# TASK 3: EXPLORATORY DATA ANALYSIS

AUTHOR: VANIPENTA BALAJI GRIP JUNE 2023 THE SPARKS FOUNDATIONS

In [1]:

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

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

## Importing the Data form Samplesuperstore

In [2]:

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

Out[2]:

| | Ship Mode | Segment | Country | City | State | Postal Code | Region | Category | Sub-Category | S |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Second Class | Consumer | United States | Henderson | Kentucky | 42420 | South | Furniture | Bookcases | 261.9 |
| 1 | Second Class | Consumer | United States | Henderson | Kentucky | 42420 | South | Furniture | Chairs | 731.9 |
| 2 | Second Class | Corporate | United States | Los Angeles | California | 90036 | West | Office Supplies | Labels | 14.6 |
| 3 | Standard Class | Consumer | United States | Fort Lauderdale | Florida | 33311 | South | Furniture | Tables | 957.5 |
| 4 | Standard Class | Consumer | United States | Fort Lauderdale | Florida | 33311 | South | Office Supplies | Storage | 22.3 |

In [3]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 13 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   Ship Mode    9994 non-null   object
 1   Segment      9994 non-null   object
 2   Country      9994 non-null   object
 3   City         9994 non-null   object
 4   State        9994 non-null   object
 5   Postal Code  9994 non-null   int64
 6   Region       9994 non-null   object
 7   Category     9994 non-null   object
 8   Sub-Category 9994 non-null   object
 9   Sales        9994 non-null   float64
 10  Quantity     9994 non-null   int64
 11  Discount     9994 non-null   float64
 12  Profit       9994 non-null   float64
dtypes: float64(3), int64(2), object(8)
```

memory usage: 1015.1+ KB

In [4]:
```python
df.describe()
```

Out[4]:

|  | 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 | 8.666500 |
| 75% | 90008.000000 | 209.940000 | 5.000000 | 0.200000 | 29.364000 |
| max | 99301.000000 | 22638.480000 | 14.000000 | 0.800000 | 8399.976000 |

In [5]:
```python
for i in df.columns:
    print(i,len(df[i].unique()))
```

```
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 [6]:
```python
df.isnull().sum()
```

