriptive statistics on the data set.

Notice, it only shows the statistics on the numerical columns. From here you can see the following statistics:

- 1.Row count, which aligns to what the shape attribute showed us.
- 2.The mean, or average.
- 3.The standard deviation, or how spread out the data is.
- 4.The minimum and maximum value of each column
- 5.The number of items that fall within the first, second, and third percentiles.

```
In [17]: df.describe().round()
```

	sales	quantity	discount	profit	shipping_cost	year
count	51290.0	51290.0	51290.0	51290.0	51290.0	51290.0
mean	246.0	3.0	0.0	29.0	26.0	2013.0
std	488.0	2.0	0.0	174.0	57.0	1.0
min	0.0	1.0	0.0	-6600.0	0.0	2011.0
25%	31.0	2.0	0.0	0.0	3.0	2012.0
50%	85.0	3.0	0.0	9.0	8.0	2013.0
75%	251.0	5.0	0.0	37.0	24.0	2014.0
max	22638.0	14.0	1.0	8400.0	934.0	2014.0

EXPLORATORY DATA ANALYSIS :

1.WHAT IS THE OVERALL SALES TREND?

```
In [18]: # the first date of our data set is...
df['order_date'].min()
```

Out[18]: Timestamp('2011-01-01 00:00:00')

```
In [19]: # the last date of our data set is...
df['order_date'].max()
```

Out[19]: Timestamp('2014-12-31 00:00:00')

```
In [21]: # Getting month and year from order_date
df['month_year'] = df['order_date'].apply(lambda x: x.strftime('%Y-%m'))
df['month_year']
```

```
Out[21]:
0      2011-01
1      2011-01
2      2011-01
3      2011-01
4      2011-01
...
51285  2014-12
51286  2014-12
51287  2014-12
51288  2014-12
51289  2014-12
Name: month_year, Length: 51290, dtype: object
```

```
In [22]: df.groupby('month_year').sum()
# grouping data based on "month_year" column and sum() gives count of sales,quantity.....per month
```

	sales	quantity	discount	profit	shipping_cost	year
month_year						
2011-01	98988.48886	1463	68.758	8321.80095	10544.78800	870763
2011-02	91152.15698	1224	52.252	12417.90688	10681.18300	760158
2011-03	145729.36736	1836	74.212	15303.56826	13096.18950	1069329
2011-04	146915.76418	2020	80.782	12902.32438	12494.52000	1134204
2011-05	146747.93810	2013	62.382	12183.82870	10443.20000	1130226
2011-06	145207.38022	2112	159.534	23415.54762	20813.19900	1164087
2011-07	151510.41912	1774	80.095	5585.92053	11844.47600	995445
2011-08	207581.49122	3026	115.162	23713.66712	22001.13900	1765658
2011-09	290214.45534	3707	137.678	35776.63840	29664.85100	2115572
2011-10	199071.26404	2727	110.192	25063.61834	21380.08200	1505614
2011-11	294495.53752	4038	178.835	32798.17772	34701.99800	2290529
2011-12	333925.73460	4493	187.230	40647.88400	37144.83100	2539903
2012-01	135780.72024	1945	74.454	10401.61764	13665.74900	1084648
2012-02	100510.21698	1473	62.784	15000.09618	11393.72000	803148
2012-03	160076.77116	2237	101.682	17962.81756	16170.78500	1331944
2012-04	161052.29952	2250	93.248	17366.96722	16767.82300	1231884
2012-05	208364.89124	2921	114.272	24676.70374	23801.61700	1690080
2012-06	256175.69842	3671	168.284	34407.15382	28158.83000	2285632
2012-07	145236.78512	2321	104.404	15685.38842	17334.43900	1125908
2012-08	203142.94238	3818	136.166	43573.87858	32038.71800	2178996
2012-09	289389.16564	4205	168.070	27776.18034	32083.17800	2460676
2012-10	252939.85020	3563	135.866	30662.88270	25085.74900	1991880
2012-11	323512.41690	5814	215.868	31620.72180	33489.74100	2937520
2012-12	338256.96660	4614	172.676	32950.75130	37563.32600	2583408
2013-01	199185.90738	2413	91.442	26610.55668	21677.43200	1427217
2013-02	167239.65040	2102	78.012	25340.02610	18268.01000	1217865
2013-03	198594.03012	2686	114.384	23433.77462	21268.01000	1541958
2013-04	177821.31684	2688	116.116	19462.03844	19133.23400	1580205
2013-05	260498.56470	3808	153.092	28486.68410	28153.21100	2127741
2013-06	396518.61190	5257	212.842	45478.41340	42814.02000	3078890
2013-07	299828.95020	3262	125.644	28863.82720	24501.84200	1860235
2013-08	226488.78936	4094	202.640	31023.68946	26573.68900	2207146
2013-09	376619.24568	5793	245.674	38905.66778	34848.40000	3385566
2013-10	293495.64298	3983	160.860	42322.22258	31174.68400	2214300
2013-11	373899.36010	5556	215.324	48062.99670	41407.16700	3213748
2013-12	405454.37802	5694	223.692	52002.87112	43183.80000	3224626
2014-01	241268.55566	3122	117.928	28091.38625	24870.80100	1848852
2014-02	184837.35556	2482	111.126	19751.69996	19638.63554	1522584
2014-03	263100.77262	3722	142.018	37267.20562	26038.63554	1505952
2014-04	242771.86130	3594	164.000	23782.30120	26272.71800	2116714
2014-05	288401.04614	4300	188.986	33953.55774	31882.58300	2595976
2014-06	403814.06310	6009	251.462	43778.60280	41994.07000	3520472
2014-07	258705.68048	3637	163.512	28035.87258	25789.33000	2189218
2014-08	456619.94236	5824	217.672	53542.89496	46759.35300	3373450
2014-09	481157.24370	6837	272.094	67979.45110	53485.43000	4064252
2014-10	422766.62916	5876	233.752	58208.93476	44622.41400	3274764
2014-11	555279.02700	7706	304.384	62056.58790	59818.35500	4324058
2014-12	503143.69348	7513	335.106	46916.52068	54853.89100	4336142

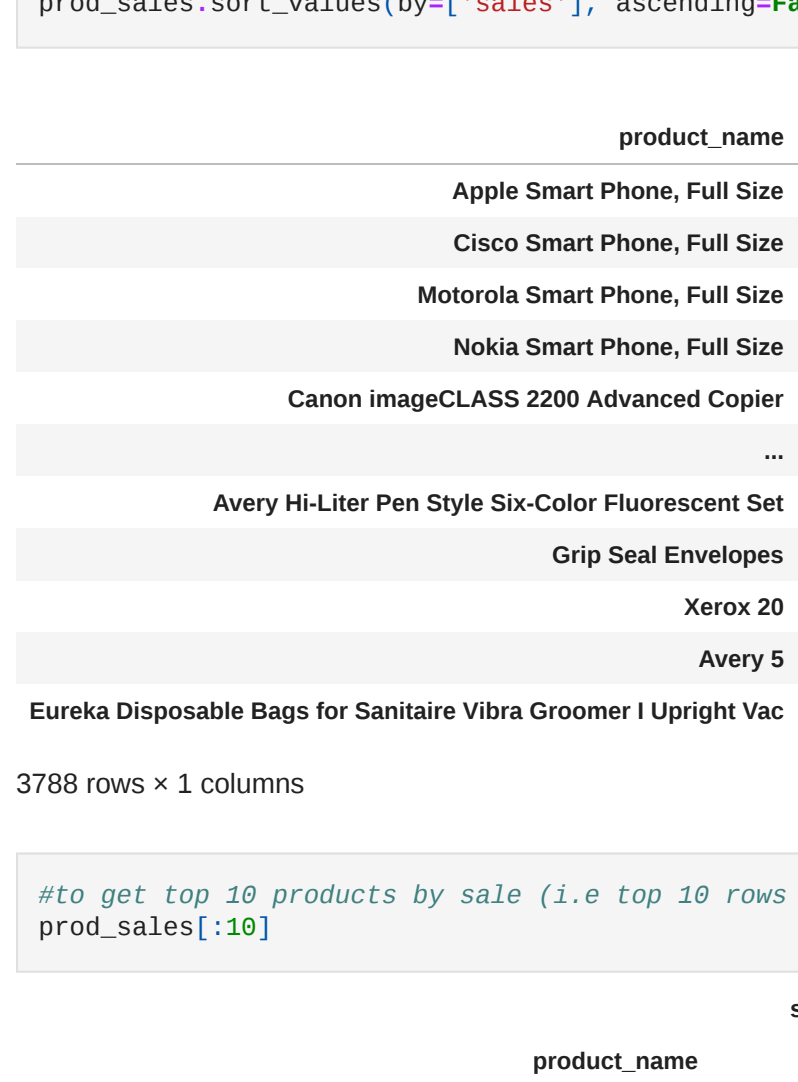
```
In [23]: #since we only need sales data we specify to sum up only sales column and reset the index values
df.groupby('month_year').sum()[['sales']].reset_index()
```

