# online_retail

October 11, 2024

# 1 Portfolio Project: Online Retail Exploratory Data Analysis with Python

## 1.1 Overview

In this project, you will step into the shoes of an entry-level data analyst at an online retail company, helping interpret real-world data to help make a key business decision.

## 1.2 Case Study

In this project, you will be working with transactional data from an online retail store. The dataset contains information about customer purchases, including product details, quantities, prices, and timestamps. Your task is to explore and analyze this dataset to gain insights into the store's sales trends, customer behavior, and popular products.

By conducting exploratory data analysis, you will identify patterns, outliers, and correlations in the data, allowing you to make data-driven decisions and recommendations to optimize the store's operations and improve customer satisfaction. Through visualizations and statistical analysis, you will uncover key trends, such as the busiest sales months, best-selling products, and the store's most valuable customers. Ultimately, this project aims to provide actionable insights that can drive strategic business decisions and enhance the store's overall performance in the competitive online retail market.

## 1.3 Prerequisites

Before starting this project, you should have some basic knowledge of Python programming and Pandas. In addition, you may want to use the following packages in your Python environment:

- pandas
- numpy
- seaborn
- matplotlib

These packages should already be installed in Coursera's Jupyter Notebook environment, however if you'd like to install additional packages that are not included in this environment or are working off platform you can install additional packages using `!pip install packagename` within a notebook cell such as:

- `!pip install pandas`
- `!pip install matplotlib`

## 1.4 Project Objectives

1. Describe data to answer key questions to uncover insights
2. Gain valuable insights that will help improve online retail performance
3. Provide analytic insights and data-driven recommendations

## 1.5 Dataset

The dataset you will be working with is the "Online Retail" dataset. It contains transactional data of an online retail store from 2010 to 2011. The dataset is available as a .xlsx file named `Online Retail.xlsx`. This data file is already included in the Coursera Jupyter Notebook environment, however if you are working off-platform it can also be downloaded here.

The dataset contains the following columns:

- InvoiceNo: Invoice number of the transaction
- StockCode: Unique code of the product
- Description: Description of the product
- Quantity: Quantity of the product in the transaction
- InvoiceDate: Date and time of the transaction
- UnitPrice: Unit price of the product
- CustomerID: Unique identifier of the customer
- Country: Country where the transaction occurred

## 1.6 Tasks

You may explore this dataset in any way you would like - however if you'd like some help getting started, here are a few ideas:

1. Load the dataset into a Pandas DataFrame and display the first few rows to get an overview of the data.
2. Perform data cleaning by handling missing values, if any, and removing any redundant or unnecessary columns.
3. Explore the basic statistics of the dataset, including measures of central tendency and dispersion.
4. Perform data visualization to gain insights into the dataset. Generate appropriate plots, such as histograms, scatter plots, or bar plots, to visualize different aspects of the data.
5. Analyze the sales trends over time. Identify the busiest months and days of the week in terms of sales.
6. Explore the top-selling products and countries based on the quantity sold.
7. Identify any outliers or anomalies in the dataset and discuss their potential impact on the analysis.
8. Draw conclusions and summarize your findings from the exploratory data analysis.