Out[6]:
```
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 [7]:
```python
df.nunique()
```

Out[7]:
```
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
dtype: int64
```
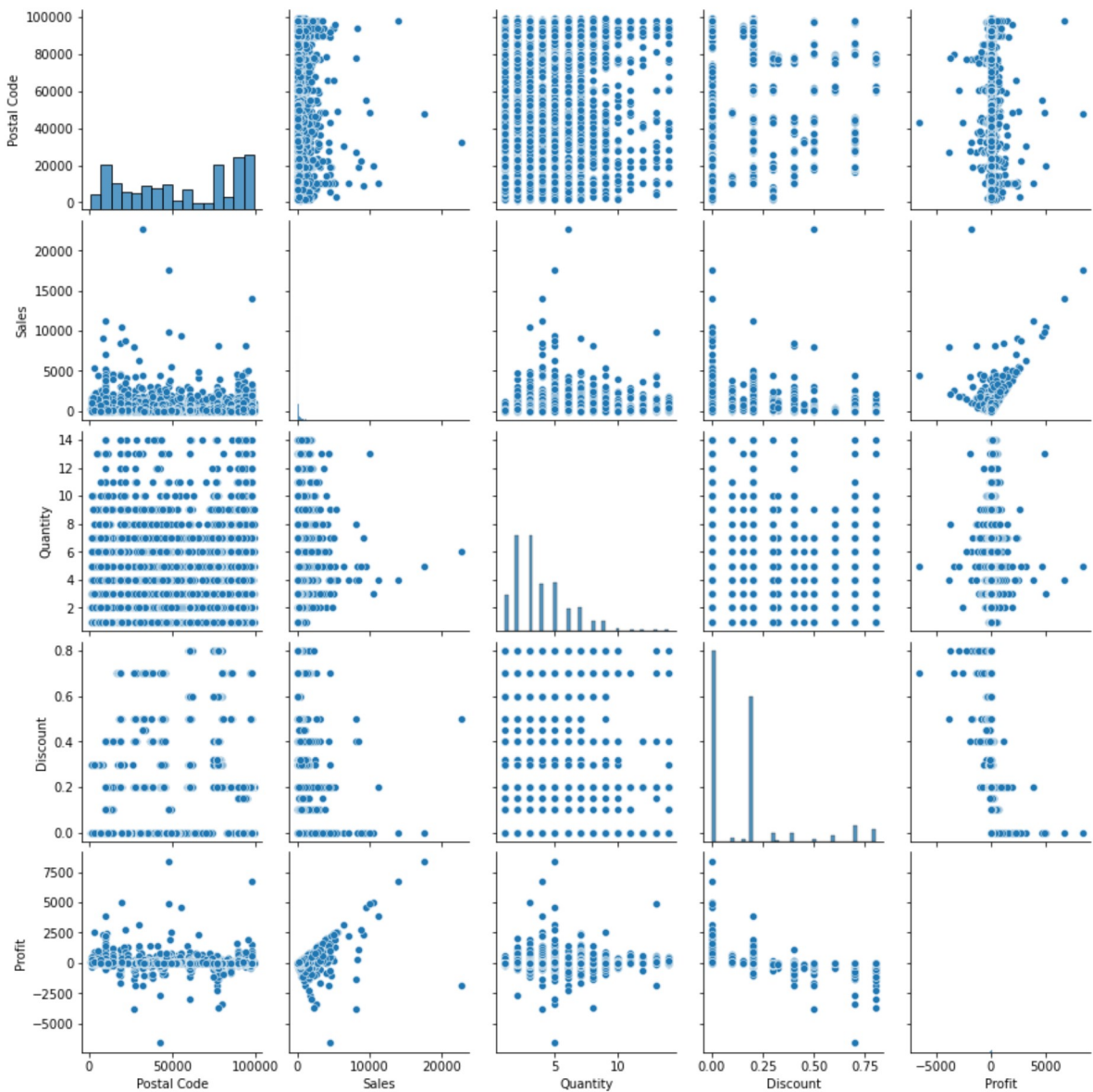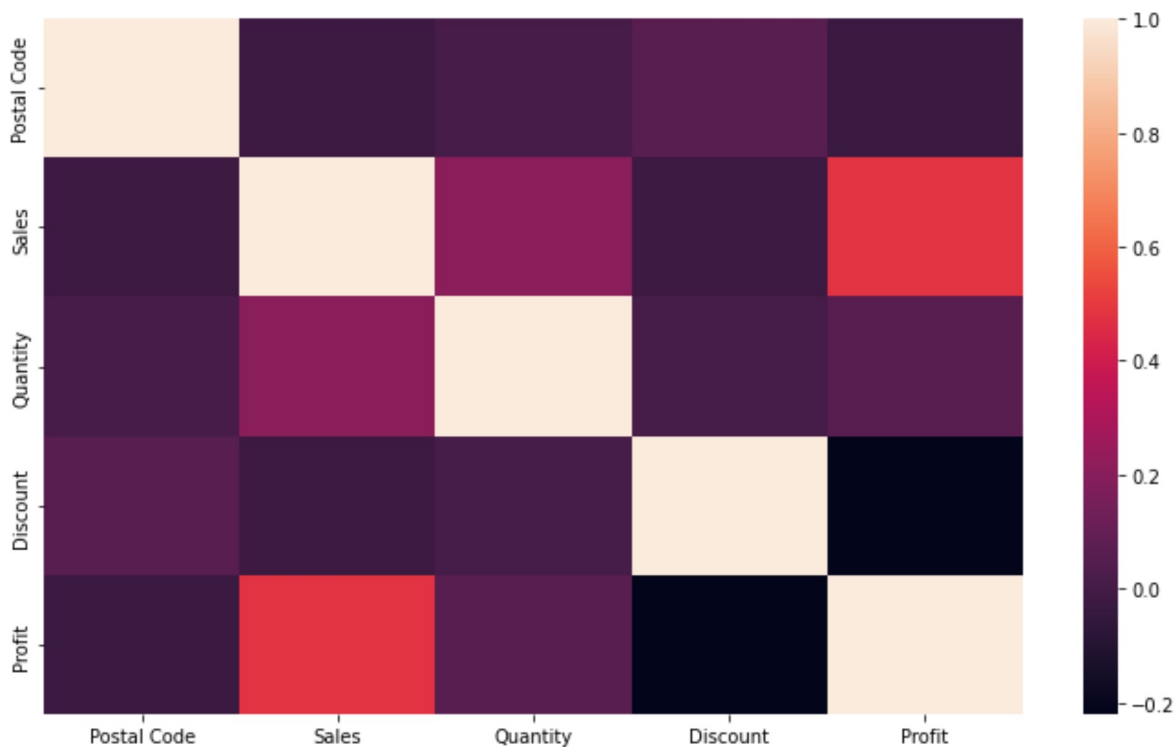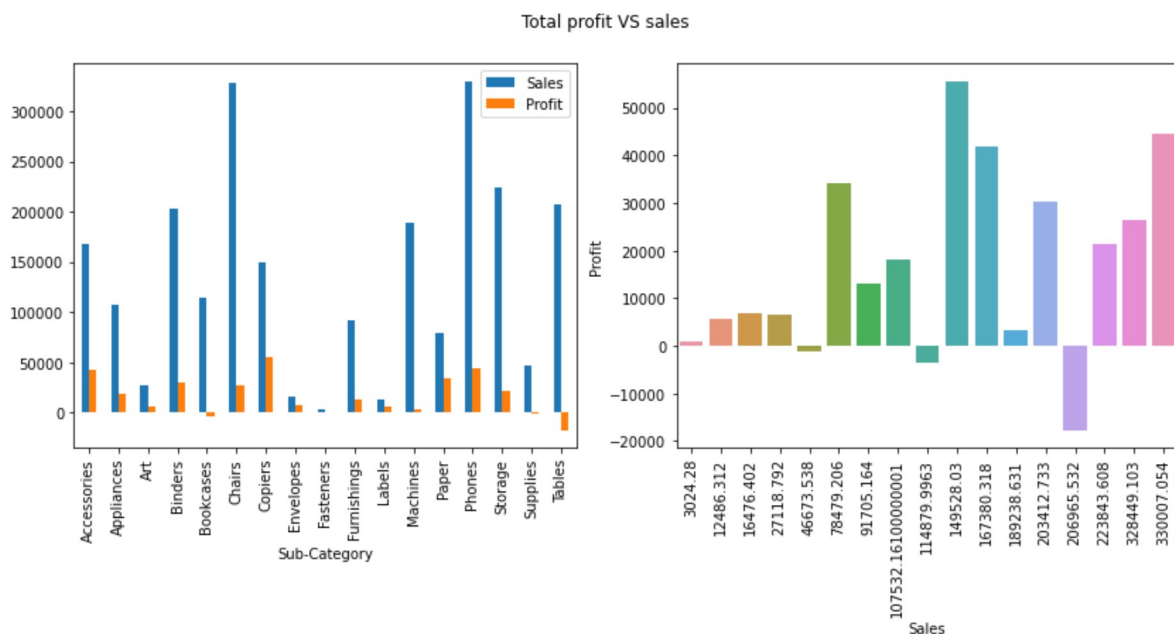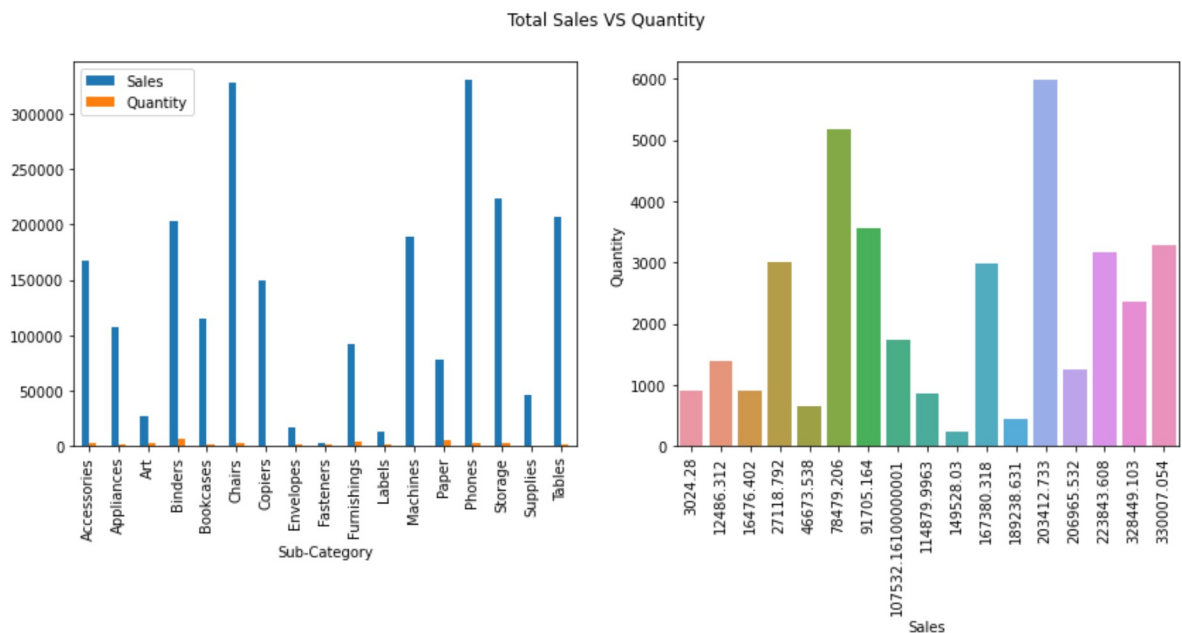
# Data Visualization

In [8]:
```python
sns.pairplot(df);
```



In [9]:
```python
fig,axes = plt.subplots(1,1,figsize=(12,7))
