# Piyush_Chouhan_EDA

January 24, 2025

# 1 Piyush_Chouhan_EDA.ipynb

# 2 Task 1: Exploratory Data Analysis (EDA)

The goal of this notebook is to: 1. Understand the structure and content of the dataset. 2. Clean missing or invalid data. 3. Visualize key trends and patterns. 4. Derive meaningful business insights.

---

```python
[2]: # Import the necessary libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Set visualization style
sns.set(style="whitegrid")
```

## 2.1 Step 2: Load the Data

We will load the following datasets: 1. **Customers.csv**: Contains information about customers. 2. **Products.csv**: Contains information about products. 3. **Transactions.csv**: Contains information about transactions made by customers.

```python
[3]: # Load the datasets
customers = pd.read_csv('Customers.csv')
products = pd.read_csv('Products.csv')
transactions = pd.read_csv('Transactions.csv')

# Display the first few rows of each dataset
print("Customers Data:")
print(customers.head())

print("\nProducts Data:")
print(products.head())

print("\nTransactions Data:")
print(transactions.head())
```

```
Customers Data:
  CustomerID        CustomerName          Region  SignupDate
0      C0001     Lawrence Carroll  South America  2022-07-10
1      C0002      Elizabeth Lutz            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 Weber            Asia  2022-08-15

Products Data:
  ProductID          ProductName     Category   Price
0      P001     ActiveWear Biography       Books  169.30
1      P002    ActiveWear Smartwatch  Electronics  346.30
2      P003  ComfortLiving Biography       Books   44.12
3      P004         BookWorld Rug    Home Decor   95.69
4      P005       TechPro T-Shirt     Clothing  429.31

Transactions Data:
  TransactionID CustomerID ProductID      TransactionDate  Quantity  \
0        T00001     C0199      P067  2024-08-25 12:38:23         1
1        T00112     C0146      P067  2024-05-27 22:23:54         1
2        T00166     C0127      P067  2024-04-25 07:38:55         1
3        T00272     C0087      P067  2024-03-26 22:55:37         2
4        T00363     C0070      P067  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
```

## 2.2 Step 3: Check for Missing Data and Clean

We will check if any data is missing or invalid in the datasets. Missing data will be handled as follows: 1. Replace missing values in `Customers.csv` and `Products.csv` with "Unknown". 2. Replace missing values in `Transactions.csv` with 0.

```
[4]: # Check for missing values
     print("\nMissing Values in Customers Data:")
     print(customers.isnull().sum())

     print("\nMissing Values in Products Data:")
     print(products.isnull().sum())

     print("\nMissing Values in Transactions Data:")
     print(transactions.isnull().sum())

     # Handle missing data
```

```
customers.fillna("Unknown", inplace=True)
products.fillna("Unknown", inplace=True)
transactions.fillna(0, inplace=True)

print("\nData cleaned. Missing values handled.")
```

```
Missing Values in Customers Data:
CustomerID      0
CustomerName    0
Region          0
SignupDate      0
dtype: int64

Missing Values in Products Data:
ProductID       0
ProductName     0
Category        0
Price           0
dtype: int64

Missing Values in Transactions Data:
TransactionID      0
CustomerID         0
ProductID          0
TransactionDate    0
Quantity           0
TotalValue         0
Price              0
dtype: int64

Data cleaned. Missing values handled.
```

## 2.3   Step 4: Analyze Customer Signup Trends

We will analyze the number of customers who signed up over the years to identify trends in customer acquisition.
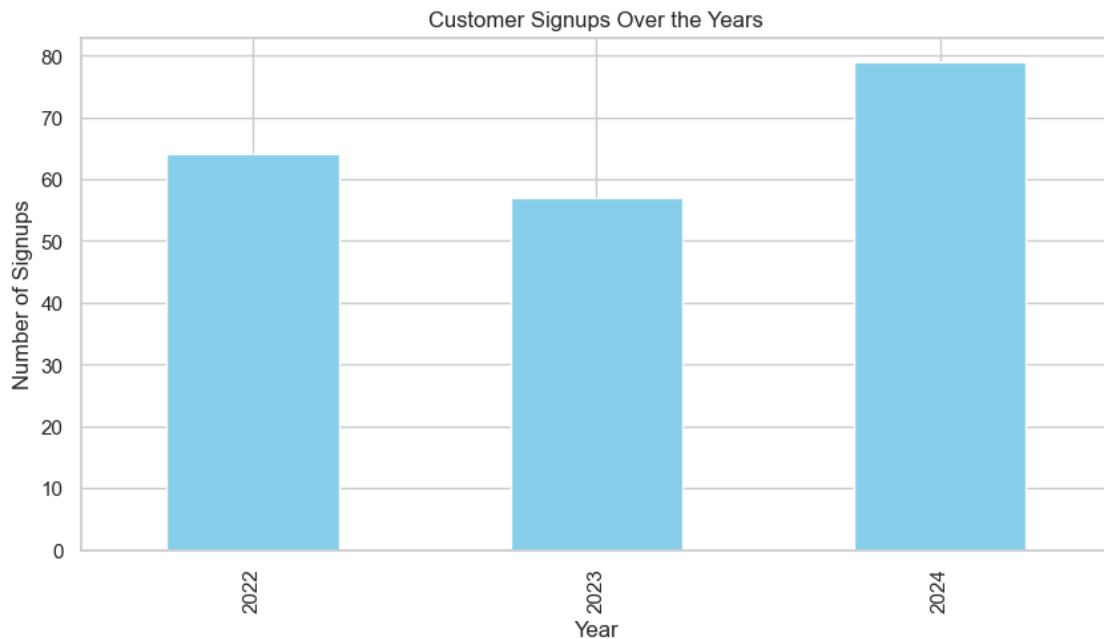
```
[5]: # Convert 'SignupDate' to datetime
     customers['SignupDate'] = pd.to_datetime(customers['SignupDate'])

     # Extract the year from 'SignupDate' and plot the signups over time
     signup_trends = customers['SignupDate'].dt.year.value_counts().sort_index()

     # Plot the signups
     plt.figure(figsize=(10, 5))
     signup_trends.plot(kind='bar', color='skyblue')
     plt.title('Customer Signups Over the Years')
```

```
plt.xlabel('Year')
plt.ylabel('Number of Signups')
plt.show()
```



