

```
In [50]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

```
In [25]: 1 File = "C:\\Users\\hitak\\Downloads\\OnlineRetail.xlsx"
2 df =pd.read_excel(File)
3 df
```

Out[25]:

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Cou
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	Ur King
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	Ur King
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
...
541904	581587	22613	PACK OF 20 SPACEBOY NAPKINS	12	2011-12-09 12:50:00	0.85	12680.0	Fr
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	6	2011-12-09 12:50:00	2.10	12680.0	Fr
541906	581587	23254	CHILDRENS CUTLERY DOLLY GIRL	4	2011-12-09 12:50:00	4.15	12680.0	Fr
541907	581587	23255	CHILDRENS CUTLERY CIRCUS PARADE	4	2011-12-09 12:50:00	4.15	12680.0	Fr
541908	581587	22138	BAKING SET 9 PIECE RETROSPOT	3	2011-12-09 12:50:00	4.95	12680.0	Fr

541909 rows × 8 columns



```
In [26]: 1 df.head()
```

Out[26]:

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom

```
In [26]: 1 df.head()
```

```
Out[26]:
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom

```
In [27]: 1 df.tail()
```

```
Out[27]:
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Cou
541904	581587	22613	PACK OF 20 SPACEBOY NAPKINS	12	2011-12-09 12:50:00	0.85	12680.0	Fre
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	6	2011-12-09 12:50:00	2.10	12680.0	Fre
541906	581587	23254	CHILDRENS CUTLERY DOLLY GIRL	4	2011-12-09 12:50:00	4.15	12680.0	Fre
541907	581587	23255	CHILDRENS CUTLERY CIRCUS PARADE	4	2011-12-09 12:50:00	4.15	12680.0	Fre
541908	581587	22138	BAKING SET 9 PIECE RETROSPOT	3	2011-12-09 12:50:00	4.95	12680.0	Fre

```
In [28]: 1 df.shape
```

```
Out[28]: (541909, 8)
```

```
In [29]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   InvoiceNo        541909 non-null object
1   StockCode        541909 non-null object
2   Description      540455 non-null object
3   Quantity         541909 non-null int64
4   InvoiceDate       541909 non-null datetime64[ns]
5   UnitPrice        541909 non-null float64
```

```
In [29]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   InvoiceNo        541909 non-null object
1   StockCode        541909 non-null object
2   Description      540455 non-null object
3   Quantity         541909 non-null int64
4   InvoiceDate      541909 non-null datetime64[ns]
5   UnitPrice        541909 non-null float64
6   CustomerID       406829 non-null float64
7   Country          541909 non-null object
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB
```

```
In [30]: 1 df.describe()
```

```
Out[30]:
```

	Quantity	UnitPrice	CustomerID
count	541909.000000	541909.000000	406829.000000
mean	9.552250	4.611114	15287.690570
std	218.081158	96.759853	1713.600303
min	-80995.000000	-11062.060000	12346.000000
25%	1.000000	1.250000	13953.000000
50%	3.000000	2.080000	15152.000000
75%	10.000000	4.130000	16791.000000
max	80995.000000	38970.000000	18287.000000

