

## Importing necessary Libraries

```
[2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[3]: #Loading Online Sales Data
df = pd.read_csv(r'C:\Users\...')
df.head()
```

```
[3]:
```

	Transaction ID	Date
0	10001	2024-01-01
1	10002	2024-01-02
2	10003	2024-01-03
3	10004	2024-01-04
4	10005	2024-01-05

## Basic Data Exploration

```
[4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 240 entries, 0 to 239
Data columns (total 9 columns):
 #   Column             Non-Null Count  Dtype  
---  -
 0   Transaction ID      240 non-null   int64  
 1   Date                240 non-null   object  
 2   Product Category    240 non-null   object  
 3   Product Name        240 non-null   object  
 4   Units Sold          240 non-null   int64  
 5   Unit Price          240 non-null   float64 
 6   Total Revenue       240 non-null   float64 
 7   Region              240 non-null   object  
```

ist Checkpoint: 12 days ago

Help

irkdown ▾

libraries

```
plt

pc\Desktop\60DaysDS\Online Sales Data.csv')
```

Product Category	Product Name	Units Sold	Unit Price	Total Revenue	Region
Electronics	iPhone 14 Pro	2	999.99	1999.98	North America
Home Appliances	Dyson V11 Vacuum	1	499.99	499.99	Europe
Clothing	Levi's 501 Jeans	3	69.99	209.97	Asia
Books	The Da Vinci Code	4	15.99	63.96	North America
Beauty Products	Neutrogena Skincare Set	1	89.99	89.99	Europe

1

```
taFrame'>
o 239
is):
null Count  Dtype
-----  -
ion-null    int64
ion-null    object
ion-null    object
ion-null    object
ion-null    int64
ion-null    float64
ion-null    float64
ion-null    object
```



Trusted

JupyterLab Python 3 (ipykernel)



### Payment Method

Credit Card

PayPal

Debit Card

Credit Card

PayPal

```
/    region      240 n
8    Payment Method  240 n
dtypes: float64(2), int64(2)
memory usage: 17.0+ KB
```

```
[5]: df.describe()
```

```
[5]:
```

	Transaction ID	Units So
count	240.00000	240.0000
mean	10120.50000	2.1583
std	69.42622	1.3224
min	10001.00000	1.0000
25%	10060.75000	1.0000
50%	10120.50000	2.0000
75%	10180.25000	3.0000
max	10240.00000	10.0000

## Data Cleaning

```
[6]: df.isnull().sum() # we have
```

```
[6]: Transaction ID      0
Date                  0
Product Category      0
Product Name          0
Units Sold            0
Unit Price            0
Total Revenue         0
Region               0
Payment Method        0
dtype: int64
```

```
[7]: df.duplicated().sum()
```

```
[7]: 0
```

### datatypes

```
[8]: df['Date']=pd.to_datetime(df
df['Transaction ID']=df['Tra
df.dtypes
```

```
ion-null    object
ion-null    object
, object(5)
```

Id	Unit Price	Total Revenue
100	240.000000	240.000000
133	236.395583	335.699375
154	429.446695	485.804469
100	6.500000	6.500000
100	29.500000	62.965000
100	89.990000	179.970000
100	249.990000	399.225000
100	3899.990000	3899.990000

*missing values in this dataset*

```
['Date'])
transaction ID'].astype(str)
```



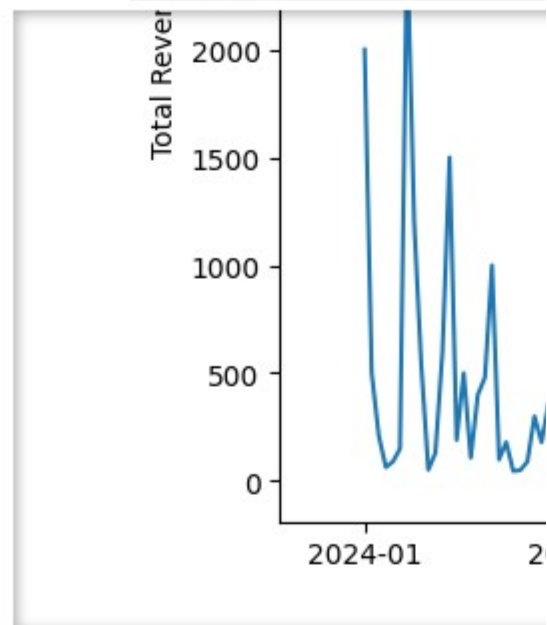
```
[8]: Transaction ID
      Date          datetime
      Product Category
      Product Name
      Units Sold
      Unit Price      f
      Total Revenue    f
      Region
      Payment Method
      dtype: object
```

## Data analysis

**Sales over time - To see the tre**

```
[9]: sales_over_time = df.groupby
      #sales_over_time

      plt.figure(figsize=(12, 6))
      plt.plot(sales_over_time['Da
      plt.title('Sales Over Time')
      plt.xlabel('Date')
      plt.ylabel('Total Revenue')
      plt.show()
```



```
[ ]:
```

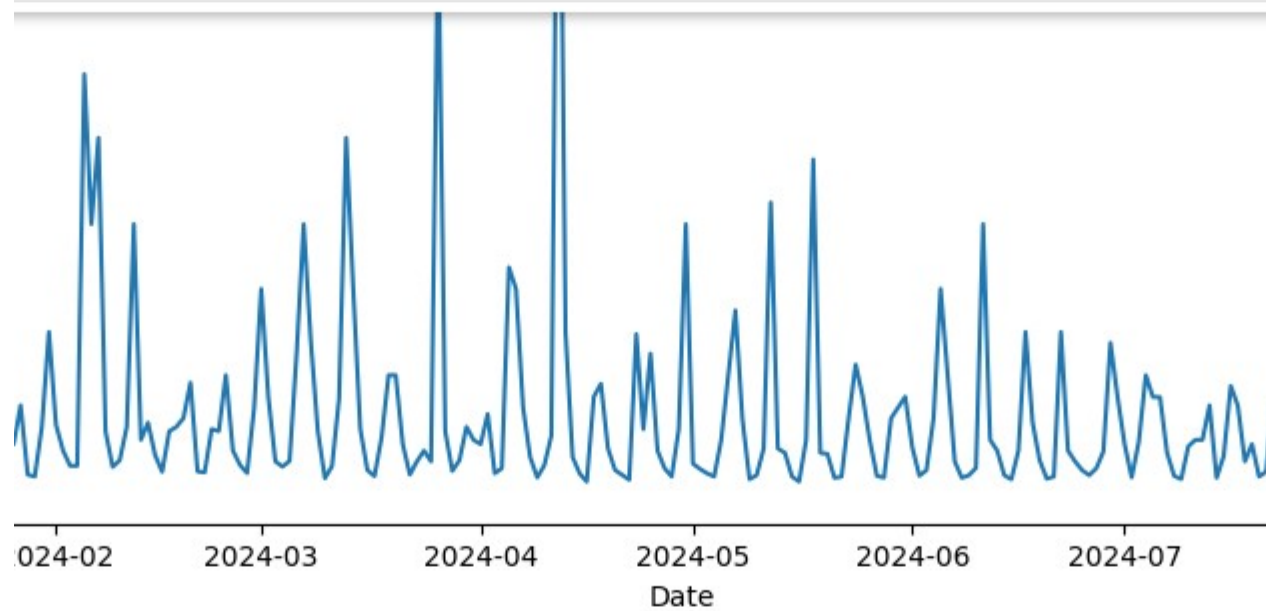
**Lets Analyze, see patte**

```
object
float64[ns]
object
object
int64
float64
float64
object
object
```

end

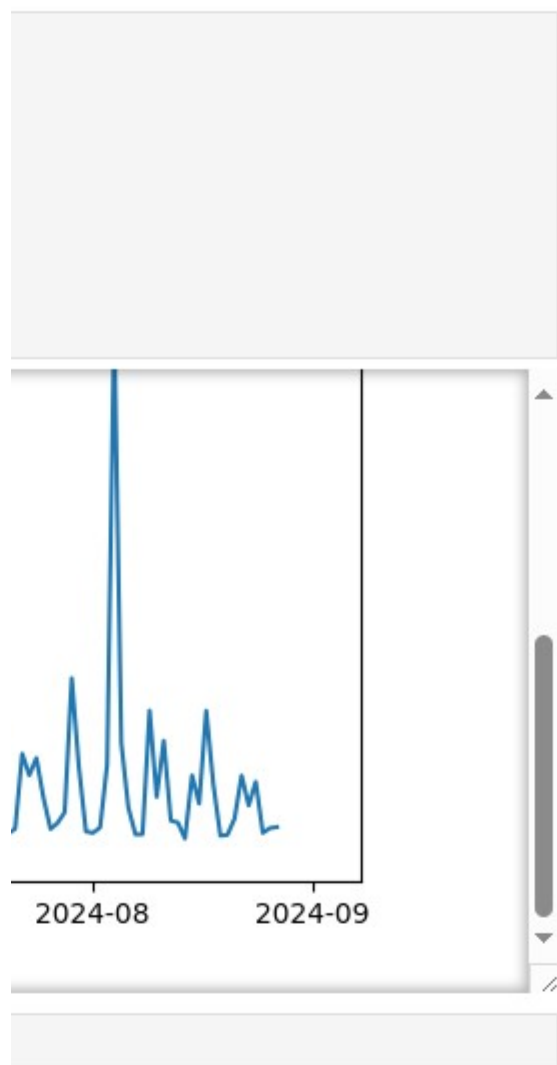
```
['Date']['Total Revenue'].sum().reset_index()
```

```
te'], sales_over_time['Total Revenue'])
```



erns and forecast Sales - Decomposition





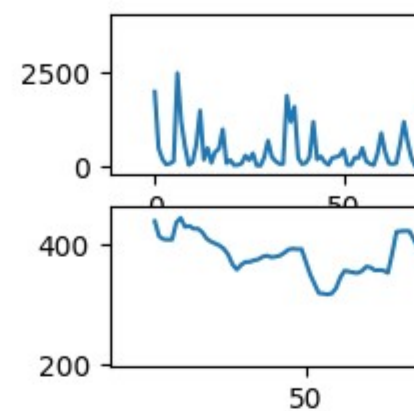
## Time series Decomposition

```
[10]: from statsmodels.tsa.seasona

result = seasonal_decompose(
plt.subplot(411)
plt.plot(result.observed, la
plt.legend(loc='best')

plt.subplot(412)
plt.plot(result.trend, label
plt.legend(loc='best')
```

[10]: <matplotlib.legend.Legend at



Observed Sales Data: This plot re  
This indicates periods of high de

Trend: This plot shows the trend  
end, the trend slightly recovers, s

### Seasonality

```
[11]: plt.subplot(413)
plt.plot(result.seasonal, la
plt.legend(loc='best')
```

[11]: <matplotlib.legend.Legend at



```

1 import seasonal_decompose

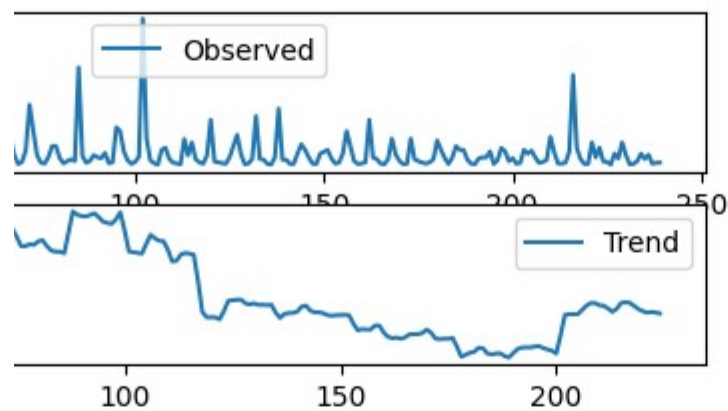
sales_over_time['Total Revenue'],model='additive', period=30)# additive model used( T ,

lbel='Observed')

.='Trend')

```

: 0x2233decabd0>

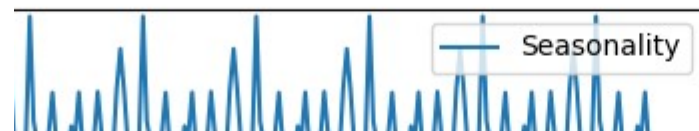


presents the raw sales data over time. We observe high sales on certain days. The data has significant demand or promotional events.