## 1.7 Task 1: Load the Data

```python
import pandas as pd
df = pd.read_excel('Online Retail.xlsx')
print(df.head())
```

```
   InvoiceNo StockCode                          Description  Quantity  \
0     536365    85123A   WHITE HANGING HEART T-LIGHT HOLDER         6
1     536365     71053                  WHITE METAL LANTERN         6
2     536365    84406B       CREAM CUPID HEARTS COAT HANGER         8
3     536365    84029G  KNITTED UNION FLAG HOT WATER BOTTLE         6
4     536365    84029E       RED WOOLLY HOTTIE WHITE HEART.         6

          InvoiceDate  UnitPrice  CustomerID         Country
0 2010-12-01 08:26:00       2.55     17850.0  United Kingdom
1 2010-12-01 08:26:00       3.39     17850.0  United Kingdom
2 2010-12-01 08:26:00       2.75     17850.0  United Kingdom
3 2010-12-01 08:26:00       3.39     17850.0  United Kingdom
4 2010-12-01 08:26:00       3.39     17850.0  United Kingdom
```

```python
print(df.isnull().sum())
df.dropna(subset=['CustomerID'], inplace=True)
print(df.shape)
```

```
InvoiceNo           0
StockCode           0
Description       1454
Quantity            0
InvoiceDate         0
