# Report on Superstore Sales Analysis

## Table of Contents

# Introduction

In today's dynamic and highly competitive retail landscape, understanding customer behavior and predicting future sales is paramount for the success of any retail business. To tackle this challenge, we delve into the world of time series analysis, a powerful tool for extracting patterns, trends, and insights from sequential data. This project revolves around a comprehensive retail dataset collected from a global superstore over a span of four years.

## Context

The dataset at hand represents a treasure trove of information from a global superstore, encompassing transactions, product details, and customer interactions. Over the years, this retail giant has amassed a vast amount of data, and our objective is clear: to harness the power of exploratory data analysis (EDA) and predictive modeling to forecast sales for the upcoming 7 days from the last date in our training dataset.

## Content

Time series analysis is a powerful tool for understanding and forecasting sequential data. While this project does not involve predicting sales for the next 7 days, it draws on the principles of EDA to extract valuable information from the dataset. EDA allows us to identify trends, correlations, and areas for improvement in the retail operation.

## Dataset

The dataset at hand is meticulously curated and self-explanatory. It encompasses a wide range of retail-related information, including transaction details, product attributes, customer demographics, and more. Its comprehensiveness and clarity make it an ideal candidate for in-depth exploratory analysis.

## Inspiration

Our inspiration for this project is twofold. First, we seek to leverage EDA to gain a comprehensive understanding of the retail dataset and extract insights that may drive business decisions. Second, by sharing our EDA results, we hope to inspire data-driven strategies and actions that can enhance the retail business's efficiency, customer experience, and profitability.

With the stage set and our objectives defined, let's embark on this data-driven journey to explore the intricacies of our retail dataset and uncover insights that can drive meaningful change.

# Data Understanding

## Dataset Description

- Size: It contains thousands of records, covering multiple years.
- Structure: The dataset includes columns such as 'Order ID,' 'Order Date,' 'Ship Date,' 'Ship Mode,' 'Customer ID,' 'Customer Name,' 'Segment,' 'Country,' 'City,' 'State,' 'Postal Code,' 'Region,' 'Product ID,' 'Category,' 'Sub-Category,' 'Product Name,' and 'Sales.'

## Data Preprocessing

- Date Format Conversion: We converted 'Order Date' and 'Ship Date' to the date format (date-month-year).
- Handling Missing Data: Identified and filled 11 missing postal code values in the city of Burlington.
- Data Sorting: Sorted the dataset based on 'Order Date' in ascending order.

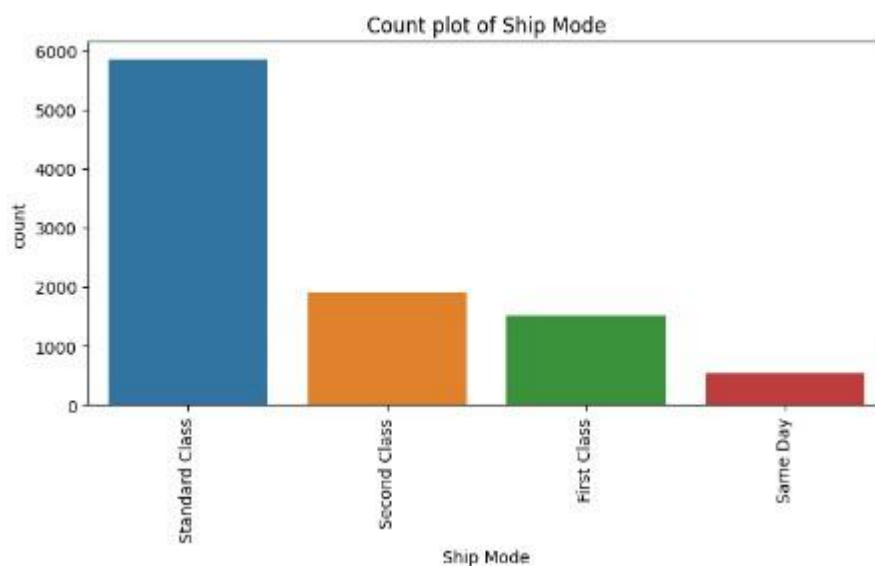# Exploratory Data Analysis (EDA)

### Ship Mode Analysis

In this analysis, we examine the distribution of ship modes and customer segments within the Superstore Sales dataset. Understanding how customers prefer to receive their orders and the different customer segments can provide valuable insights for optimizing shipping strategies and tailoring marketing efforts.

### Ship Mode Analysis

We start by analyzing the distribution of ship modes chosen by customers for their orders. Ship mode represents the method of shipping, and it is essential to ensure efficient and timely deliveries based on customer preferences.

