# MINI PROJECT

## 1.   PROBLEM STATEMENT

Designing a mini project focused on a superstore sales dataset, where  aim to perform data preprocessing and visualization tasks. The project will involve cleaning and transforming the data, handling missing values and outliers, and then creating insightful visualizations to uncover patterns, trends, and anomalies in the sales data. The goal is to provide actionable insights for improving the store's performance and optimizing sales strategies. The project will require Python programming skills, data analysis libraries (e.g., pandas, matplotlib, and  seaborn ) and will culminate in a comprehensive report and presentation summarizing the findings and recommendations.

## 2.   DATA PREPROCESSING

 Superstore sale dataset preprocessing  involves the following steps:

1. Data Cleaning: Removing duplicates, handling missing values, and correcting any erroneous entries to ensure data integrity.

2. Data Transformation: Standardizing data types, converting categorical variables into numerical representations (encoding), and scaling features if needed.
3. Feature Selection: Identifying and selecting relevant features to reduce dimensionality and improve model performance.
4. Outlier Detection: Identifying and handling outliers that may skew the analysis or modeling results.
5. Data Splitting: Dividing the dataset into training, validation, and test sets to evaluate and validate models effectively.

## PROGRAM:

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

```python
df=pd.read_csv('/content/SuperStoreOrders.csv')
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 | AG-2011-2040 | 1/1/2011 | 6/1/2011 | Standard Class | Toby Braunhardt | Consumer | Constantine | Algeria | Africa | Africa | ... | Office Supplies | Storage | Tenex Lockers, Blue | 408 | 2 | 0.0 | 106.140 |
| 1 | IN-2011-47883 | 1/1/2011 | 8/1/2011 | Standard Class | Joseph Holt | Consumer | New South Wales | Australia | APAC | Oceania | ... | Office Supplies | Supplies | Acme Trimmer, High Speed | 120 | 3 | 0.1 | 36.036 |
| 2 | HU-2011-1220 | 1/1/2011 | 5/1/2011 | Second Class | Annie Thurman | Consumer | Budapest | Hungary | EMEA | EMEA | ... | Office Supplies | Storage | Tenex Box, Single Width | 66 | 4 | 0.0 | 29.640 |
| 3 | IT-2011-3647632 | 1/1/2011 | 5/1/2011 | Second Class | Eugene Moren | Home Office | Stockholm | Sweden | EU | North | ... | Office Supplies | Paper | Enermax Note Cards, Premium | 45 | 3 | 0.5 | -26.055 |
| 4 | IN-2011-47883 | 1/1/2011 | 8/1/2011 | Standard Class | Joseph Holt | Consumer | New South Wales | Australia | APAC | Oceania | ... | Furniture | Furnishings | Eldon Light Bulb, Duo Pack | 114 | 5 | 0.1 | 37.770 |

5 rows × 21 columns

```
df.shape
```

```
(19911, 21)
```

```
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'], dtype='object')
```

```
df.dtypes
```

| | |
|---|---|
| vorder_id | object |
| order_date | object |
| ship_date | object |
| ship_mode | object |
| customer_name | object |
| segment | object |
| state | object |
| country | object |
| market | object |
| region | object |
| product_id | object |
| category | object |
| sub_category | object |

```
product_name       object
sales              object
quantity          float64
discount          float64
profit            float64
shipping_cost     float64
order_priority     object
year              float64
dtype: object
```

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

```
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           1
sub_category       1
product_name       1
sales              1
quantity           1
```

```
discount          1
profit            1
shipping_cost     1
order_priority    1
year              1
dtype: int64
```

```
df['country'].value_counts()
```

```
United States     4084
Mexico            1057
Australia         1050
France            1032
Germany            868
                  ...
Tajikistan           1
Macedonia            1
Mauritania           1
South Sudan          1
Sri Lanka            1
Name: country, Length: 141, dtype: int64
```

```
df['category'].unique()
```

```
array(['Office Supplies', 'Furniture', 'Technology',
nan], dtype=object)
```

```python
df['category'].value_counts()
```

```
Office Supplies    12115
Technology          3994
Furniture           3801
Name: category, dtype: int64
```

```python
df['sub_category'].nunique()
```

```
17
```

```python
df['sub_category'].value_counts()
```

```
Binders       2337
Storage       1989
Art           1907
Paper         1364
Phones        1340
Chairs        1303
Furnishings   1236
Accessories   1201
Labels        1018
Supplies       933
Envelopes      931
Fasteners      929
Bookcases      927
Copiers        855
Appliances     707
Machines       598
```

