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


# Load datasets
customers = pd.read_csv("/content/Customers.csv")
products = pd.read_csv("/content/Products.csv")
transactions = pd.read_csv("/content/Transactions.csv")
```

## EDA (Exploratory Data Analysis)

### Overview of datasets

```
# Overview of data
print(customers.info())
print(products.info())
print(transactions.info())
```

```
<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    object
dtypes: object(4)
memory usage: 6.4+ KB
None
<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: 3.3+ KB
None
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 7 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
dtypes: float64(2), int64(1), object(4)
memory usage: 54.8+ KB
None
```

```
# Display basic statistics
print(customers.describe())
print(products.describe())
print(transactions.describe())
```

|        | CustomerID | CustomerName     | Region        | SignupDate |
|--------|------------|------------------|---------------|------------|
| count  | 200        | 200              | 200           | 200        |
| unique | 200        | 200              | 4             | 179        |
| top    | C0001      | Lawrence Carroll | South America | 2024-11-11 |
| freq   | 1          | 1                | 59            | 3          |

|       | Price      |
|-------|------------|
| count | 100.000000 |
| mean  | 267.551700 |
| std   | 143.219383 |
| min   | 16.080000  |
| 25%   | 147.767500 |

```
50%       292.875000
75%       397.090000
max       497.760000
          Quantity    TotalValue        Price
count   1000.000000  1000.000000  1000.00000
mean       2.537000    689.995560   272.55407
std        1.117981    493.144478   140.73639
min        1.000000     16.080000    16.08000
25%        2.000000    295.295000   147.95000
50%        3.000000    588.880000   299.93000
75%        4.000000   1011.660000   404.40000
max        4.000000   1991.040000   497.76000
```

## Missing values checking

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

```
CustomerID        0
CustomerName      0
Region            0
SignupDate        0
dtype: int64
ProductID         0
ProductName       0
Category          0
Price             0
dtype: int64
TransactionID     0
CustomerID        0
ProductID         0
TransactionDate   0
Quantity          0
TotalValue        0
Price             0
dtype: int64
```
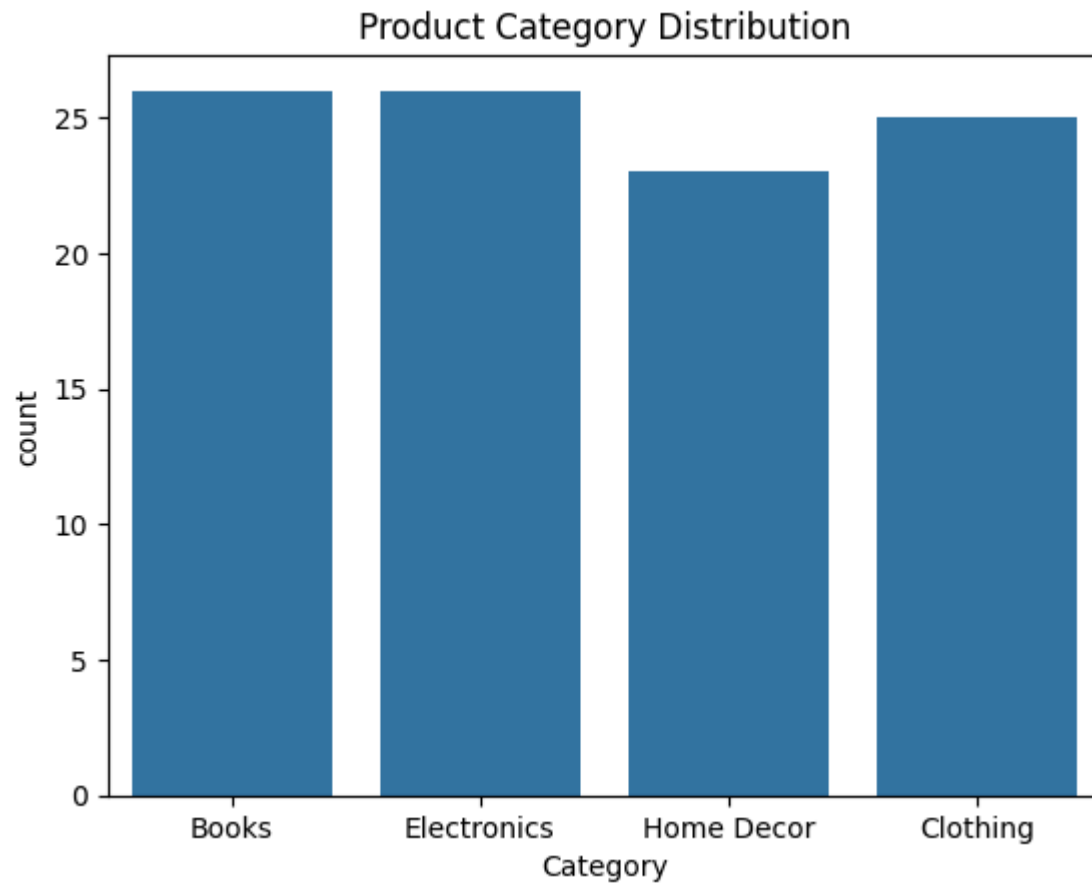
## Visualization