```
In [31]: 1 df.duplicated().sum()
```

```
Out[31]: 5268
```

```
In [32]: 1 df[df.duplicated()]
```

```
Out[32]:
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID
517	536409	21866	UNION JACK FLAG LUGGAGE TAG	1	2010-12-01 11:45:00	1.25	17908.0
527	536409	22866	HAND WARMER SCOTTY DOG DESIGN	1	2010-12-01 11:45:00	2.10	17908.0
537	536409	22900	SET 2 TEA TOWELS I LOVE LONDON	1	2010-12-01 11:45:00	2.95	17908.0
539	536409	22111	SCOTTIE DOG HOT WATER	1	2010-12-01 11:45:00	4.95	17908.0

```
1 df.isnull().sum()
```

```
Out[33]: InvoiceNo      0
StockCode      0
Description    1454
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID    135080
Country        0
dtype: int64
```

```
539 530409 22111 WATER 1 11:45:00 4.95 17908.0
1 df.isnull().sum()
```

```
Out[33]: InvoiceNo      0
StockCode      0
Description    1454
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID    135080
Country        0
dtype: int64
```

```
In [41]: 1 df["Description"].fillna("No Description", inplace=True)
2 df["CustomerID"].fillna(-1, inplace=True)
3 df
```

```
Out[41]:
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Cou
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	Ur King
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	Ur King
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	Ur King
...
541904	581587	22613	PACK OF 20 SPACEBOY NAPKINS	12	2011-12-09 12:50:00	0.85	12680.0	Fr
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	6	2011-12-09 12:50:00	2.10	12680.0	Fr
541906	581587	23254	CHILDRENS CUTLERY DOLLY GIRL	4	2011-12-09 12:50:00	4.15	12680.0	Fr
541907	581587	23255	CHILDRENS CUTLERY CIRCUS PARADE	4	2011-12-09 12:50:00	4.15	12680.0	Fr
541908	581587	22138	BAKING SET 9 PIECE RETROSPOT	3	2011-12-09 12:50:00	4.95	12680.0	Fr

541909 rows × 8 columns

```
In [43]: df.info
```

```
Out[42]: <bound method DataFrame.info of
Out[42]: Description Quantity \
0 Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate',
1 'UnitPrice', 'CustomerID', 'Country'], dtype=object)
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
... ..
541904 581587 22613 PACK OF 20 SPACEBOY NAPKINS 12
541905 581587 22899 CHILDREN'S APRON DOLLY GIRL 6
541906 581587 23254 CHILDRENS CUTLERY DOLLY GIRL 4
```

```
In [43]: df.info
```

```
Out[43]: <bound method DataFrame.info of
Description Quantity \
Index(['InvoiceNo', 'StockCode', 'Description', 'Quantity', 'InvoiceDate',
1      'UnitPrice', 'CustomerID', 'Country'],
2      dtype=object)
3      536365      84029G  KNITTED UNION FLAG HOT WATER BOTTLE      6
4      536365      84029E  RED WOOLLY HOTTIE WHITE HEART.      6
...      ...      ...      ...      ...
541904      581587      22613      PACK OF 20 SPACEBOY NAPKINS      12
541905      581587      22899      CHILDREN'S APRON DOLLY GIRL      6
541906      581587      23254      CHILDRENS CUTLERY DOLLY GIRL      4
541907      581587      23255      CHILDRENS CUTLERY CIRCUS PARADE      4
541908      581587      22138      BAKING SET 9 PIECE RETROSPOT      3

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
...      ...      ...      ...      ...
541904 2011-12-09 12:50:00      0.85      12680.0      France
541905 2011-12-09 12:50:00      2.10      12680.0      France
541906 2011-12-09 12:50:00      4.15      12680.0      France
541907 2011-12-09 12:50:00      4.15      12680.0      France
541908 2011-12-09 12:50:00      4.95      12680.0      France