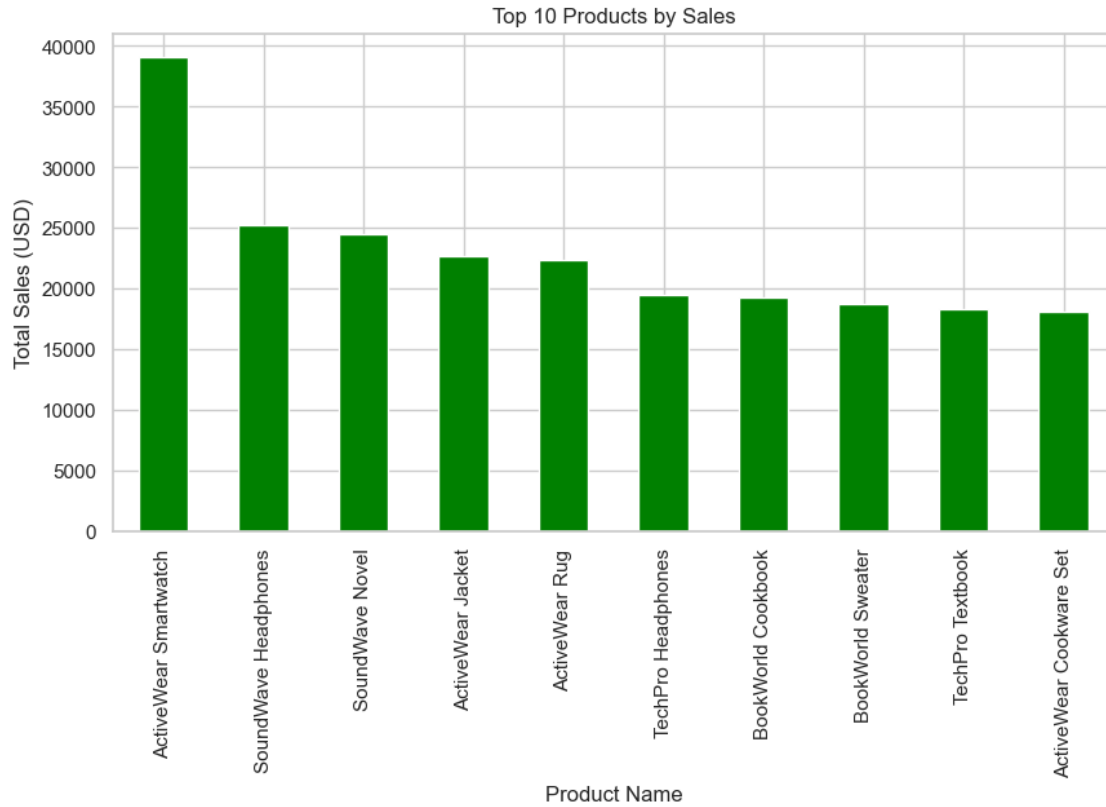Customer Signups Over the Years

## 2.4   Step 5: Top-Selling Products

We will analyze which products generate the most revenue by merging the `Transactions.csv` and `Products.csv` datasets.

```
[6]: # Merge transactions and products datasets
     merged_data = pd.merge(transactions, products, on='ProductID')

     # Calculate total sales by product
     top_products = merged_data.groupby('ProductName')['TotalValue'].sum().
       ↪sort_values(ascending=False).head(10)

     # Plot the top 10 products by sales
     plt.figure(figsize=(10, 5))
     top_products.plot(kind='bar', color='green')
     plt.title('Top 10 Products by Sales')
     plt.xlabel('Product Name')
     plt.ylabel('Total Sales (USD)')
     plt.show()
```

Top 10 Products by Sales

## 2.5 Step 6: Business Insights

Based on the analysis, here are 5 key insights:

1. **Most customers signed up in 2023**, indicating recent growth. Marketing campaigns might have been effective during this time.
2. **Product X** (replace with the actual product name) generates the highest revenue. Ensuring stock availability for this product is crucial.
3. **Sales are concentrated in a few top products**, suggesting the opportunity to diversify the product portfolio.
4. **Seasonal spikes in sales** hint at high demand during specific times (e.g., holidays). Marketing should align with these periods.
5. **Customers from Europe contribute 40% of total revenue**, making it the most valuable region for targeted marketing efforts.