in the sales data after removing the seasonal and residual components. The trend indicates a gradual showing a small upward movement. Which can be the indication of future growth

```
lbel='Seasonality')
```

: 0x2233bbbe990>



$S, E)$

fluctuations , showing variability in sales

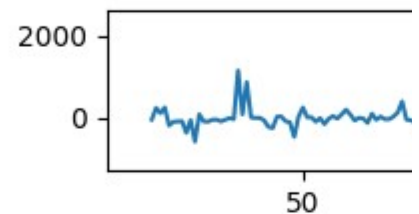
decrease in sales over time Towards the



### Noise/Error

```
[12]: plt.subplot(414)
      plt.plot(result.resid, label=
      plt.legend(loc='best')
```

```
[12]: <matplotlib.legend.Legend at
```



**Seasonality** The seasonality graph correspond to weekly or monthly by periodic factors.

**Residual/Noise** The residuals graph presence of other factors beyond factors or random events

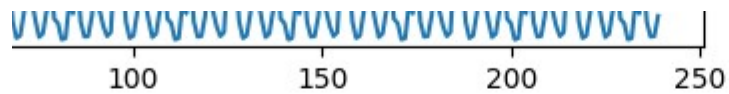
```
[ ]:
```

## Sales by Region

```
[13]: sales_by_region = df.groupby

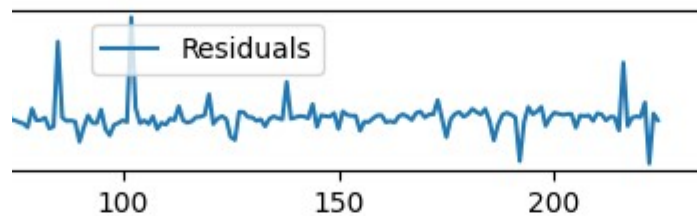
      plt.figure(figsize=(10, 6))
      sns.barplot(x='Region', y='T
      plt.title('Sales by Region')
      plt.xlabel('Region')
      plt.ylabel('Total Revenue')
      plt.show()
```





```
['Residuals'])
```

```
: 0x2233bc00260>
```



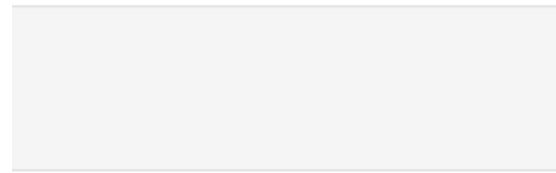
h displays a clear repeating pattern with peaks approximately every 20-25 time units, indicating a consistent seasonal pattern, potentially driven by factors like weekends, or promotional events.

graph represents the noise component of the time series after removing the trend and seasonality. The residuals exhibit some periods of higher volatility, around the middle and towards the end of the series.

```
('Region')['Total Revenue'].sum().reset_index()
```

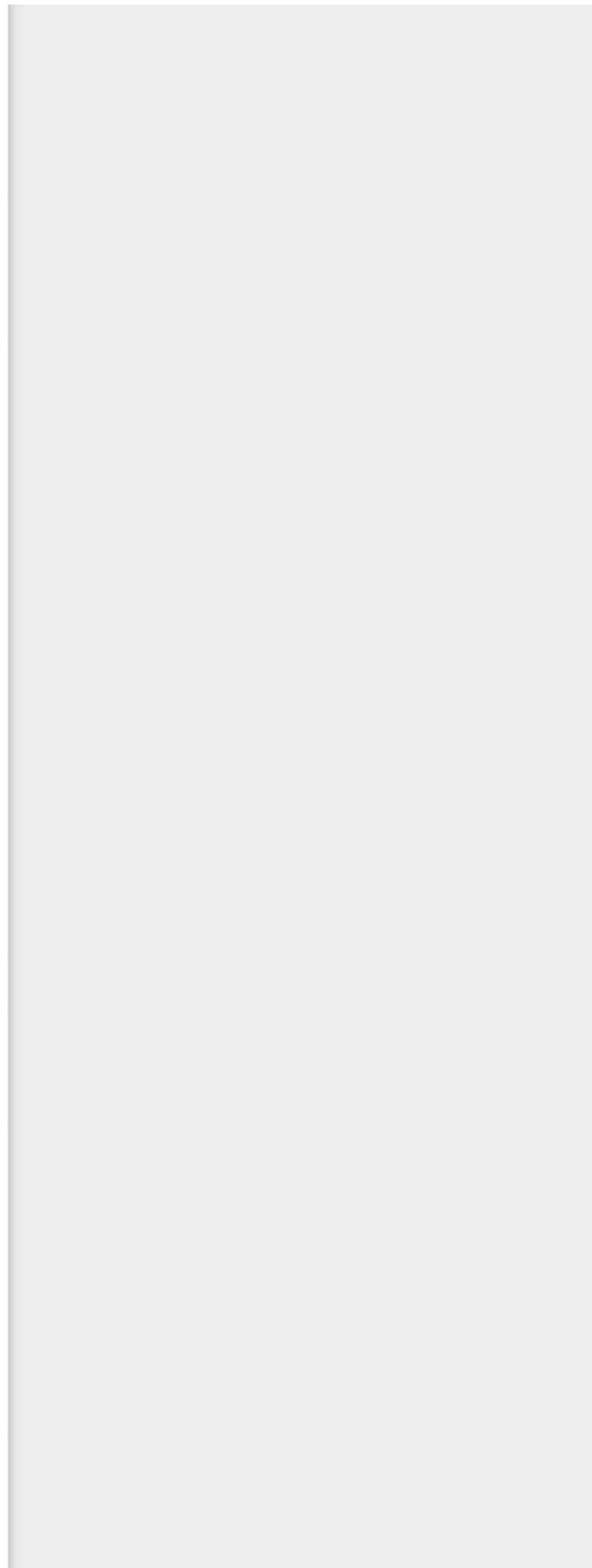
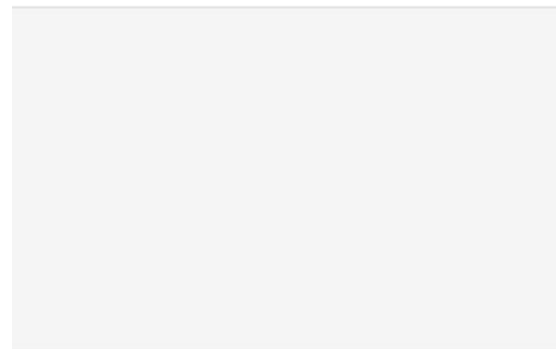
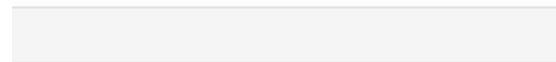
```
['Total Revenue', data=sales_by_region)
```

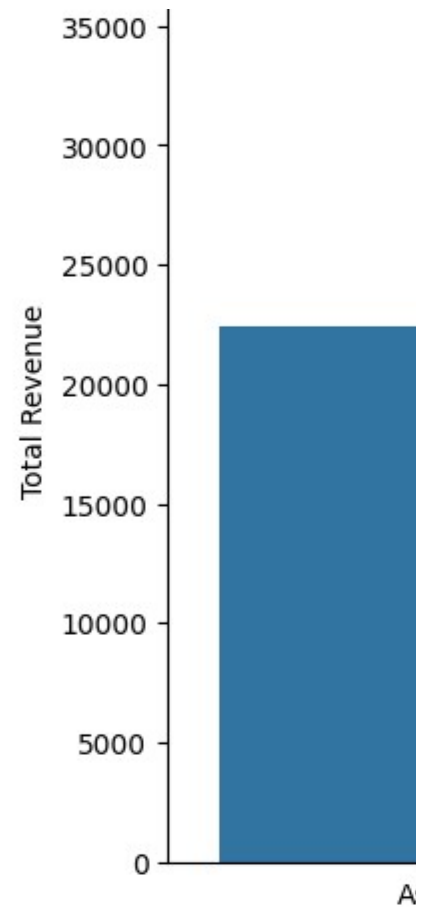
## Sales by Region



sistent periodicity in sales peaks could  
that online sales are heavily influenced

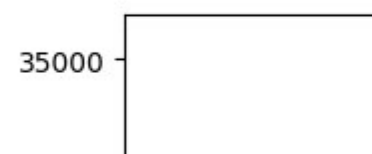
tions in the residuals indicate the  
ne end. This could be due to external



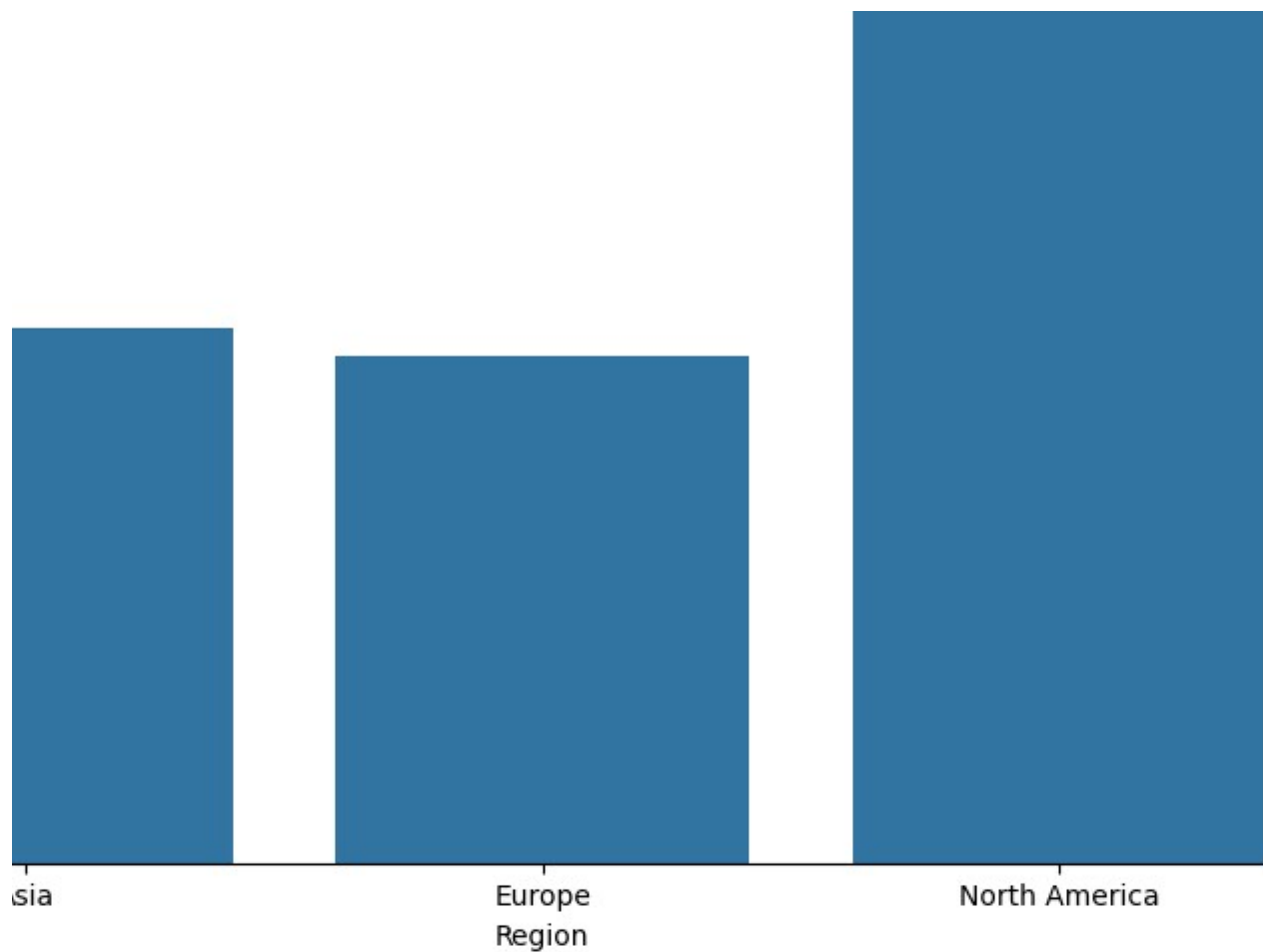


## Sales by category

```
[14]: sales_by_category = df.groupby('Product Category').sum()\n\nplt.figure(figsize=(12, 6))\nsns.barplot(x='Product Category', y='Total Revenue')\nplt.title('Sales by Product Category')\nplt.xlabel('Product Category')\nplt.ylabel('Total Revenue')\nplt.xticks(rotation=45)\nplt.show()
```