sns.heatmap(df.corr())
plt.show()
```
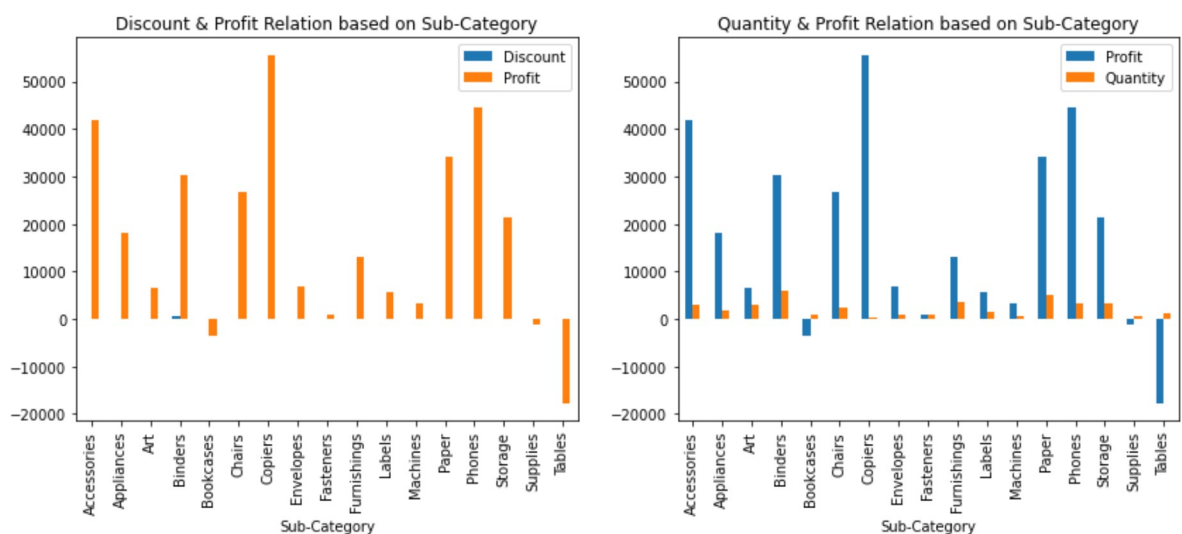
```
In [10]:  fig,axes = plt.subplots(1,2,figsize=(14,5))
          fig.suptitle("Total profit VS sales ")
          sns.barplot(data=df.groupby('Sub-Category')['Sales','Profit'].agg(sum),x='Sales',y
          df.groupby('Sub-Category')['Sales','Profit'].agg(sum).plot(kind='bar',ax=axes[0])
          plt.xticks(rotation=90)
          plt.show()
```