**Key Findings:**

- Ship Mode Preferences: The dataset contains four ship modes: "Standard Class," "Second Class," "First Class," and "Same Day."
- Distribution: The most commonly chosen ship mode is "Standard Class," followed by "Second Class," "First Class," and "Same Day."
- Recommendation: Given the preference for "Standard Class," the superstore should focus on optimizing the efficiency and cost-effectiveness of standard shipping to meet customer expectations.
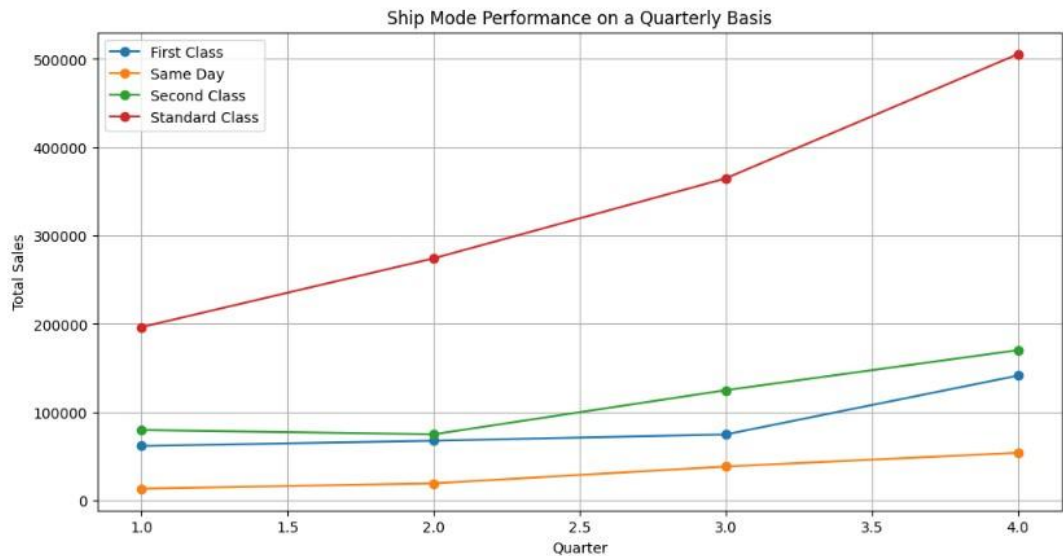


**Ship Mode Performance on a Quarterly Basis**

Next, we'll analyze how the ship modes perform over time on a quarterly basis.

**Key Findings:**
- The peak shipping period is from Q3- Q4 i.e. during holiday season
- Insights from this analysis can help in adjusting staffing and resources during peak shipping periods.

Ship Mode Performance on a Quarterly Basis

## Time to Ship Analysis

To understand the relationship between ship modes and the time it takes to ship orders, we conducted an analysis of "Time to Ship."

## Key Findings

- An analysis of variance (ANOVA) test was performed to assess the significance of differences in "Time to Ship" across ship modes.
- The ANOVA results indicate a significant relationship between "Ship Mode" and "Time to Ship" ($p < 0.05$).

## Tukey HSD Test

- To further explore the significant relationship between ship modes and time to ship, we conducted a Tukey HSD test to identify specific differences.

## Key Findings

- There are significant differences in "Time to Ship" between all pairs of ship modes. All comparisons have p-adj values less than 0.05, indicating that the mean times to ship are significantly different between these ship modes.
- Specifically, "Same Day" and "Standard Class" have the largest mean difference in "Time to Ship" ("meandiff" = 4.9638), and all comparisons involving these two ship modes are significant.
- "First Class" and "Second Class" also have a significant difference in "Time to Ship" ("meandiff" = 1.07).

## Conclusion

- Analyzing ship modes provides valuable insights into customer preferences and shipping trends.
- The superstore can use this information to optimize its shipping strategies, improve customer satisfaction, and potentially reduce shipping costs.
- Further analysis may involve investigating the relationship between ship mode and other factors like sales, region, or product category.

**Recommendations**

- Consider offering promotions or discounts for customers who choose more cost-effective ship modes like "Standard Class" to incentivize this shipping method.
- Continuously monitor ship mode performance to adapt to changing customer preferences and market trends.
- Explore the impact of ship mode on customer satisfaction and retention, as this can have a direct influence on sales and profitability.

## Customer Segment Analysis

The "Segment" column in the dataset provides information about different customer segments, such as "Consumer," "Corporate," and "Home Office."
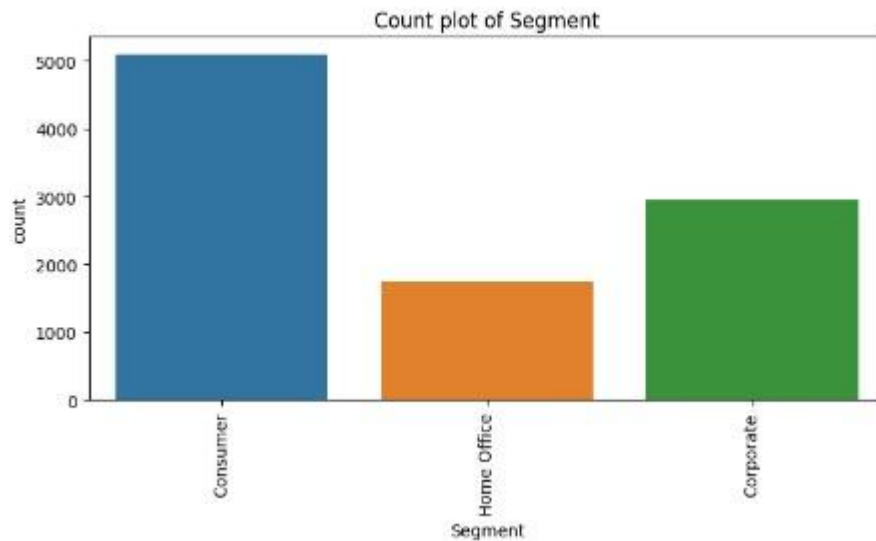
Analyzing customer segments helps the superstore understand its customer base and tailor marketing and sales strategies accordingly.