```python
# Distribution of regions
sns.countplot(data=customers, x="Region")
plt.title("Customer Distribution by Region")
plt.show()
```
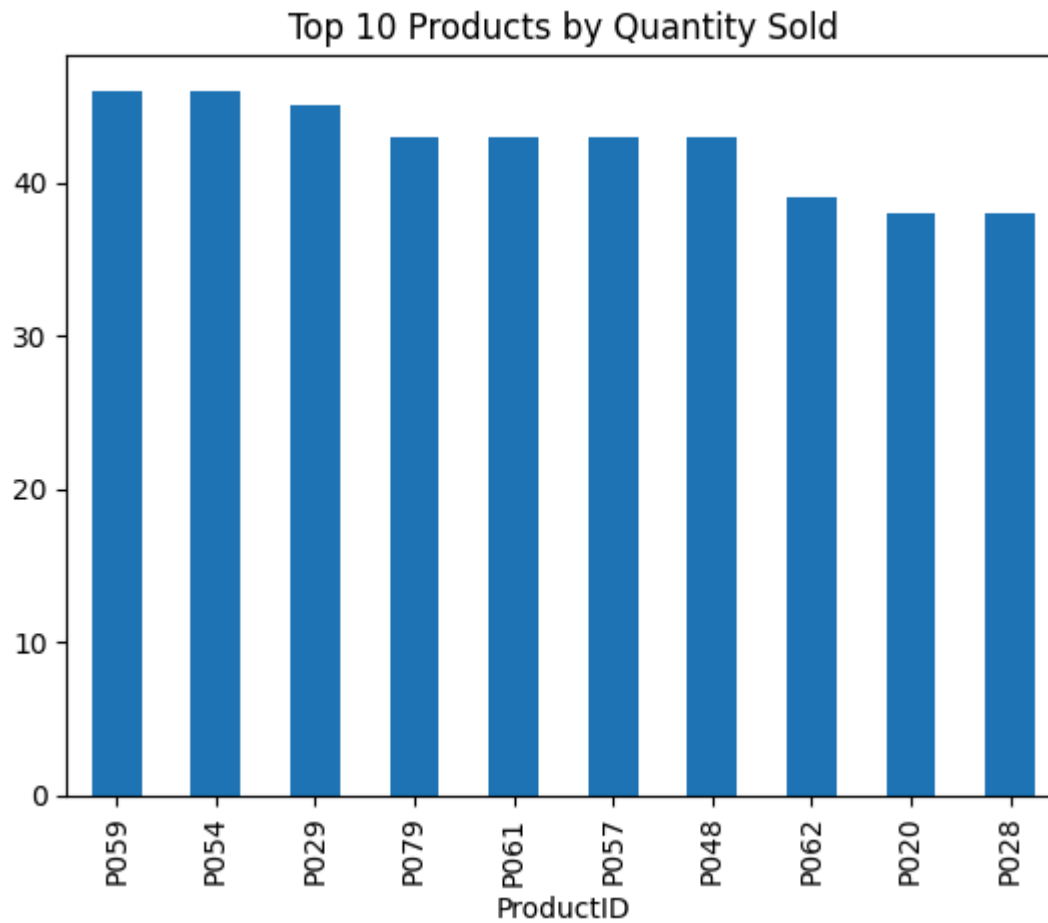
### Customer Distribution by Region



```python
# Product category distribution
sns.countplot(data=products, x="Category")
```

```
plt.title("Product Category Distribution")
plt.show()
```

Product Category Distribution



```
# Top products by transaction quantity
top_products = transactions.groupby("ProductID")["Quantity"].sum().sort_values(ascending=False).head(10)
top_products.plot(kind="bar", title="Top 10 Products by Quantity Sold")
plt.show()
```

Top 10 Products by Quantity Sold

1. The majority of customers are from the "Asia" region, indicating a strong presence in that market.

2. Products in the "Electronics" category have the highest sales, suggesting their popularity.

3. A small number of products contribute to the majority of the sales (Pareto principle).

Start coding or generate with AI.