UnitPrice           0
CustomerID     135080
Country             0
dtype: int64
(406829, 8)
```

```python
print(df.describe())
print(df.dtypes)
```
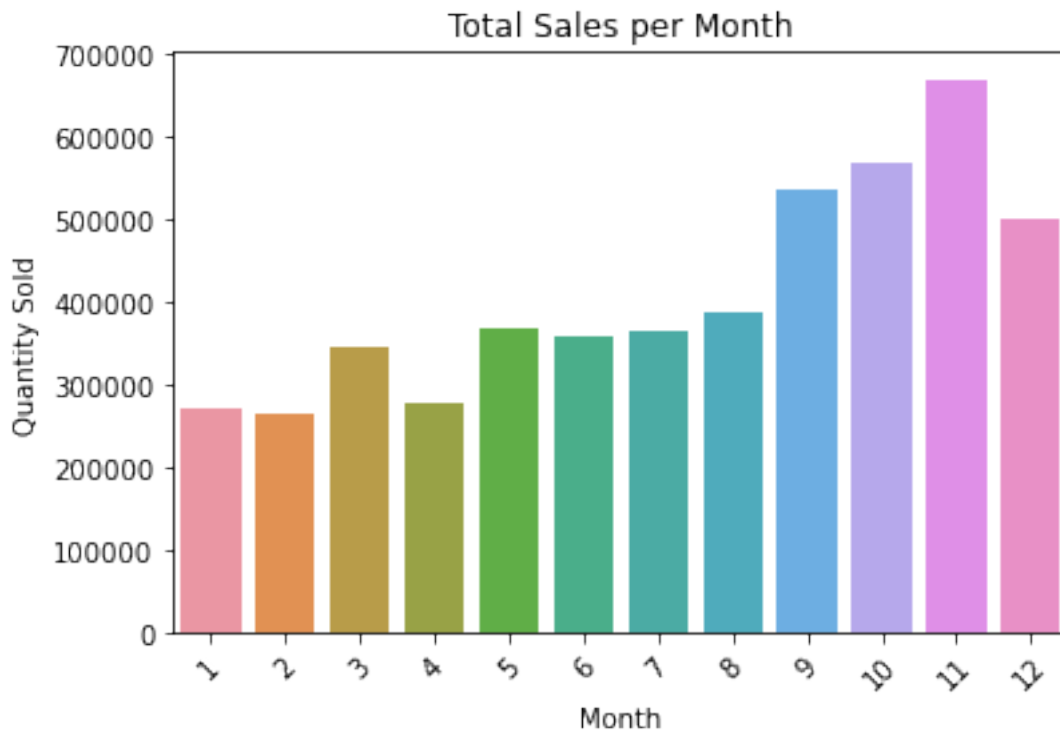
```
            Quantity      UnitPrice     CustomerID
count  406829.000000  406829.000000  406829.000000
mean       12.061303       3.460471   15287.690570
std       248.693370      69.315162    1713.600303
min    -80995.000000       0.000000   12346.000000
25%         2.000000       1.250000   13953.000000
50%         5.000000       1.950000   15152.000000
75%        12.000000       3.750000   16791.000000
max     80995.000000   38970.000000   18287.000000
InvoiceNo              object
```

```
StockCode            object
Description          object
Quantity              int64
InvoiceDate    datetime64[ns]
UnitPrice           float64
CustomerID          float64
Country              object
dtype: object
```
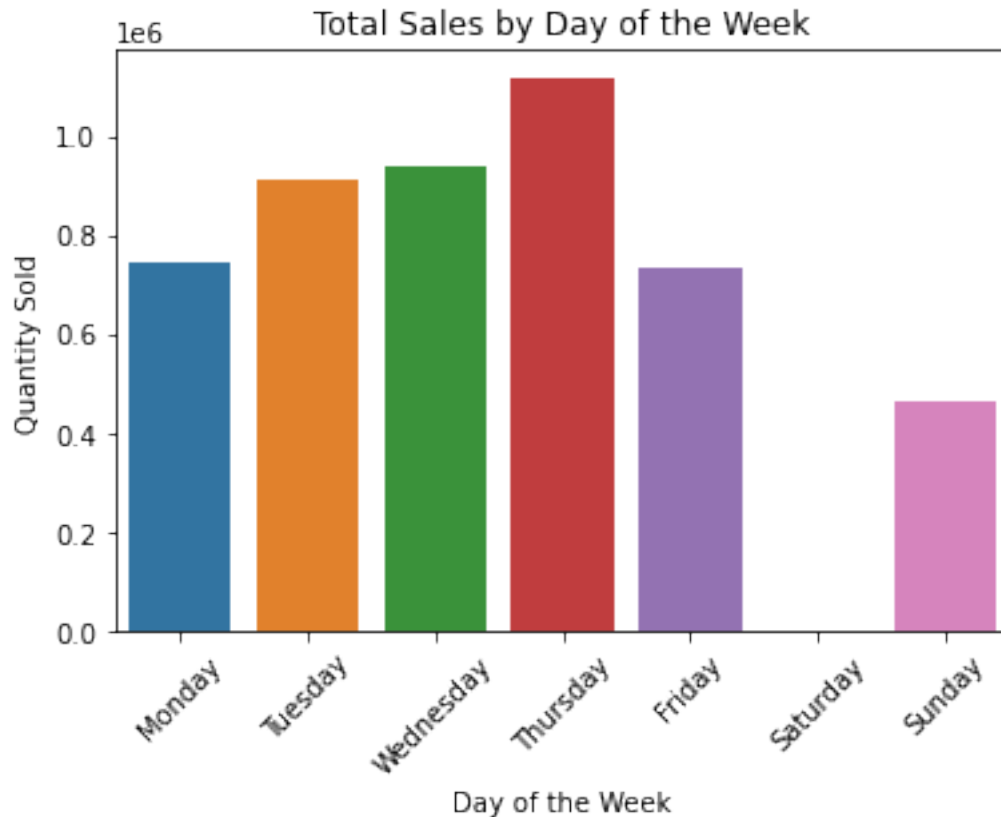
[4]:
```python
import matplotlib.pyplot as plt
import seaborn as sns

df['InvoiceDate'] = pd.to_datetime(df['InvoiceDate'])
df['Month'] = df['InvoiceDate'].dt.month