```
In [11]:  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),x='Sales
          df.groupby('Sub-Category')['Sales','Quantity'].agg(sum).plot(kind='bar',ax=axes[0]
          plt.xticks(rotation=90)
          plt.show()
```
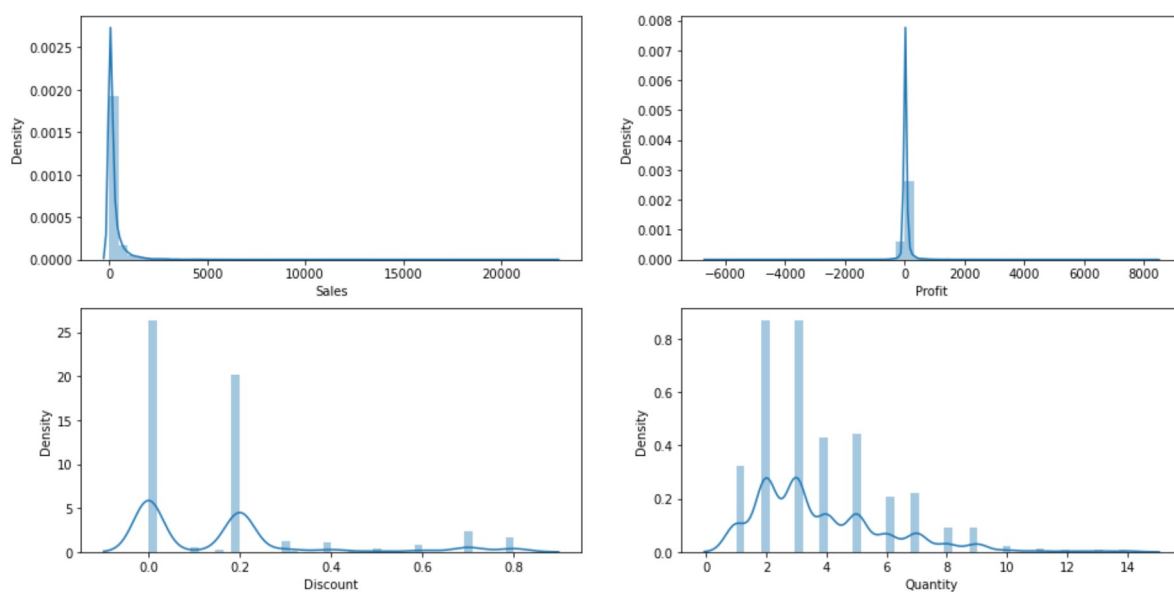
Total Sales VS Quantity