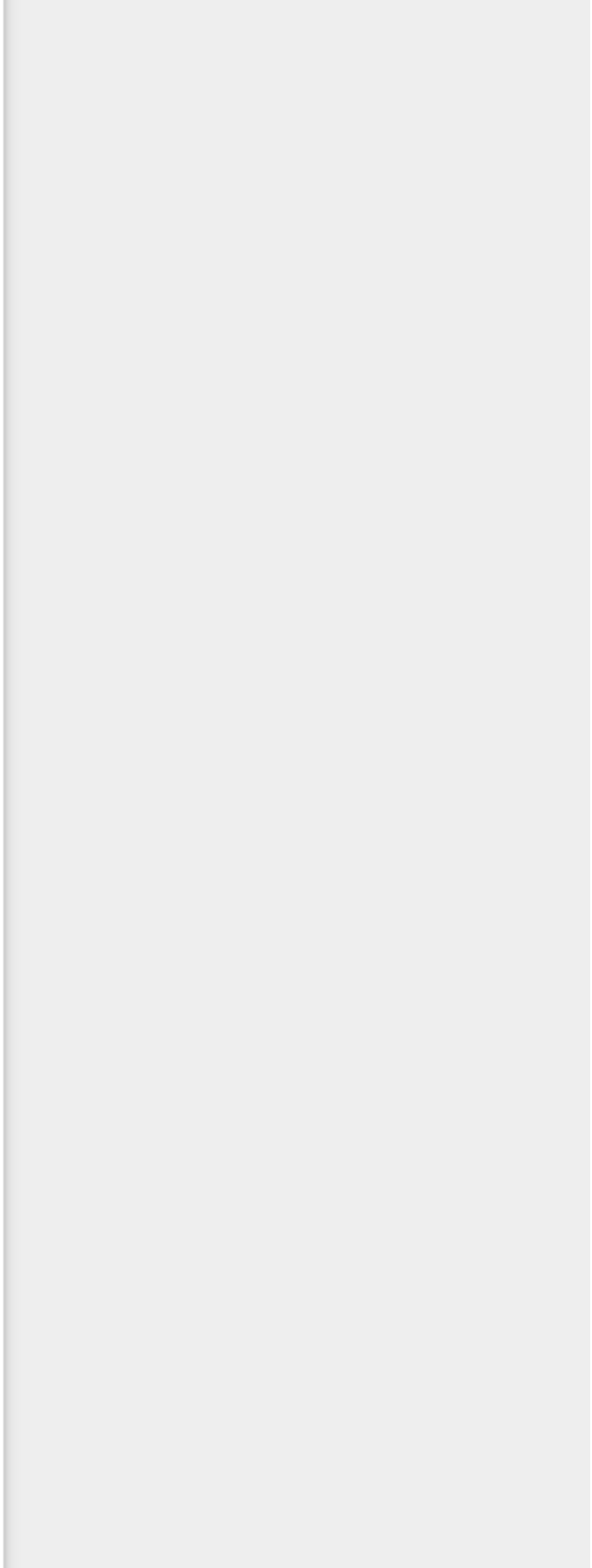
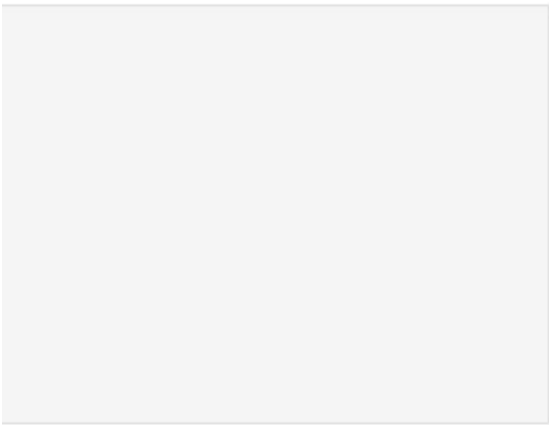


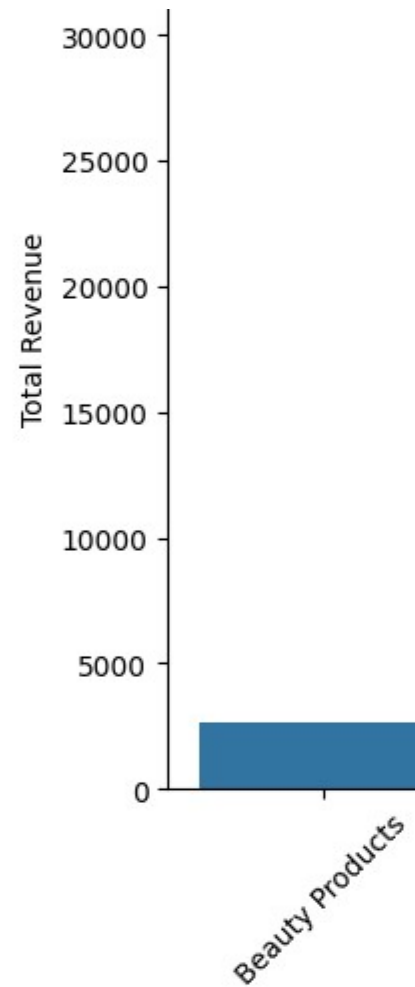
```
by('Product Category')['Total Revenue'].sum().reset_index()

ory', y='Total Revenue', data=sales_by_category)
Category')
')
```

Sales by Product Category







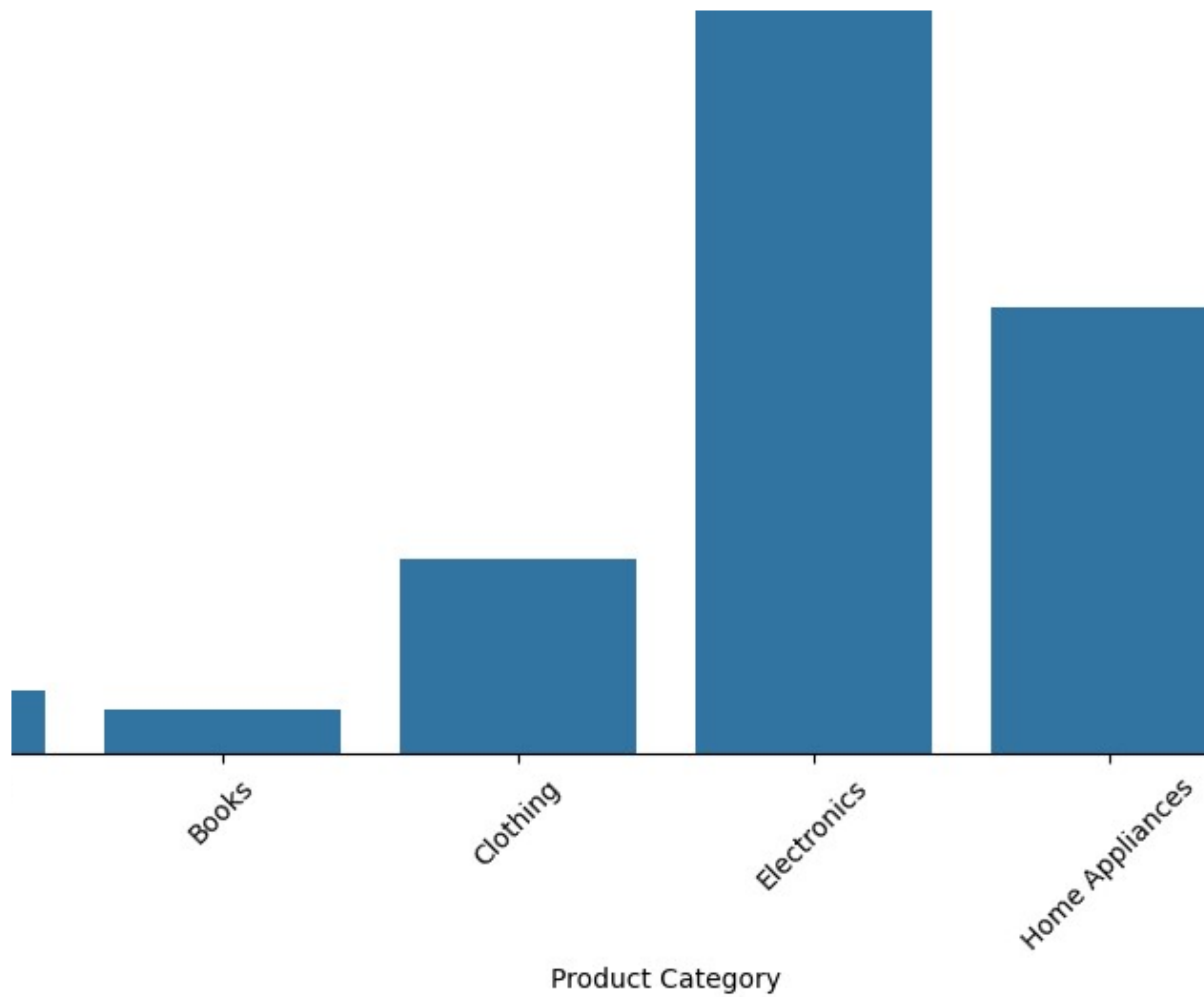
## payment methods

```
[15]: #df.head()
payment_methods = df['Payment Methods']

plt.figure(figsize=(8, 6))
payment_methods.plot.pie(autopct='%1.1f%%')
plt.title('Payment Methods Distribution')
plt.ylabel('')
plt.show()
```