	month_year	sales
0	2011-01	98988.48886
1	2011-02	91152.15698
2	2011-03	145729.36736
3	2011-04	146915.76418
4	2011-05	146747.93810
5	2011-06	145207.38022
6	2011-07	151510.41912
7	2011-08	207581.49122
8	2011-09	290214.45534
9	2011-10	199071.26404
10	2011-11	294495.53752
11	2011-12	333925.73460
12	2012-01	135780.72024
13	2012-02	100510.21698
14	2012-03	160076.77116
15	2012-04	161052.29952
16	2012-05	208364.89124
17	2012-06	256175.69842
18	2012-07	145236.78512
19	2012-08	203142.94238
20	2012-09	289389.16564
21	2012-10	252939.85020
22	2012-11	323512.41690
23	2012-12	338256.96660
24	2013-01	199185.90738
25	2013-02	167239.65040
26	2013-03	198594.03012
27	2013-04	177821.31684
28	2013-05	260498.56470
29	2013-06	396518.61190
30	2013-07	229928.95020
31	2013-08	326488.78936
32	2013-09	376619.24568
33	2013-10	293495.64298
34	2013-11	373899.36010
35	2013-12	405454.37802
36	2014-01	241268.55566
37	2014-02	184837.35556
38	2014-03	263100.77262
39	2014-04	242771.86130
40	2014-05	288401.04614
41	2014-06	403814.06310
42	2014-07	258705.68048
43	2014-08	456619.94236
44	2014-09	481157.24370
45	2014-10	422766.62916
46	2014-11	555279.02700
47	2014-12	503143.69348

```
In [26]: #rename the new data set that we made now as df_MonthYearSales
df_MonthYearSales=df.groupby('month_year').sum()[['sales']].reset_index()
```

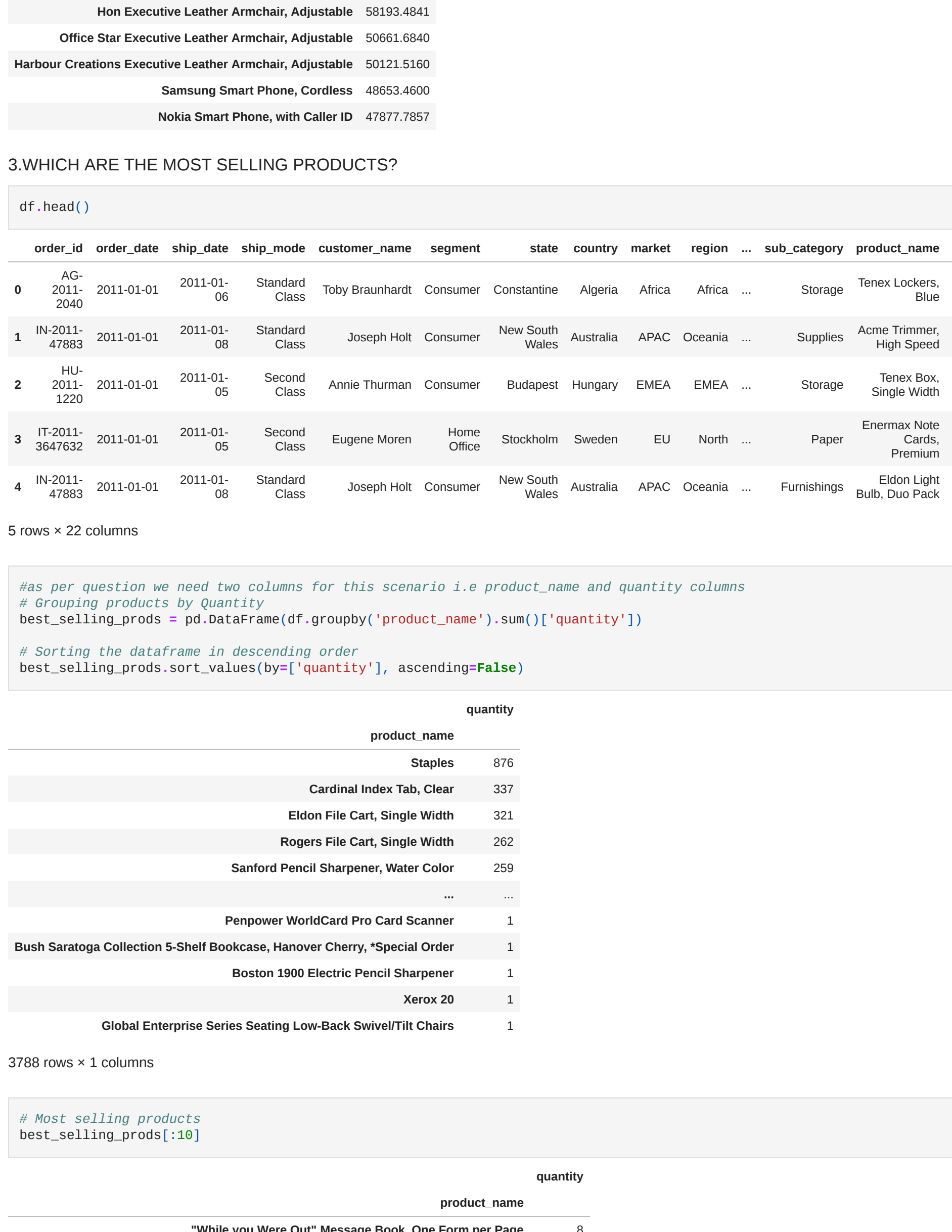
```
In [27]: #plotting: plt.plot(x,y)
plt.plot(df_MonthYearSales['month_year'],df_MonthYearSales['sales'])
```

```
Out[27]: [<matplotlib.lines.Line2D at 0x3217ed4fe>]
```



```
In [29]: #we cannot see x axis labels clearly and the graph is aligned too close(in the above line graph)so set the figure size
plt.figure(figsize=(16, 5))
plt.plot(df_MonthYearSales['month_year'],df_MonthYearSales['sales'])
```

```
Out[29]: [<matplotlib.lines.Line2D at 0x322d6c4fe>]
```



```
In [31]: #the graph is aligned properly but since x axis is still not clear hence we tilt the labels to 90 degrees vertical by using xticks(rotation)
plt.figure(figsize=(16, 5))
plt.xticks(rotation='vertical', size=8)
plt.plot(df_MonthYearSales['month_year'],df_MonthYearSales['sales'])
plt.show()
```

2.WHICH ARE THE TOP 10 PRODUCTS BY SALES?

```
In [32]: df.head()
```

	order_id	order_date	ship_date	ship_mode	customer_name	segment	state	country	market	region	...	sub_category	product_name	sales	quantity	discount	profit	shipping_cost
0	AC-2011-47883	2011-01-01	2011-01-06	Standard Class	Roby Braunhardt	Consumer	Constantine	Algeria	Africa	Africa	Storage	Tenex Lockers, Blue	408.300	2	0.0	106.14	
1	IN-2011-37913	2011-01-01	2011-01-08	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC	Oceania	Supplies	Acme Trimmer, High Speed	120.366	3	0.1	36.06	
2	HU-2011-1220	2011-01-01	2011-01-05	Second Class	Annie Thurman	Consumer	Budapest	Hungary	EMEA	EMEA	Storage	Tenex Box, Single Width	66.120	4	0.0	29.64	
3	IT-2011-3647632	2011-01-																