[541909 rows x 8 columns]>
```

```
In [44]: 1 df['Description'].value_counts()
```

```
Out[44]: WHITE HANGING HEART T-LIGHT HOLDER      2369
REGENCY CAKESTAND 3 TIER      2200
JUMBO BAG RED RETROSPOT      2159
PARTY BUNTING      1727
LUNCH BAG RED RETROSPOT      1638
...
Missing      1
historic computer difference?...se      1
DUSTY PINK CHRISTMAS TREE 30CM      1
WRAP BLUE RUSSIAN FOLKART      1
PINK BERTIE MOBILE PHONE CHARM      1
Name: Description, Length: 4224, dtype: int64
```

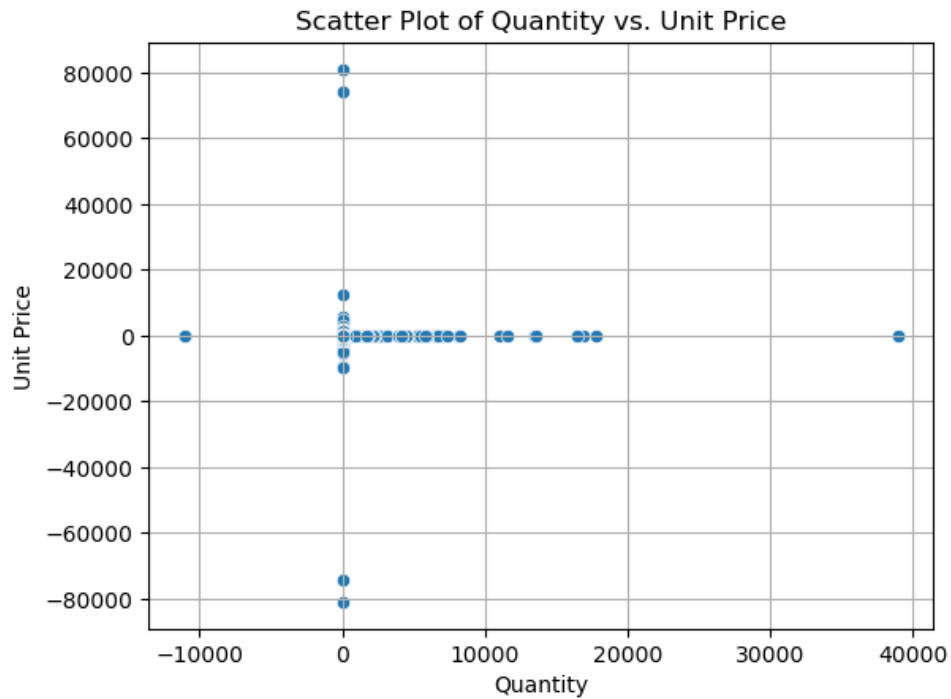
visualization

```
In [56]: 1 #scatter plot
2
3 sns.scatterplot(x='UnitPrice', y='Quantity',data=df)
4 plt.title('Scatter Plot of Quantity vs. Unit Price')
5 plt.xlabel('Quantity')
6 plt.ylabel('Unit Price')
7 plt.grid(True)
```

Scatter Plot of Quantity vs. Unit Price

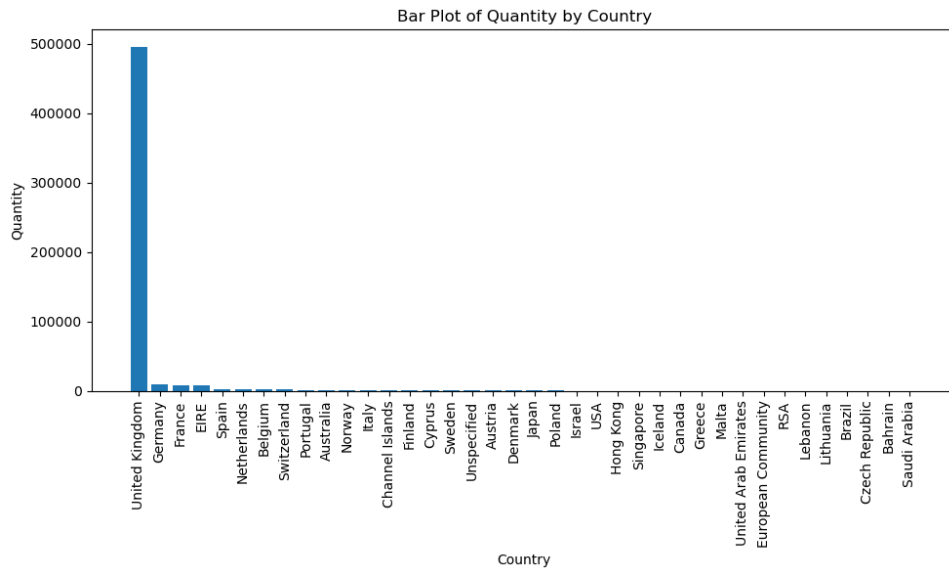
visualization