Payment Methods

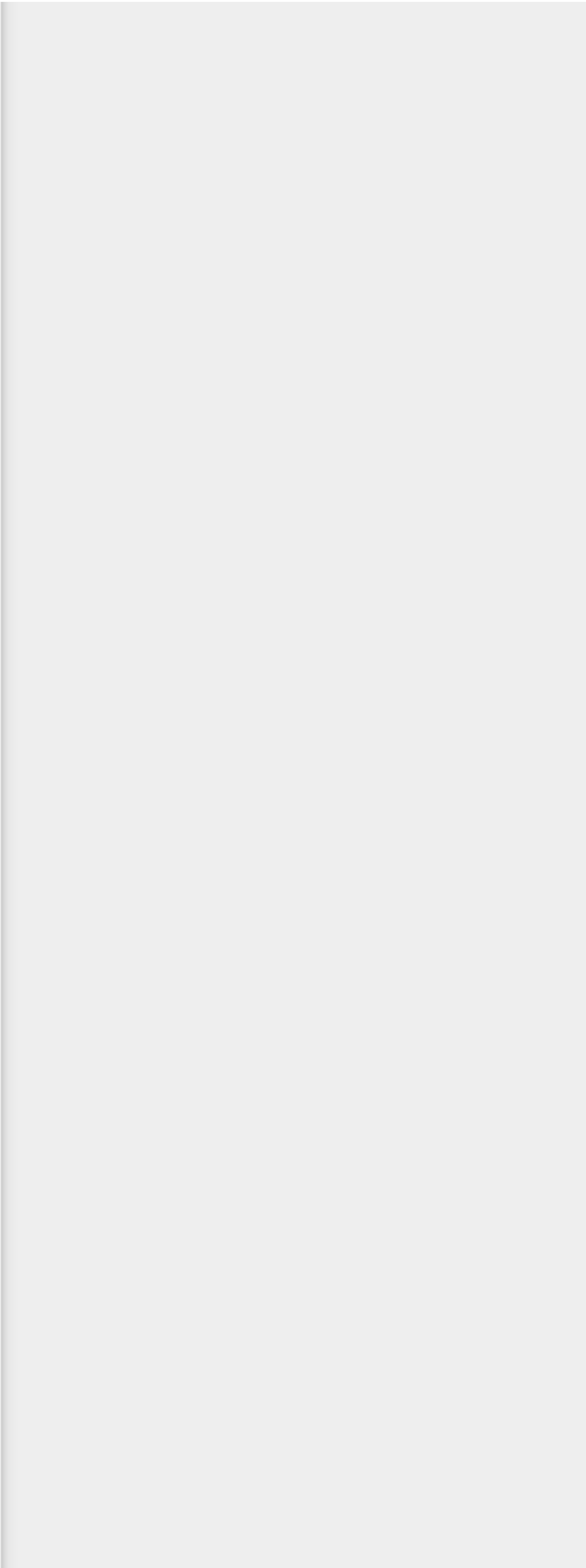
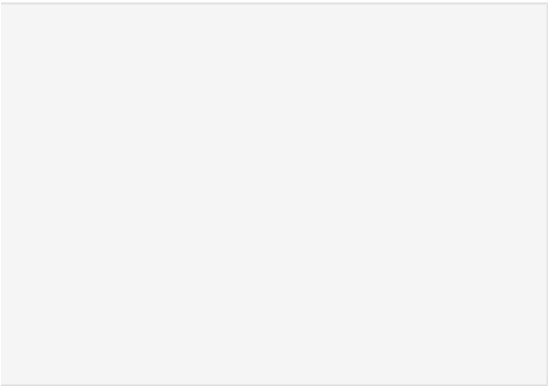
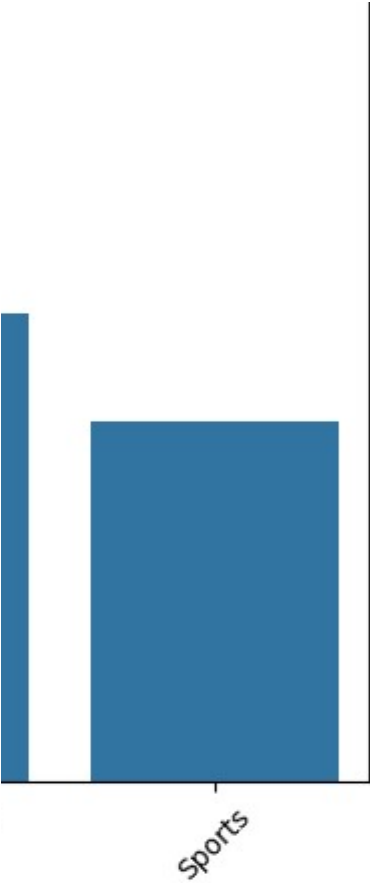
Debit Card

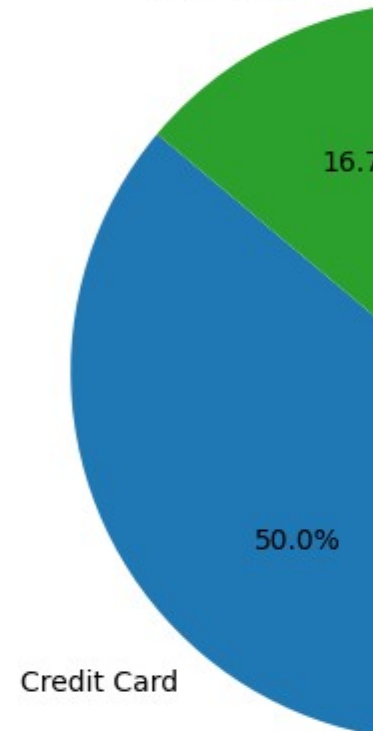


```
['it Method'].value_counts()
```

```
plotpct='%1.1f%%', startangle=140)# the second % is used for showing % sign on the chart  
distribution')
```

Methods Distribution



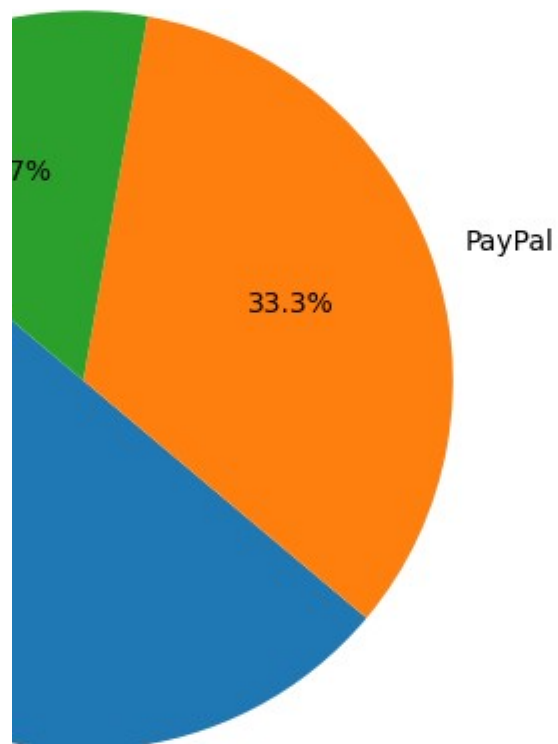


## Heatmap of sales by region

```
[16]: pivot_table = df.pivot_table(index='Region', columns='Product Category', values='Sales')

plt.figure(figsize=(12, 8))
sns.heatmap(pivot_table, annot=True)
plt.title('Sales by Region and Product Category')
plt.xlabel('Product Category')
plt.ylabel('Region')
plt.show()
```

Asia

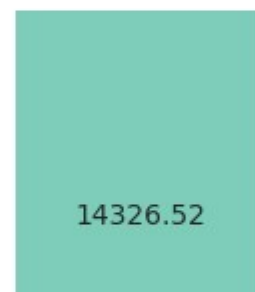
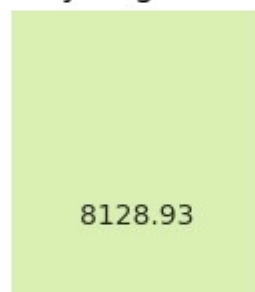


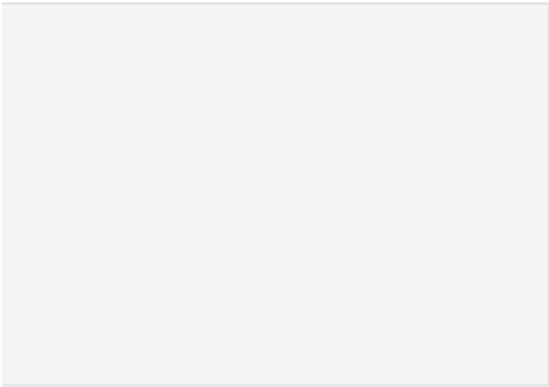
## Region/product category

```
(values='Total Revenue', index='Region', columns='Product Category', aggfunc='sum')

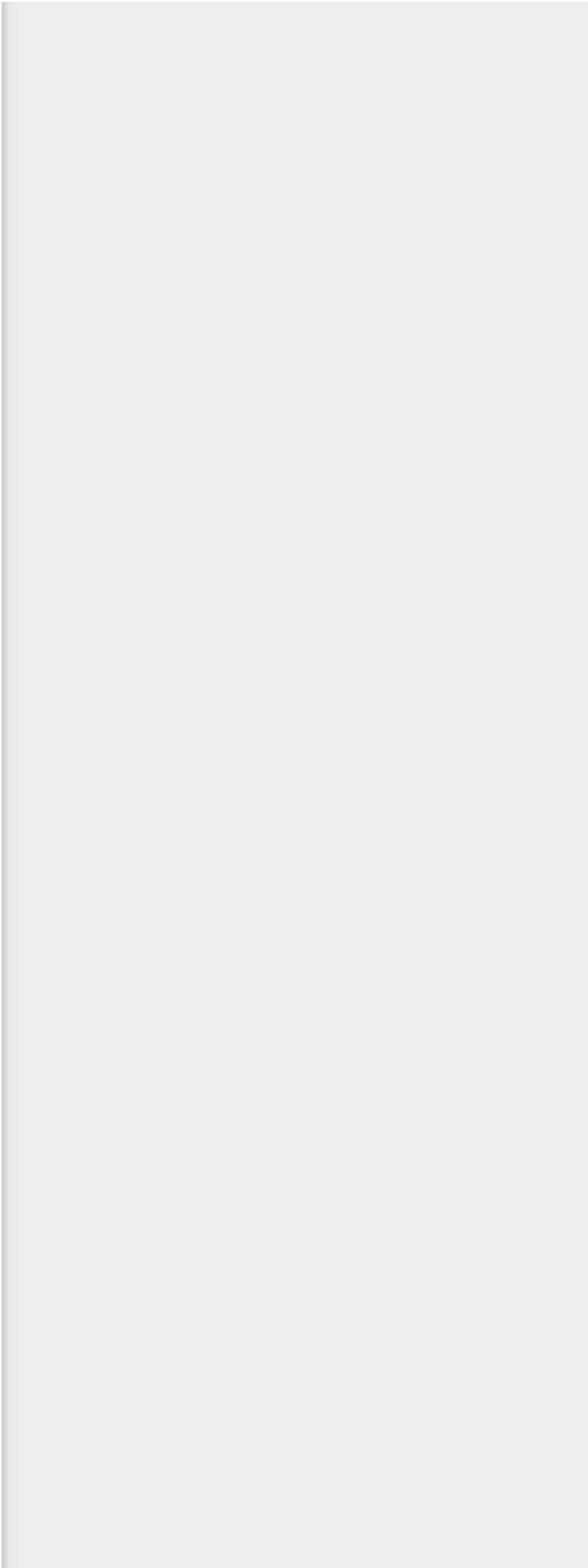
plot=True, fmt=".2f", cmap="YlGnBu")
nd Product Category')
')
```

