

```
In [7]: import pandas as pd

# Reading the CSV files
products = pd.read_csv('Products.csv.csv')

# Display the first few rows of each dataset
print(products.head())

+-----+
| ProductID | ProductName | Category | Price |
|-----|-----|-----|-----|
0 | P001 | ActiveWear Biography | Books | 149.30
1 | P002 | ActiveWear Smartwatch | Electronics | 446.30
2 | P003 | ComfortLiving Biography | Books | 46.12
3 | P004 | BookWorld Rug | Home Decor | 95.49
4 | P005 | TechPro T-Shirt | Clothing | 429.31
+-----+

In [10]: # Inspect the first few rows of each dataset
print(products.head())
print(transactions.head())
print(transactions.info())

# Check for missing values
print(customers.isnull().sum())
print(products.isnull().sum())
print(transactions.isnull().sum())

+-----+
| CustomerID | CustomerName | Region | SignupDate |
|-----|-----|-----|-----|
0 | C0001 | Lawrence Carroll | South America | 2022-07-10
1 | C0002 | Elizabeth Tate | Asia | 2022-02-13
2 | C0003 | Michael Rivera | South America | 2024-03-07
3 | C0004 | Kathleen Rodriguez | South America | 2022-10-09
4 | C0005 | Laura Baker | Asia | 2022-08-15
+-----+
| ProductID | ProductName | Category | Price |
|-----|-----|-----|-----|
0 | P001 | ActiveWear Biography | Books | 149.30
1 | P002 | ActiveWear Smartwatch | Electronics | 446.30
2 | P003 | ComfortLiving Biography | Books | 46.12
3 | P004 | BookWorld Rug | Home Decor | 95.49
4 | P005 | TechPro T-Shirt | Clothing | 429.31
+-----+
| TransactionID | CustomerID | ProductID | TransactionDate | Quantity |
|-----|-----|-----|-----|-----|
0 | T00001 | C0199 | P047 | 2024-08-25 12:18:23 | 1
1 | T00112 | C0146 | P047 | 2024-03-27 22:23:54 | 1
2 | T00146 | C0127 | P047 | 2024-04-26 07:38:55 | 1
3 | T00272 | C0087 | P047 | 2024-03-24 22:55:37 | 2
4 | T00363 | C0070 | P047 | 2024-03-21 15:10:10 | 3
+-----+
| TotalValue | Price |
|-----|-----|
0 | 300.68 | 300.68
1 | 300.68 | 300.68
2 | 300.68 | 300.68
3 | 601.36 | 300.68
4 | 902.04 | 300.68
+-----+
| CustomerID | 0 |
|-----|-----|
0 | CustomerID | 0 |
1 | Region | 0 |
2 | SignupDate | 0 |
3 | ProductID | int64 |
4 | TransactionID | 0 |
5 | Category | 0 |
6 | Price | 0 |
7 | Quantity | 0 |
8 | TotalValue | 0 |
9 | Price | 0 |
+-----+
| dtype: int64 |
+-----+
| TransactionID | 0 |
|-----|-----|
0 | CustomerID | 0 |
1 | ProductID | 0 |
2 | TransactionDate | 0 |
3 | Quantity | 0 |
4 | TotalValue | 0 |
5 | Price | 0 |
+-----+
| dtype: int64 |
+-----+
```

```
In [13]: # Fill or drop missing values based on your analysis
customers.fillna('Unknown', inplace=True) # Example: filling missing values with 'Unknown'
```

```
In [12]: # Group by product name and sum total sales
products.groupby('Product Name').sum()
products.groupby('Product Name').sum().reset_index(inplace=True)
```

```
In [11]: import matplotlib.pyplot as plt
import seaborn as sns

# Count customers by region
region_counts = customers['Region'].value_counts()

# Bar chart
plt.figure(figsize=(8, 5))
region_counts.plot(kind='bar', color='skyblue')
plt.title('Customers by Region')
plt.xlabel('Region')
plt.ylabel('Number of Customers')
plt.xticks(rotation=45)
plt.show()
```



```
In [13]: # Merge transactions with products
merged_data = transactions.merge(products, on='ProductID')

# Total sales by product
top_products = merged_data.groupby('ProductName')['TotalValue'].sum().nlargest(5)

# Bar chart
plt.figure(figsize=(8, 5))
top_products.plot(kind='bar', figsize=(8, 5), color='orange')
plt.title('Top 5 Products by Total Sales')
plt.xlabel('Product Name')
plt.ylabel('Total Sales (USD)')
plt.xticks(rotation=45)
plt.show()
```



```
In [20]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load the datasets (update file names if necessary)