```
In [56]: 1 #scatter plot
2
3 sns.scatterplot(x='UnitPrice', y='Quantity',data=df)
4 plt.title('Scatter Plot of Quantity vs. Unit Price')
5 plt.xlabel('Quantity')
6 plt.ylabel('Unit Price')
7 plt.grid(True)
```

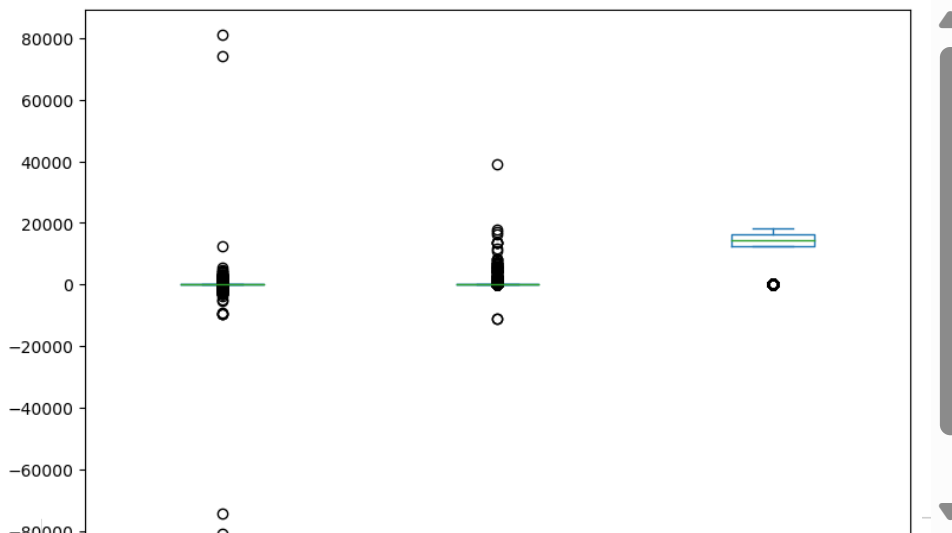


```
In [58]: 1 #Barplot
2
3 country_counts = df['Country'].value_counts()
4 countries = country_counts.index
5 quantity_values = country_counts.values
6
7 plt.figure(figsize=(10, 6))
8 plt.bar(countries, quantity_values)
9 plt.title('Bar Plot of Quantity by Country')
10 plt.xlabel('Country')
11 plt.ylabel('Quantity')
12 plt.xticks(rotation=90)
13 plt.tight_layout()
```

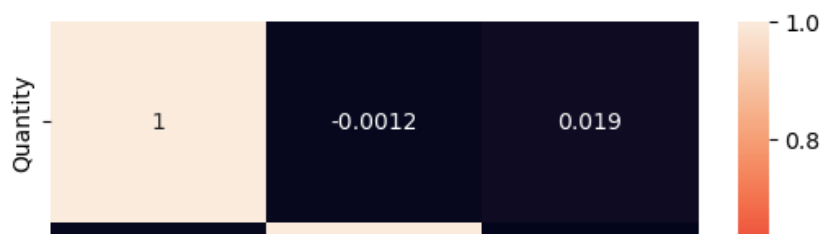
```
In [58]: 1 #BarPlot
2
3 country_counts = df['Country'].value_counts()
4 countries = country_counts.index
5 quantity_values = country_counts.values
6
7 plt.figure(figsize=(10, 6))
8 plt.bar(countries, quantity_values)
9 plt.title('Bar Plot of Quantity by Country')
10 plt.xlabel('Country')
11 plt.ylabel('Quantity')
12 plt.xticks(rotation=90)
13 plt.tight_layout()
14 plt.show()
15
```



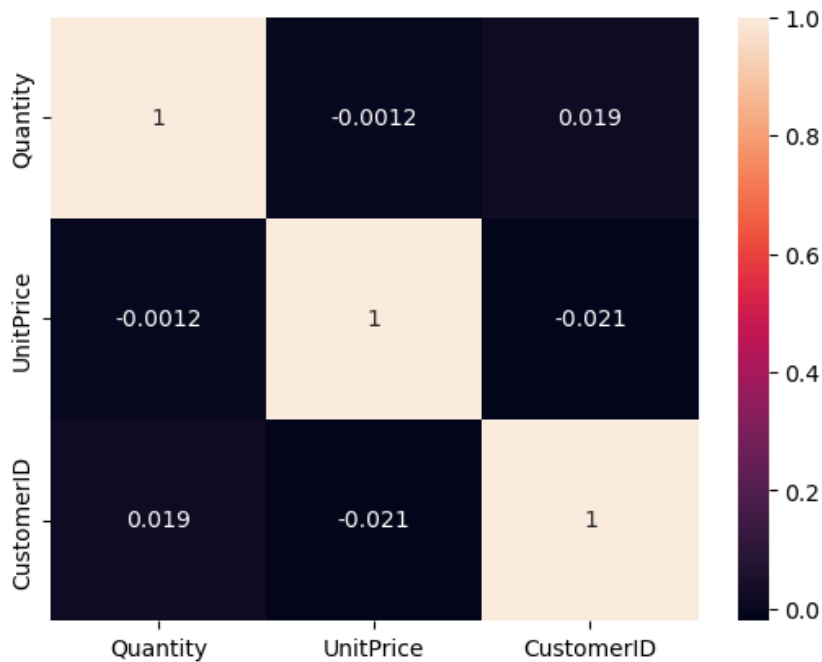
```
In [52]: 1 #box plot for outlier detection
2 df.plot(kind='box', figsize=(9,6))
3 plt.show()
```