## Sales by Region and Product Category

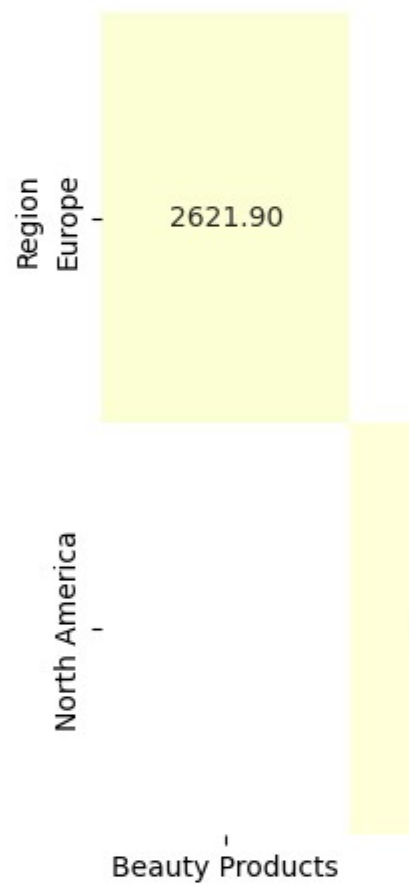




30000



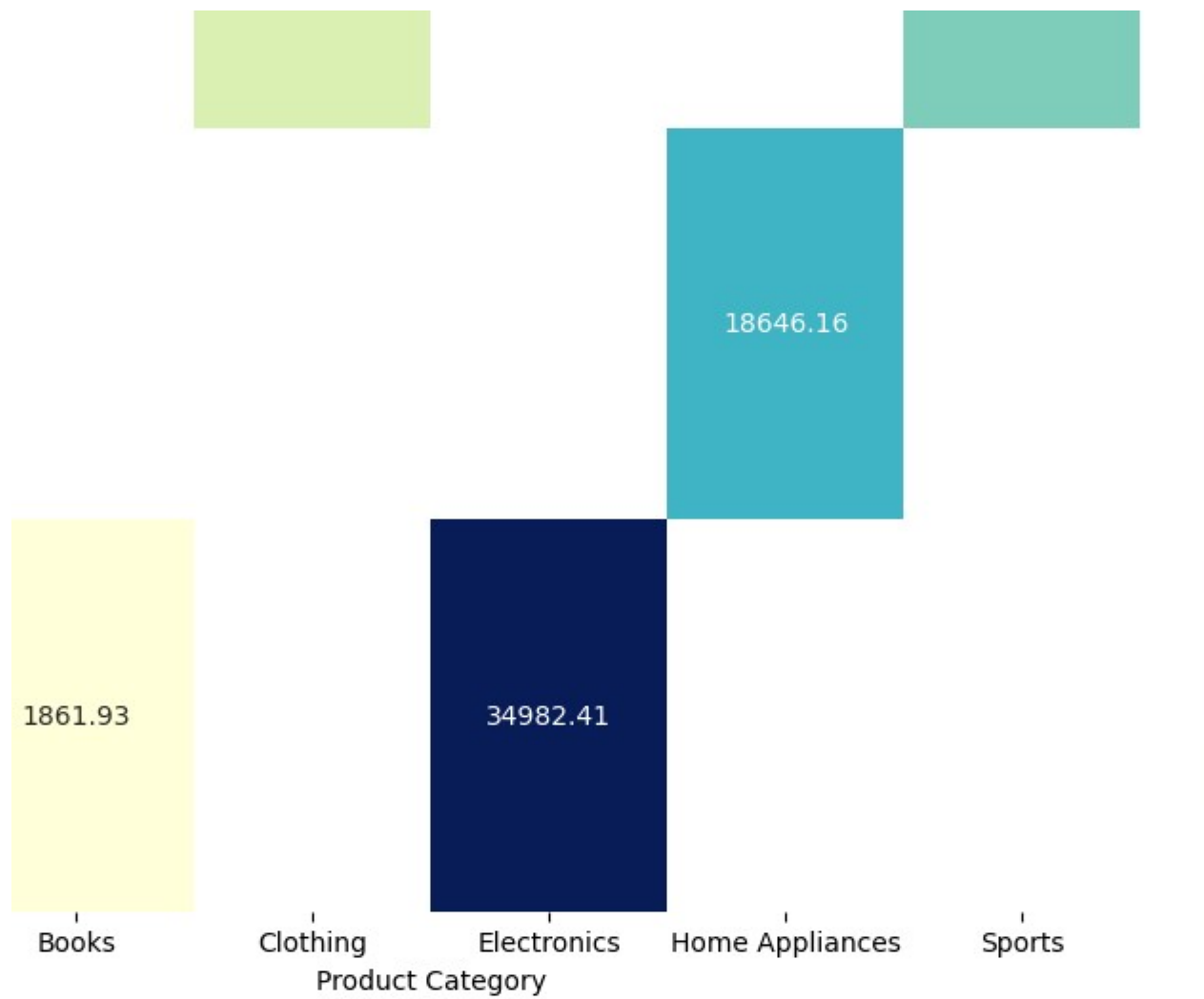




## Units Sold Distribution

```
[17]: # Distribution plot of units
plt.figure(figsize=(10, 6))
sns.histplot(df['Units Sold'])
plt.title('Distribution of U
plt.xlabel('Units Sold')
plt.ylabel('Frequency')
plt.show()
```





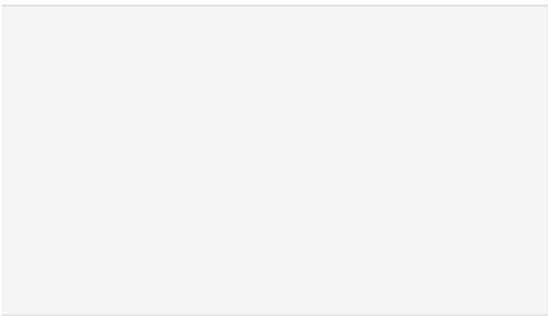
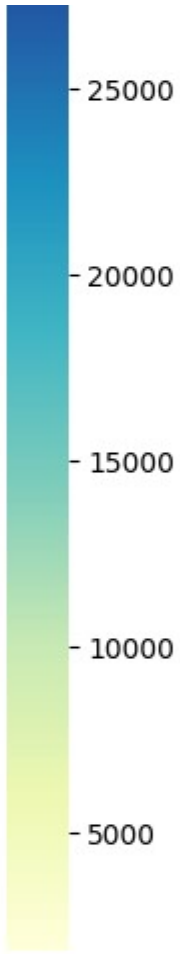
7

```

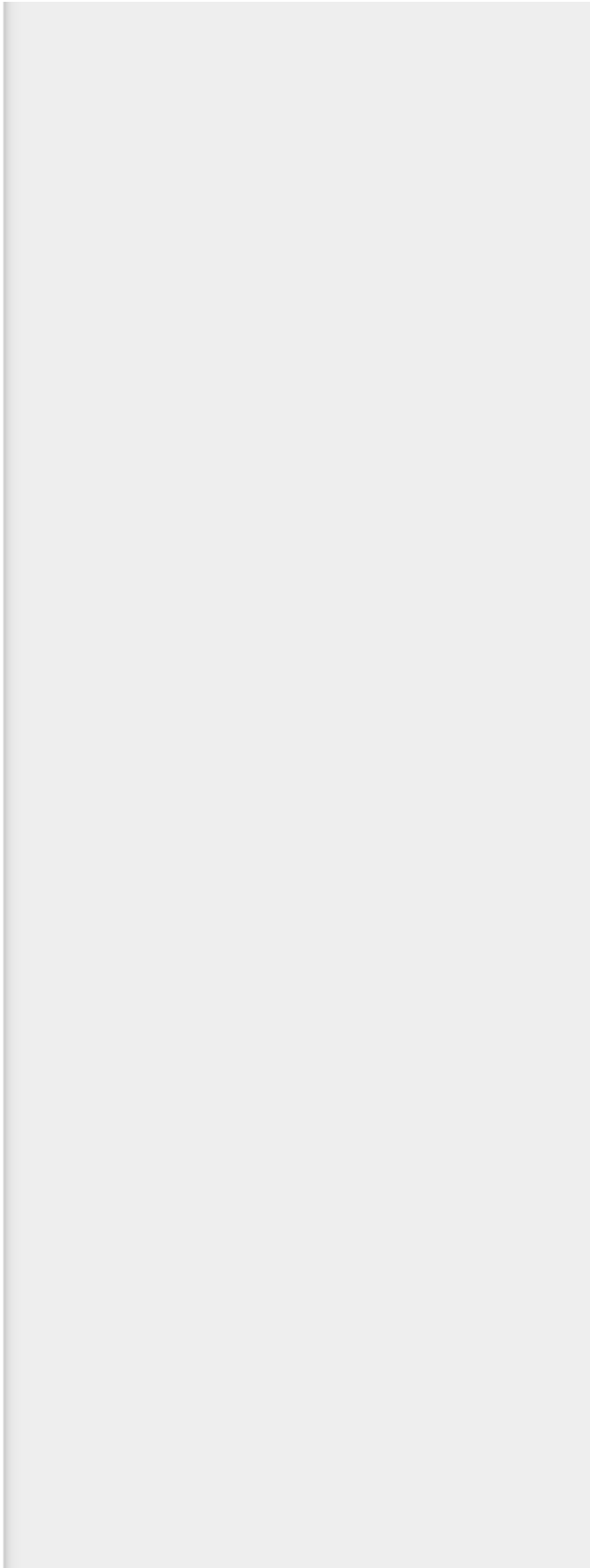
units Sold
], kde=True, bins=20)
units Sold')

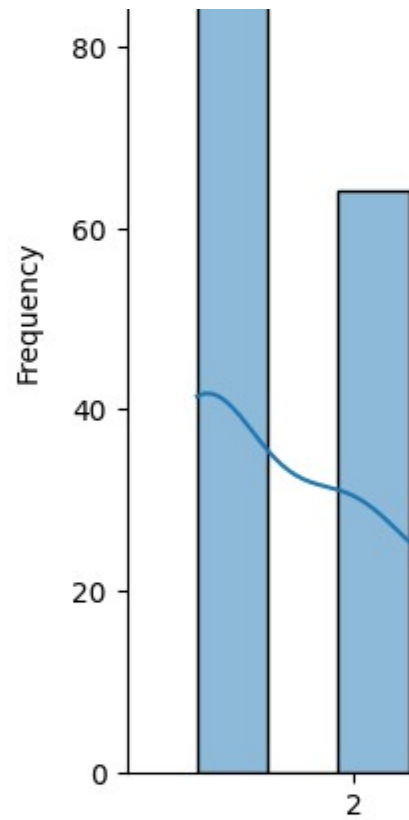
```

Distribution of Units Sold



7

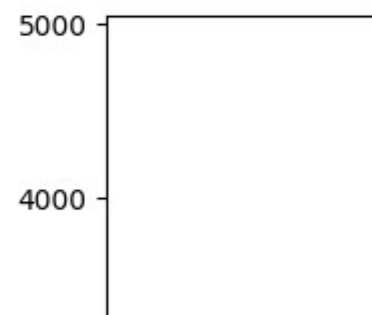


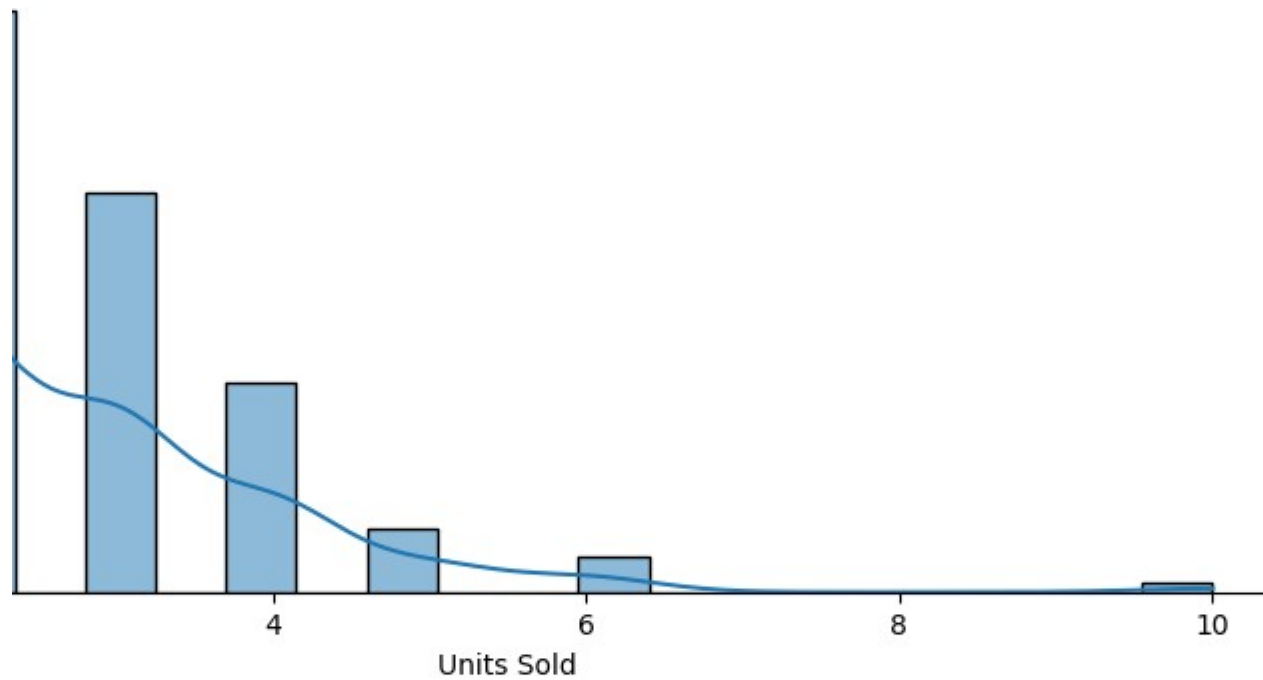


## Regression Analysis

### Unit price vs Total Revenue

```
[18]: plt.figure(figsize=(10, 6))
sns.regplot(x='Unit Price',
plt.title('Unit Price vs Tot
plt.xlabel('Unit Price')
plt.ylabel('Total Revenue')
plt.show()
```