customers_df = pd.read_csv('Customers.csv.csv')
products_df = pd.read_csv('Products.csv.csv')
transactions_df = pd.read_csv('Transactions.csv.csv')

# Convert date columns to datetime format
customers_df['SignupDate'] = pd.to_datetime(customers_df['SignupDate'])
transactions_df['TransactionDate'] = pd.to_datetime(transactions_df['TransactionDate'])

# Display summary of Customers dataset
print('Customers Dataset Summary:')
print(customers_df.info())
print(customers_df.describe(include='all'))

# Optionally, display the first few rows to inspect the data
print(customers_df.head())
print(transactions_df.head())
```

```
Customers Dataset Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   CustomerID  200 non-null    object
 1   CustomerName  200 non-null    object
 2   Region      200 non-null    object
 3   SignupDate   200 non-null    datetime64[ns]
dtypes: datetime64[ns](1), object(3)
memory usage: 6.4+ KB
None
```

```
count      CustomerID      CustomerName      Region      SignupDate
unique      200              200              4              NaN
top         C0001      Lawrence Carroll      South America      2022-07-10 08:31:12
mean        NaN              NaN              NaN              2023-07-19 08:31:12
std         NaN              NaN              NaN              2022-01-22 00:00:00
min         NaN              NaN              NaN              2022-01-22 00:00:00
max         NaN              NaN              NaN              2024-04-12 12:00:00
NaN         NaN              NaN              NaN              2024-12-28 00:00:00
```

```
ProductID      ProductName      Category      Price
unique          5              5              4              NaN
top             P001      ActiveWear Smartwatch      Electronics      446.30
mean           NaN              NaN              NaN              267.557000
std            NaN              NaN              NaN              143.233900
min            NaN              NaN              NaN              46.120000
max            NaN              NaN              NaN              147.785000
NaN            NaN              NaN              NaN              147.785000
NaN            NaN              NaN              NaN              479.760000
```

```
In [1]: # Summary of Products dataset
print('Products Dataset Summary:')
print(products_df.info())
print(products_df.describe(include='all'))
```

```
Products Dataset Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   ProductID   100 non-null    object
 1   ProductName  100 non-null    object
 2   Category     100 non-null    object
 3   Price        100 non-null    float64
dtypes: float64(1), object(3)
memory usage: 1.3+ KB
None
```

```
count      ProductID      ProductName      Category      Price
unique      100              100              100              100.000000
top         P001      ActiveWear Smartwatch      Books      NaN
mean        NaN              NaN              NaN              267.557000
std         NaN              NaN              NaN              143.233900
min         NaN              NaN              NaN              46.120000
max         NaN              NaN              NaN              147.785000
NaN         NaN              NaN              NaN              147.785000
NaN         NaN              NaN              NaN              479.760000
```

```
In [2]: # Load the Transactions dataset
transactions_df = pd.read_csv('Transactions.csv.csv')

# Convert any necessary columns (e.g., dates) to datetime format if required
transactions_df['TransactionDate'] = pd.to_datetime(transactions_df['TransactionDate'])

# Display summary of Transactions dataset
print('Transactions Dataset Summary:')
print(transactions_df.info())