```
In [12]: fig,axes = plt.subplots(1,2,figsize=(14,5))
         df.groupby('Sub-Category')['Discount','Profit'].agg(sum).plot(kind='bar',ax=axes[0
         df.groupby('Sub-Category')['Profit','Quantity'].agg(sum).plot(kind='bar',ax=axes[1
         plt.xticks(rotation=90)
         plt.show()
```



```
In [13]: 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()
```
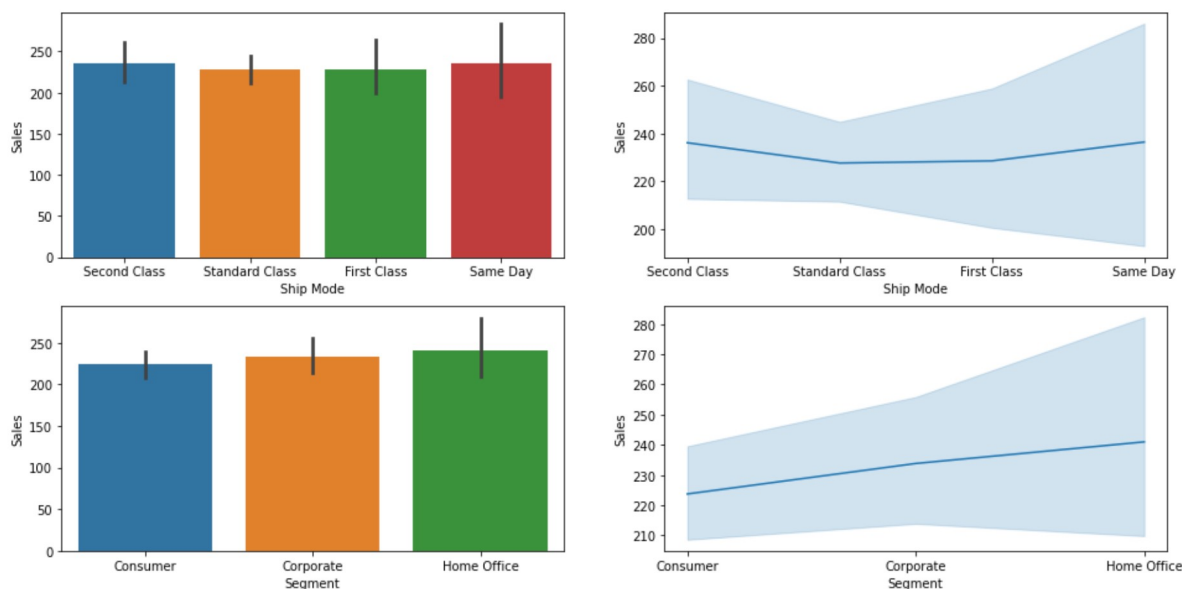
Distribution plots