df['Day'] = df['InvoiceDate'].dt.day_name()

monthly_sales = df.groupby('Month').agg({'Quantity': 'sum'}).reset_index()
sns.barplot(x='Month', y='Quantity', data=monthly_sales)
plt.title('Total Sales per Month')
plt.xlabel('Month')
plt.ylabel('Quantity Sold')
plt.xticks(rotation=45)
plt.show()
```
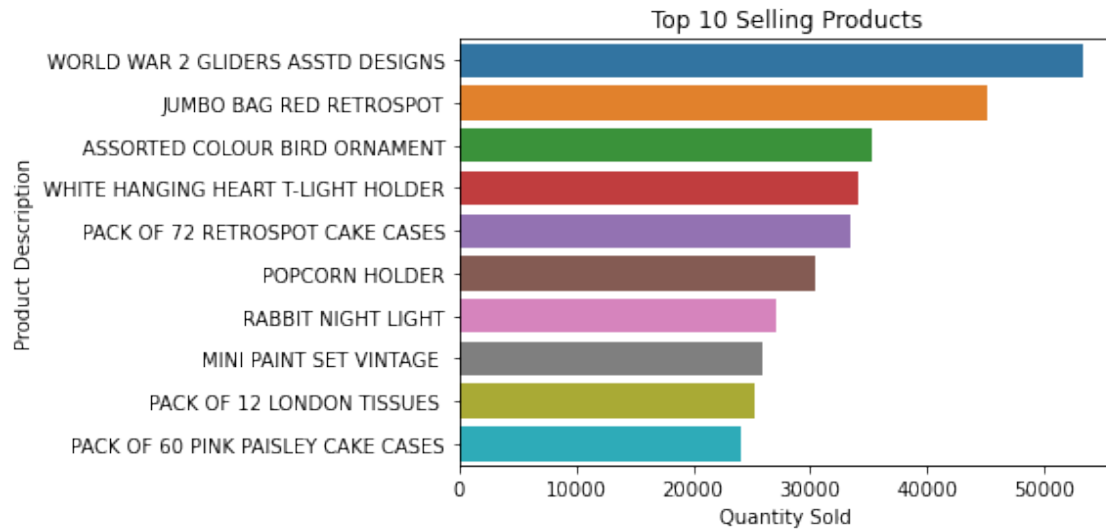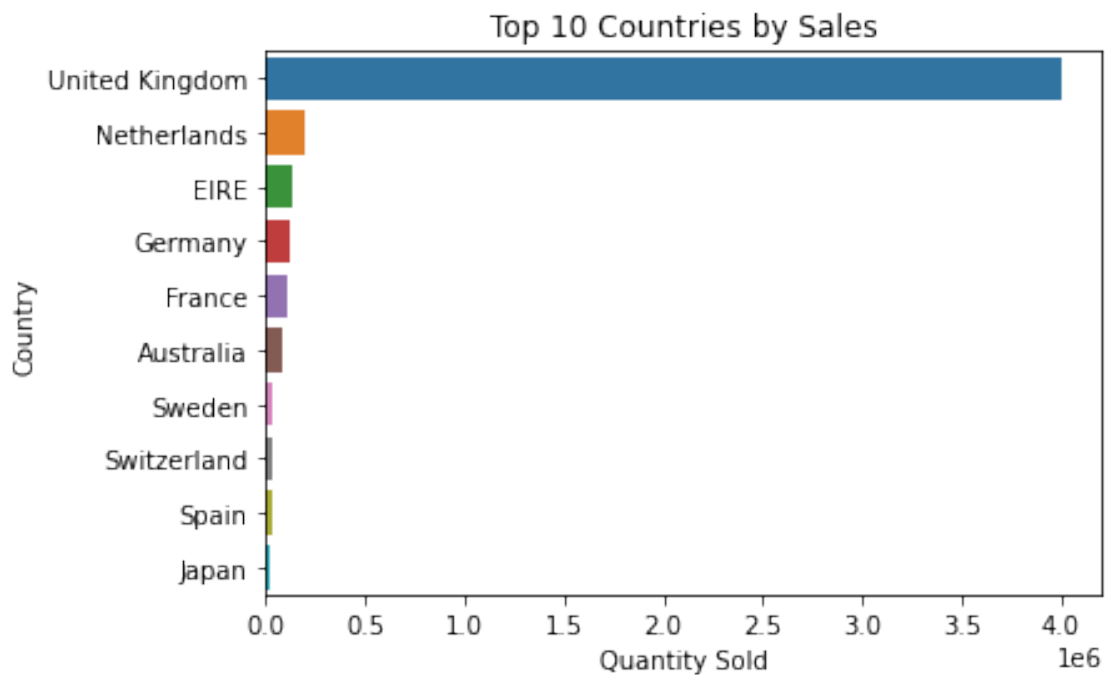
```
[5]: weekly_sales = df.groupby('Day').agg({'Quantity': 'sum'}).reindex(['Monday',␣
     ↪'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday'])
     sns.barplot(x=weekly_sales.index, y='Quantity', data=weekly_sales.reset_index())
     plt.title('Total Sales by Day of the Week')
     plt.xlabel('Day of the Week')
     plt.ylabel('Quantity Sold')
     plt.xticks(rotation=45)
     plt.show()
```
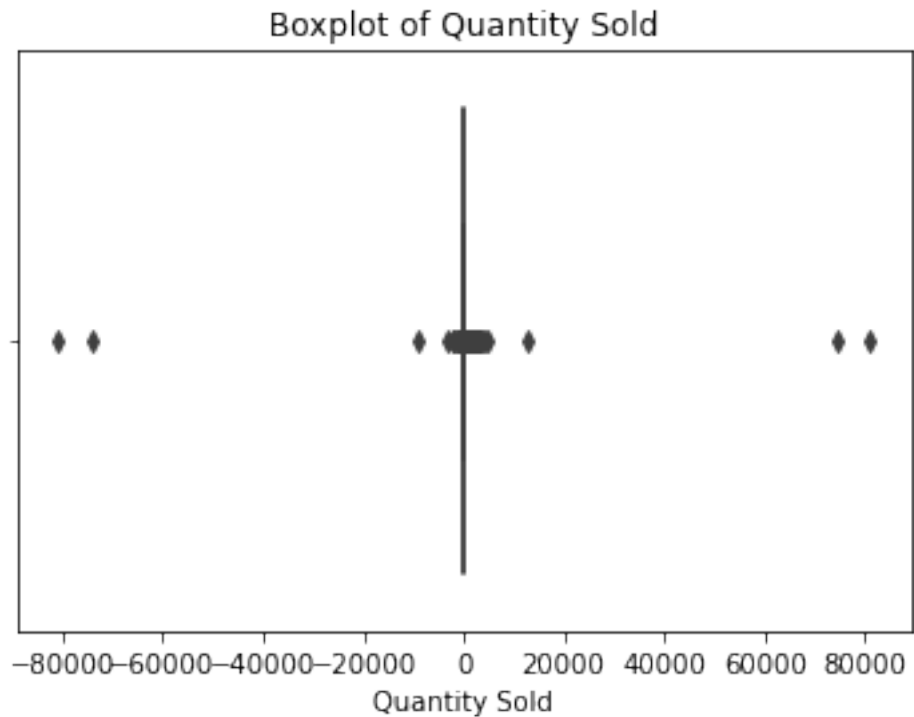


```
[6]: top_products = df.groupby('Description').agg({'Quantity': 'sum'}).nlargest(10,␣
     ↪'Quantity')
     top_products.reset_index(inplace=True)
     sns.barplot(x='Quantity', y='Description', data=top_products)
     plt.title('Top 10 Selling Products')
     plt.xlabel('Quantity Sold')
