

MINI PROJECT TO CLEAN AND BUILD A DASH BOARD

In [113]: 1 *#The aim of this project is to clean and preprocess the data, uncover the insights in*

In [114]: 1 *# importing the data to notebook*
 2 **import** pandas **as** pd
 3 **import** numpy **as** np
 4 **import** os
 5 **import** matplotlib.pyplot **as** plt
 6 **import** seaborn **as** sns

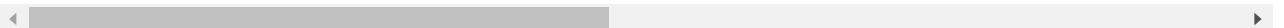
In [115]: 1 df = pd.read_csv(r"C:\Users\USER\Documents\WORKSPACE\Visualization\Data\superstore_sales.csv")

In [116]: 1 df

Out[116]:

	order_id	order_date	ship_date	ship_mode	customer_name	segment	state	country	market
0	AG-2011-2040	1/1/2011	1/6/2011	Standard Class	Toby Braunhardt	Consumer	Constantine	Algeria	Africa
1	IN-2011-47883	1/1/2011	1/8/2011	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC C
2	HU-2011-1220	1/1/2011	1/5/2011	Second Class	Annie Thurman	Consumer	Budapest	Hungary	EMEA
3	IT-2011-3647632	1/1/2011	1/5/2011	Second Class	Eugene Moren	Home Office	Stockholm	Sweden	EU
4	IN-2011-47883	1/1/2011	1/8/2011	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC C
...
51399	CA-2014-115427	12/31/2014	1/4/2015	Standard Class	Erica Bern	Corporate	California	United States	US
51400	MO-2014-2560	12/31/2014	1/5/2015	Standard Class	Liz Preis	Consumer	Souss-Massa-Draâ	Morocco	Africa
51401	MX-2014-110527	12/31/2014	1/2/2015	Second Class	Charlotte Melton	Consumer	Managua	Nicaragua	LATAM
51402	MX-2014-114783	12/31/2014	1/6/2015	Standard Class	Tamara Dahlen	Consumer	Chihuahua	Mexico	LATAM
51403	CA-2014-156720	12/31/2014	1/4/2015	Standard Class	Jill Matthias	Consumer	Colorado	United States	US

51404 rows × 21 columns



```
In [117]: 1 #inspecting the structure and properties of our data .
          2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51404 entries, 0 to 51403
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   order_id              51404 non-null  object
1   order_date            51404 non-null  object
2   ship_date             51404 non-null  object
3   ship_mode             51404 non-null  object
4   customer_name         51404 non-null  object
5   segment               51404 non-null  object
6   state                 51404 non-null  object
7   country               51404 non-null  object
8   market                51404 non-null  object
9   region                51404 non-null  object
10  product_id            51404 non-null  object
11  category              51404 non-null  object
12  sub_category          51404 non-null  object
13  product_name          51404 non-null  object
14  sales                 51399 non-null  object
15  quantity              51401 non-null  float64
16  discount              51400 non-null  float64
17  profit                51401 non-null  float64
18  shipping_cost         51402 non-null  float64
19  order_priority        51404 non-null  object
20  year                  51404 non-null  int64
dtypes: float64(4), int64(1), object(16)
memory usage: 8.2+ MB
```

```
In [118]: 1 # To get the List of my columns to know what is going on in my data
          2 df.columns
```

```
Out[118]: Index(['order_id', 'order_date', 'ship_date', 'ship_mode', 'customer_name',
                  'segment', 'state', 'country', 'market', 'region', 'product_id',
                  'category', 'sub_category', 'product_name', 'sales', 'quantity',
                  'discount', 'profit', 'shipping_cost', 'order_priority', 'year'],
                  dtype='object')
```

```
In [119]: 1 # Getting quantity columns and checking the nature of the data in that column
          2 df['quantity']
```

```
Out[119]: 0      2.0
          1      3.0
          2      4.0
          3      3.0
          4      5.0
          ...
        51399    2.0
        51400    1.0
        51401    3.0
        51402    1.0
        51403    3.0
Name: quantity, Length: 51404, dtype: float64
```

```
In [120]: 1 # Getting discount columns and checking the nature of the data in that column
          2 df['discount']
```

```
Out[120]: 0      0.0
          1      0.1
          2      0.0
          3      0.5
          4      0.1
          ...
          51399  0.2
          51400  0.0
          51401  0.0
          51402  0.0
          51403  0.2
          Name: discount, Length: 51404, dtype: float64
```

```
In [121]: 1 # Getting year columns and checking the nature of the data in that column
          2 df['year']
```

```
Out[121]: 0      2011
          1      2011
          2      2011
          3      2011
          4      2011
          ...
          51399  2014
          51400  2014
          51401  2014
          51402  2014
          51403  2014
          Name: year, Length: 51404, dtype: int64
```

```
In [122]: 1 # Inspecting the column for their unique values
          2 df['profit']
```

```
Out[122]: 0      106.1400
          1      36.0360
          2      29.6400
          3     -26.0550
          4      37.7700
          ...
          51399    4.5188
          51400    0.4200
          51401    12.3600
          51402    0.5600
          51403   -0.6048
          Name: profit, Length: 51404, dtype: float64
```

```
In [158]: 1 #checking for unique value for sales
          2 df['sales'].unique()
```

```
Out[158]: array([ 408., 120., 66., ..., 1763., 1821., 1831.])
```

```
In [125]: 1 # converting the sales colounm from strings to float
```

```
In [126]: 1 df.sales.unique()
```

```
Out[126]: array(['408', '120', '66', ..., '1,763', '1,821', '1,831'], dtype=object)
```

```
In [127]: 1 df['sales'] = df['sales'].str.replace(',', '').astype(float)
```

```
In [128]: 1 # checking to see that the sales column is converted to float
          2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51404 entries, 0 to 51403
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   order_id              51404 non-null  object
1   order_date            51404 non-null  object
2   ship_date             51404 non-null  object
3   ship_mode             51404 non-null  object
4   customer_name         51404 non-null  object
5   segment               51404 non-null  object
6   state                 51404 non-null  object
7   country               51404 non-null  object
8   market                51404 non-null  object
9   region                51404 non-null  object
10  product_id            51404 non-null  object
11  category              51404 non-null  object
12  sub_category          51404 non-null  object
13  product_name          51404 non-null  object
14  sales                  51399 non-null  float64
15  quantity              51401 non-null  float64
16  discount               51400 non-null  float64
17  profit                51401 non-null  float64
18  shipping_cost          51402 non-null  float64
19  order_priority         51404 non-null  object
20  year                  51404 non-null  int64
dtypes: float64(5), int64(1), object(15)
memory usage: 8.2+ MB
```

```
In [129]: 1 #checking for missing values
          2 df.isna().sum()
```

```
Out[129]: 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         5
quantity      3
discount      4
profit        3
shipping_cost 2
order_priority 0
year          0
dtype: int64
```

```
In [130]: 1 #dropping the missing values because it will affect the outcome of my work
          2 df=df.dropna()
```

```
In [131]: 1 df.isna().sum()
```

```
Out[131]: 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 [132]: 1 #checking for duplicate
          2 df.duplicated().sum()
```

```
Out[132]: 114
```

```
In [133]: 1 #dropping duplicate
          2 df = df.drop_duplicates()
```

```
In [136]: 1 #Grouping the data for total sales
          2 df["Total_sales"] = df["sales"] * df["quantity"]
```

```
In [137]: 1 #checking to see if the total sales have been added to the coloum
          2 print(df.columns)
```

```
Index(['order_id', 'order_date', 'ship_date', 'ship_mode', 'customer_name',
      'segment', 'state', 'country', 'market', 'region', 'product_id',
      'category', 'sub_category', 'product_name', 'sales', 'quantity',
      'discount', 'profit', 'shipping_cost', 'order_priority', 'year',
      'Total_sales'],
      dtype='object')
```