print(transactions_df.describe(include='all'))

# Optionally, display the first few rows to inspect the data
print(transactions_df.head())
```

```
Transactions Dataset Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   TransactionID  1000 non-null    object
 1   CustomerID     1000 non-null    object
 2   ProductID      1000 non-null    object
 3   TransactionDate  1000 non-null    object
 4   Quantity        1000 non-null    int64
 5   TotalValue      1000 non-null    float64
dtypes: float64(1), int64(1), object(3)
memory usage: 14.8+ KB
None
```

```
count      TransactionID      CustomerID      ProductID      TransactionDate      Quantity \
unique      1000              1000              1000              1000              1000.000000
top         T00992      C0199      P047      2024-04-21 10:10:14      NaN
mean        NaN              NaN              NaN              NaN              1.117990
std         NaN              NaN              NaN              NaN              2.337000
min         NaN              NaN              NaN              NaN              1.000000
max         NaN              NaN              NaN              NaN              4.000000
NaN         NaN              NaN              NaN              NaN              4.000000
```

```
count      TotalValue      Price
unique      1000              1000.000000
top         NaN              NaN
mean        493.993500      272.548700
std         493.144678      140.738330
min         46.120000      46.120000
max         285.235000      147.785000
NaN         285.235000      147.785000
NaN         1011.600000      404.400000
NaN         1991.000000      497.000000
```

```
TransactionID      CustomerID      ProductID      TransactionDate      Quantity \
0 | T00001 | C0199 | P047 | 2024-08-25 12:18:23 | 1
1 | T00112 | C0146 | P047 | 2024-03-27 22:23:54 | 1
2 | T00146 | C0127 | P047 | 2024-03-26 07:38:55 | 1
3 | T00272 | C0087 | P047 | 2024-03-24 22:55:37 | 2
4 | T00363 | C0070 | P047 | 2024-03-21 15:10:10 | 3
+-----+
| TotalValue | Price |
|-----|-----|
0 | 300.68 | 300.68
1 | 300.68 | 300.68
2 | 300.68 | 300.68
3 | 601.36 | 300.68
4 | 902.04 | 300.68
+-----+
```

```
In [24]: # EDA - Customers dataset
print('Customers Region Distribution:')
# Display value counts for the 'Region' column
print(customers_df['Region'].value_counts())

# Plot the customer distribution by region using Seaborn's countplot
plt.figure(figsize=(8, 5))
sns.countplot(data=customers_df, x='Region', palette='viridis')
plt.title('Customer Distribution by Region')
plt.xlabel('Number of Customers')
plt.ylabel('Region')
plt.xticks(rotation=45)
plt.show()
```

```
Customers Region Distribution:
Region:
South America    50
Europe           48
North America    45
Asia             42
dtype: int64

C:\Users\WP\AppData\Local\Temp\ipykernel_13340\208497260.py:8: FutureWarning:
    Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.
sns.countplot(data=customers_df, x='Region', palette='viridis')
```



```
In [25]: # EDA - Products dataset
print('Products Category Distribution:')
# Display value counts for the 'Category' column
print(products_df['Category'].value_counts())

# Plot the product distribution by category using Seaborn's countplot
plt.figure(figsize=(8, 5))
sns.countplot(data=products_df, x='Category', palette='crest')
plt.title('Product Distribution by Category')
plt.xlabel('Number of Products')
plt.ylabel('Category')
plt.xticks(rotation=45)
plt.show()
```

```
Products Category Distribution:
Category:
Books           26
Electronics    14
Home Decor     23
Clothing       25
dtype: int64

C:\Users\WP\AppData\Local\Temp\ipykernel_13340\208497260.py:8: FutureWarning:
    Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.