**Customer Segment Distribution**

Let's start by examining the distribution of customer segments within the dataset.

**Key Findings**

- The dataset contains three customer segments: "Consumer," "Corporate," and "Home Office."
- "Consumer" is the most common segment, followed by "Corporate" and "Home Office."
- Understanding the distribution of customer segments provides insights into the composition of the customer base.
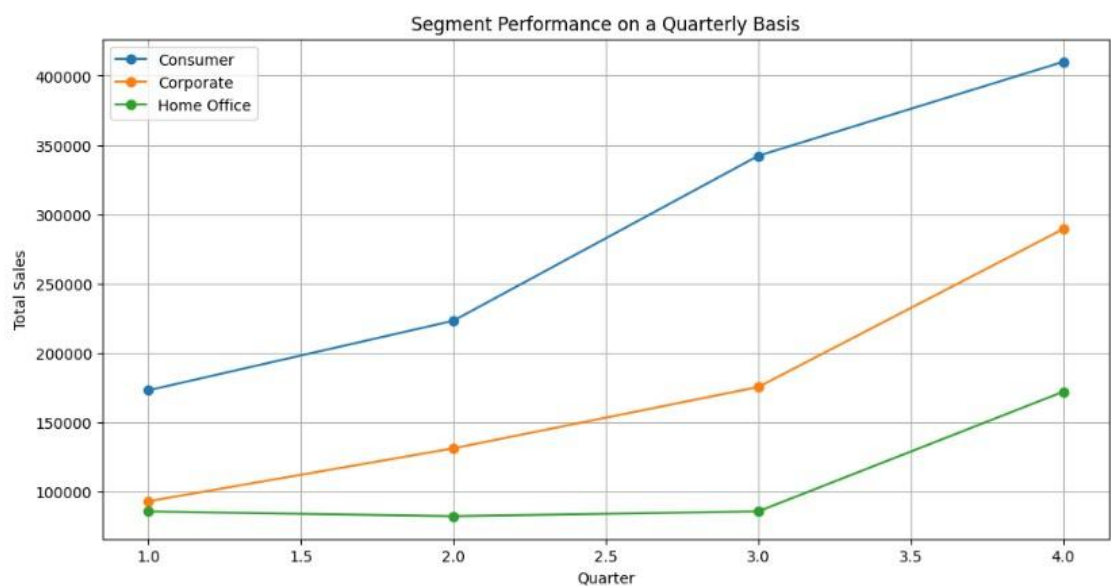
Count plot of Segment

## Segment Performance on a Quarterly Basis

Next, we'll analyze how each customer segment performs over time on a quarterly basis.

## Key Findings

- The analysis reveals the performance of customer segments over the years and allows for identifying any seasonal or long-term trends.
- Insights from this analysis can help in tailoring marketing campaigns and product offerings to specific segments.



Segment Performance on a Quarterly Basis

## Conclusion

- Analyzing customer segments provides insights into the superstore's customer base and their purchasing behavior.
- Tailoring strategies and offerings to different segments can lead to improved customer satisfaction and increased sales.

## Recommendations

- Develop customized marketing campaigns for each segment based on their unique preferences and behaviors.
- Consider loyalty programs or incentives to encourage high-value segments to make repeat purchases.
- Continuously monitor segment performance and adjust strategies as needed to stay competitive in the market.
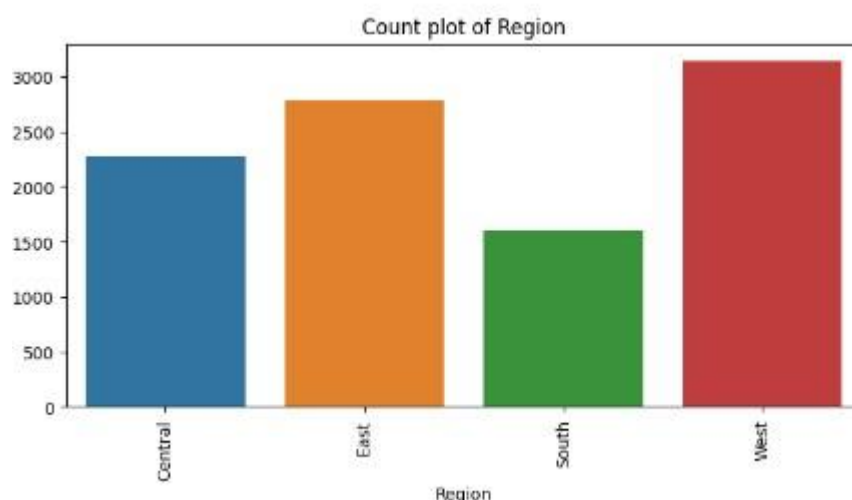
# Regional Sales Analysis

The "Region" column in the dataset provides information about different regions where the superstore operates. Analyzing regional sales can help in understanding which regions contribute the most to overall sales and identifying potential growth opportunities.