```
Tables            335
Name: sub_category, dtype: int64
```

## 3. DATA VISUALIZATION

A superstore sale dataset visualization is a graphical representation of sales data from a large retail store. It typically uses charts, graphs, or other visual elements to illustrate key sales metrics, such as revenue, product categories, and trends over time. These visualizations help businesses and analysts gain insights into sales performance, identify patterns, and make informed decisions for inventory management, marketing strategies, and more. They provide a clear and concise way to communicate complex sales information, aiding in data-driven decision-making and strategic planning.

### PROGRAM:

```
plt.figure(figsize=(12,10))
df['sub_category'].value_counts().plot.pie(auto
pct="%1.1f%%")
plt.show()
```

```
df.groupby('sub_category')['profit','sales'].ag
g(['sum']).plot.bar()
plt.title('Total Profit and Sales per Sub-
Category')
plt.legend('Profit')
plt.legend('Sales')
plt.show()
```
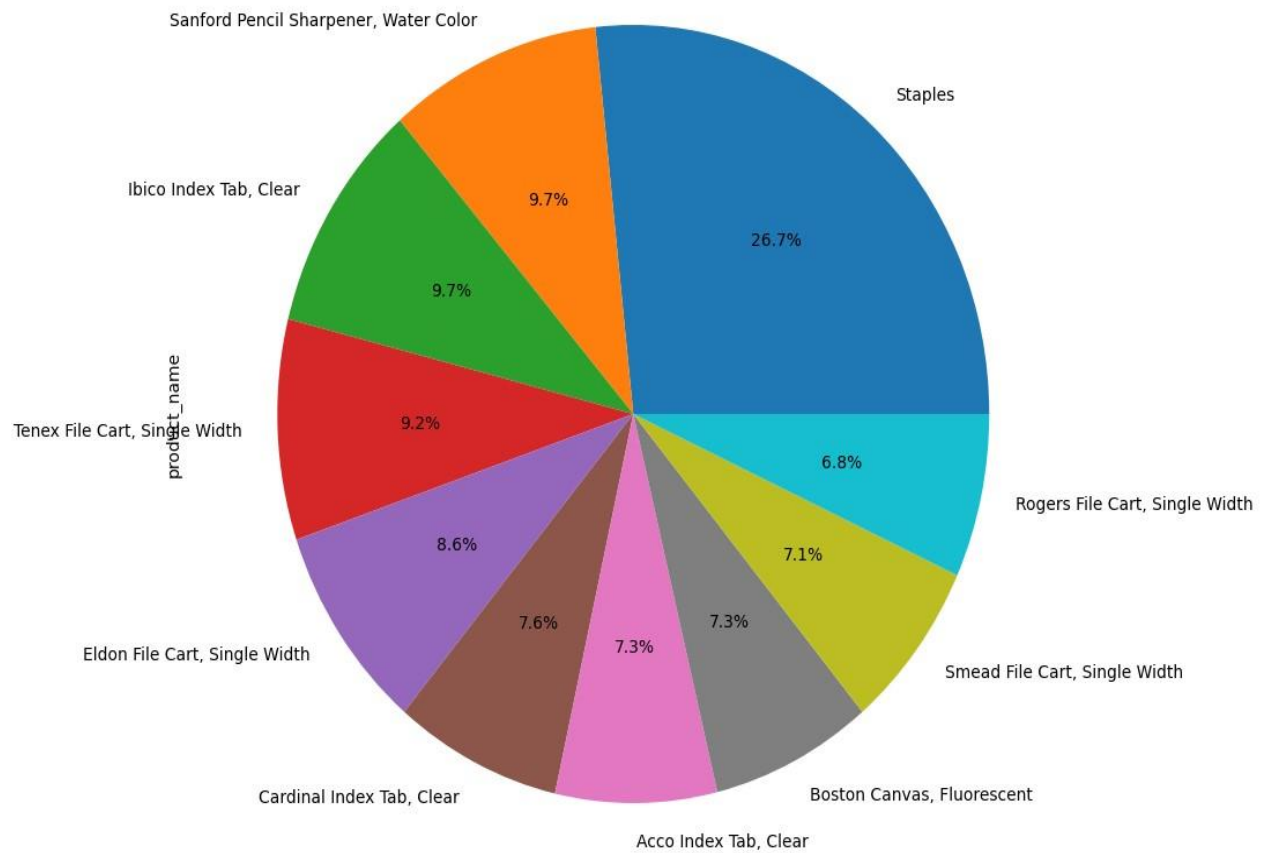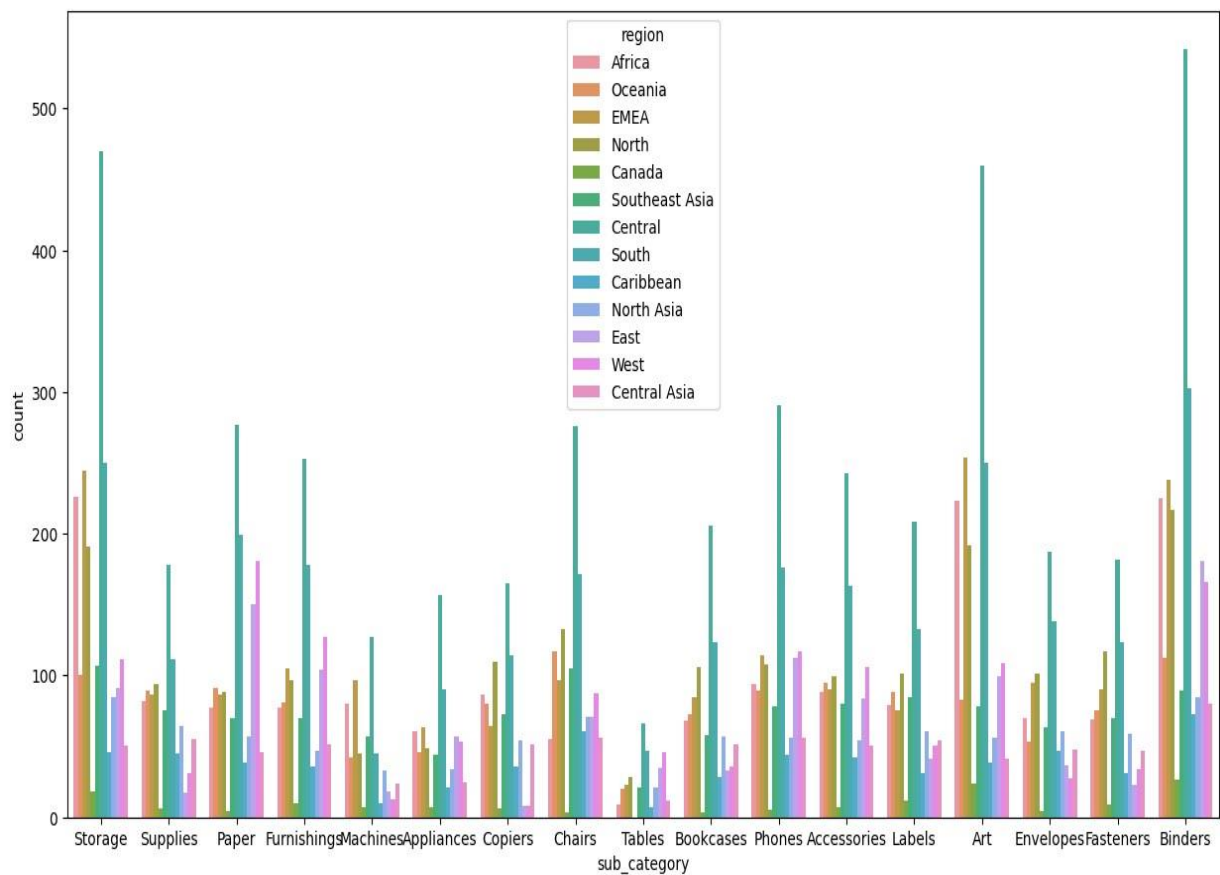
## Total Profit and Sales per Sub-Category