sns.countplot(data=products_df, x='Category', palette='crest')
```



```
In [26]: # EDA - Transactions dataset
print('Transaction Value Statistics:')
# Display statistics for the 'TotalValue' column
print(transactions_df['TotalValue'].describe())

# Display statistics for the 'Quantity' column
print(transactions_df['Quantity'].describe())
```

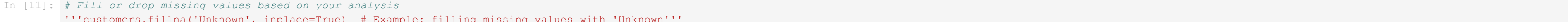
```
Transaction Value Statistics:
count      1000.000000
mean        493.993500
std         493.144678
min         46.120000
25%        285.235000
50%        388.480000
75%        1011.600000
max         1991.000000
Name: TotalValue, dtype: float64

Transaction Quantity Statistics:
count      1000.000000
mean         2.337000
std          1.117990
min           1.000000
25%           2.000000
50%           2.000000
75%           4.000000
max           4.000000
Name: Quantity, dtype: float64
```

```
In [27]: # EDA - Transactions dataset: Distribution of Transaction Value
plt.figure(figsize=(8, 5))
sns.kdeplot(transactions_df['TotalValue'], bins=30, kde=True, color='blue')
plt.title('Transaction Value Distribution')
plt.xlabel('Total Value (USD)')
plt.ylabel('Frequency')
plt.show()
```



```
In [28]: # EDA - Transactions dataset: Distribution of Transaction Quantity
plt.figure(figsize=(8, 5))
sns.kdeplot(transactions_df['Quantity'], bins=30, kde=True, color='green')
plt.title('Transaction Quantity Distribution')
plt.xlabel('Quantity')
plt.ylabel('Frequency')
plt.show()
```



```
In [29]: # Merge datasets for further analysis
merged_df = transactions_df.merge(customers_df, on='CustomerID').merge(products_df, on='ProductID')

# Insights from merged dataset
print('Merged Dataset Summary:')
print(merged_df.info())
```

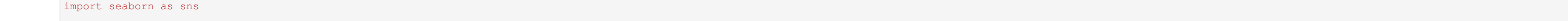
```
Merged Dataset Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 13 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   TransactionID  1000 non-null    object
 1   CustomerID     1000 non-null    object
 2   ProductID      1000 non-null    object
 3   TransactionDate  1000 non-null    object
 4   Quantity        1000 non-null    int64
 5   TotalValue      1000 non-null    float64
 6   Price          1000 non-null    float64
 7   CustomerName   1000 non-null    object
 8   Region         1000 non-null    object
 9   SignupDate     1000 non-null    datetime64[ns]
10  ProductName    1000 non-null    object
11  Category       1000 non-null    object
12  Price_y        1000 non-null    float64
dtypes: datetime64[ns](1), float64(3), int64(1), object(8)
memory usage: 101.7+ KB
None
```

```
In [30]: # Sample Insights: Sales by Region
region_sales = merged_df.groupby('Region')['TotalValue'].sum().sort_values(ascending=False)

# Print the total sales by region
print('Total Sales by Region:')
print(region_sales)

# Visualize the sales by region
plt.figure(figsize=(8, 5))
region_sales.plot(kind='bar', color='teal')
plt.title('Total Sales by Region')
plt.xlabel('Region')
plt.ylabel('Total Sales (USD)')
plt.xticks(rotation=45)
plt.show()
```

```
Total Sales by Region:
Region:
South America    21932.56
Europe           14624.43
North America    15231.40
Asia             15274.97
Name: TotalValue, dtype: float64
```

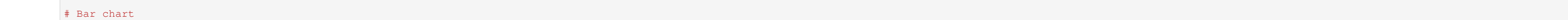


```
In [31]: # Sample Insights: Top 10 Products by Sales
product_sales = merged_df.groupby('ProductName')['TotalValue'].sum().sort_values(ascending=False).head(10)

# Print the top 10 products by sales
print('Top 10 Products by Sales:')
print(product_sales)

# Visualize the top 10 products by sales
plt.figure(figsize=(10, 6))
product_sales.plot(kind='bar', color='orange')
plt.title('Top 10 Products by Sales')
plt.xlabel('Product Name')
plt.ylabel('Total Sales (USD)')
plt.xticks(rotation=45)
plt.show()
```

```
Top 10 Products by Sales:
ProductName
ActiveWear Smartwatch    38096.97
SoundWave Headphones    23211.64
ActiveWear Jacket        24507.90
ActiveWear Rug           22314.43
BookWorld Headphones     19815.80
ActiveWear Backpack       19221.49
BookWorld Rug            18743.79
TechPro T-shirt          18267.94
ActiveWear Cookware Set   18069.73
Name: TotalValue, dtype: float64
```



```
In [32]: # Save the merged dataset to the current working directory
merged_df.to_csv('MergedDataset.csv', index=False)
```

```
In [33]: import os
os.getcwd()
```