## Regional Sales Distribution

Let's begin by examining the distribution of sales across different regions.

## Key Findings
- The dataset contains sales data for three regions: "Central," "East," and "West."
- "West" is the region with the highest sales, followed by "East" and "Central."
- Understanding the distribution of sales by region provides insights into where the superstore's primary markets are located.

**Time to Ship Analysis**

To understand the relationship between regional distribution and "Time to Ship," we conducted an analysis of variance (ANOVA) test.

**Key Findings**

The ANOVA test results indicate a significant relationship between "Region" and "Time to Ship" ($p < 0.05$).

**Tukey HSD Test**

To explore the specific differences in "Time to Ship" among the regions, we performed a Tukey HSD test.

**Key Findings**

- There are significant differences in "Time to Ship" between the Central and East regions, with the East region having a shorter time to ship ("meandiff" = -0.1556, p-adj < 0.05).
- There is a significant difference in "Time to Ship" between the Central and West regions, with the West region having a shorter time to ship ("meandiff" = -0.1356, p-adj < 0.05).
- There is no significant difference in "Time to Ship" between the East and South regions ("meandiff" = 0.051, p-adj > 0.05).
- There is no significant difference in "Time to Ship" between the East and West regions ("meandiff" = 0.02, p-adj > 0.05).
- There is no significant difference in "Time to Ship" between the South and West regions ("meandiff" = -0.0309, p-adj > 0.05).

**Conclusion**

- Analyzing regional sales helps identify the superstore's top-performing regions and regions that may require additional attention.
- It provides insights into sales trends and can inform strategies for regional growth and expansion.

**Recommendations**

- Allocate marketing and promotional efforts to regions with potential for growth.
- Consider optimizing supply chain and logistics operations in regions with high sales volumes.
- Continuously monitor regional sales trends and adjust strategies to capitalize on market dynamics.

# Product Category Analysis

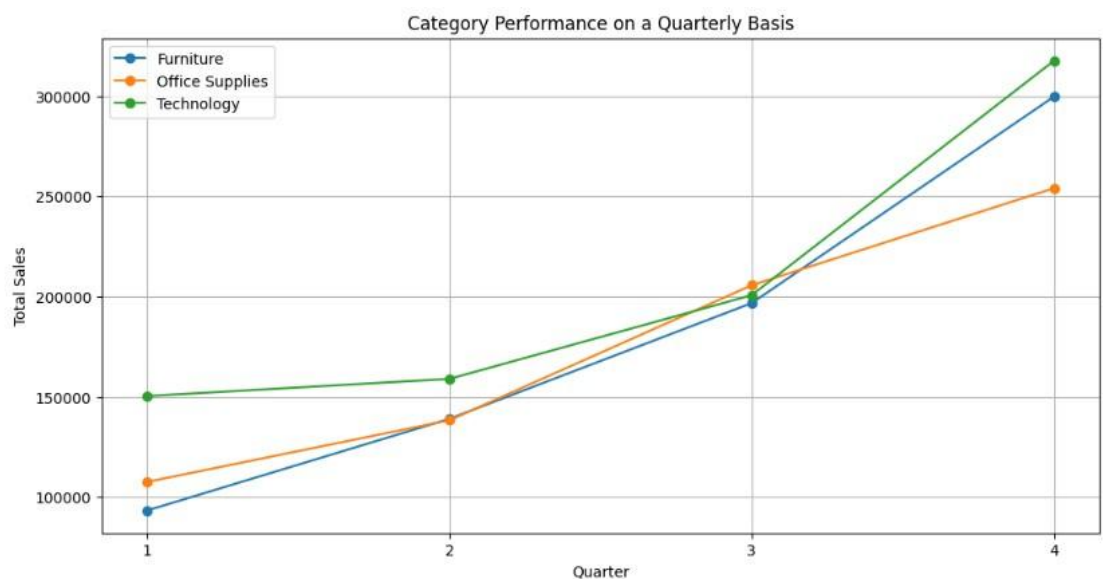In this section, we delve into the performance of product categories over different quarters of the year.

We aim to identify trends, changes, and seasonality in the sales of Furniture, Office Supplies, and Technology categories.

**Quarterly Sales Trends**

Analyzing sales trends on a quarterly basis helps us understand how each product category performs throughout the year.