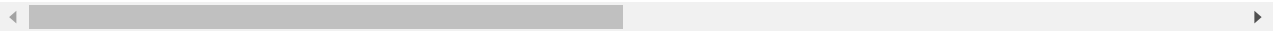
```
In [138]: 1 #Grouping the total sales by years
          2 sales_df = df.groupby("year")["Total_sales"].sum().reset_index()
          3 profit_df = df.groupby("year")["profit"].sum().reset_index()
```

In [139]: 1 df.head()

Out[139]:

	order_id	order_date	ship_date	ship_mode	customer_name	segment	state	country	market	region
0	AG-2011-2040	1/1/2011	1/6/2011	Standard Class	Toby Braunhardt	Consumer	Constantine	Algeria	Africa	Africa
1	IN-2011-47883	1/1/2011	1/8/2011	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC	Oceania
2	HU-2011-1220	1/1/2011	1/5/2011	Second Class	Annie Thurman	Consumer	Budapest	Hungary	EMEA	EMEA
3	IT-2011-3647632	1/1/2011	1/5/2011	Second Class	Eugene Moren	Home Office	Stockholm	Sweden	EU	North America
4	IN-2011-47883	1/1/2011	1/8/2011	Standard Class	Joseph Holt	Consumer	New South Wales	Australia	APAC	Oceania