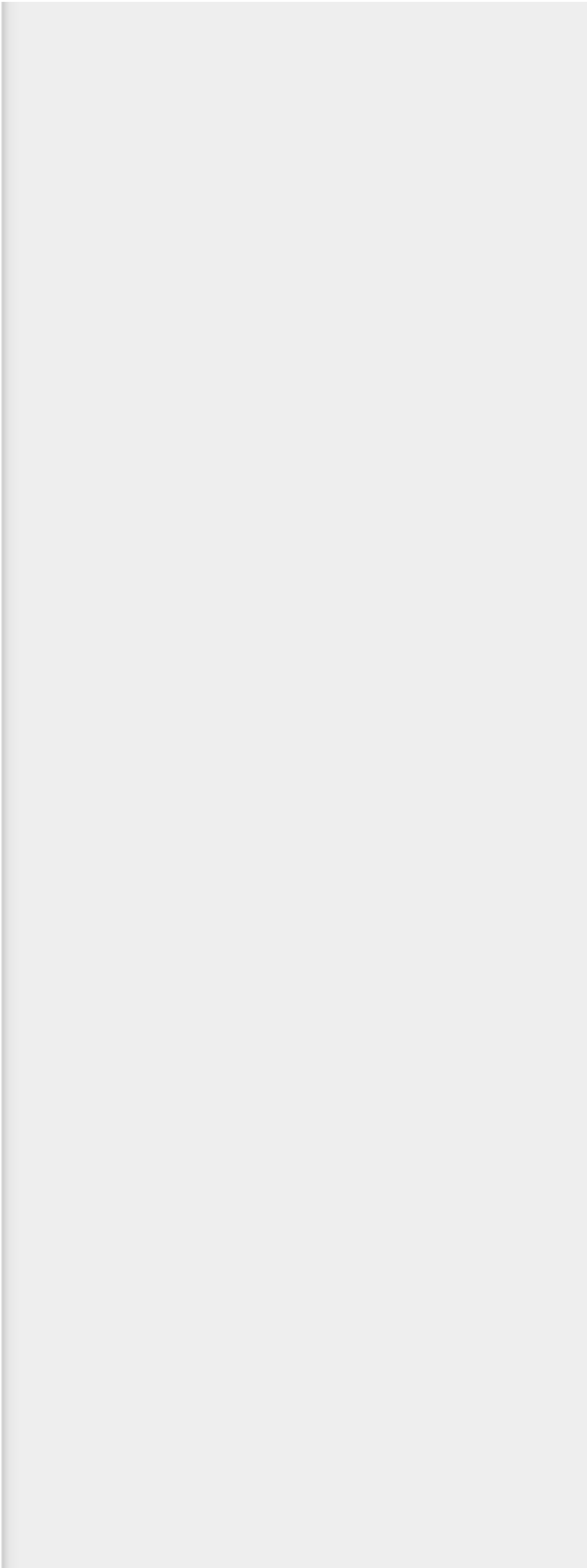
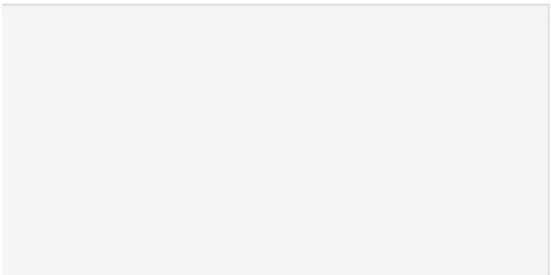
```
y='Total Revenue', data=df)# scatter with regression line  
al Revenue')
```

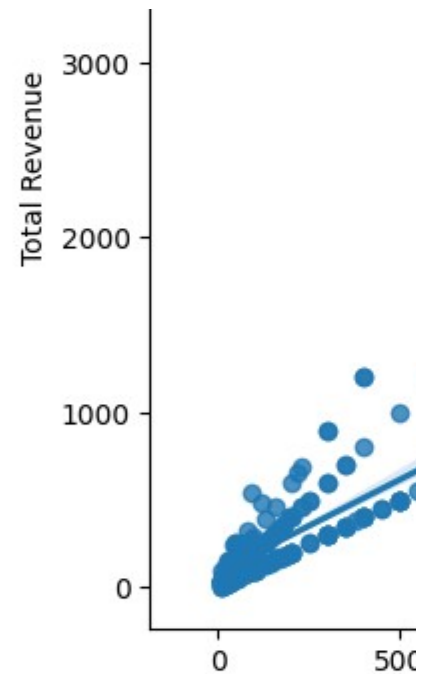
---

Unit Price vs Total Revenue

---







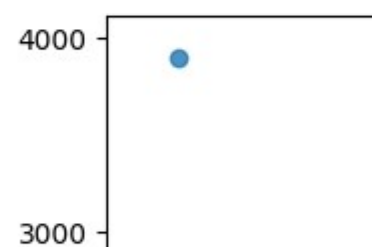
The positive slope of the Reg line

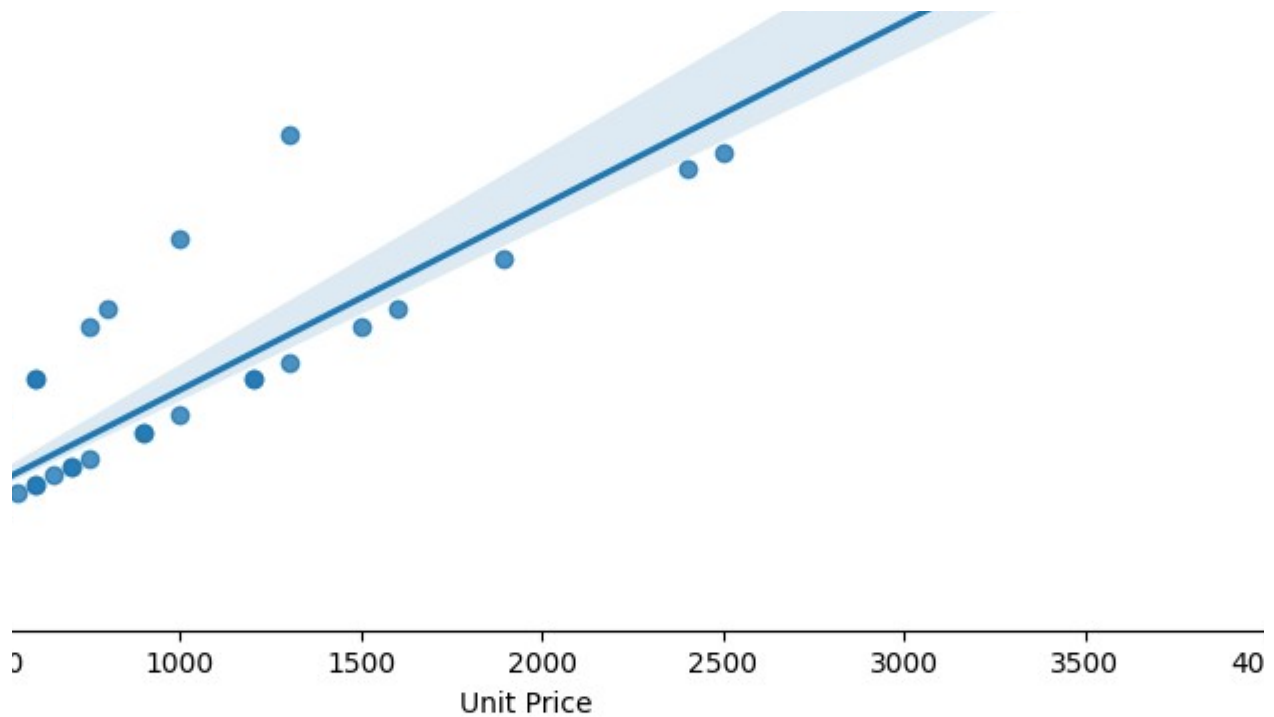
The scatter plot shows a clear up

There are some dots that are dist  
low for a given unit price. Outlier

### Units Sold vs Total Revenue

```
[19]: plt.figure(figsize=(10, 6))
sns.regplot(x='Units Sold',
plt.title('Units Sold vs Tot
plt.xlabel('Units Sold')
plt.ylabel('Total Revenue')
plt.show()
```





It indicates a +ive relationship between unit price and total revenue.

Upward trend, suggesting a strong positive correlation between unit price and total revenue.

Points far from the regression line, indicating outliers. These outliers could represent exceptional cases where products can provide valuable insights, such as identifying products with exceptional performance.

```
df['Total Revenue'] = df['Unit Price'] * df['Units Sold']
df['Total Revenue']
```

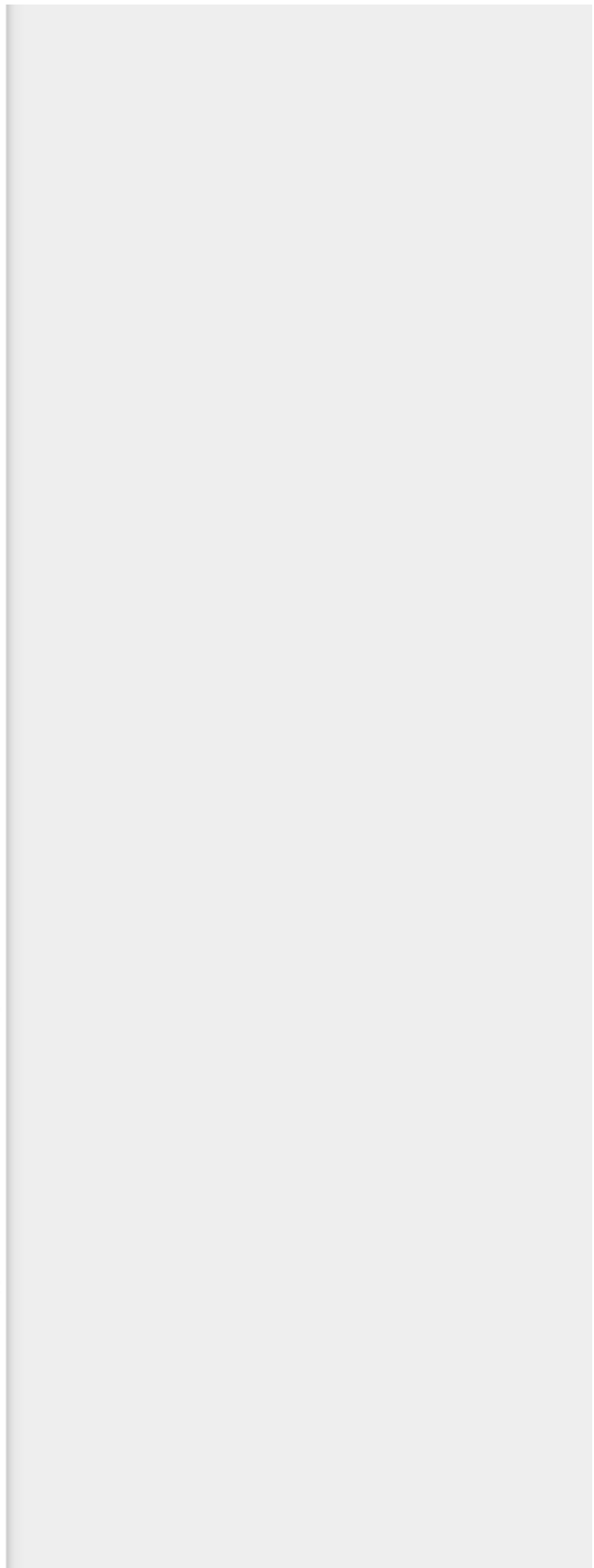
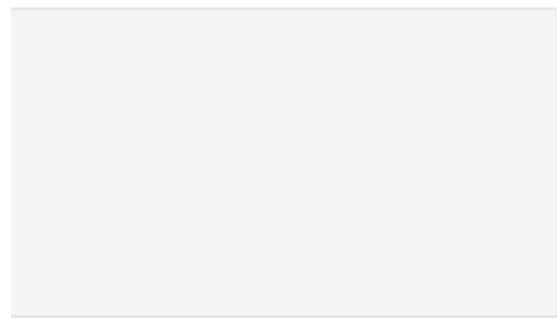
---

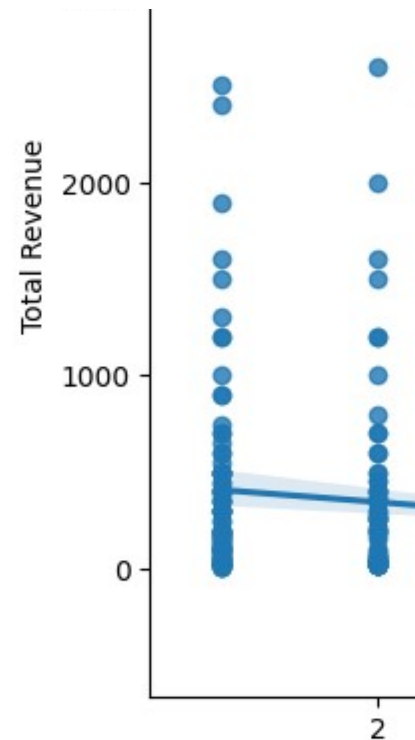
Units Sold vs Total Revenue



100

are the total revenue is unusually high or





There is a slight negative relation  
total revenue There is significant  
result in a wide range of total rev

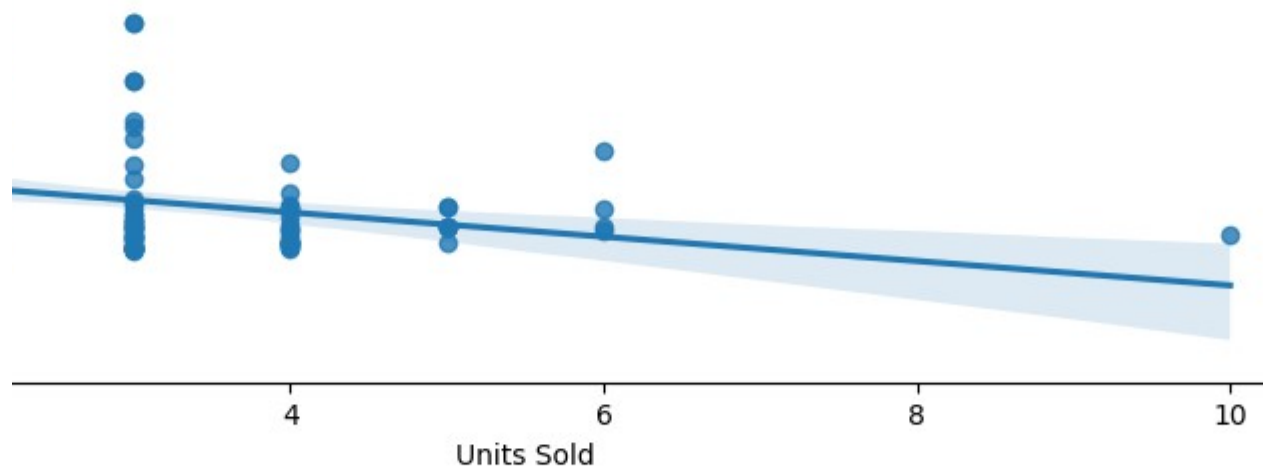
Outliers: A few outliers exist, espe

The trend suggests that as more

**transaction with the h**

```
[20]: transaction = df.loc[df['Tot
transaction
```

```
[20]: Transaction ID
Date                2024-04-
Product Category
Product Name        Canon EO
Units Sold
Unit Price
Total Revenue
Region              No
Payment Method
Name: 102, dtype: object
```



relationship between the number of units sold and total revenue. This suggests that selling more units does not necessarily result in higher total revenue. There is significant variability in total revenue when fewer units are sold (particularly between 1 and 3 units). This indicates that total revenue is influenced by unit prices.

especially at higher total revenue levels for low units sold. These could represent high-value items sold in small quantities. As the number of units sold increases, the total revenue tends to decrease slightly, which could be indicative of discounts or bulk pricing.

## Identifying the highest total revenue

```
data['Total Revenue'].idxmax()
```

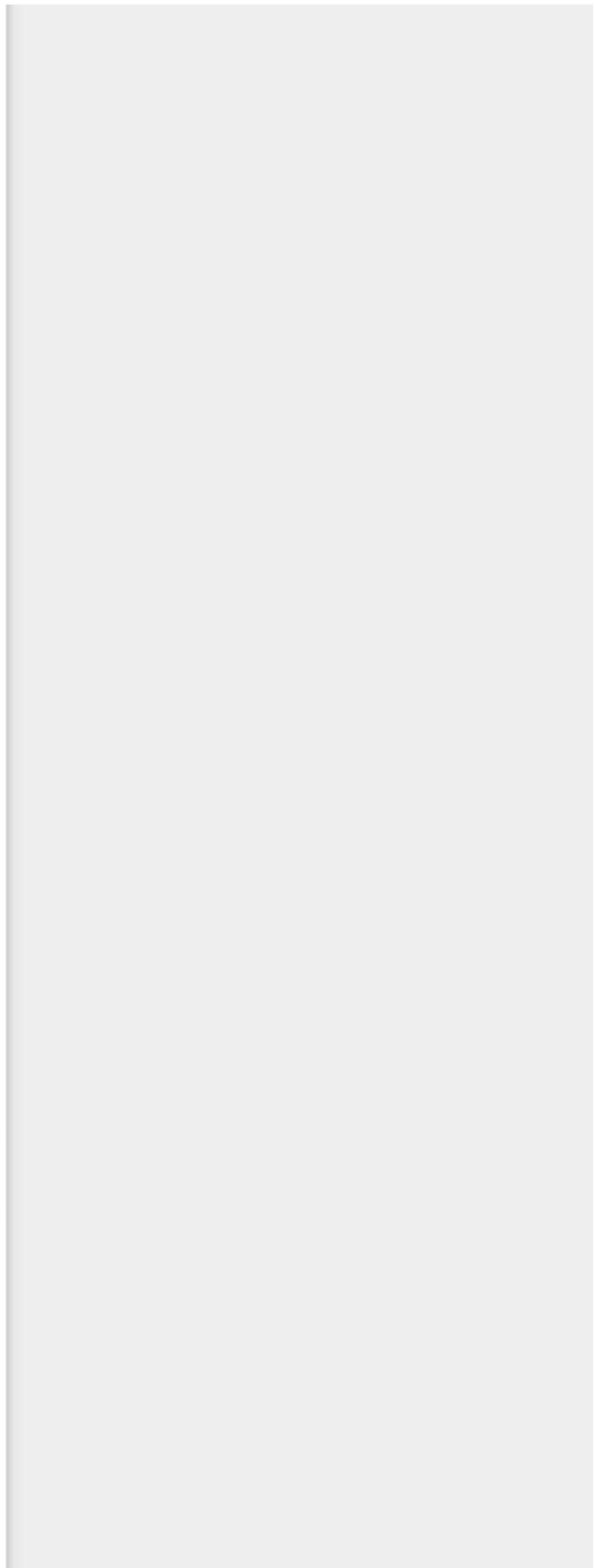
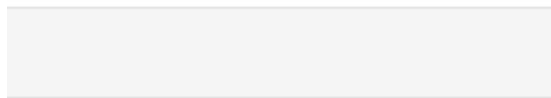
```
10103
12 00:00:00
Electronics
DS R5 Camera
1
3899.99
3899.99
North America
Credit Card
```



not necessarily correspond to higher  
es that a small number of units sold can

small quantities.

lower-priced items being sold in bulk.



## product with the high

```
[21]: top_product_units = df.loc[d
top_product_units
```

```
[21]: Transaction ID
Date 20
Product Category
Product Name Hanes Co
Units Sold
Unit Price
Total Revenue
Region
Payment Method
Name: 62, dtype: object
```

## product with most rev

```
[22]: prod= df.groupby('Product Na
prod
```

```
[22]: Product Name
Canon EOS R5 Camera
LG OLED TV
MacBook Pro 16-inch
Apple MacBook Pro 16-inch
iPhone 14 Pro
```

```
Neutrogena Hydro Boost Water
Biore UV Aqua Rich Watery Es
The Ordinary Hyaluronic Acid
The Ordinary Caffeine Soluti
The Ordinary Niacinamide Ser
Name: Total Revenue, Length:
```

```
[23]: prod= df.groupby('Product Na
top5=prod.head(5)
plt.figure(figsize=(12, 6))
sns.barplot(x=top5.index, y=
plt.title('Total Revenue by
plt.xlabel('Product Name')
plt.ylabel('Total Revenue')
plt.xticks(rotation=90)
plt.show()
```

## est units sold

```
lf['Units Sold'].idxmax()]
```

```
10063
24-03-03 00:00:00
Clothing
importSoft T-Shirt
10
9.99
99.9
Asia
Debit Card
```

## revenue produced

```
me')['Total Revenue'].sum().sort_values(ascending = False)
```

```
3899.99
2599.98
2499.99
2399.00
1999.98
...
Gel 16.99
sence Sunscreen 15.00
l Serum 6.80
on 5% + EGCG 6.70
um 6.50
232, dtype: float64
```

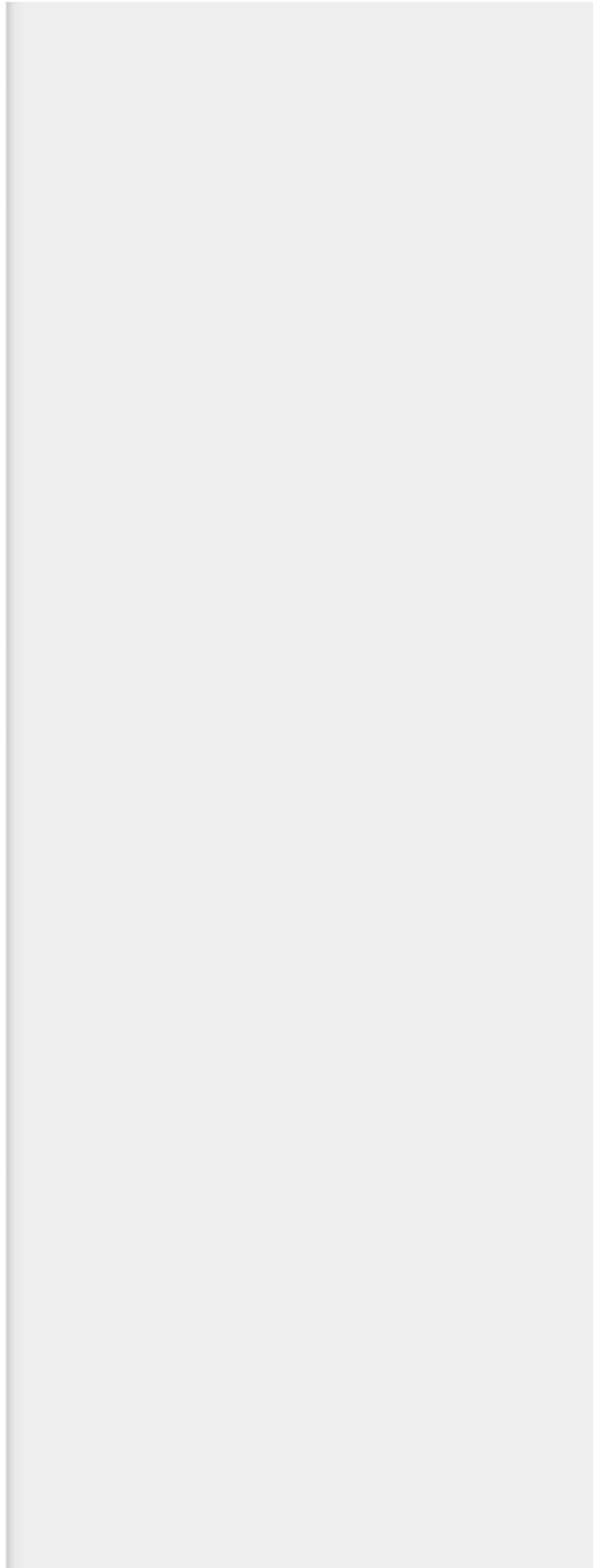
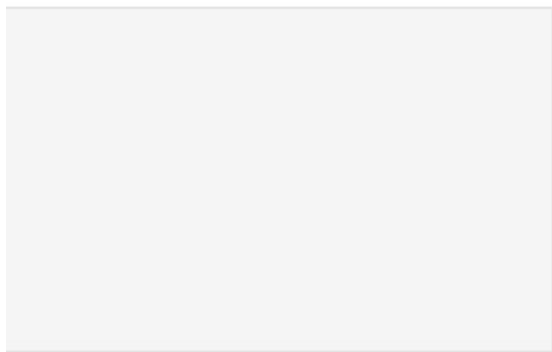
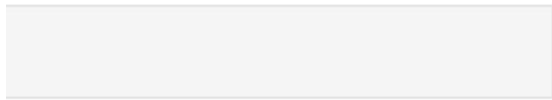
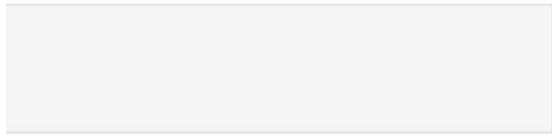
```
me')['Total Revenue'].sum().sort_values(ascending = False)
```