**Key Findings**

- We've plotted quarterly sales trends for each product category: Furniture, Office Supplies, and Technology.
- This visualization highlights any seasonality or variations in sales over the four-year period.
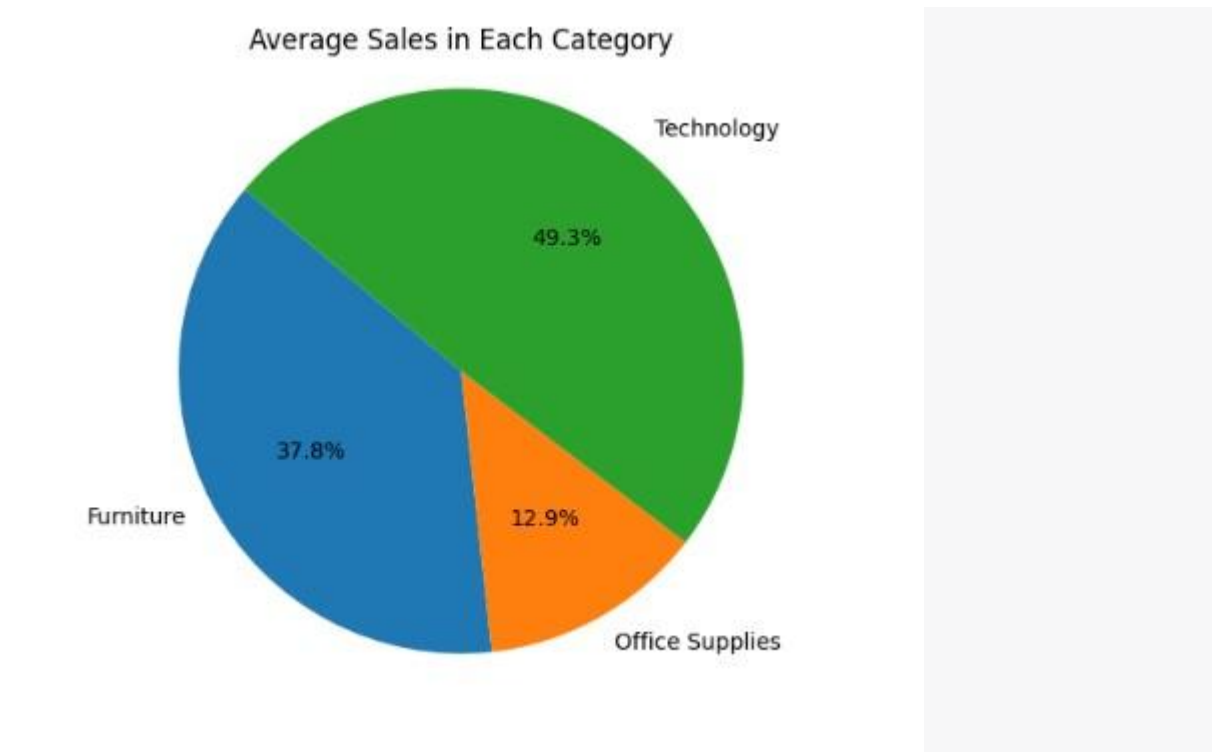


**Average Sales in Each Category**

Calculating average sales figures provides insights into the overall performance of product categories over time.

**Key Findings**

- We've computed and visualized the average sales in each category (Furniture, Office Supplies, and Technology).
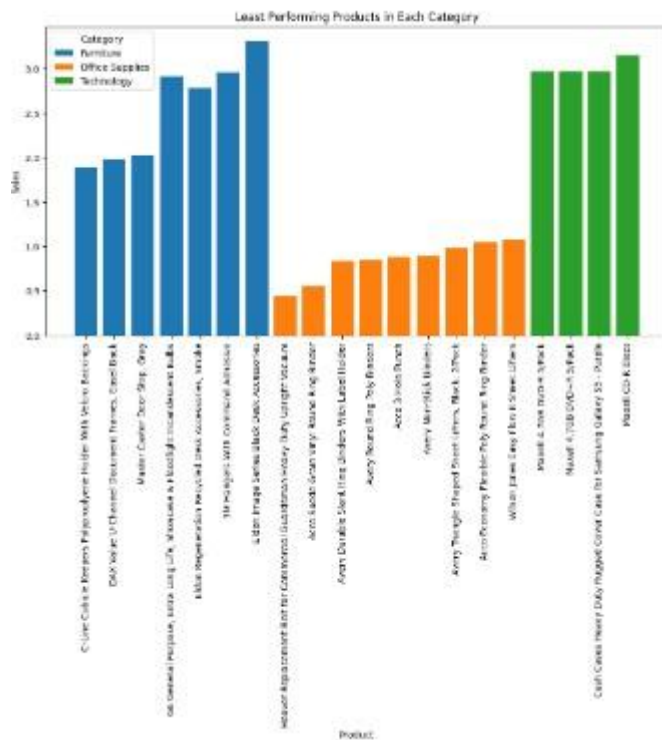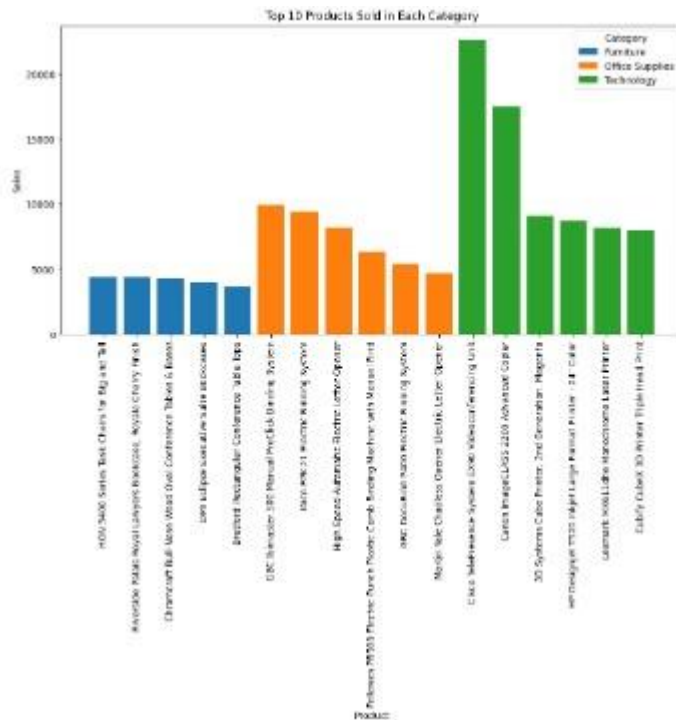
## Average Sales in Each Category

Technology 49.3%

Furniture 37.8%

Office Supplies 12.9%

**Seasonal Effects**

Analyzing seasonal effects can reveal when certain product categories experience higher or lower demand.

**Key Findings**

We've identified any significant seasonal effects or trends within each product category, shedding light on the seasonality of product sales.

Top 10 Products Sold in Each Category



Least Performing Products in Each Category

**Conclusion**

- This analysis provides valuable insights into how product categories perform on a quarterly basis.
- It highlights any seasonality or changes in sales patterns that can inform marketing and inventory management decisions.

**Recommendations**

- Use insights from quarterly sales trends to plan marketing campaigns and promotions strategically.
- Prepare for seasonality by optimizing inventory levels and marketing efforts in anticipation of high-demand quarters.
- Continuously monitor and analyze sales patterns to adapt to changing market conditions.

## Product Performance Analysis

- In this section, we explore the performance of individual products or product IDs within the Superstore Sales dataset.
- Our goal is to identify the best-selling and least-selling products and provide insights into their sales patterns.

**Best-Selling Products**

- Identifying the best-selling products is crucial for understanding which items drive the majority of sales.
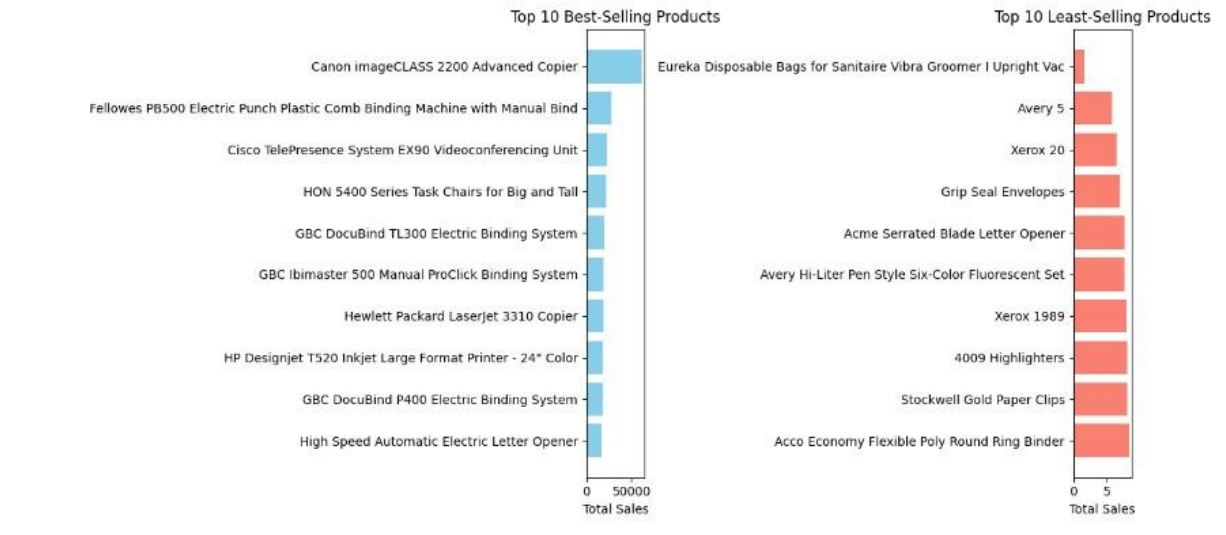
**Key Findings**

- We've identified the top 10 best-selling products in the dataset based on their total sales.
- Visualizations and tables highlight these high-performing products, allowing for easy identification.

**Least-Performing Products**

- Identifying the least-selling products can help pinpoint areas for improvement or potential discontinuation.

**Key Findings**

- We've identified the bottom 10 least-selling products in the dataset based on their total sales.
- Visualizations and tables showcase these underperforming products for further analysis.