```python
df['product_name'].nunique()
```

```
3532
```

```python
plt.figure(figsize=(12,10))
df['product_name'].value_counts().head(10).plot.pie(autopct="%1.1f%%")
```
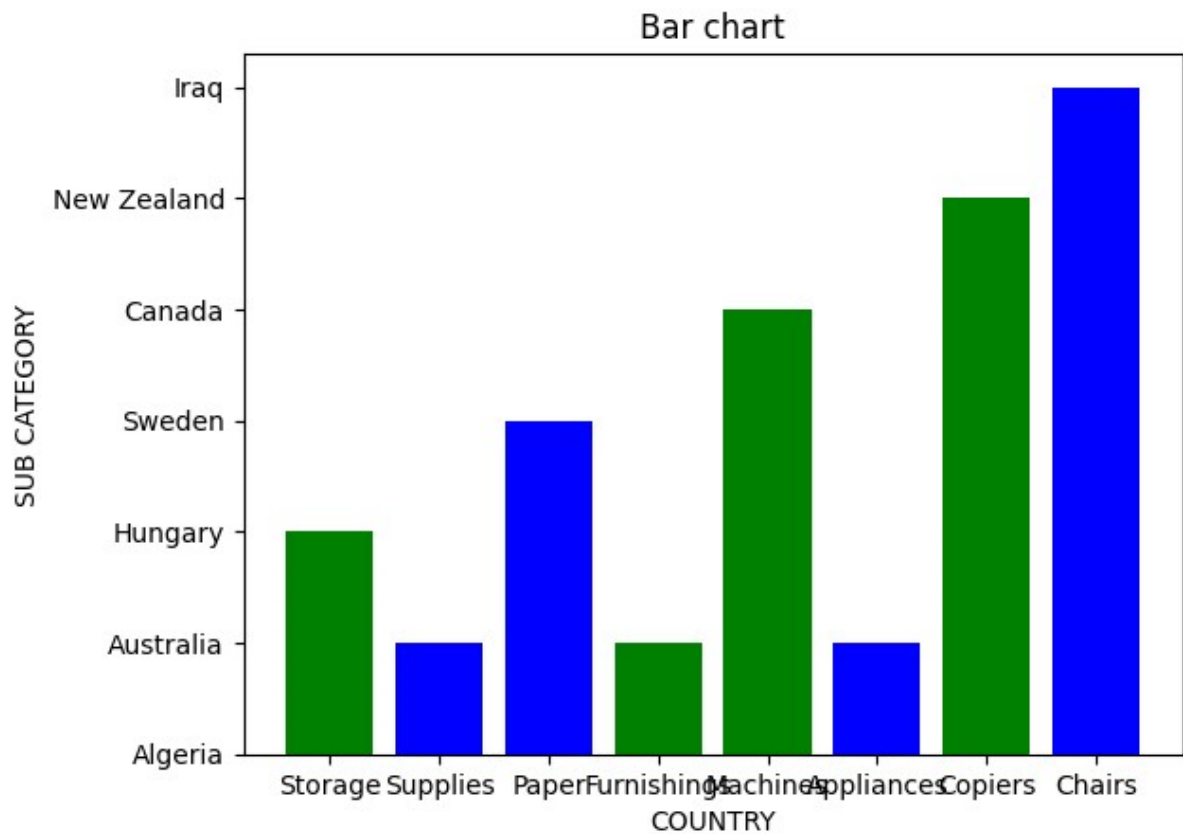
Pie chart titled with product_name showing: Staples 26.7%, Sanford Pencil Sharpener, Water Color 9.7%, Ibico Index Tab, Clear 9.7%, Tenex File Cart, Single Width 9.2%, Eldon File Cart, Single Width 8.6%, Cardinal Index Tab, Clear 7.6%, Acco Index Tab, Clear 7.3%, Boston Canvas, Fluorescent 7.3%, Smead File Cart, Single Width 7.1%, Rogers File Cart, Single Width 6.8%

```
plt.figure(figsize=(15,8))
sns.countplot(x="sub_category", hue="region",
data=df)
plt.show()
```

```
x=df['sub_category'].head(10)

y=df['country'].head(10)

plt.bar(df['country'],df['sub_category'],color=
['green','blue'])
plt.title("Bar chart")
plt.xlabel('COUNTRY')
plt.ylabel('SUB CATEGORY')
plt.show()
```
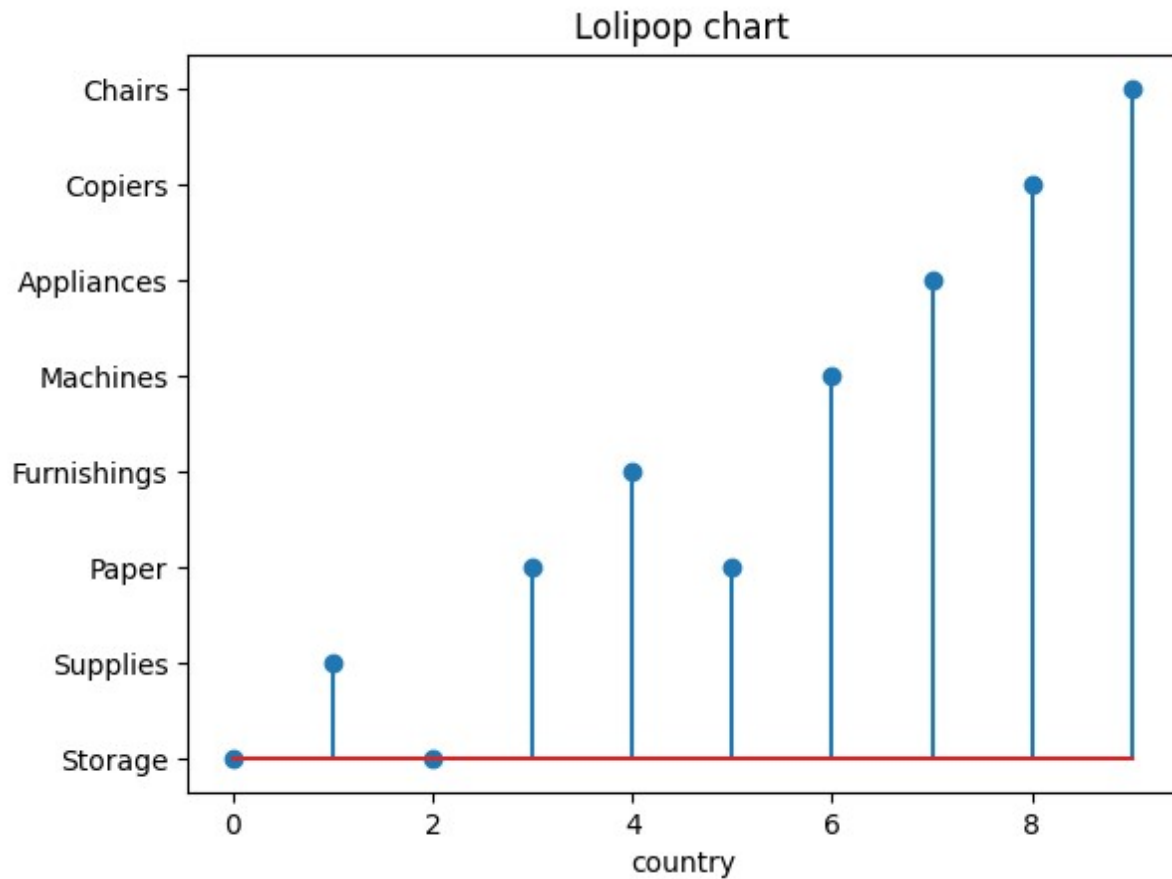
Bar chart

```
plt.hist(df['country'])
plt.title("Histrogram")
plt.xlabel('country')
plt.show()
```
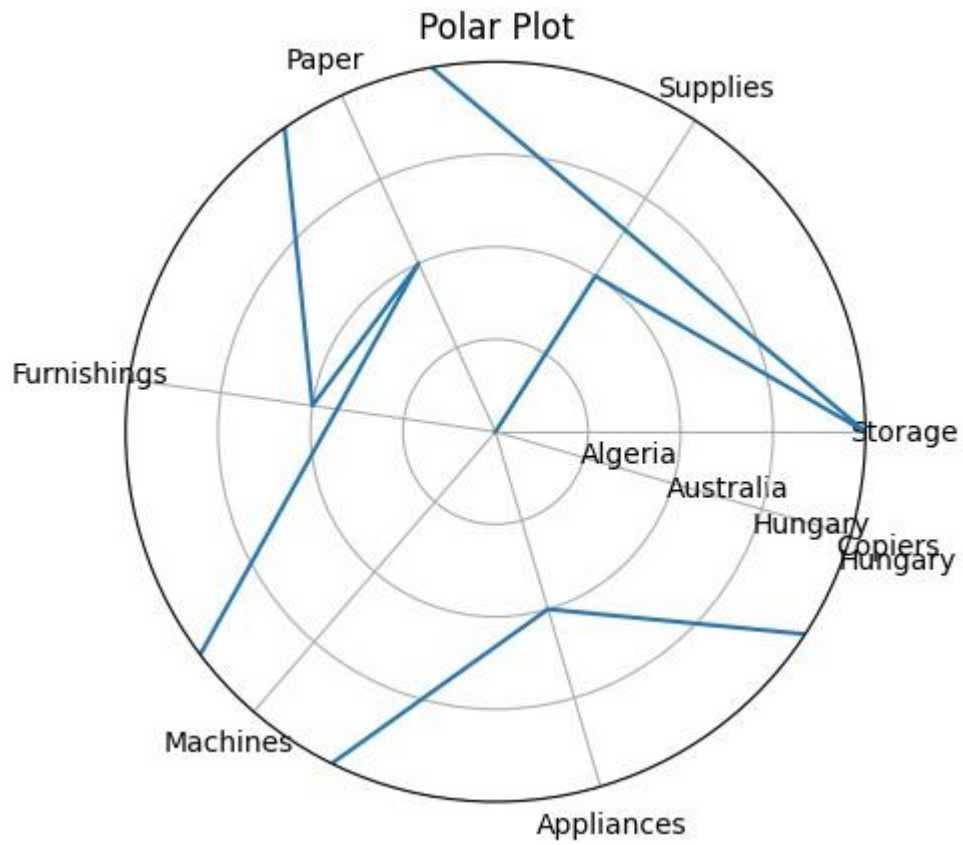


Histrogram

```
plt.plot(x,y)
plt.title("Line chart")
plt.xlabel('COUNTRY')
plt.ylabel('SUB CATEGORY')
plt.show()
```
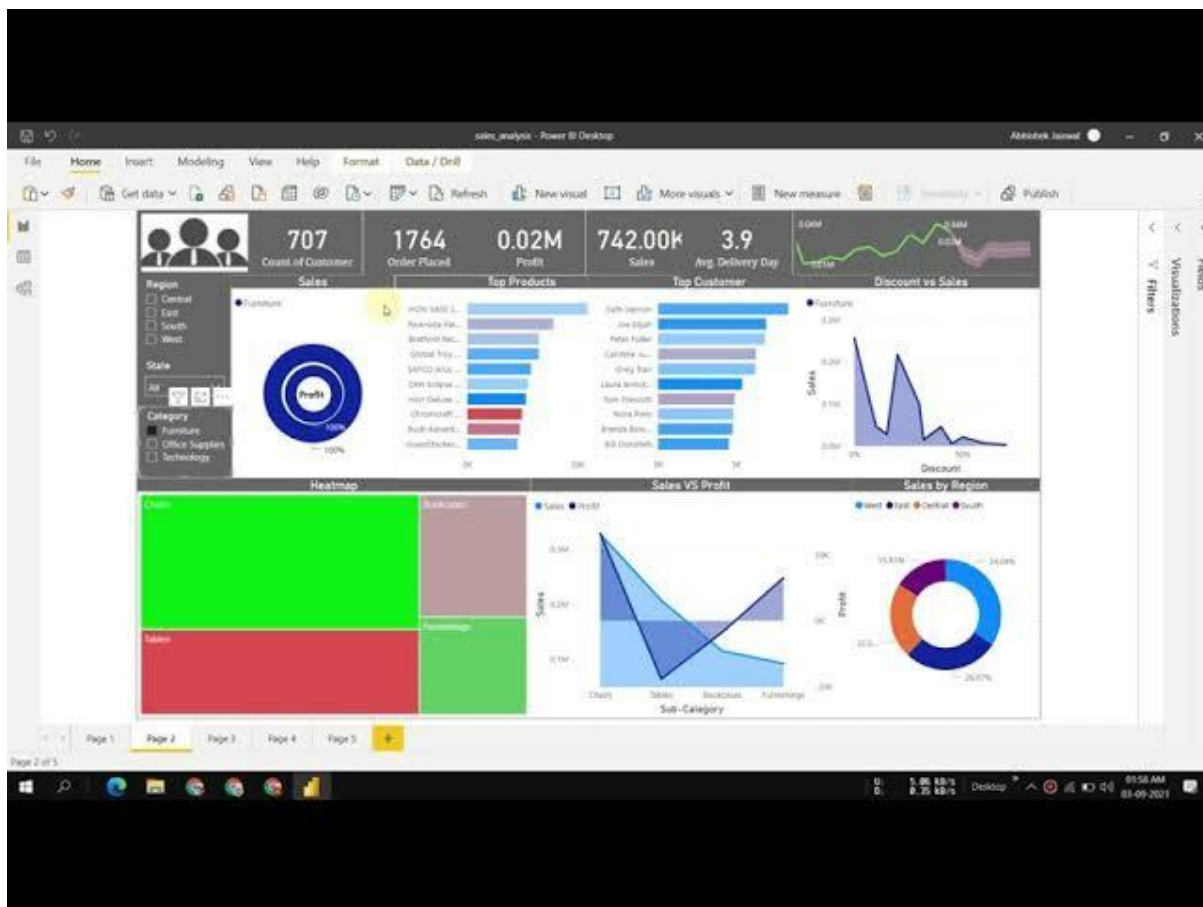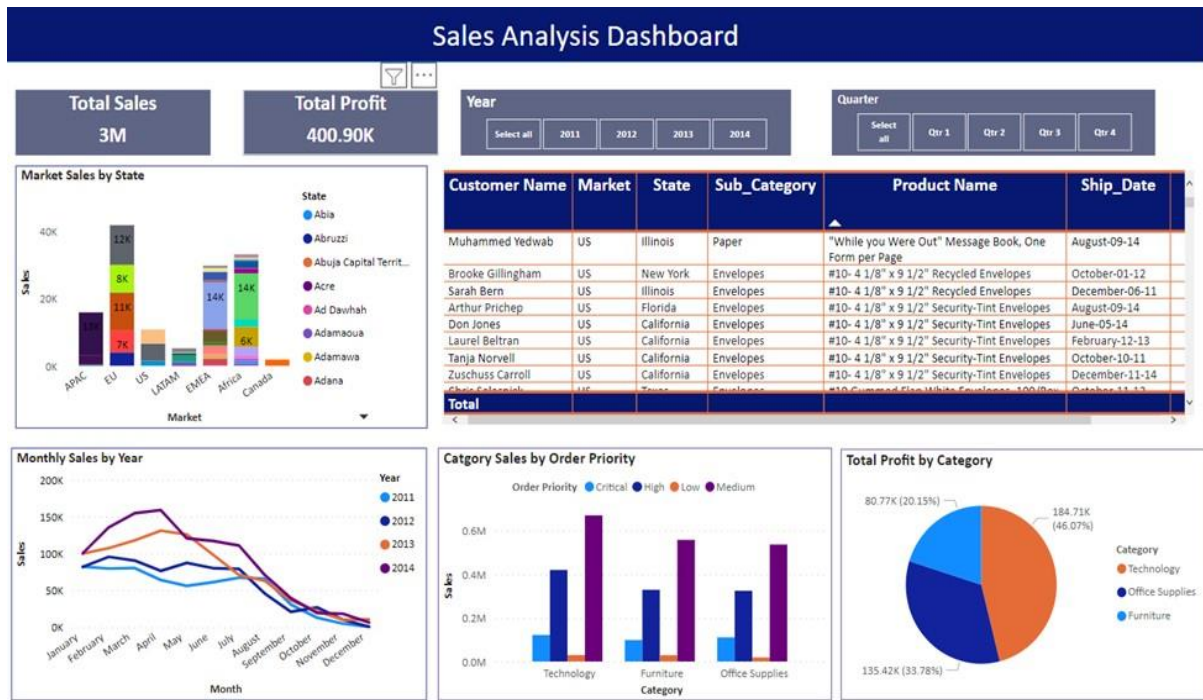


Line chart

```
plt.stem(x)
plt.title("Lolipop chart")
plt.xlabel('country')
plt.show()
```

Lolipop chart

```
fig, ax = plt.subplots(subplot_kw={'projection':
'polar'})
ax.plot(x, y)
ax.set_rmax(2)
ax.set_rticks([0.5, 1, 1.5, 2])
ax.set_rlabel_position(-22.5)
ax.grid(True)
ax.set_title("Polar Plot", va='bottom')
plt.show()
```
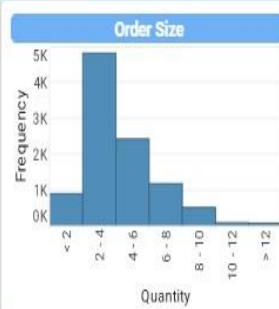
Polar Plot

**DASHBOARD**

Sales Analysis Dashboard

# REPORT