5 rows × 22 columns

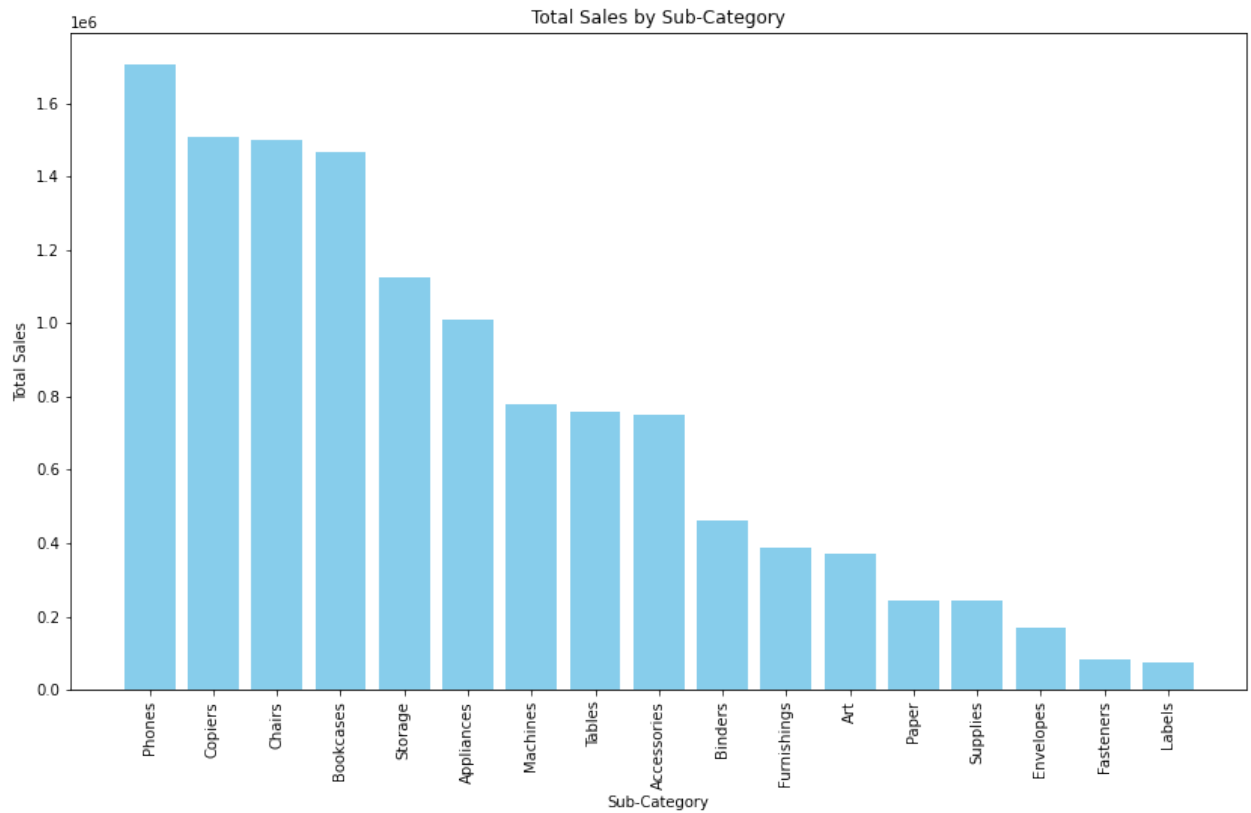


In [27]:

```

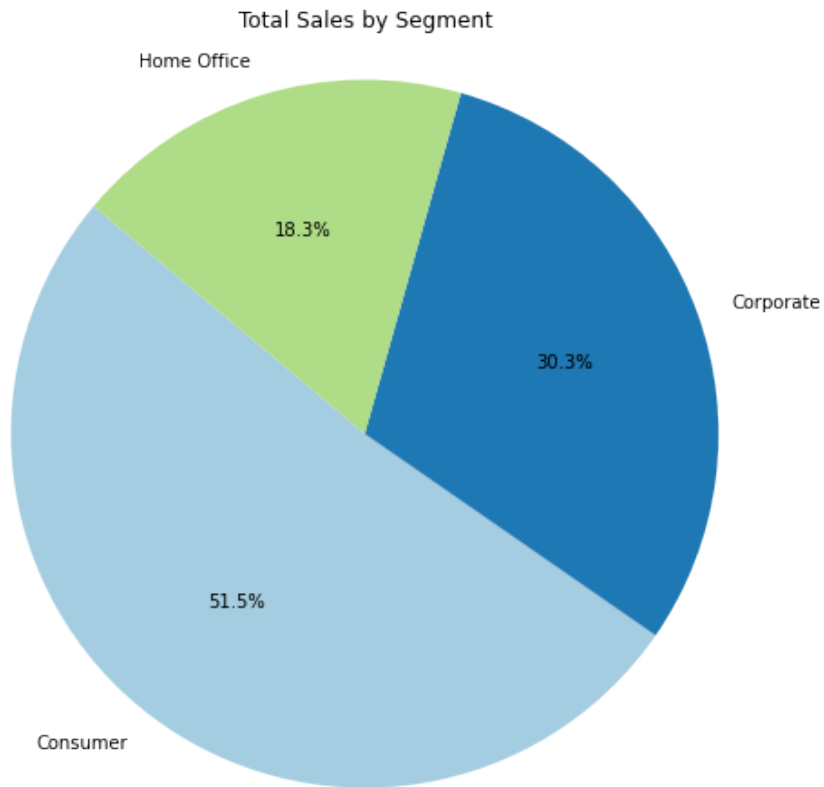
1 # Group by sub-category and sum the sales
2 total_sales_by_sub_category = df.groupby('sub_category')['sales'].sum().reset_index()
3
4 # Sort the data by total sales for better visualization
5 total_sales_by_sub_category = total_sales_by_sub_category.sort_values(by='sales', ascending=False)
```

```
In [140]: 1 # plotting the total sales by sub-category
2 plt.figure(figsize=(14, 8))
3 plt.bar(total_sales_by_sub_category['sub_category'], total_sales_by_sub_category['sales'])
4 plt.title('Total Sales by Sub-Category')
5 plt.xlabel('Sub-Category')
6 plt.ylabel('Total Sales')
7 plt.xticks(rotation=90)
8 plt.show()
```



```
In [141]: 1 # Group by segment and sum the sales
2 total_sales_by_segment = df.groupby('segment')['sales'].sum().reset_index()
```

```
In [142]: 1 # Plotting the data
2 plt.figure(figsize=(10, 8))
3 plt.pie(total_sales_by_segment['sales'], labels=total_sales_by_segment['segment'], autopct='%1.1f%%')
4 plt.title('Total Sales by Segment')
5 plt.axis('equal')
6 plt.show()
```



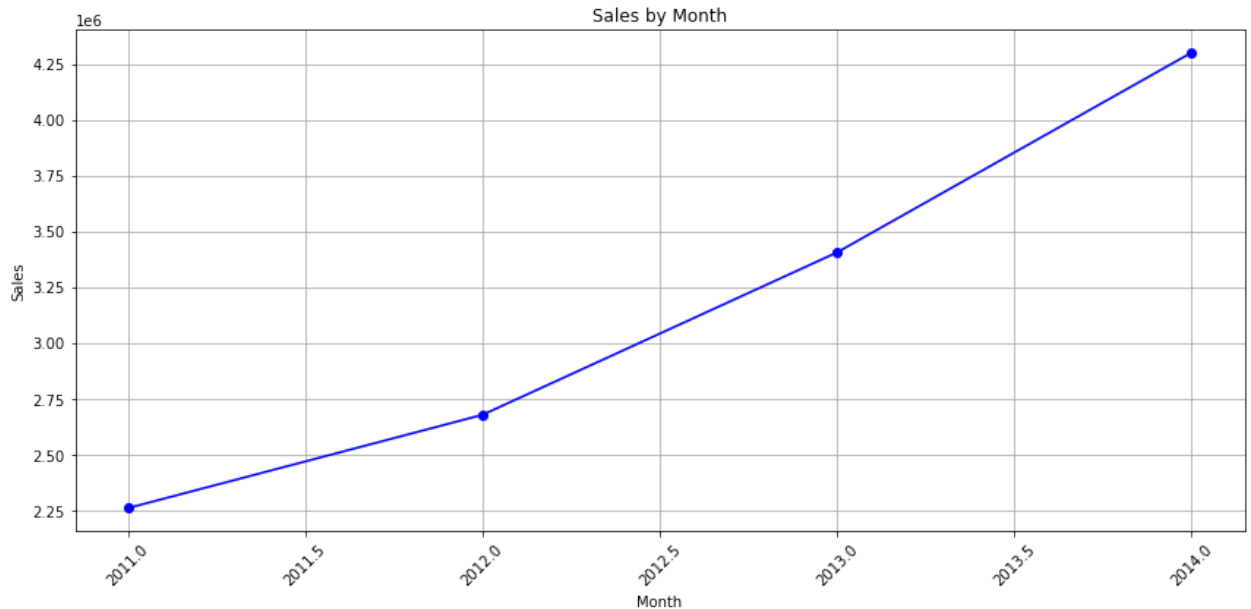
```
In [143]: 1 # Aggregate sales data by month
2 monthly_sales = df.groupby('year')['sales'].sum().reset_index()
```