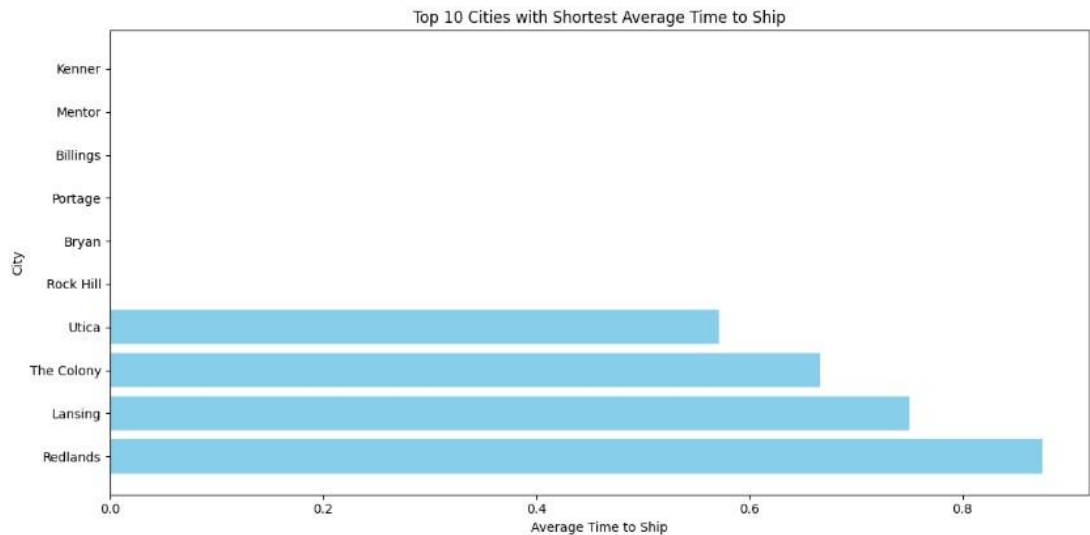Top 10 Best-Selling Products | Top 10 Least-Selling Products

## Time to Ship

- Efficient shipping and timely delivery are critical factors in customer satisfaction.
- In this section, we identify the top 10 cities with the shortest average time to ship orders, highlighting areas where the superstore excels in fast order processing.

## Key Findings

- We've computed and ranked the cities based on their average time to ship orders.
- The top 10 cities with the shortest average time to ship are as follows:
    Kenner - Average Time to Ship: 0.00 days
    Mentor - Average Time to Ship: 0.00 days
    Billings - Average Time to Ship: 0.00 days
    Portage - Average Time to Ship: 0.00 days
    Bryan - Average Time to Ship: 0.00 days
    Rock Hill - Average Time to Ship: 0.00 days
    Utica - Average Time to Ship: 0.57 days
    The Colony - Average Time to Ship: 0.67 days
    Lansing - Average Time to Ship: 0.75 days
    Redlands - Average Time to Ship: 0.88 days

## Implications

- These cities showcase excellent order processing and shipping efficiency.
- Identifying the reasons behind their success can help other regions improve their time-to-ship metrics.
- Customer satisfaction is likely higher in these areas due to quick order fulfillment.

Top 10 Cities with Shortest Average Time to Ship

## Conclusion

- This analysis sheds light on the performance of specific products within the Superstore Sales dataset.
- It identifies both the best-selling and least-selling products, enabling businesses to strategize accordingly.

## Recommendations

- Focus marketing and promotional efforts on the top-performing products to maximize revenue.
- Consider evaluating and potentially discontinuing the least-performing products to optimize inventory.
- Monitor and analyze sales trends for individual products to adapt to changing market conditions and customer preferences.

# Customer Analysis

Understanding customer behavior and preferences is vital for tailoring marketing strategies and improving overall customer satisfaction.

In this section, we delve into customer-related insights, including order frequency, churn rate, and Customer Lifetime Value (CLV).

## Order Frequency

Order frequency measures how often customers place orders with the superstore.

We calculated the order frequency for each customer by analyzing their purchase history.

**Churn Rate**

Churn rate is a key metric that quantifies the percentage of customers who stop purchasing from the superstore within a specific time frame.

We calculated the churn rate to assess customer retention and loyalty.

**Customer Lifetime Value (CLV)**

CLV represents the predicted total revenue a customer will generate throughout their entire relationship with the superstore.

We computed the CLV to estimate the long-term value of each customer.

**Key Findings**

**Order Frequency:**

We identified segments of customers with high, medium, and low order frequencies.