```
In [14]:   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()
```
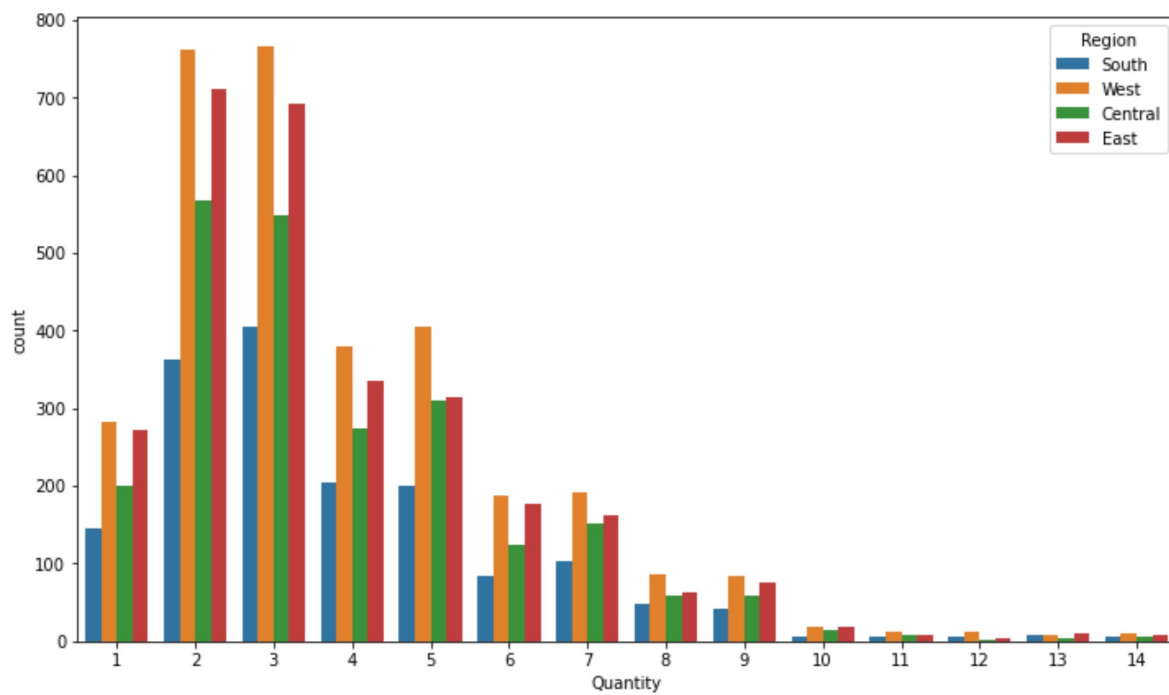
Sales with different shipping modes and Segments



```
In [15]:   fig,ax= plt.subplots(1,1,figsize=(12,7))
           sns.countplot(df['Quantity'],hue=df['Region'])
           plt.show()
```

In [ ]: