

GRIP: The Sparks Foundation

Data Science and Business Analytics Internship Task 3: Exploratory Data Analysis - Retail

Perform 'Exploratory Data Analysis' on dataset 'SampleSuperstore' Python 3 Jupyter Notebook Name-Chetana Thorat

In [7]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

import warnings
warnings.filterwarnings('ignore')
```

In [8]:

```
df = pd.read_csv("SampleSuperstore.csv")
```

In [13]:

```
df.sample(5)
```

Out [13]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
906	Standard Class	Consumer	United States	New York City	New York	10009	East	Furniture	Bookcases	323.136	4	0.2	12.1176
847	Standard Class	Consumer	United States	Louisville	Kentucky	40214	South	Furniture	Chairs	287.940	3	0.0	77.7438
9393	First Class	Consumer	United States	New York City	New York	10011	East	Technology	Phones	629.950	5	0.0	157.4875
9952	Standard Class	Consumer	United States	Jackson	Tennessee	38301	South	Office Supplies	Art	23.128	7	0.2	2.8910
6167	Standard Class	Consumer	United States	Houston	Texas	77036	Central	Office Supplies	Paper	117.456	3	0.2	44.0460

In [14]:

```
df.head()
```

Out [14]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
0	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	261.9600	2	0.00	41.9136
1	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	731.9400	3	0.00	219.5820
2	Second Class	Corporate	United States	Fort Los Angeles	California	90036	West	Office Supplies	Labels	14.6200	2	0.00	6.8714
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	957.5775	5	0.45	-383.0310
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	22.3680	2	0.20	2.5164

In [15]:

```
df.tail()
```

Out [15]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
9989	Second Class	Consumer	United States	Miami	Florida	33180	South	Furniture	Furnishings	25.246	3	0.2	4.1028
9990	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Furniture	Furnishings	91.9690	2	0.0	15.6332
9991	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Technology	Phones	258.576	2	0.0	19.3932
9992	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Office Supplies	Paper	29.6890	4	0.0	13.3200
9993	Second Class	Consumer	United States	Westminster	California	92683	West	Office Supplies	Appliances	243.1600	2	0.0	72.9480

In [16]:

```
df.shape
```

Out [16]:

```
(9994, 13)
```

In [17]:

```
df.info
```

Out [17]:

```
<bound method DataFrame.info of
0      Second Class  Consumer  United States  Henderson  Kentucky  42420  South  Furniture  Bookcases  261.9600  2  0.00  41.9136
1      Second Class  Consumer  United States  Henderson  Kentucky  42420  South  Furniture  Chairs      731.9400  3  0.00  219.5820
2      Second Class  Corporate United States  Fort Los Angeles  California  90036  West  Office Supplies  Labels      14.6200  2  0.00  6.8714
3      Standard Class Consumer  United States  Fort Lauderdale  Florida  33311  South  Furniture  Tables      957.5775  5  0.45 -383.0310
4      Standard Class Consumer  United States  Fort Lauderdale  Florida  33311  South  Office Supplies  Storage      22.3680  2  0.20  2.5164

...
9989  Second Class  Consumer  United States  Miami          Florida  33180  South  Furniture  Furnishings  25.246    3  0.20  4.1028
9990  Standard Class Consumer  United States  Costa Mesa      California  92627  West  Furniture  Furnishings  91.9690   2  0.00  15.6332
9991  Standard Class Consumer  United States  Costa Mesa      California  92627  West  Technology  Phones      258.576   2  0.00  19.3932
9992  Standard Class Consumer  United States  Costa Mesa      California  92627  West  Office Supplies  Paper       29.6890   4  0.00  13.3200
9993  Second Class  Consumer  United States  Westminster     California  92683  West  Office Supplies  Appliances  243.1600   2  0.00  72.9480

Postal Code Region Category Sub-Category Sales Quantity \
0      42420  South  Furniture  Bookcases  261.9600  2
1      42420  South  Furniture  Chairs      731.9400  3
2      90036  West  Office Supplies  Labels      14.6200  2
3      33311  South  Furniture  Tables      957.5775  5
4      33311  South  Office Supplies  Storage      22.3680  2
...
9989  33180  South  Furniture  Furnishings  25.246    3
9990  92627  West  Furniture  Furnishings  91.9690   2
9991  92627  West  Technology  Phones      258.576   2
9992  92627  West  Office Supplies  Paper       29.6890   4
9993  92683  West  Office Supplies  Appliances  243.1600   2

Discount Profit
0      0.00  41.9136
1      0.00  219.5820
2      0.00  6.8714
3      0.45 -383.0310
4      0.20  2.5164

9989  0.20  4.1028
9990  0.00  15.6332
9991  0.00  19.3932
9992  0.00  13.3200
9993  0.00  72.9480

[9994 rows x 13 columns]>
```

In [18]:

```
df.describe()
```

Out [18]:

	Postal Code	Sales	Quantity	Discount	Profit
count	9994.000000	9994.000000	9994.000000	9994.000000	9994.000000
mean	55190.379428	229.858001	3.789574	0.156203	28.656896
std	32063.693350	623.245101	2.225110	0.206452	234.260108
min	1040.000000	0.444000	1.000000	0.000000	-6599.978000
25%	23223.000000	17.280000	2.000000	0.000000	1.728750
50%	56430.500000	54.490000	3.000000	0.200000	8.666650
75%	90006.000000	209.940000	5.000000	0.200000	29.364000
max	99320.000000	22638.480000	14.000000	0.800000	8399.976000

In [19]:

```
for i in df.columns:
    print(i, len(df[i].unique()))
```

Out [19]:

```
Ship Mode 4
Segment 3
Country 1
City 531
State 49
Postal Code 631
Region 4
Category 3
Sub-Category 17
Sales 5825
Quantity 14
Discount 12
Profit 7287
```

In [20]:

```
df.isnull().sum()
```

Out [20]:

```
Ship Mode 0
Segment 0
Country 0
City 0
State 0
Postal Code 0
Region 0
Category 0
Sub-Category 0
Sales 0
Quantity 0
Discount 0
Profit 0
dtype: int64
```

In [21]:

```
sns.pairplot(df)
```

Out [21]:

In [22]:

```
fig, axes = plt.subplots(1,1,figsize=(12,7))
sns.heatmap(df.corr())
plt.show()
```

Out [22]:

In [23]:

```
fig, axes = plt.subplots(1,2,figsize=(14,5))
sns.barplot(data=df.groupby(['Sales','Profit']).agg(sum),xs='Sales',ys='Profit',ax=axes[1])
df.groupby('Sub-Category')['Sales','Profit'].agg(sum).plot(kind='bar',ax=axes[0])
plt.xticks(rotation=90)
plt.show()
```

Out [23]:

In [24]:

```
fig, axes = plt.subplots(1,2,figsize=(14,5))
fig.suptitle('Total Sales VS Quantity')
sns.barplot(data=df.groupby('Sub-Category')['Sales','Quantity'].agg(sum),xs='Sales',ys='Quantity',ax=axes[1])
df.groupby('Sub-Category')['Sales','Quantity'].agg(sum).plot(kind='bar',ax=axes[0])
plt.xticks(rotation=90)
plt.show()
```

Out [24]:

In [25]:

```
# computing top categories in terms of sales from first 100 observations
top_category_s = df.groupby('Category').Sales.sum().nlargest(n=100)
# computing top categories in terms of profit from first 100 observations
top_category_p = df.groupby('Category').Profit.sum().nlargest(n=100)

# plotting to see it visually
plt.style.use('seaborn')
top_category_s.plot(kind='bar',figsize=(10,5),fontsize=14)
top_category_p.plot(kind='bar',figsize=(10,5),fontsize=14,color='red')
plt.xlabel('Category',fontsize=15)
plt.ylabel('Total Sales/Profits',fontsize=15)
plt.title('Top Category Sales vs Profit',fontsize=15)
plt.show()
```

Out [25]:

In [26]:

```
fig, axes = plt.subplots(1,2,figsize=(14,5))
fig.suptitle('Discount & Profit Relation based on Sub-Category')
df.groupby('Sub-Category')['Discount','Profit'].agg(sum).plot(kind='bar',ax=axes[0]).set_title('Discount & Profit Relation based on Sub-Category')
df.groupby('Sub-Category')['Profit','Quantity'].agg(sum).plot(kind='bar',ax=axes[1]).set_title('Quantity & Profit Relation based on Sub-Category')
plt.xticks(rotation=90)
plt.show()
```

Out [26]:

In [27]:

```
# computing top sub-categories in terms of sales from first 100 observations
top_subcategory_s = df.groupby('Sub-Category').Sales.sum().nlargest(n=100)
# computing top sub-categories in terms of profit from first 100 observations
top_subcategory_p = df.groupby('Sub-Category').Profit.sum().nlargest(n=100)

# plotting to see it visually
plt.style.use('seaborn')
top_subcategory_s.plot(kind='bar',figsize=(10,5),fontsize=14)
top_subcategory_p.plot(kind='bar',figsize=(10,5),fontsize=14,color='red')
plt.xlabel('Sub-Category',fontsize=15)
plt.ylabel('Total Sales/Profits',fontsize=15)
plt.title('Top Sub-Category Sales vs Profit',fontsize=15)
plt.show()
```

Out [27]:

In [28]:

```
fig, axes = plt.subplots(2,2,figsize=(16,8))
fig.suptitle('Distribution plots',fontsize=16)
sns.distplot(df['Sales'],ax=axes[0,0])
sns.distplot(df['Profit'],ax=axes[0,1])
sns.distplot(df['Discount'],ax=axes[1,0])
sns.distplot(df['Quantity'],ax=axes[1,1])
plt.show()
```

Out [28]:

In [29]:

```
plt.figure(figsize=(6,6))
plt.title('Region')
plt.pie(df['Region'].value_counts(), labels=df['Region'].value_counts().index,autopct='%1.1f%%')
plt.show()
```

Out [29]:

In [30]:

```
fig, axes = plt.subplots(2,2,figsize=(16,8))
fig.suptitle('Sales with different shipping modes and Segments',fontsize=16)
sns.barplot(df['Ship Mode'],df['Sales'],ax=axes[0,0])
sns.lineplot(df['Ship Mode'],df['Sales'],ax=axes[0,1])
sns.barplot(df['Segment'],df['Sales'],ax=axes[1,0])
sns.lineplot(df['Segment'],df['Sales'],ax=axes[1,1])
plt.show()
```

Out [30]:

In [31]:

```
fig, ax = plt.subplots(1,1,figsize=(12,7))
sns.countplot(df['Quantity'],hue=df['Region'])
plt.show()
```

Out [31]:

The weak areas where one can work to make more profit are : ##### We should limit sales of furniture and increase that of technology and office suppliers as furniture has very less profit as compared to sales. ##### Considering the sub-categories sales of tables should be minimized, as it increase sales more in the east as profit is more. ##### We should concentrate on the states like New York and California to make more profits. ##### The features Profit and Discount are highly related ##### Over Less quantity of products also the sales were high. ##### The maximum quantity of product in demand was in range 2-4. ##### The mode of shipping doesn't affect much to the sales. ##### The Home Office provides highest sales followed by Corporate by a slight variation

In []: