



```
In [1]: 1 import pandas as pd
        2 import numpy as np
```

```
In [2]: 1 df = pd.read_csv(r'C:\Users\LENOVO\Downloads\Online Retail Sample.csv')
        2 print("Initial shape:", df.shape)
```

Initial shape: (3066, 8)

```
In [3]: 1 df.head()
```

Out[3]:

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A WHITE HANGING HEART T-LIGHT HOLDER	6	1/12/2010 8:26	2.55	17850.0	United Kingdom
1	536365	71053 WHITE METAL LANTERN	6	1/12/2010 8:26	3.39	17850.0	United Kingdom
2	536365	84406B CREAM CUPID HEARTS COAT HANGER	8	1/12/2010 8:26	2.75	17850.0	United Kingdom
3	536365	84029G KNITTED UNION FLAG HOT WATER BOTTLE	6	1/12/2010 8:26	3.39	17850.0	United Kingdom
4	536365	84029E RED WOOLLY HOTTIE WHITE HEART.	6	1/12/2010 8:26	3.39	17850.0	United Kingdom



Data Cleaning

```
In [4]: 1 # Data types and missing values
        2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3066 entries, 0 to 3065
Data columns (total 8 columns):
 #   Column        Non-Null Count  Dtype  
---  --
 0   InvoiceNo      3066 non-null   object  
 1   StockCode     3066 non-null   object  
 2   Description    3056 non-null   object  
 3   Quantity      3066 non-null   int64   
 4   InvoiceDate    3066 non-null   object  
 5   UnitPrice     3066 non-null   float64  
 6   CustomerID    1932 non-null   float64  
 7   Country       3066 non-null   object  
dtypes: float64(2), int64(1), object(5)
memory usage: 191.8+ KB
```

```
In [5]: 1 # Count missing values
        2 print("\nMissing values:\n", df.isnull().sum())
```

```
Missing values:
InvoiceNo      0
StockCode      0
Description    10
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID    1134
Country        0
dtype: int64
```

```
In [6]: 1 # Drop rows with missing CustomerID or Description
        2 df_cleaned = df.dropna(subset=['CustomerID', 'Description'])
        3 print("After dropping missing values:", df_cleaned.shape)
```

After dropping missing values: (1932, 8)

```
In [7]: 1 # Cancelled orders have InvoiceNo starting with 'C'
        2 df_cleaned = df_cleaned[~df_cleaned['InvoiceNo'].astype(str).str.startswith('C')]
        3 print("After removing cancelled invoices:", df_cleaned.shape)
```

After removing cancelled invoices: (1906, 8)

```
In [8]: 1 #Remove Negative or Zero Quantities and Unit Prices
        2 df_cleaned = df_cleaned[(df_cleaned['Quantity'] > 0) & (df_cleaned['UnitPrice'] > 0)]
        3 print("After removing negative/zero quantity & price:", df_cleaned.shape)
```

After removing negative/zero quantity & price: (1906, 8)

```
In [9]: 1 # Convert InvoiceDate to Datetime...
        2 df_cleaned['InvoiceDate'] = pd.to_datetime(df_cleaned['InvoiceDate'])
        3 df_cleaned['InvoiceDate'].dtype
```

Out[9]: dtype('<M8[ns]')

```
In [10]: 1 #Add Total Amount Column
        2 df_cleaned['TotalAmount'] = df_cleaned['Quantity'] * df_cleaned['UnitPrice']
        3 df_cleaned['TotalAmount'].describe()
```

```
Out[10]: count    1906.000000
         mean      24.234208
         std       66.378259
         min        0.290000
         25%        4.200000
         50%       11.800000
         75%       20.280000
         max      1627.200000
         Name: TotalAmount, dtype: float64
```



```
In [11]: 1 # Final check for nulls or duplicates
          2 print("Final missing values:\n", df_cleaned.isnull().sum())
          3 print("Number of duplicate rows:", df_cleaned.duplicated().sum())
```

Final missing values:

InvoiceNo 0

StockCode 0

Description 0

Quantity 0

InvoiceDate 0

UnitPrice 0

CustomerID 0

Country 0

TotalAmount 0

dtype: int64

Number of duplicate rows: 44

```
In [12]: 1 # Drop exact duplicate rows
          2 df_cleaned = df_cleaned.drop_duplicates()
          3 print("Duplicates after dropping:", df_cleaned.duplicated().sum())
```

Duplicates after dropping: 0



Exploratory Data Analysis (EDA)

```
In [13]: 1 import matplotlib.pyplot as plt
          2 import seaborn as sns
          3
          4 # Set plot style
          5 plt.style.use('seaborn-v0_8-whitegrid')
          6 sns.set_palette('muted')
```

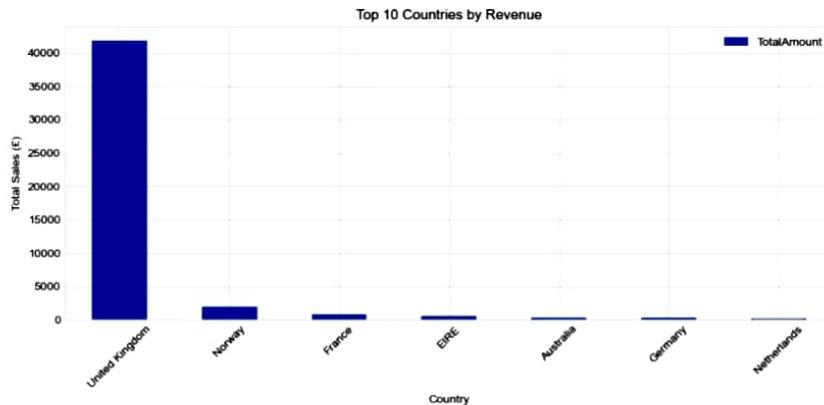
```
In [14]: 1 # Check basic stats
          2 df_cleaned.describe()
          3
          4 # Number of unique products, customers, and countries
          5 print("Unique Products:", df_cleaned['Description'].nunique())
          6 print("Unique Customers:", df_cleaned['CustomerID'].nunique())
          7 print("Countries:", df_cleaned['Country'].nunique())
```

```
Unique Products: 927
Unique Customers: 94
Countries: 7
```

```
In [15]: 1 print("Min date:", df_cleaned['InvoiceDate'].min())
          2 print("Max date:", df_cleaned['InvoiceDate'].max())
          3
```

```
Min date: 2010-01-12 08:26:00
Max date: 2010-01-12 17:24:00
```

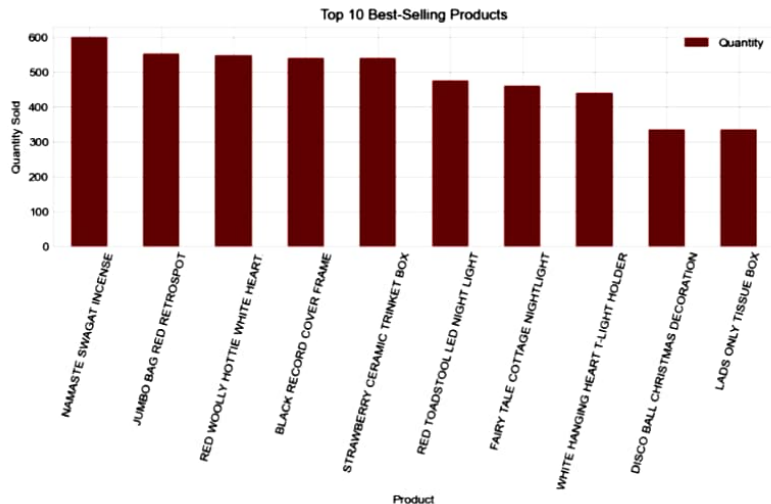
```
In [30]: 1 # Top 10 Countries by Revenue
2 country_sales = df_cleaned.groupby('Country')['TotalAmount'].sum().sort_values(ascending=False).head(10)
3
4 # Plot
5 plt.figure(figsize=(10,5))
6 country_sales.plot(kind='bar', color='darkblue')
7 plt.title('Top 10 Countries by Revenue')
8 plt.ylabel('Total Sales (£)')
9 plt.xlabel('Country')
10 plt.xticks(rotation=45)
11 plt.legend()
12 plt.tight_layout()
13 plt.savefig('Top_10_Countries_by_Revenue.png',dpi=300, bbox_inches='tight')
14 plt.show()
```



```

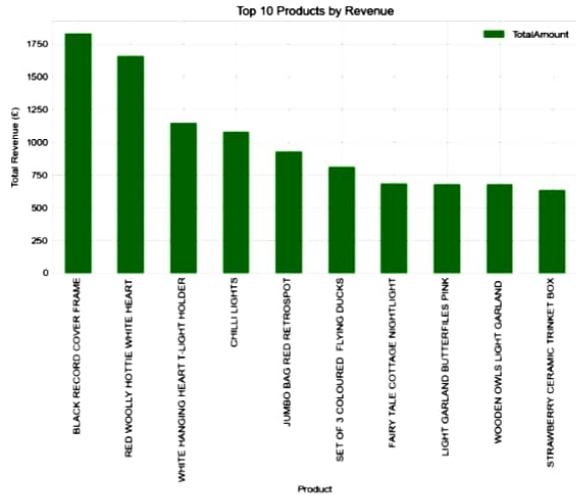
In [31]: 1 #Best-Selling Products
2 product_sales = df_cleaned.groupby('Description')['Quantity'].sum().sort_values(ascending=False).head(10)
3
4 # Plot
5 plt.figure(figsize=(9,6))
6 product_sales.plot(kind='bar', color='darkred')
7 plt.title('Top 10 Best-Selling Products')
8 plt.ylabel('Quantity Sold')
9 plt.xlabel('Product')
10 plt.xticks(rotation=75)
11 plt.legend()
12 plt.tight_layout()
13 plt.savefig('Top_10_Best-Selling_Products',dpi=300,bbox_inches='tight')
14 plt.show()

```

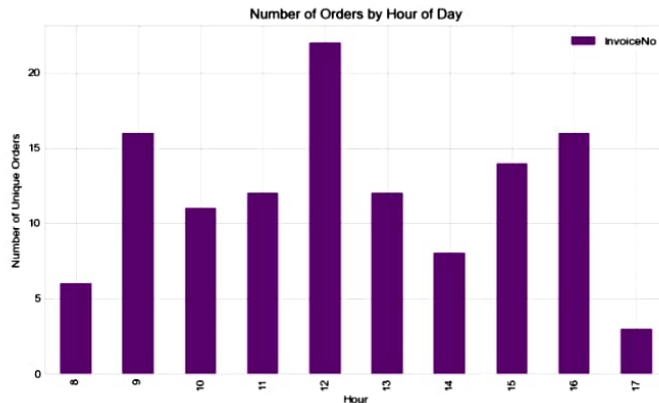


In [33]:

```
1 #High Revenue Products
2 top_revenue_products = df_cleaned.groupby('Description')['TotalAmount'].sum().sort_values(ascending=False).head(10)
3
4 # Plot
5 plt.figure(figsize=(8,7))
6 top_revenue_products.plot(kind='bar', color='green')
7 plt.title('Top 10 Products by Revenue')
8 plt.ylabel('Total Revenue (£)')
9 plt.xlabel('Product')
10 plt.xticks(rotation=90)
11 plt.legend()
12 plt.tight_layout()
13 plt.savefig('Top_10_Products_by_Revenue', dpi=300, bbox_inches='tight')
14 plt.show()
15
```




```
In [34]: 1 #Number of Orders per Hour
2 df_cleaned['Hour'] = df_cleaned['InvoiceDate'].dt.hour
3 hourly_orders = df_cleaned.groupby('Hour')['InvoiceNo'].nunique()
4
5 # Plot
6 plt.figure(figsize=(8,5))
7 hourly_orders.plot(kind='bar', color='purple')
8 plt.title('Number of Orders by Hour of Day')
9 plt.xlabel('Hour')
10 plt.ylabel('Number of Unique Orders')
11 plt.legend()
12 plt.tight_layout()
13 plt.savefig('Number_of_Orders_by_Hour_of_Day',dpi=300,bbbox_inches='tight')
14 plt.show()
15
```



```
In [42]: 1 # Returns Analysis
2 # Look at negative quantities
3 if not most_returned.empty:
4     most_returned.plot(kind='bar', color='red')
5     plt.title('Most Frequently Returned Products')
6     plt.ylabel('Number of Returns')
7     plt.xlabel('Product')
8     plt.xticks(rotation=75)
9     plt.tight_layout()
10    plt.show()
11 else:
12     print("No returned products found in the dataset.")
13
```

No returned products found in the dataset.

```
In [35]: 1 #Top Customers by Spend
2 top_customers = df_cleaned.groupby('CustomerID')['TotalAmount'].sum().sort_values(ascending=False).head(10)
3
4 # Plot
5 plt.figure(figsize=(8,6))
6 top_customers.plot(kind='bar', color='brown')
7 plt.title('Top 10 Customers by Revenue')
8 plt.ylabel('Total Spend ($)')
9 plt.xlabel('Customer ID')
10 plt.legend()
11 plt.tight_layout()
12 plt.savefig('Top_10_Customers_by_Revenue',dpi=100,bbox_inches='tight')
13 plt.show()
14
```