```
In [53]: 1 #HeatMap
2 sns.heatmap(df.corr(),annot=True)
3 plt.show()
```



```
In [53]: 1 #HeatMap
2 sns.heatmap(df.corr(),annot=True)
3 plt.show()
```

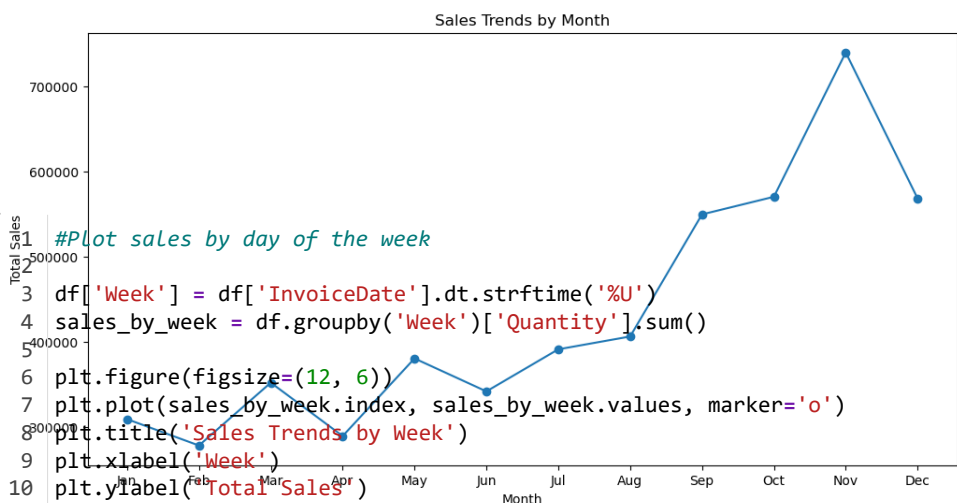


Analyze sales trends over time

```
In [70]: 1 df['InvoiceDate'] = pd.to_datetime(df['InvoiceDate'])
```

```
In [78]: 1 df['Month'] = df['InvoiceDate'].dt.month
2 df['Hour'] = df['InvoiceDate'].dt.hour
3
4 sales_by_month = df.groupby('Month')['Quantity'].sum()
```

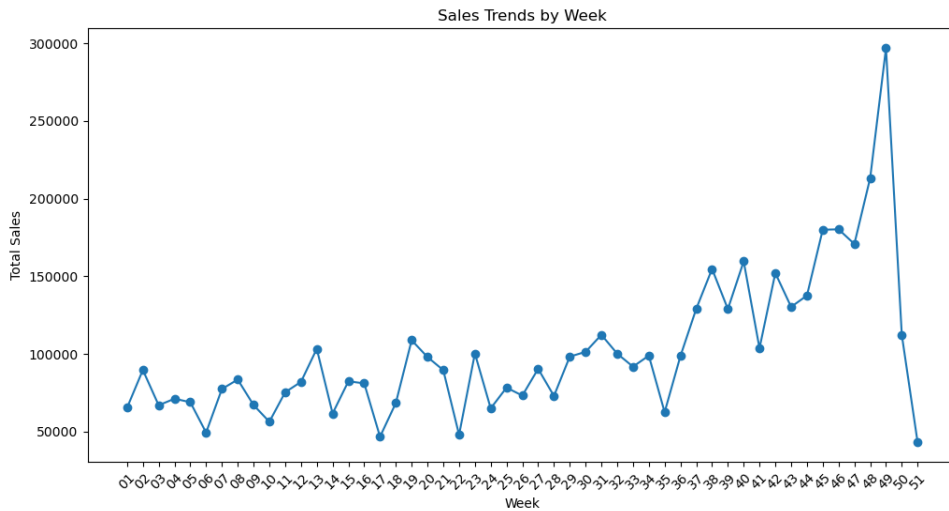
```
In [79]: 1 # Plot sales by month
2
3 plt.figure(figsize=(12, 6))
4 plt.plot(sales_by_month.index, sales_by_month.values, marker='o')
5 plt.title('Sales Trends by Month')
6 plt.xlabel('Month')
7 plt.ylabel('Total Sales')
8 plt.xticks(range(1, 13), ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul',
9 'Aug', 'Sep', 'Oct', 'Nov', 'Dec'])
9 plt.show()
10
```



```
In [80]: 1 #Plot sales by day of the week
2
3 df['Week'] = df['InvoiceDate'].dt.strftime('%U')
4 sales_by_week = df.groupby('Week')['Quantity'].sum()
5
6 plt.figure(figsize=(12, 6))
7 plt.plot(sales_by_week.index, sales_by_week.values, marker='o')
8 plt.title('Sales Trends by Week')
9 plt.xlabel('Week')
10 plt.ylabel('Total Sales')
11 plt.xticks(rotation=45)
12 plt.show()
```


In [80]:

```
1 #Plot sales by day of the week
2
3 df['Week'] = df['InvoiceDate'].dt.strftime('%U')
4 sales_by_week = df.groupby('Week')['Quantity'].sum()
5
6 plt.figure(figsize=(12, 6))
7 plt.plot(sales_by_week.index, sales_by_week.values, marker='o')
8 plt.title('Sales Trends by Week')
9 plt.xlabel('Week')
10 plt.ylabel('Total Sales')
11 plt.xticks(rotation=45)
12 plt.show()
```



Top selling product and countries

In [85]:

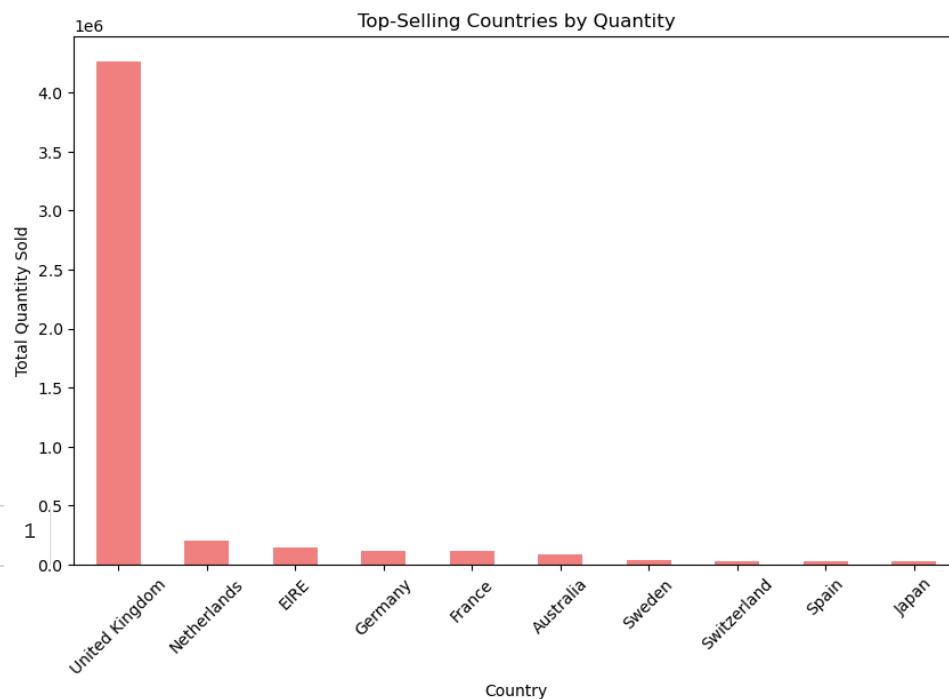
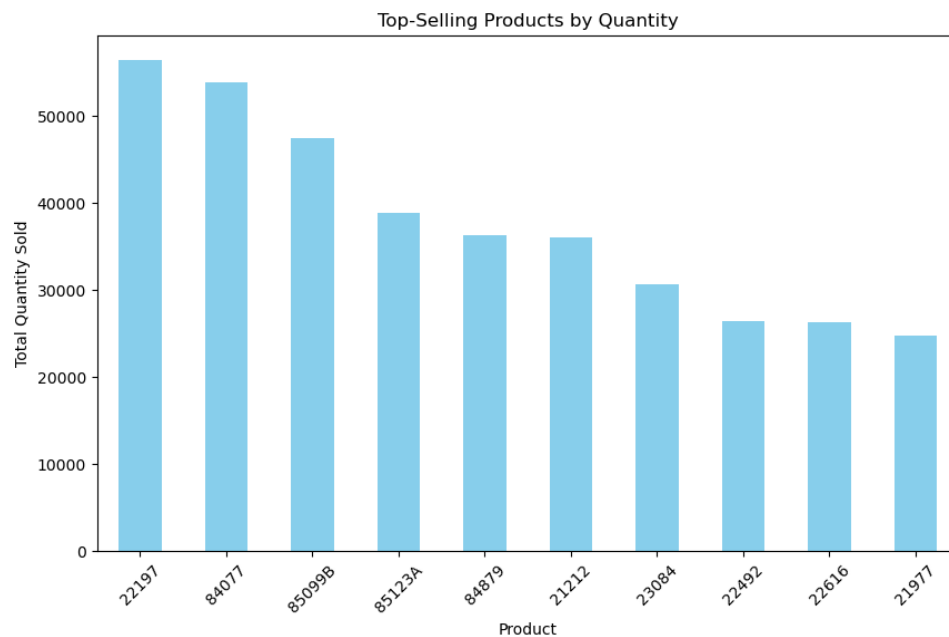
```
1 product_country_sales = df.groupby(['StockCode', 'Country'])['Quantity'].
2
3 top_selling_products = product_country_sales.groupby('StockCode')['Quantity'].
4
5 top_selling_countries = product_country_sales.groupby('Country')['Quantity'].
```

In [86]:

```
1 # Plot top-selling products
2 plt.figure(figsize=(10, 6))
3 top_selling_products.plot(kind='bar', color='skyblue')
4 plt.title('Top-Selling Products by Quantity')
5 plt.xlabel('Product')
6 plt.ylabel('Total Quantity Sold')
7 plt.xticks(rotation=45)
8 plt.show()
9
10 # Plot top-selling countries
11 plt.figure(figsize=(10, 6))
12 top_selling_countries.plot(kind='bar', color='lightcoral')
13 plt.title('Top-Selling Countries by Quantity')
```

In [86]:

```
1 # Plot top-selling products
2 plt.figure(figsize=(10, 6))
3 top_selling_products.plot(kind='bar', color='skyblue')
4 plt.title('Top-Selling Products by Quantity')
5 plt.xlabel('Product')
6 plt.ylabel('Total Quantity Sold')
7 plt.xticks(rotation=45)
8 plt.show()
9
10 # Plot top-selling countries
11 plt.figure(figsize=(10, 6))
12 top_selling_countries.plot(kind='bar', color='lightcoral')
13 plt.title('Top-Selling Countries by Quantity')
14 plt.xlabel('Country')
15 plt.ylabel('Total Quantity Sold')
16 plt.xticks(rotation=45)
17 plt.show()
18
```



In []:

