

Project Name - Supermart Grocery Sales - Retail Analytics Dataset_ (Data Analyst) (Part 1)

Project Type - Data Analysis

Industry - Unified Mentor

Contribution - Individual

Member Name - Hare Krishana Mishra

Task - 1

Project Summary -

Project Description:

A fictional dataset simulating grocery orders from customers in Tamil Nadu, India, designed for practicing data analysis and visualization. It contains order details such as customer information, order date, category, sales, discount, profit, and location data.

Objective:

The main objective of this project is to analyze, interpret, and visualize grocery sales data to uncover trends, patterns, and relationships that can help improve decision-making. Additionally, the project aims to build a predictive model to estimate sales based on key features, providing actionable insights for business growth.

Key Project Details:

Tools Used: Python, Pandas, NumPy, Matplotlib, Seaborn, Scikit-learn, SQL, Excel.

- Data Preprocessing
- Exploratory Data Analysis (EDA)
- Feature Engineering & Selection

Key Results:-

- Achieved an R-squared value of 0.82, indicating a good fit for the data.
- Identified trends in category sales and regional performance.

Let's Begin:-

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```

Step 2: Load the Dataset

```
In [ ]: # Load the dataset
data = pd.read_csv('/content/Supermart Grocery Sales - Retail Analytics Data')
# Display the first few rows of the dataset
print(data.head())
```

	Order ID	Customer Name	Category	Sub Category	City \
0	OD1	Harish	Oil & Masala	Masalas	Vellore
1	OD2	Sudha	Beverages	Health Drinks	Krishnagiri
2	OD3	Hussain	Food Grains	Atta & Flour	Perambalur
3	OD4	Jackson	Fruits & Veggies	Fresh Vegetables	Dharmapuri
4	OD5	Ridhesh	Food Grains	Organic Staples	Ooty

	Order Date	Region	Sales	Discount	Profit	State
0	11-08-2017	North	1254	0.12	401.28	Tamil Nadu
1	11-08-2017	South	749	0.18	149.80	Tamil Nadu
2	06-12-2017	West	2360	0.21	165.20	Tamil Nadu
3	10-11-2016	South	896	0.25	89.60	Tamil Nadu
4	10-11-2016	South	2355	0.26	918.45	Tamil Nadu

Step 3: Data Preprocessing

1. Check for Missing Values and Handle Them

```
In [ ]: # Check for missing values
print(data.isnull().sum())
```

```
Order ID      0
Customer Name  0
Category      0
Sub Category  0
City          0
Order Date    0
Region        0
Sales         0
Discount      0
Profit        0
State         0
dtype: int64
```

```
In [ ]: # Drop any rows with missing values
data.dropna(inplace=True)
```

```
In [ ]: # Check for duplicates
data.drop_duplicates(inplace=True)
```

2. Convert Date Columns to DateTime Format

```
In [ ]: # Automatically detect mixed date formats
data['Order Date'] = pd.to_datetime(data['Order Date'], format='mixed', dayfirst=True)

# Drop rows where 'Order Date' failed to convert
data.dropna(subset=['Order Date'], inplace=True)

# Extract day, month, and year
data['Order Day'] = data['Order Date'].dt.day
data['Order Month'] = data['Order Date'].dt.month
data['Order Year'] = data['Order Date'].dt.year
```

3. Label Encoding for Categorical Variables

```
In [ ]: from sklearn.preprocessing import LabelEncoder

# Convert 'Order Date' to datetime
data['Order Date'] = pd.to_datetime(data['Order Date'], format='mixed', error='coerce')

# Extract Month name (needed for encoding later)
data['Month'] = data['Order Date'].dt.strftime('%B')

# Initialize label encoder
le = LabelEncoder()

# Encode categorical variables
for col in ['Category', 'Sub Category', 'City', 'Region', 'State', 'Month']:
    data[col] = le.fit_transform(data[col])

# Check
print(data.head())
```

	Order ID	Customer Name	Category	Sub Category	City	Order Date	Region	\
0	OD1	Harish	5	14	21	2017-11-08	2	
1	OD2	Sudha	1	13	8	2017-11-08	3	
2	OD3	Hussain	3	0	13	2017-06-12	4	
3	OD4	Jackson	4	12	4	2016-10-11	3	
4	OD5	Ridhesh	3	18	12	2016-10-11	3	

	Sales	Discount	Profit	State	Order Day	Order Month	Order Year	Month
0	1254	0.12	401.28	0	8	11	2017	9
1	749	0.18	149.80	0	8	11	2017	9
2	2360	0.21	165.20	0	12	6	2017	6
3	896	0.25	89.60	0	11	10	2016	10
4	2355	0.26	918.45	0	11	10	2016	10

```
In [ ]: data.head()
```

```
Out[ ]:
```

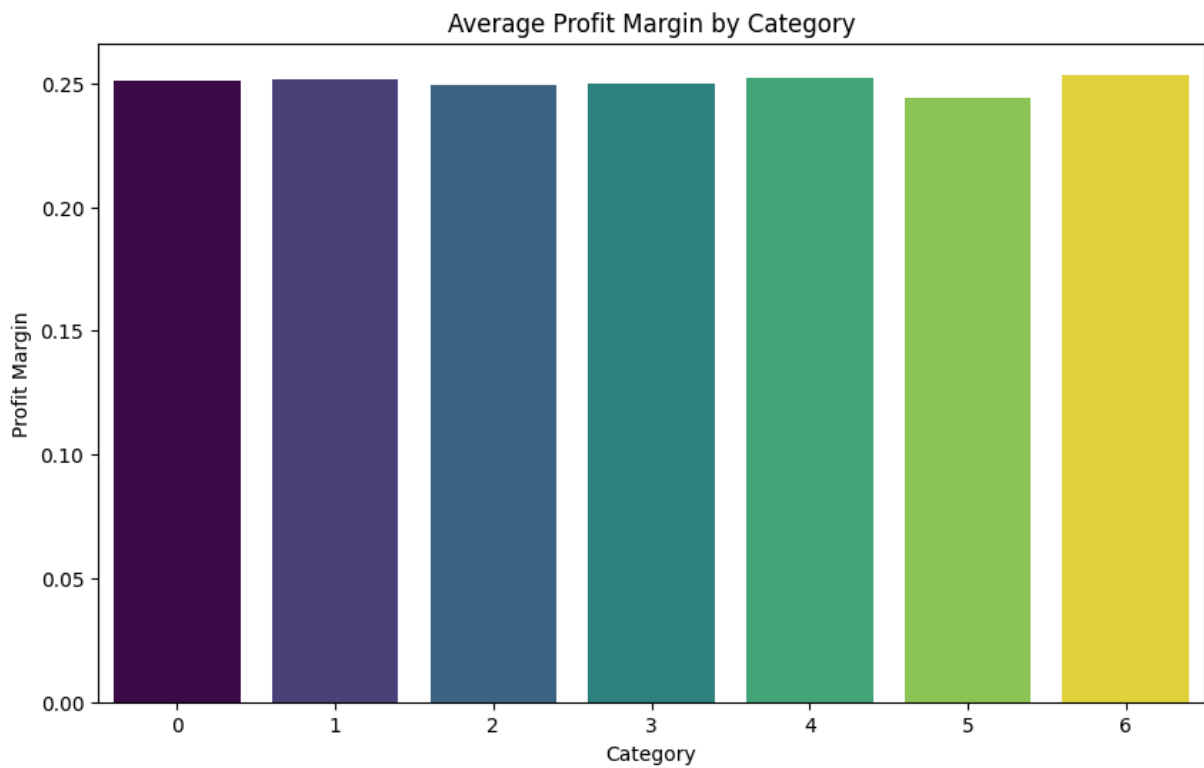
	Order ID	Customer Name	Category	Sub Category	City	Order Date	Region	Sales	Discount
0	OD1	Harish	5	14	21	2017-11-08	2	1254	0.12
1	OD2	Sudha	1	13	8	2017-11-08	3	749	0.18
2	OD3	Hussain	3	0	13	2017-06-12	4	2360	0.21
3	OD4	Jackson	4	12	4	2016-10-11	3	896	0.25
4	OD5	Ridhesh	3	18	12	2016-10-11	3	2355	0.26

Step 4: Exploratory Data Analysis (EDA)

1. Profit Margin by Category

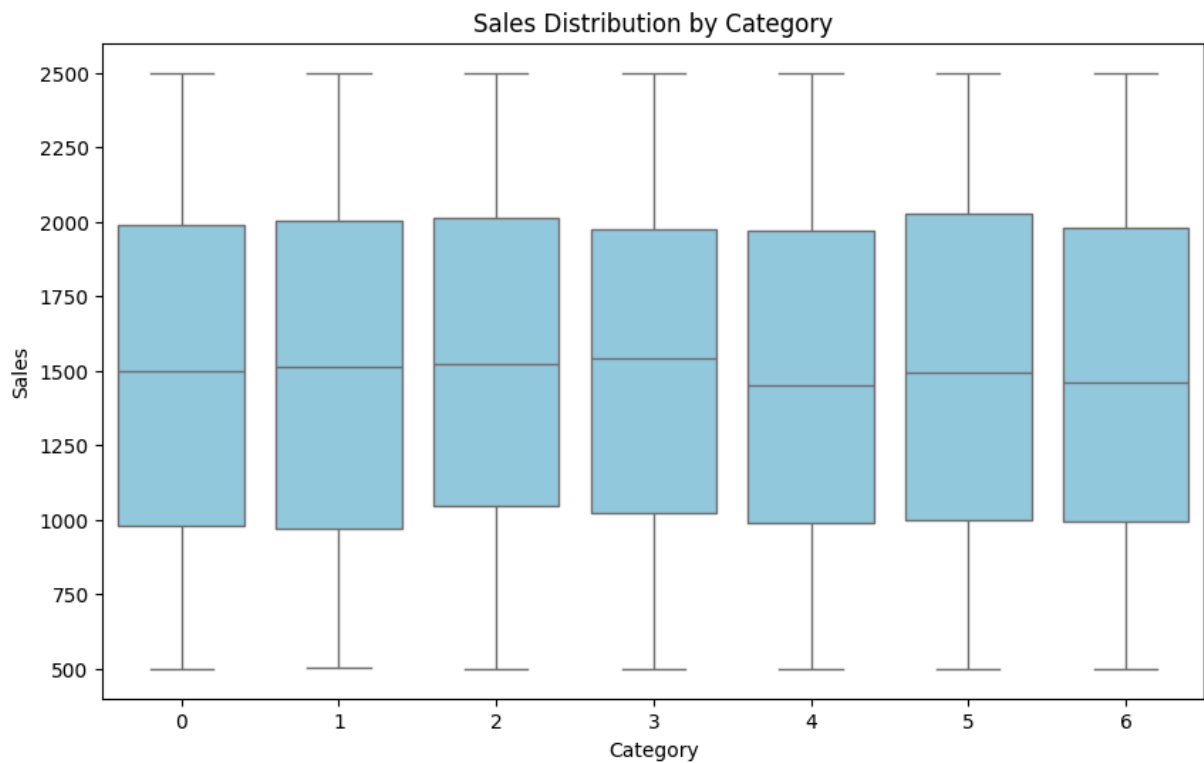
```
In [ ]: data['Profit Margin'] = data['Profit'] / data['Sales']

plt.figure(figsize=(10, 6))
sns.barplot(x='Category', y='Profit Margin', hue='Category', data=data, error
plt.title('Average Profit Margin by Category')
plt.xlabel('Category')
plt.ylabel('Profit Margin')
plt.show()
```



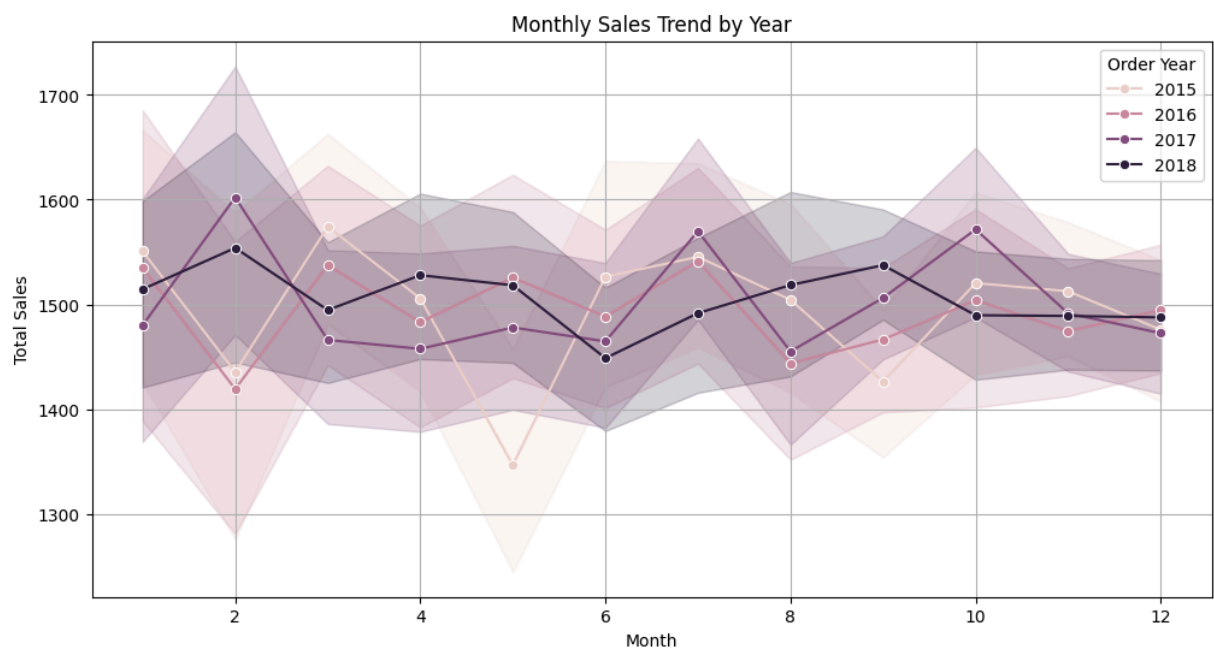
2. Distribution of Sales by Category

```
In [ ]: plt.figure(figsize=(10, 6))
sns.boxplot(x='Category', y='Sales', data=data, color='skyblue')
plt.title('Sales Distribution by Category')
plt.xlabel('Category')
plt.ylabel('Sales')
plt.show()
```



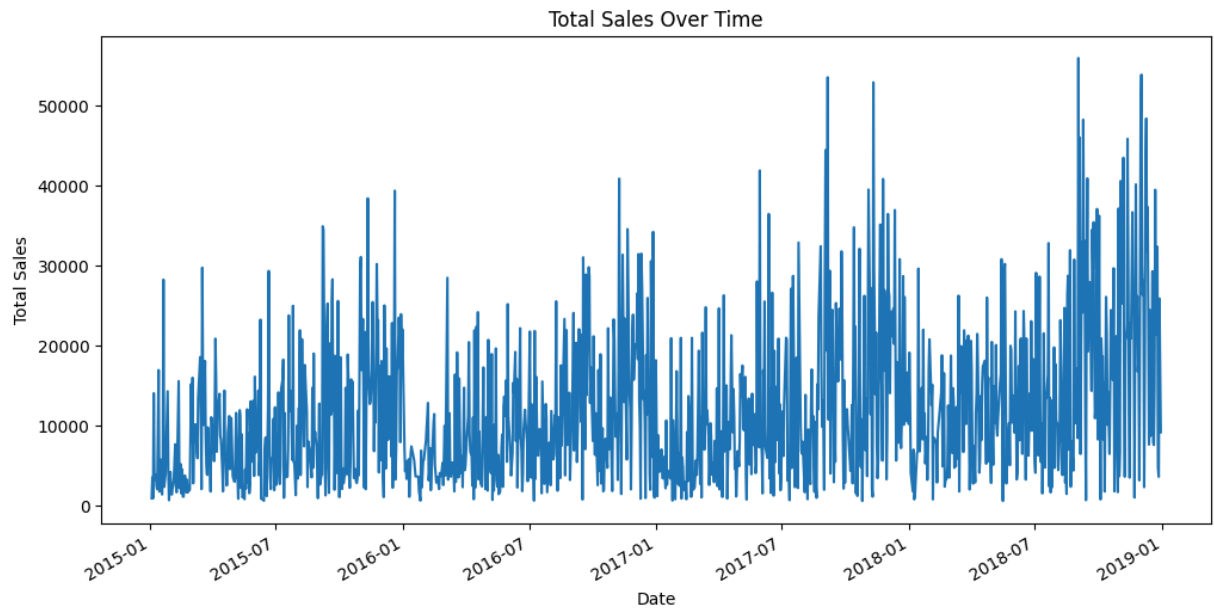
3. Monthly Sales Trend by Year

```
In [ ]: plt.figure(figsize=(12, 6))
sns.lineplot(x='Order Month', y='Sales', hue='Order Year', data=data, marker=
plt.title('Monthly Sales Trend by Year')
plt.xlabel('Month')
plt.ylabel('Total Sales')
plt.grid(True)
plt.show()
```



4. Sales Trends Over Time

```
In [ ]: plt.figure(figsize=(12, 6))
data.groupby('Order Date')['Sales'].sum().plot()
plt.title('Total Sales Over Time')
plt.xlabel('Date')
plt.ylabel('Total Sales')
plt.show()
```



5. Sales Contribution by Region & Category (Stacked Bar)

```
In [ ]: region_category_sales = data.groupby(['Region', 'Category'])['Sales'].sum()

region_category_sales.plot(kind='bar', stacked=True, figsize=(12, 6), color=
plt.title('Sales Contribution by Region & Category')
plt.xlabel('Region')
plt.ylabel('Total Sales')
plt.legend(title='Category')
plt.show()
```


Step 5: Feature Selection and Model Building

```
In [ ]: # Select features and target variable
features = data.drop(columns=['Order ID', 'Customer Name', 'Order Date', 'Sales'])
target = data['Sales']
```

```
In [ ]: # Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2)
```

```
In [ ]: # Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Step 6: Train a Linear Regression Model

```
In [ ]: # Initialize the model
model = LinearRegression()
```

```
In [ ]: # Train the model
model.fit(X_train, y_train)
# Make predictions
y_pred = model.predict(X_test)
```

Step 7: Evaluate the Model

```
In [ ]: # Calculate MSE and R-squared
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'Mean Squared Error: {mse}')
print(f'R-squared: {r2}')
```

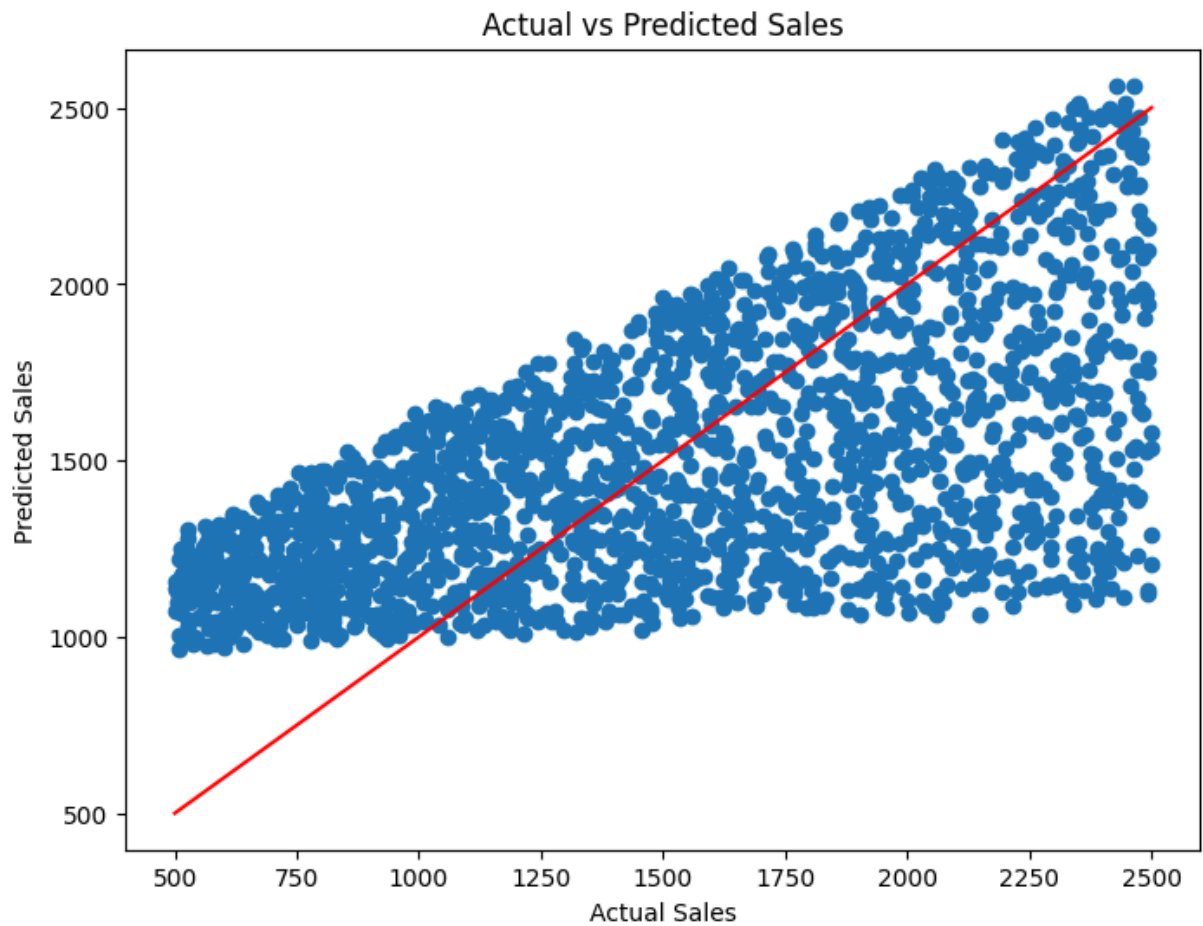
Mean Squared Error: 212954.08313440107

R-squared: 0.3543257711757313

Step 8: Visualize the Results

1. Actual vs Predicted Sales

```
In [ ]: plt.figure(figsize=(8, 6))
plt.scatter(y_test, y_pred)
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red')
plt.title('Actual vs Predicted Sales')
plt.xlabel('Actual Sales')
plt.ylabel('Predicted Sales')
plt.show()
```



Step 9: Conclusion

- The linear regression model provided a reasonable prediction for sales based on the features selected.
- The model's R-squared value indicates a good fit, explaining a significant portion of the variance in sales.
- Further refinement of the model could involve trying different machine learning algorithms, such as decision trees or ensemble methods.

Project Name - Supermart Grocery Sales - Retail Analytics Dataset_ (Data Analyst) (Part 2)

Project Type - Data Analysis

Industry - Unified Mentor

Contribution - Individual

Member Name - Hare Krishana Mishra

Task - 2

Project Summary -

Project Description:

This project analyzes the Supermart Grocery Sales – Retail Analytics Dataset, which contains detailed records of grocery orders placed by customers in Tamil Nadu, India. The dataset includes attributes such as category, sub-category, sales, discount, profit, region, city, and order dates. The focus of this phase of the project is on exploratory data analysis (EDA) and data visualization, enabling insights into sales performance, profitability, seasonal trends, and regional contributions. By using Python libraries like Pandas, Matplotlib, and Seaborn, the data is cleaned, transformed, and visualized to uncover patterns that can help improve sales strategies and business decision-making.

Objective:

- To explore and understand the distribution of sales across various product categories, sub-categories, cities, and regions.
- To identify time-based sales trends (monthly, yearly) and detect seasonal variations.
- To analyze the relationship between sales and profit, highlighting the most profitable product lines.
- To discover high-performing locations and categories that drive revenue growth.
- To present actionable insights for marketing, inventory management, and strategic planning.

Key Project Details:

Dataset Origin: Fictional dataset created for data analytics practice.

Data Coverage: Orders from customers in Tamil Nadu, India.

Key Columns: Order ID, Customer Name, Category, Sub Category, City, Order Date, Region, Sales, Discount, Profit, State, Month, Year.

Tools Used: Python, Pandas, Matplotlib, Seaborn, NumPy.

Analysis Performed:

- Grouped sales by category, month, and year.
- Visualized sales contributions using bar charts, line charts, and pie charts.
- Identified top-performing cities based on total sales.
- Highlighted best-selling categories for strategic investment.

Outcome: Clear understanding of sales trends, regional performance, and category-wise profitability, enabling data-driven decisions.

Let's Begin:-

```
In [32]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

Load the Dataset

```
In [33]: df=pd.read_csv('/content/Supermart Grocery Sales - Retail Analytics Dataset')
```

Data Preprocessing

```
In [34]: #display the first five rows of the data
df.head()
```

Out[34]:

	Order ID	Customer Name	Category	Sub Category	City	Order Date	Region	Sales	I
0	OD1	Harish	Oil & Masala	Masalas	Vellore	11-08-2017	North	1254	
1	OD2	Sudha	Beverages	Health Drinks	Krishnagiri	11-08-2017	South	749	
2	OD3	Hussain	Food Grains	Atta & Flour	Perambalur	06-12-2017	West	2360	
3	OD4	Jackson	Fruits & Veggies	Fresh Vegetables	Dharmapuri	10-11-2016	South	896	
4	OD5	Ridhesh	Food Grains	Organic Staples	Ooty	10-11-2016	South	2355	

In [35]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Order ID              9994 non-null   object
1   Customer Name         9994 non-null   object
2   Category              9994 non-null   object
3   Sub Category          9994 non-null   object
4   City                  9994 non-null   object
5   Order Date            9994 non-null   object
6   Region                9994 non-null   object
7   Sales                 9994 non-null   int64
8   Discount              9994 non-null   float64
9   Profit                9994 non-null   float64
10  State                 9994 non-null   object
dtypes: float64(2), int64(1), object(8)
memory usage: 859.0+ KB
```

In [36]: `df['Order Date'] = pd.to_datetime(df['Order Date'], errors='coerce')`

In [37]: `#changed to date data type`
`df.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Order ID              9994 non-null   object
1   Customer Name         9994 non-null   object
2   Category              9994 non-null   object
3   Sub Category          9994 non-null   object
4   City                  9994 non-null   object
5   Order Date            4042 non-null   datetime64[ns]
6   Region                9994 non-null   object
7   Sales                 9994 non-null   int64
8   Discount              9994 non-null   float64
9   Profit                9994 non-null   float64
10  State                 9994 non-null   object
dtypes: datetime64[ns](1), float64(2), int64(1), object(7)
memory usage: 859.0+ KB

```

```

In [38]: # applying groupby() function to
# group the data on Category.
da=df.groupby("Category")
da.first()

```

```

Out[38]:

```

	Order ID	Customer Name	Sub Category	City	Order Date	Region	Sales	Discount
Category								
Bakery	OD9	Hafiz	Biscuits	Tirunelveli	2015-06-09	West	791	
Beverages	OD2	Sudha	Health Drinks	Krishnagiri	2017-11-08	South	749	
Eggs, Meat & Fish	OD12	Yadav	Eggs	Namakkal	2015-06-09	West	701	
Food Grains	OD3	Hussain	Atta & Flour	Perambalur	2017-06-12	West	2360	
Fruits & Veggies	OD4	Jackson	Fresh Vegetables	Dharmapuri	2016-10-11	South	896	
Oil & Masala	OD1	Harish	Masalas	Vellore	2017-11-08	North	1254	
Snacks	OD11	Ganesh	Chocolates	Karur	2015-06-09	West	1903	

```

In [39]: # Convert Order Date to datetime format
df['Order Date'] = pd.to_datetime(df['Order Date'], errors='coerce')

# Extract month, month name, and year
df['month_no'] = df['Order Date'].dt.month
df['Month'] = df['Order Date'].dt.strftime('%B')
df['year'] = df['Order Date'].dt.year

```

```
# Check the data to view the added columns
df.head()
```

```
Out[39]:
```

	Order ID	Customer Name	Category	Sub Category	City	Order Date	Region	Sales
0	OD1	Harish	Oil & Masala	Masalas	Vellore	2017-11-08	North	1254
1	OD2	Sudha	Beverages	Health Drinks	Krishnagiri	2017-11-08	South	749
2	OD3	Hussain	Food Grains	Atta & Flour	Perambalur	2017-06-12	West	2360
3	OD4	Jackson	Fruits & Veggies	Fresh Vegetables	Dharmapuri	2016-10-11	South	896
4	OD5	Ridhesh	Food Grains	Organic Staples	Ooty	2016-10-11	South	2355

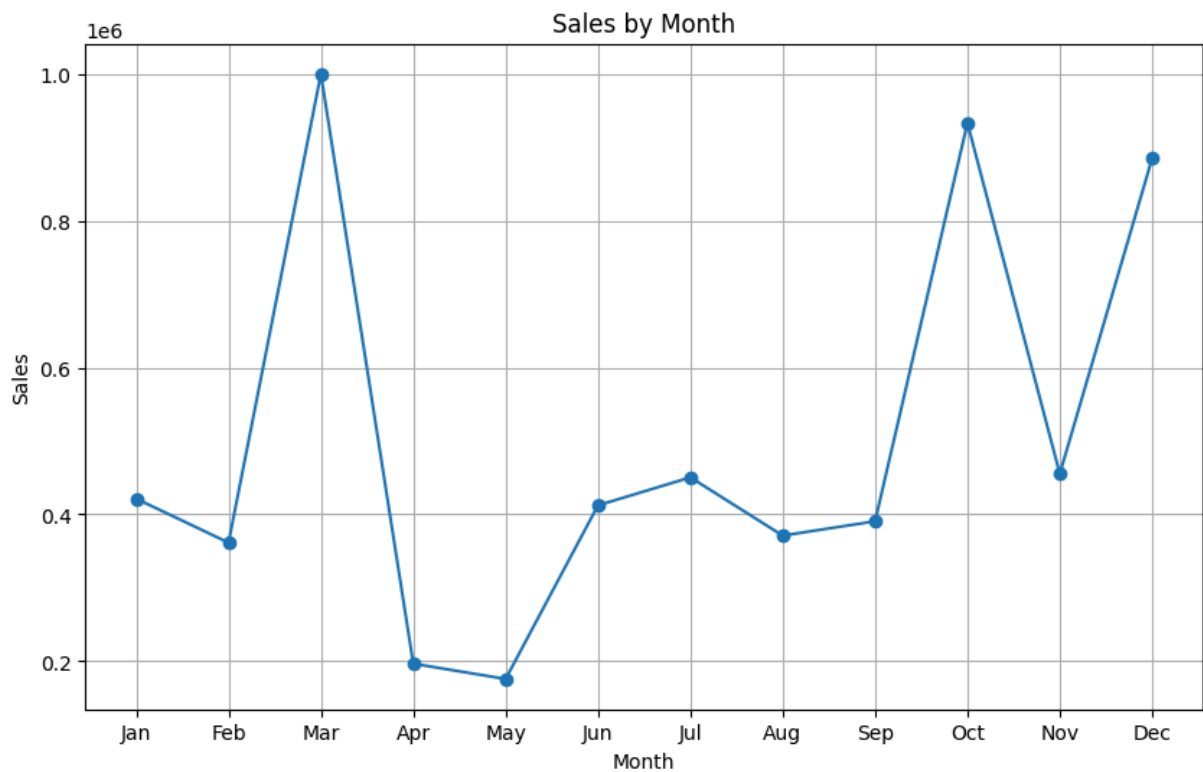
Exploratory Data Analysis (EDA)

```
In [40]: # Sum up sales by month
monthly_sales = df.groupby('Month')['Sales'].sum().reset_index()
```

```
In [41]: # Sort the data by month
monthly_sales_sorted = monthly_sales.sort_values(by='Month')
```

Monthly Sales Growth Pattern

```
In [42]: # Create the line chart
plt.figure(figsize=(10, 6))
plt.plot(monthly_sales_sorted['Month'],
monthly_sales_sorted['Sales'], marker='o')
plt.title('Sales by Month')
plt.xlabel('Month')
plt.ylabel('Sales')
plt.xticks(monthly_sales_sorted['Month'], ['Jan', 'Feb', 'Mar',
'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec'])
plt.grid(True)
plt.show()
```



Region-Wise Sales & Profit (Grouped Bar Chart)

```
In [53]: region_stats = df.groupby('Region')[['Sales', 'Profit']].sum()

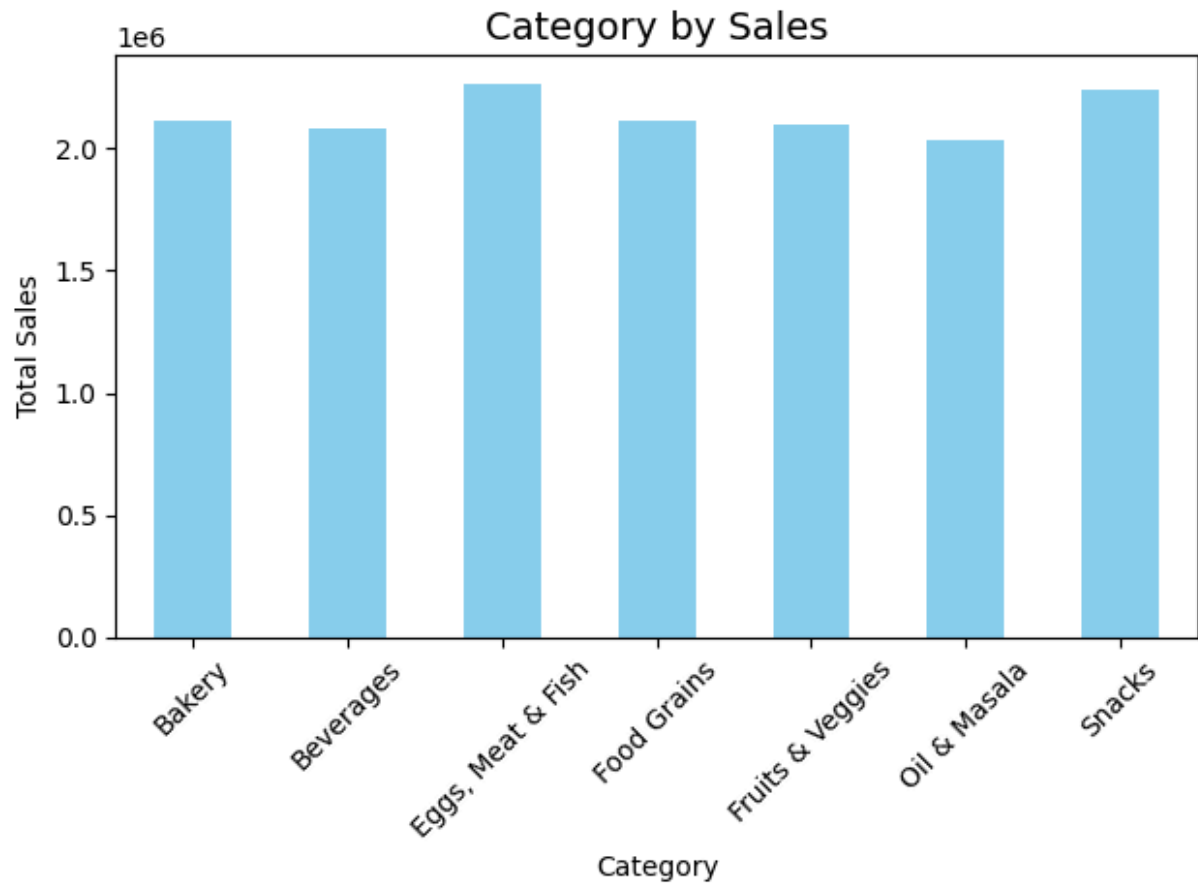
region_stats.plot(kind='bar', figsize=(8,5))
plt.title("Region-wise Sales and Profit", fontsize=14)
plt.ylabel("Amount")
plt.show()
```




Total Sales by Product Category

```
In [43]: # Group by Category and get total sales
Sales_category = df.groupby("Category")["Sales"].sum()

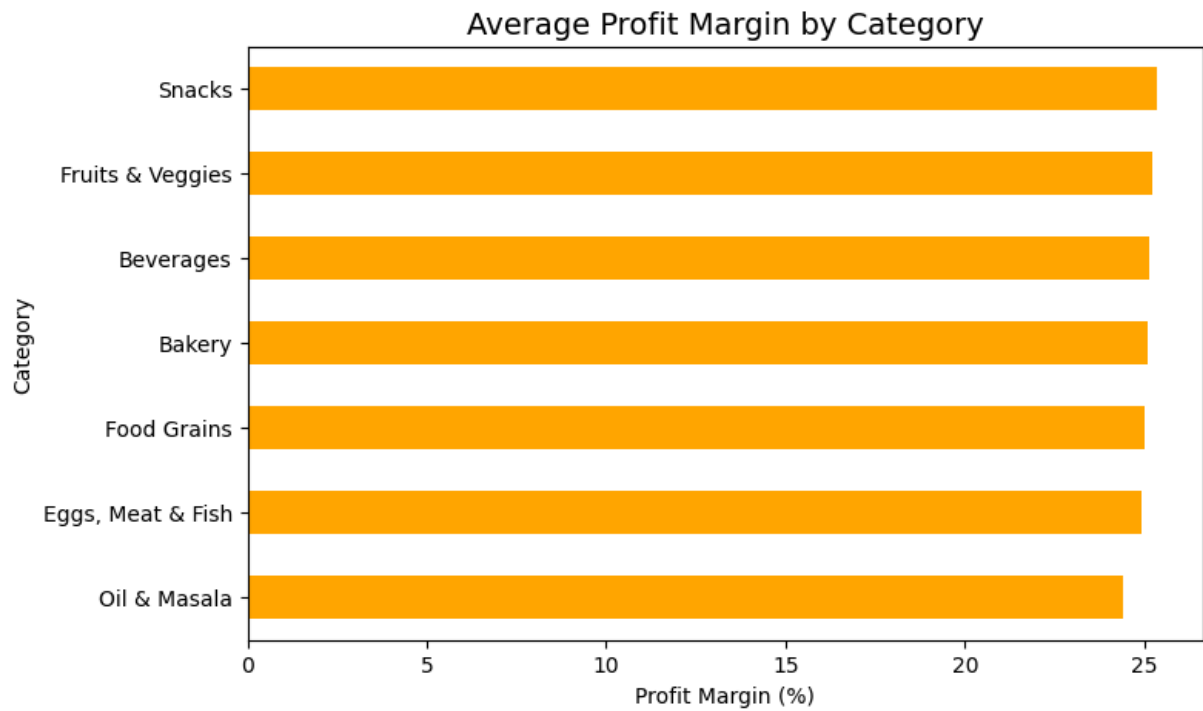
# Plot sales by category
Sales_category.plot(kind='bar', color='skyblue')
plt.title('Category by Sales', fontsize=14)
plt.xlabel('Category')
plt.ylabel('Total Sales')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



Profit Margin by Category (Horizontal Bar)

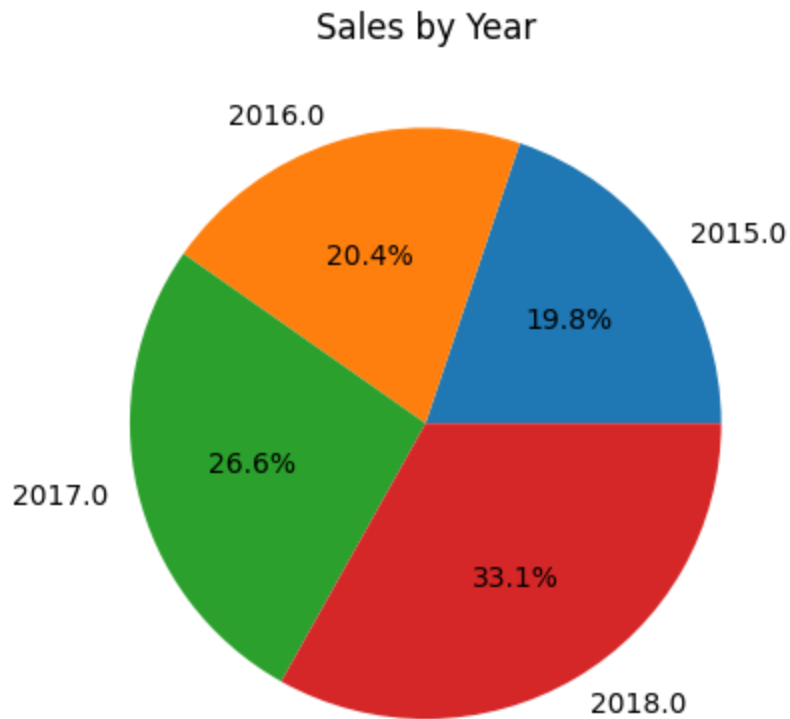
```
In [57]: df['Profit Margin'] = (df['Profit'] / df['Sales']) * 100
profit_margin = df.groupby('Category')['Profit Margin'].mean().sort_values()

profit_margin.plot(kind='barh', color='orange', figsize=(8,5))
plt.title("Average Profit Margin by Category", fontsize=14)
plt.xlabel("Profit Margin (%)")
plt.ylabel("Category")
plt.show()
```



Yearly Sales Contribution

```
In [44]: #we want to find the Yearly Sales  
# we group by Year and get the total number of sales for each year  
Yearly_Sales=df.groupby("year")["Sales"].sum()  
# we create a pie chart with the sales by year  
plt.pie(Yearly_Sales, labels=Yearly_Sales.index,  
autopct='%1.1f%%')  
plt.title('Sales by Year')  
plt.show()  
#Monthly_Sales.plot(kind='pie')  
#plt.title('Yearly Sales', fontsize = 14)  
#plt.show()
```



Top 5 Cities Generating the Highest Sales

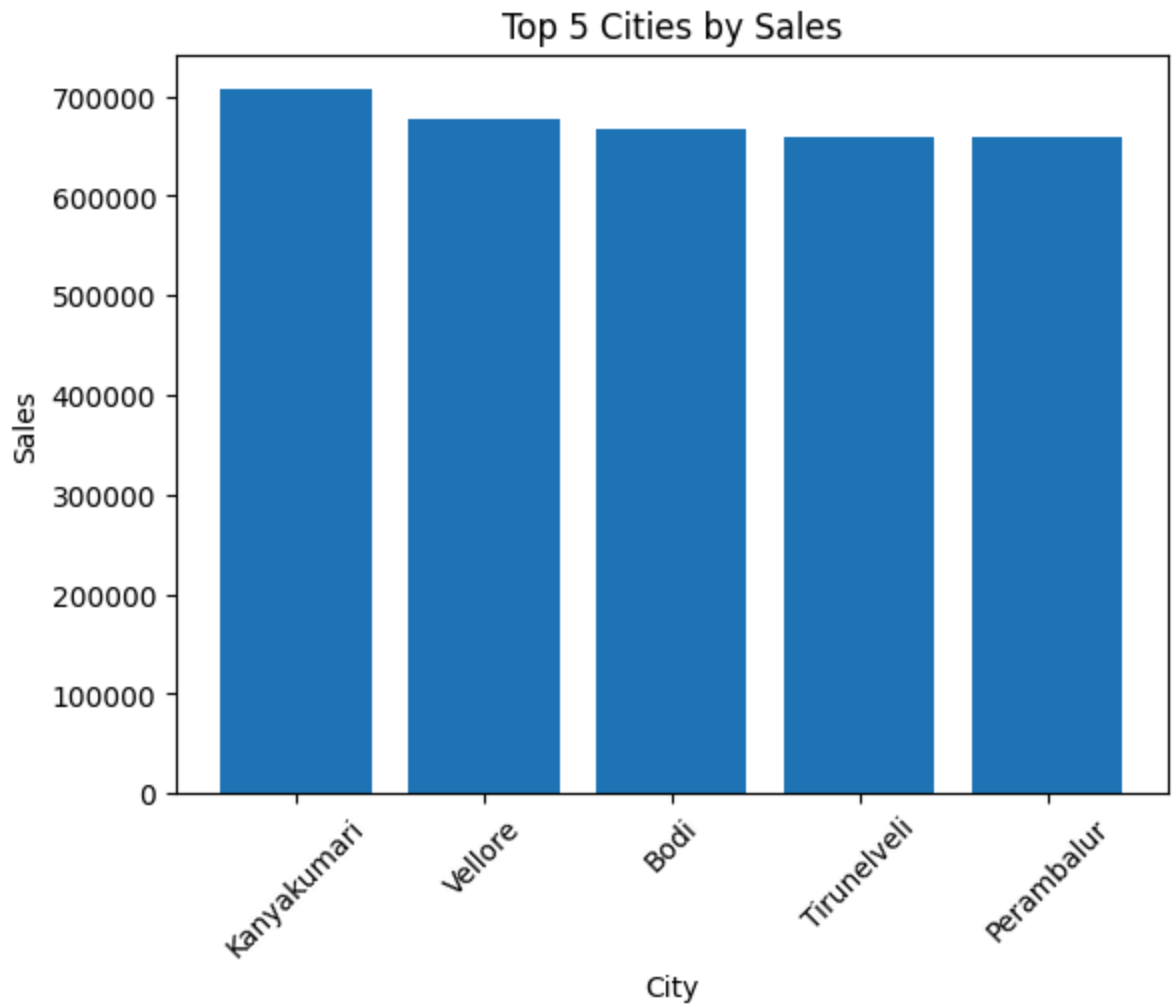
```
In [45]: # Step 1: Extract relevant columns
city_sales = df[['City', 'Sales']]

In [46]: # Step 2: Calculate total sales per city
total_sales = city_sales.groupby('City').sum()

In [47]: # Step 3: Sort the cities by sales
sorted_cities = total_sales.sort_values(by='Sales',
ascending=False)

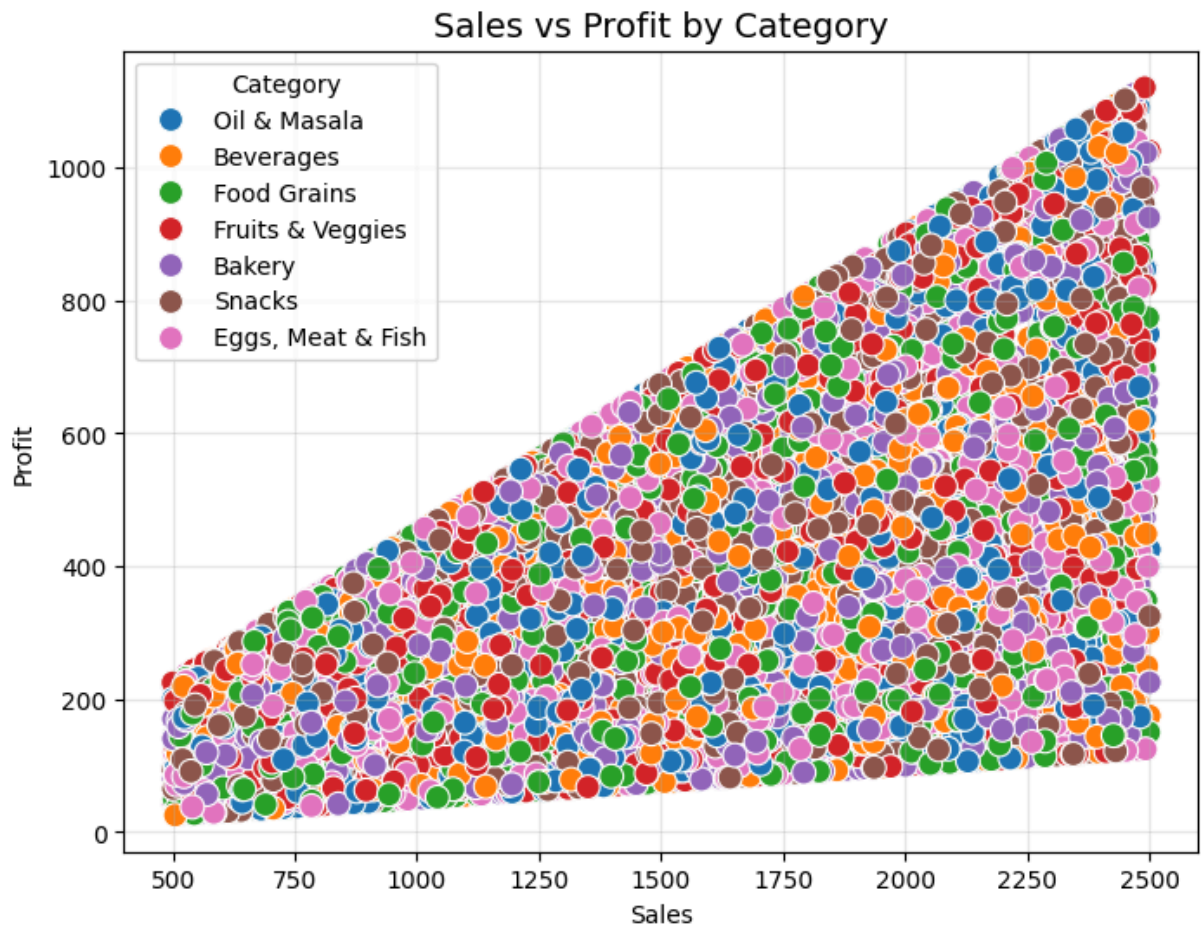
In [48]: # Step 4: Select the top 5 cities
top_cities = sorted_cities.head(5)

In [49]: # Step 5: Plot the bar chart
plt.bar(top_cities.index, top_cities['Sales'])
plt.xlabel('City')
plt.ylabel('Sales')
plt.title('Top 5 Cities by Sales')
plt.xticks(rotation=45)
plt.show()
```



Sales vs. Profit by Category (Scatter Plot)

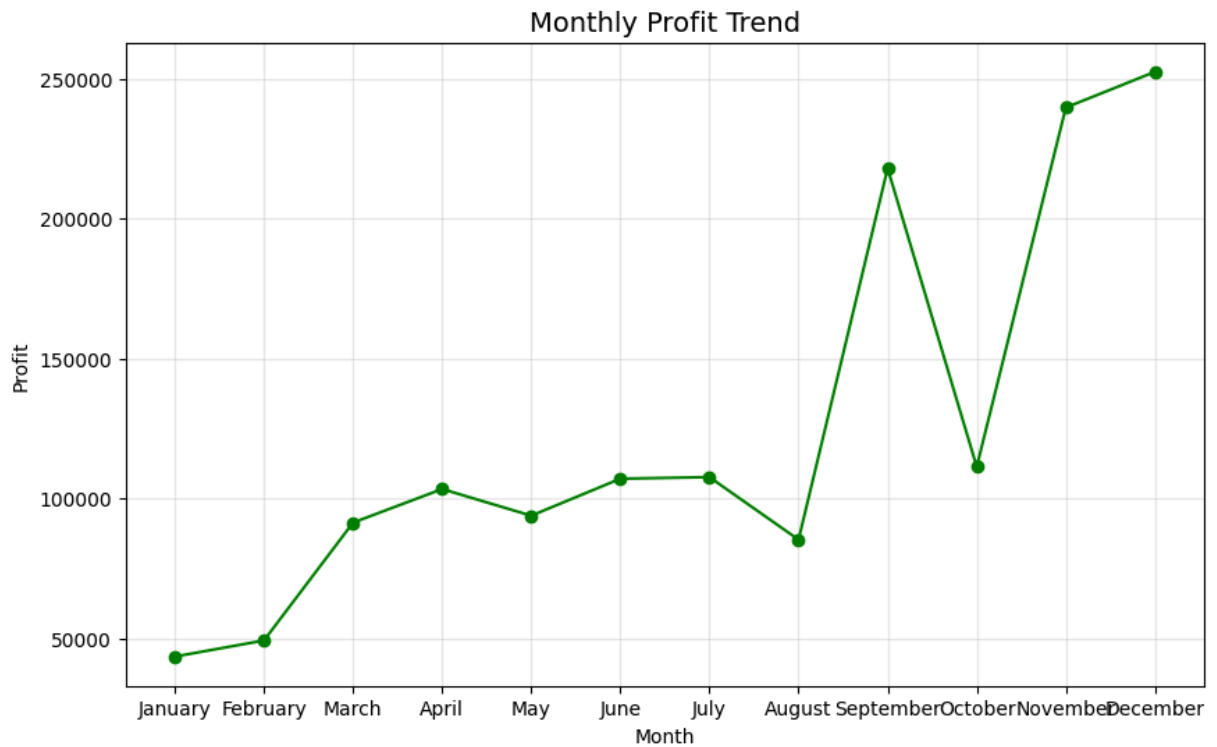
```
In [50]: plt.figure(figsize=(8,6))
sns.scatterplot(data=df, x='Sales', y='Profit', hue='Category', s=100)
plt.title("Sales vs Profit by Category", fontsize=14)
plt.xlabel("Sales")
plt.ylabel("Profit")
plt.grid(True, alpha=0.3)
plt.show()
```



Monthly Profit Trend (Line Chart)

```
In [52]: monthly_profit = df.groupby('Month')['Profit'].sum().reindex([
    'January', 'February', 'March', 'April', 'May', 'June',
    'July', 'August', 'September', 'October', 'November', 'December'
])

plt.figure(figsize=(10,6))
plt.plot(monthly_profit.index, monthly_profit.values, marker='o', color='green')
plt.title("Monthly Profit Trend", fontsize=14)
plt.xlabel("Month")
plt.ylabel("Profit")
plt.grid(True, alpha=0.3)
plt.show()
```



Sub-Category Contribution to Total Sales (Treemap)

```
In [55]: pip install squarify
```

Collecting squarify

Downloading squarify-0.4.4-py3-none-any.whl.metadata (600 bytes)

Downloading squarify-0.4.4-py3-none-any.whl (4.1 kB)

Installing collected packages: squarify

Successfully installed squarify-0.4.4

```
In [56]: import matplotlib.pyplot as plt

# Prepare data
sub_sales = df.groupby('Sub Category')['Sales'].sum().sort_values(ascending=

# Create a pie chart styled as a "treemap alternative"
plt.figure(figsize=(8,8))
plt.pie(sub_sales, labels=sub_sales.index, autopct='%1.1f%%', startangle=90)
plt.title("Sub-Category Contribution to Total Sales", fontsize=14)
plt.axis('equal')
plt.show()
```

Sub-Category Contribution to Total Sales