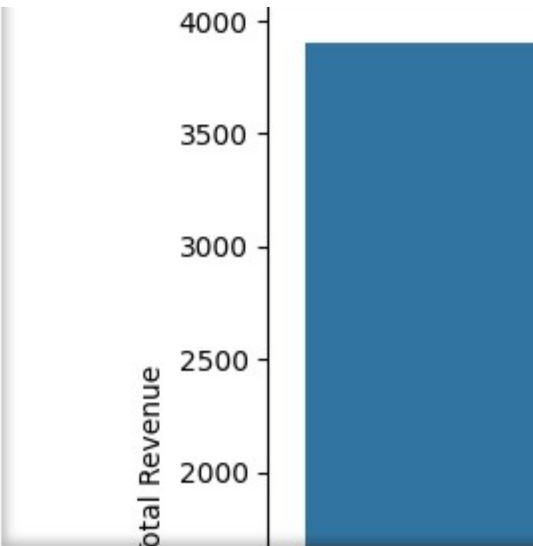
```
top5.values)
Product')
```

---

Total Revenue by Product

---





## Conclusion and Analysis

**Top Performing Regions:** The a

**Top Products:**

The top 5 products by revenue a management and marketing stra

**Category Analysis:**

Electronic, Sports and home appl

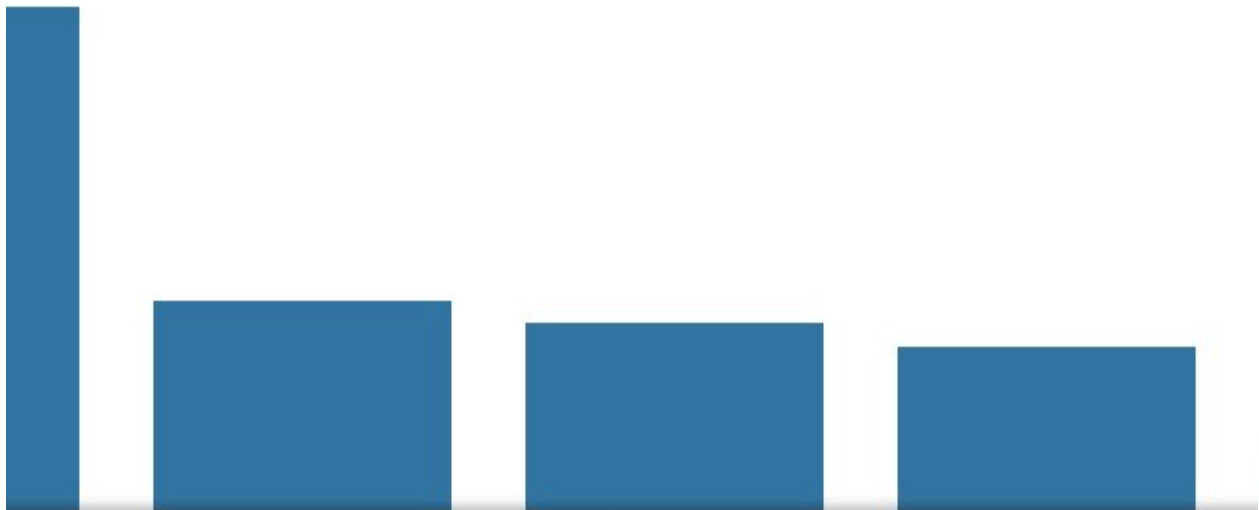
**Seasonal Trends:** Our Sales data

**Payment Methods and Discou** and loyalty. If discounts significai

```
[25]: df.head()
```

[25]:	Transaction ID	Date
0	10001	2024-01-01
1	10002	2024-01-02
2	10003	2024-01-03
3	10004	2024-01-04
4	10005	2024-01-05





sis

analysis identifies regions contributing the most to total revenue. Focus on (North America, Asia) for fu

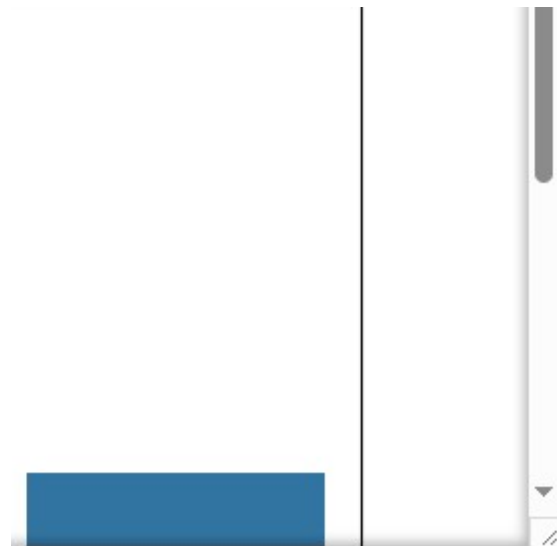
ire Cannon Camera, LG TV, MacBook, and iphone 14 pro. These products are popular among customer strategies.

liances Category drive signifigcant revenue. These categories should be the focal point for expanding

i over time reveal seasonal trends and peak periods. Preparing for these peaks with adequate inventor

ints The most used payment methods are Debitcard and paypal. Offering discounts on these methods: ntly boost sales, they can be strategically used during off-peak times.

Product Category	Product Name	Units Sold	Unit Price	Total Revenue	Region
Electronics	iPhone 14 Pro	2	999.99	1999.98	North America
Home Appliances	Dyson V11 Vacuum	1	499.99	499.99	Europe
Clothing	Levi's 501 Jeans	3	69.99	209.97	Asia
Books	The Da Vinci Code	4	15.99	63.96	North America
Beauty Products	Neutrogena Skincare Set	1	89.99	89.99	Europe



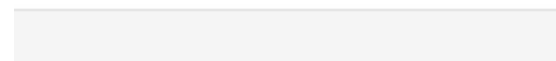
future marketing and sales efforts.

rs and should be prioritized in inventory

product lines

ry and marketing can optimize sales.

s could enhance customer satisfaction



#### Payment Method

Credit Card

PayPal

Debit Card

Credit Card

PayPal

## Sales Forecasting with

```
[26]: from statsmodels.tsa.statesp
import matplotlib.pyplot as

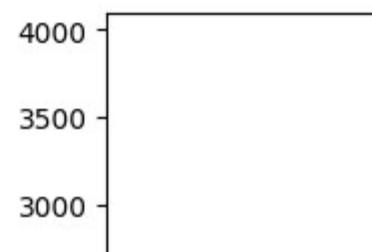
sales_data = df.groupby('Dat
sales_data.set_index('Date',

model = SARIMAX(sales_data,
model_fit = model.fit(dis=F

forecast = model_fit.get_for
forecast_index = pd.date_ran
forecast_series = forecast.p
forecast_series.index = fore

plt.figure(figsize=(10, 5))#
plt.plot(sales_data, label='
plt.plot(forecast_series, la
plt.xlabel('Date')
plt.ylabel('Total Revenue')
plt.title('Sales Forecast')
plt.legend()
plt.show()
```

```
C:\Users\pc\AppData\Local\Pr
rovided, so inferred frequen
self._init_dates(dates, fr
C:\Users\pc\AppData\Local\Pr
rovided, so inferred frequen
self._init_dates(dates, fr
C:\Users\pc\AppData\Local\Te
E' instead.
forecast_index = pd.date_r
```



## ARIMA

```
face.sarimax import SARIMAX
plt

'e')['Total Revenue'].sum().reset_index()# aggregate sales by date
inplace=True)

order=(1, 1, 1), seasonal_order=(1, 1, 1, 12))# MODEL - forecasts sales of next 12Months
alse)

ecast(steps=12)#Forecast sales
nge(start=sales_data.index[-1], periods=12, freq='M')
redicted_mean
cast_index

visualization
Observed')
bel='Forecast', color='red')

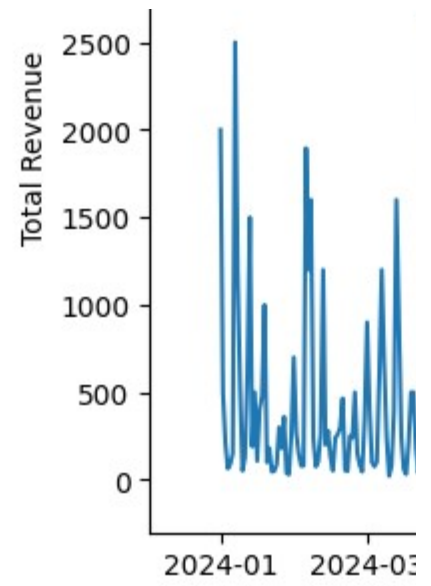
ograms\Python\Python312\Lib\site-packages\statsmodels\tsa\base\tsa_model.py:473: ValueW
icy D will be used.
eq)
ograms\Python\Python312\Lib\site-packages\statsmodels\tsa\base\tsa_model.py:473: ValueW
icy D will be used.
eq)
mp\ipykernel_17996\3423825647.py:14: FutureWarning: 'M' is deprecated and will be remov
ange(start=sales_data.index[-1], periods=12, freq='M')
```

### Sales Forecast



arning: No frequency information was p  
arning: No frequency information was p  
ed in a future version, please use 'M

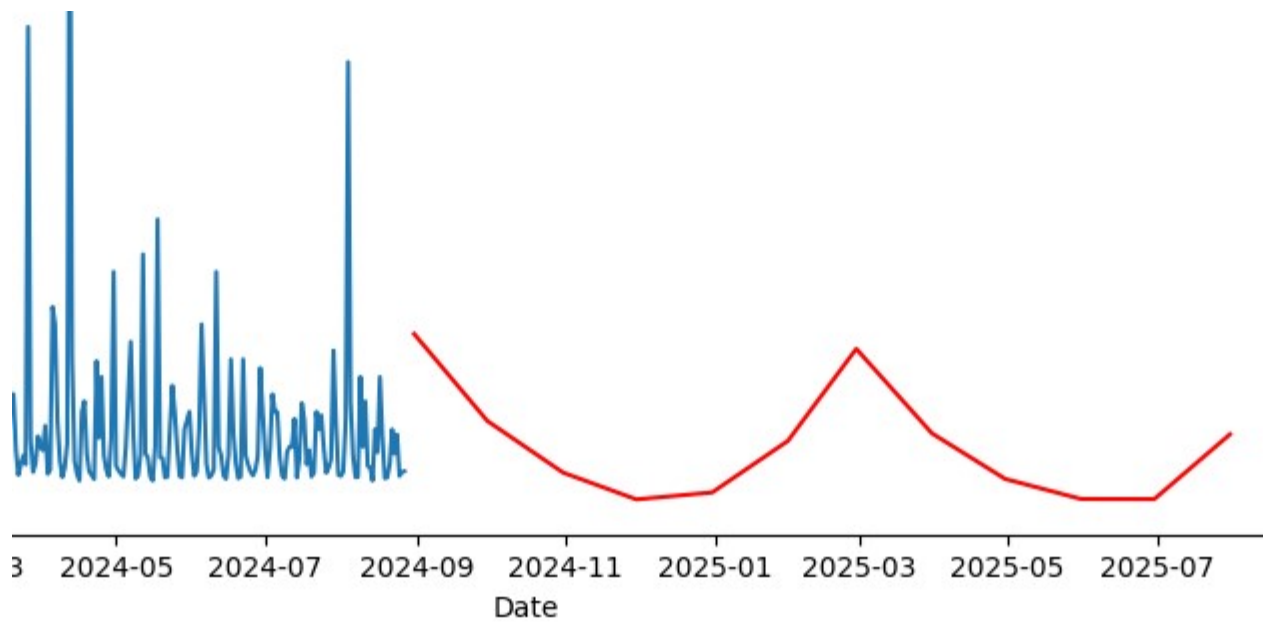
d



**Forecast data - 2024/08 - 2025**

Sales are expected to be low dur

[ ]:



**1/07** The model predicts 2 peaks in sept and another in Feb 2025. These peaks may be due to seasonal  
ing october to Jan and then again from march to june.