In [144]:

```

1 # Plotting the data
2 plt.figure(figsize=(12, 6))
3 plt.plot(monthly_sales['year'], monthly_sales['sales'], marker='o', linestyle='-', color='blue')
4 plt.title('Sales by Month')
5 plt.xlabel('Month')
6 plt.ylabel('Sales')
7 plt.xticks(rotation=45)
8 plt.grid(True)
9 plt.tight_layout()
10 plt.show()

```



In [145]:

```

1 # Convert the order_date column to datetime
2 df['order_date'] = pd.to_datetime(df['order_date'])

```

In [146]:

```

1 # Extract the month and year from the date column
2 df['year_month'] = df['order_date'].dt.to_period('M')

```

In [147]:

```

1 # Aggregate sales data by month
2 monthly_sales = df.groupby('year_month')['sales'].sum().reset_index()

```

In [148]:

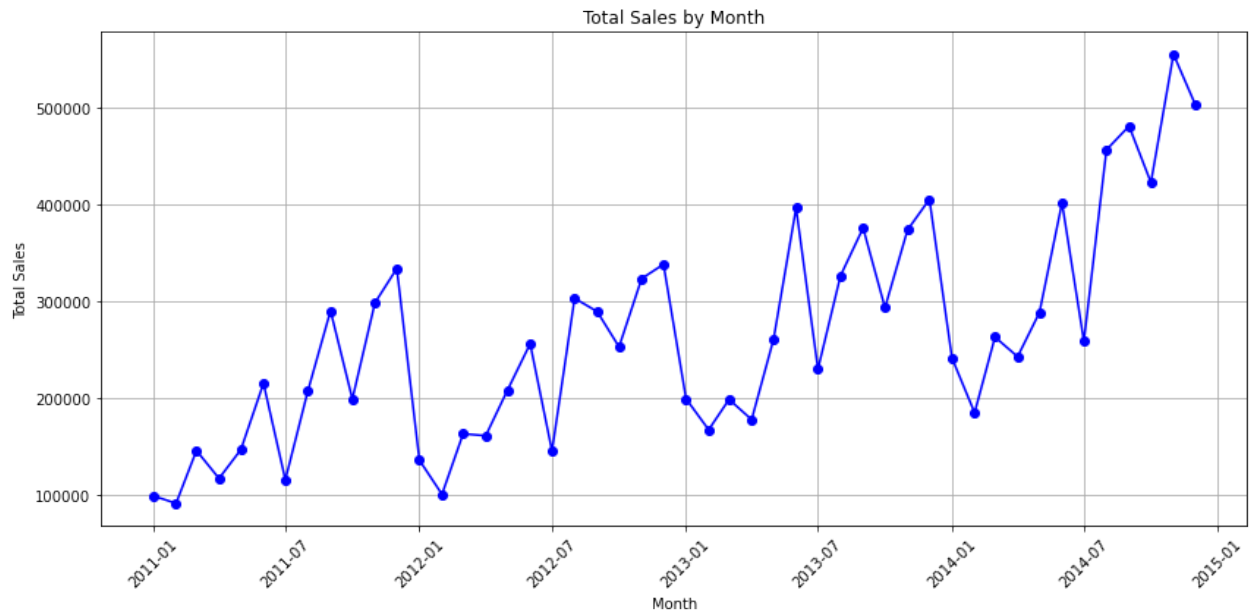
```

1 # Convert the 'year_month' to a datetime object for plotting
2 monthly_sales['year_month'] = monthly_sales['year_month'].dt.to_timestamp()