     plt.ylabel('Product Description')
     plt.show()
```

## Top 10 Selling Products

| Product Description | Quantity Sold |
|---|---|



```
[7]: country_sales = df.groupby('Country').agg({'Quantity': 'sum'}).nlargest(10,␣
     ↪'Quantity')
     country_sales.reset_index(inplace=True)
     sns.barplot(x='Quantity', y='Country', data=country_sales)
     plt.title('Top 10 Countries by Sales')
     plt.xlabel('Quantity Sold')
     plt.ylabel('Country')
     plt.show()
```

```
[8]: sns.boxplot(x=df['Quantity'])
     plt.title('Boxplot of Quantity Sold')
     plt.xlabel('Quantity Sold')
     plt.show()
```

Boxplot of Quantity Sold



```
[9]: summary = {
         "Busiest Month": monthly_sales['Month'][monthly_sales['Quantity'].idxmax()],
         "Top Selling Product": top_products['Description'].iloc[0],
         "Most Active Country": country_sales['Country'].iloc[0]
     }

     print("Summary of Findings:")
     for key, value in summary.items():
         print(f"{key}: {value}")
```

```
Summary of Findings:
Busiest Month: 11
Top Selling Product: WORLD WAR 2 GLIDERS ASSTD DESIGNS
Most Active Country: United Kingdom
```