```

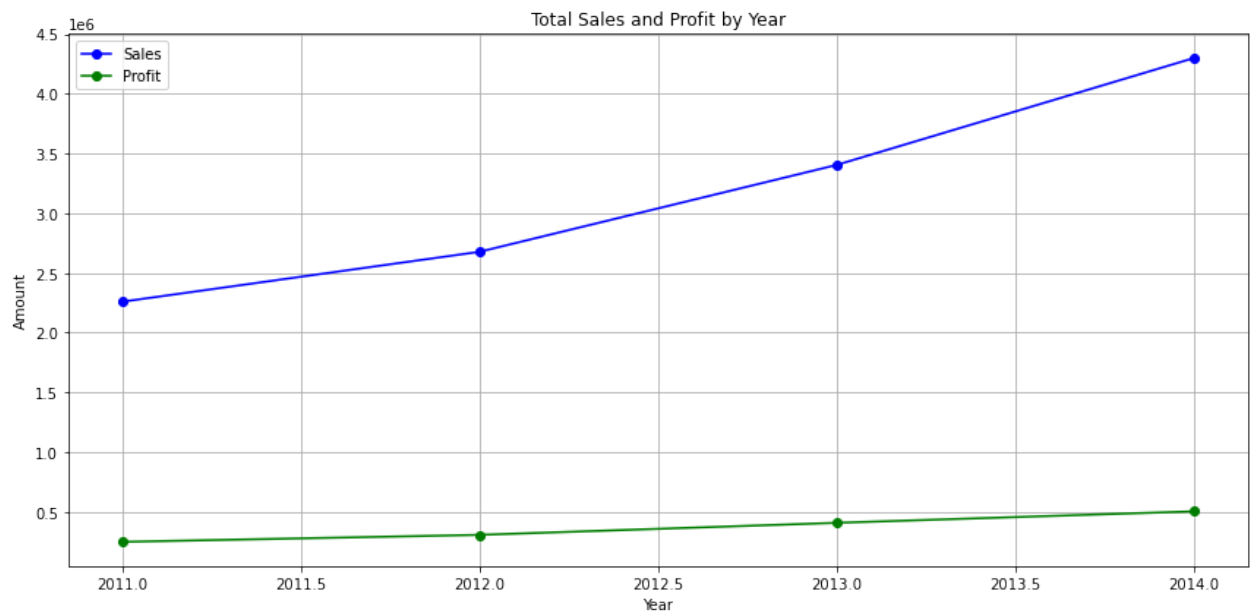
In [149]:

```
1 # Plotting the data
2 plt.figure(figsize=(12, 6))
3 plt.plot(monthly_sales['year_month'], monthly_sales['sales'], marker='o', linestyle='--')
4 plt.title('Total Sales by Month')
5 plt.xlabel('Month')
6 plt.ylabel('Total Sales')
7 plt.xticks(rotation=45)
8 plt.grid(True)
9 plt.tight_layout()
10 plt.show()
```

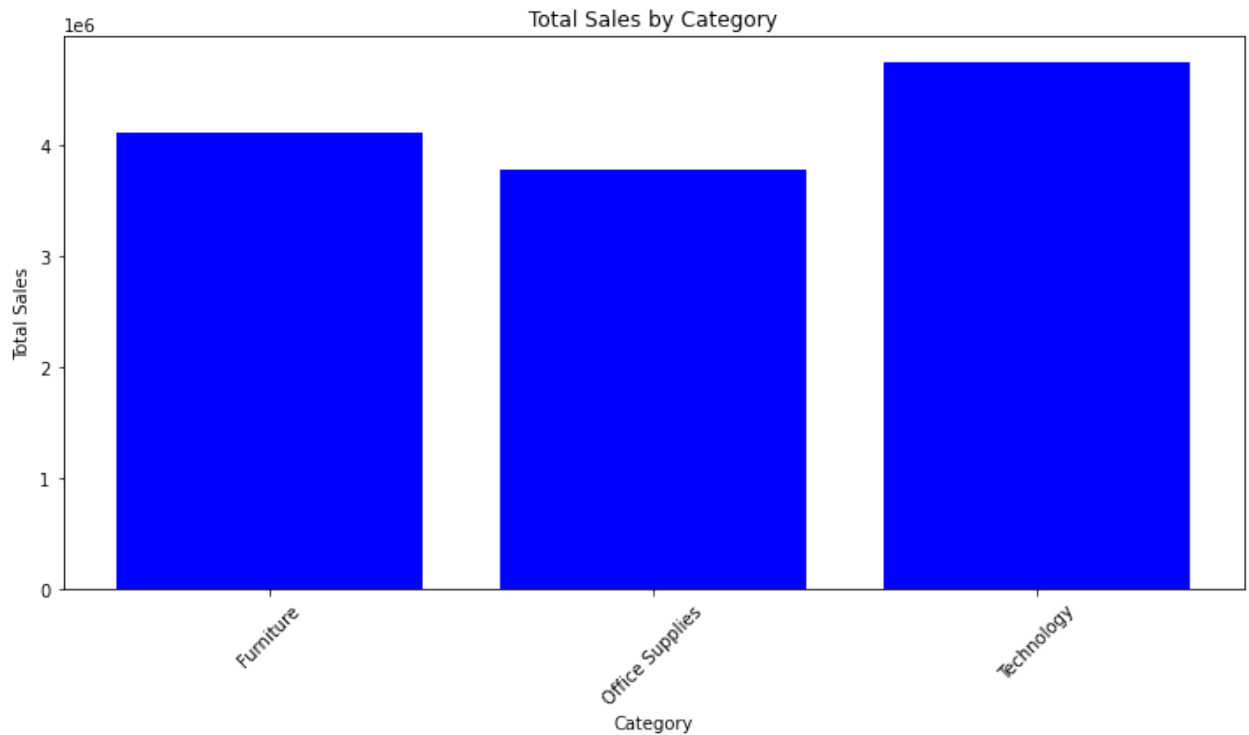


In [150]:

```
1 # Convert the order_date column to datetime
2 df['order_date'] = pd.to_datetime(df['order_date'])
3
4 # Extract the year from the date column
5 df['year'] = df['order_date'].dt.year
6
7 # Aggregate sales and profit data by year
8 yearly_data = df.groupby('year')[['sales', 'profit']].sum().reset_index()
9
10 # Plot the data
11 plt.figure(figsize=(12, 6))
12 plt.plot(yearly_data['year'], yearly_data['sales'], marker='o', linestyle='-', color='blue')
13 plt.plot(yearly_data['year'], yearly_data['profit'], marker='o', linestyle='-', color='green')
14 plt.title('Total Sales and Profit by Year')
15 plt.xlabel('Year')
16 plt.ylabel('Amount')
17 plt.legend()
18 plt.grid(True)
19 plt.tight_layout()
20 plt.show()
```

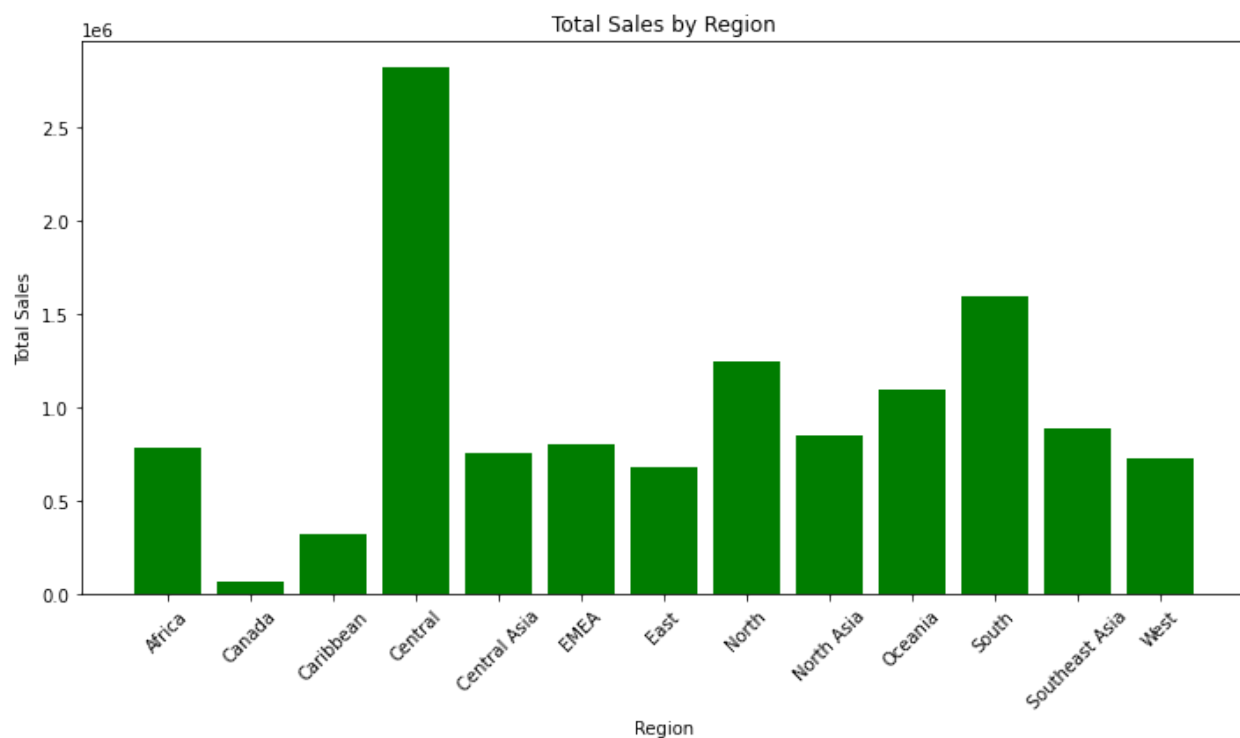


```
In [151]: 1 #Aggregate total sales by category
2 category_sales = df.groupby('category')['sales'].sum().reset_index()
3
4 # Plotting the data
5 plt.figure(figsize=(10, 6))
6 plt.bar(category_sales['category'], category_sales['sales'], color='blue')
7 plt.title('Total Sales by Category')
8 plt.xlabel('Category')
9 plt.ylabel('Total Sales')
10 plt.xticks(rotation=45)
11 plt.tight_layout()
12 plt.show()
```

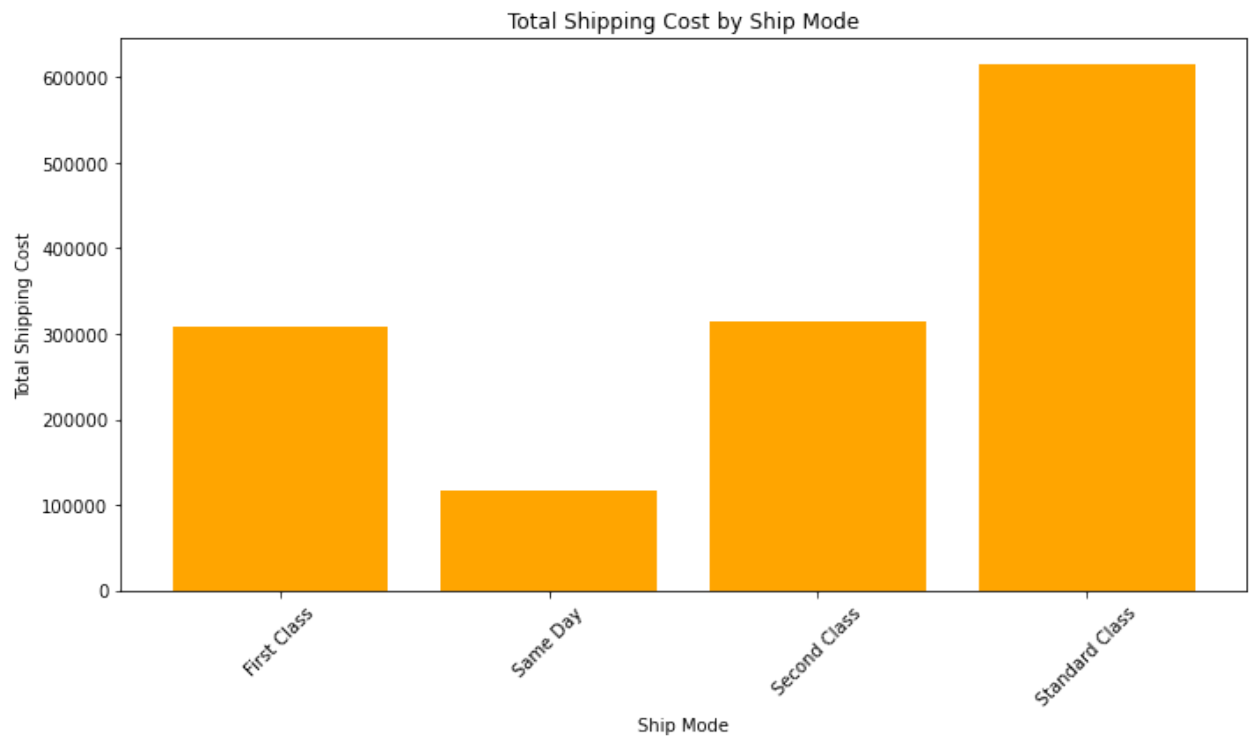


In [152]:

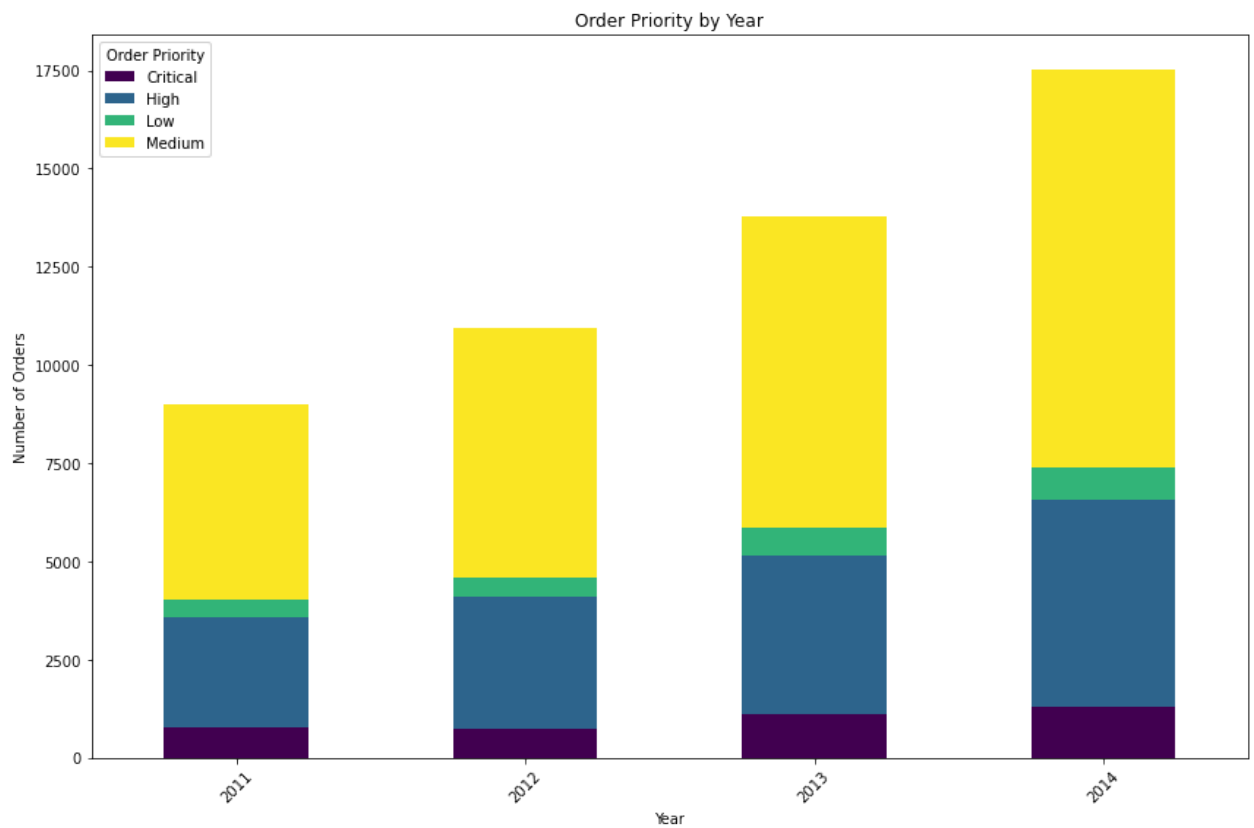
```
1 # Aggregate total sales by region
2 region_sales = df.groupby('region')['sales'].sum().reset_index()
3
4 # Plot the data
5 plt.figure(figsize=(10, 6))
6 plt.bar(region_sales['region'], region_sales['sales'], color='green')
7 plt.title('Total Sales by Region')
8 plt.xlabel('Region')
9 plt.ylabel('Total Sales')
10 plt.xticks(rotation=45)
11 plt.tight_layout()
12 plt.show()
```



```
In [153]: 1 # Aggregate total shipping cost by ship mode
2 ship_mode_cost = df.groupby('ship_mode')['shipping_cost'].sum().reset_index()
3
4 # Plot the data
5 plt.figure(figsize=(10, 6))
6 plt.bar(ship_mode_cost['ship_mode'], ship_mode_cost['shipping_cost'], color='orange')
7 plt.title('Total Shipping Cost by Ship Mode')
8 plt.xlabel('Ship Mode')
9 plt.ylabel('Total Shipping Cost')
10 plt.xticks(rotation=45)
11 plt.tight_layout()
12 plt.show()
```



```
In [154]: 1 # Aggregate the count of orders by year and order priority
2 order_priority_by_year = df.groupby(['year', 'order_priority']).size().unstack(fill_value=0)
3
4 # Plot the data
5 order_priority_by_year.plot(kind='bar', stacked=True, figsize=(12, 8), colormap='viridis')
6 plt.title('Order Priority by Year')
7 plt.xlabel('Year')
8 plt.ylabel('Number of Orders')
9 plt.legend(title='Order Priority')
10 plt.xticks(rotation=45)
11 plt.tight_layout()
12 plt.show()
```



```
In [ ]: 1 # CREATING DASHBOARD
```

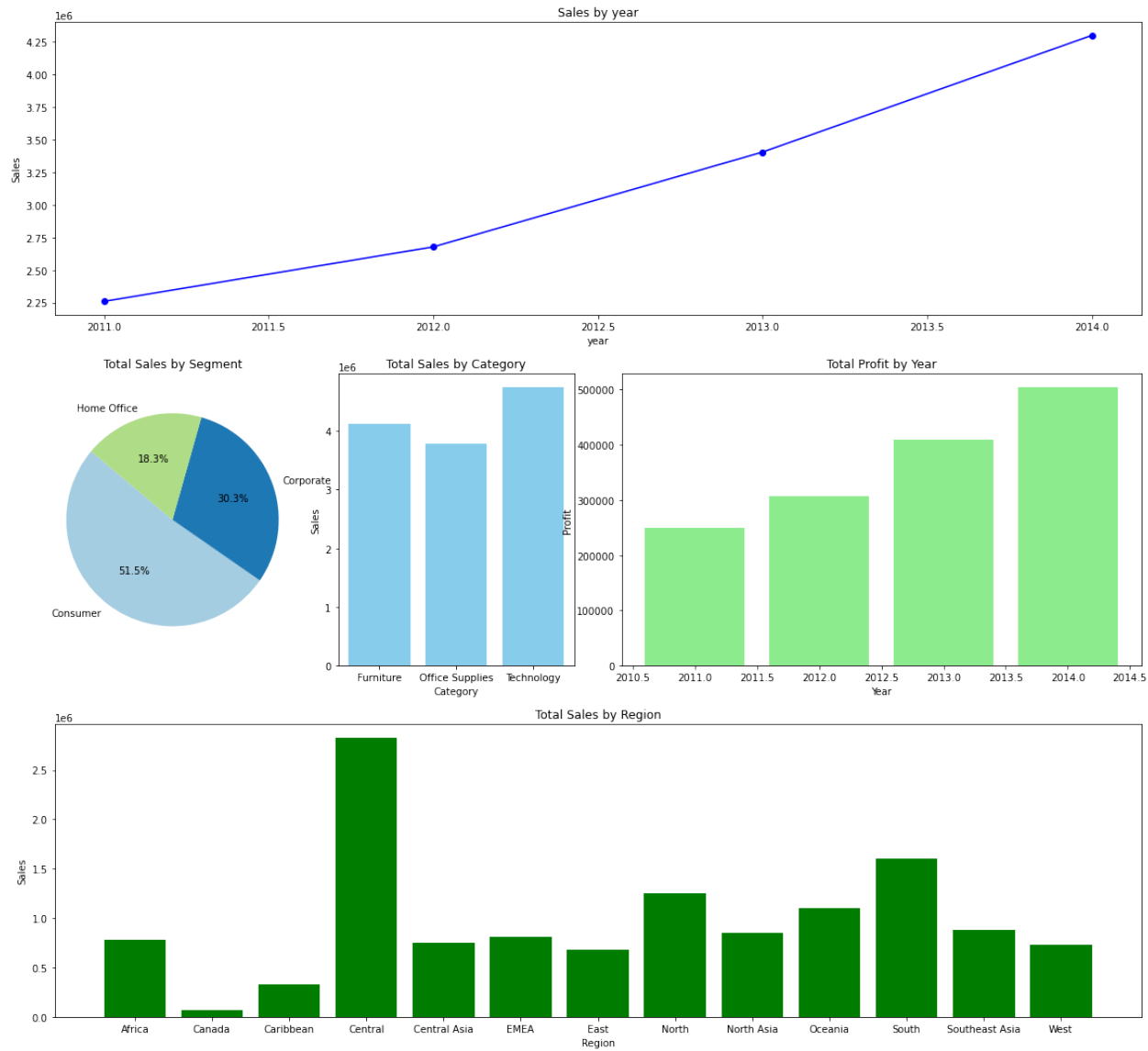
In [156]:

```

1  fig = plt.figure(figsize=(30,25))
2
3  ax = [None for _ in range(8)]
4
5  ax[0] = plt.subplot2grid((4,6), (0,0), colspan=4)
6  ax[1] = plt.subplot2grid((4,6), (1,0), colspan=1)
7  ax[2] = plt.subplot2grid((4,6), (1,1), colspan=1)
8  ax[3] = plt.subplot2grid((4,6), (1,2), colspan=2)
9  ax[4] = plt.subplot2grid((4,6), (2,3), colspan=3)
10 ax[5] = plt.subplot2grid((4,6), (2,0), colspan=4)
11
12
13
14 # to Load content into the content
15 sales_by_month = df.groupby('year')['sales'].sum().reset_index()
16 ax[0].plot(sales_by_month['year'], sales_by_month['sales'], marker='o', color='b')
17 ax[0].set_title('Sales by year')
18 ax[0].set_xlabel('year')
19 ax[0].set_ylabel('Sales')
20
21
22 ax[1].pie(total_sales_by_segment['sales'], labels=total_sales_by_segment['segment'], ax
23 ax[1].set_title('Total Sales by Segment')
24 ax[1].axis('equal')
25
26 ax[2].bar(total_sales_by_category['category'], total_sales_by_category['sales'], color=
27 ax[2].set_title('Total Sales by Category')
28 ax[2].set_xlabel('Category')
29 ax[2].set_ylabel('Sales')
30
31 ax[3].bar(profit_by_year['year'], profit_by_year['profit'], color='lightgreen')
32 ax[3].set_title('Total Profit by Year')
33 ax[3].set_xlabel('Year')
34 ax[3].set_ylabel('Profit')
35
36 ax[5].bar(total_sales_by_region['region'], total_sales_by_region['sales'], color='green
37 ax[5].set_title('Total Sales by Region')
38 ax[5].set_xlabel('Region')
39 ax[5].set_ylabel('Sales')
40
41
42

```

Out[156]: Text(0, 0.5, 'Sales')



In []:

1