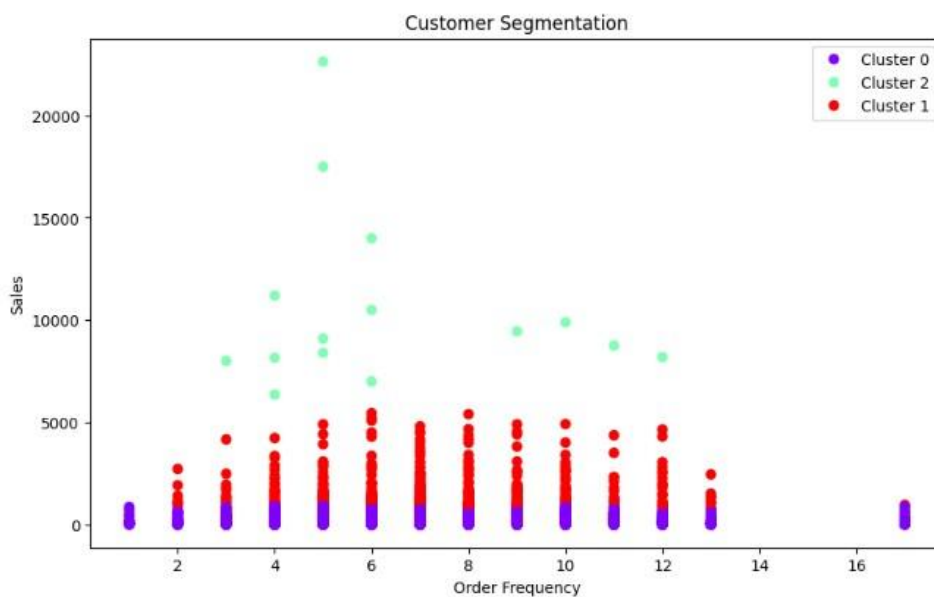| | Cluster | Cluster Label | Order Frequency | Sales | Count |
|---|---|---|---|---|---|
| 0 | 0 | High Frequency, High Sales | 7.214938 | 134.537164 | 9305 |
| 1 | 1 | Low Frequency, Low Sales | 6.333333 | 10608.891067 | 15 |
| 2 | 2 | Medium Frequency, Medium Sales | 7.189583 | 1771.948131 | 480 |

Based on the K-Means clustering analysis of customer data, we have identified three distinct customer segments:

- Cluster 0 - High Frequency, High Sales: Customers in this segment have a high order frequency (approximately 7.21 orders) and high sales (approximately $134.54). This cluster is the largest, containing 9,305 customers.
- Cluster 1 - Low Frequency, Low Sales: Customers in this segment have a low order frequency (approximately 6.33 orders) and low sales (approximately $10,608.89). This cluster is the smallest, with only 15 customers.
- Cluster 2 - Medium Frequency, Medium Sales: Customers in this segment exhibit medium order frequency (approximately 7.19 orders) and medium sales (approximately $1,771.95).

This cluster includes 480 customers. These findings provide valuable insights into the different customer behaviors within the dataset.

- Cluster 0 represents a large group of customers who frequently make purchases with relatively high sales. These customers are likely loyal and contribute significantly to revenue.
- Cluster 1 includes a small group of customers who make infrequent and low-value purchases. It may be essential to engage with this segment to increase their activity or explore reasons for their limited engagement.
- Cluster 2 comprises customers with moderate purchase frequency and spending. Understanding their needs and preferences can help tailor marketing strategies.



Overall, these customer segments allow for more targeted marketing, product recommendations, and customer retention efforts, ultimately contributing to business growth and optimization.

**Churn Rate:**

The churn rate for the superstore is 1.6393442622950838, indicating good customer retention for the superstore. It's important to continue monitoring and engaging with customers to maintain or improve this retention rate.

**Customer Lifetime Value (CLV):**

The average CLV for all customers is 1739.6436789999982. This means that, on average, a customer is expected to generate around $1,739.64 in revenue over the one-year CLV period.

**Recommendations:**

Understanding customer behavior and their relationship with the superstore is crucial for tailoring marketing efforts, enhancing customer experiences, and optimizing business strategies. The insights from this customer analysis can guide data-driven decisions and help build stronger customer relationships.

## Customer Segmentation

Customer segmentation helps categorize customers based on shared characteristics, enabling targeted marketing strategies and personalized services.

In this section, we present an analysis of customer segments within the superstore dataset, including their count, mean, standard deviation, and quartile statistics.

### Key Statistics by Customer Segment

### Consumer Segment

- Count: 5,101 customers
- Mean Order Frequency: 225.07
- Standard Deviation of Order Frequency: 588.93
- Minimum Order Frequency: 0.44
- 25th Percentile (Q1): 17.28
- Median (Q2): 53.98
- 75th Percentile (Q3): 208.56
- Maximum Order Frequency: 13,999.96

### Corporate Segment

- Count: 2,953 customers
- Mean Order Frequency: 233.15
- Standard Deviation of Order Frequency: 600.11
- Minimum Order Frequency: 0.56
- 25th Percentile (Q1): 17.38
- Median (Q2): 56.30
- 75th Percentile (Q3): 213.08
- Maximum Order Frequency: 17,499.95

### Home Office Segment

- Count: 1,746 customers
- Mean Order Frequency: 243.40
- Standard Deviation of Order Frequency: 762.86
- Minimum Order Frequency: 0.99
- 25th Percentile (Q1): 17.07

- Median (Q2): 52.56
- 75th Percentile (Q3): 211.69
- Maximum Order Frequency: 22,638.48

**Findings**

- The customer segments, namely Consumer, Corporate, and Home Office, exhibit variations in order frequency.
- The Consumer segment has the highest count of customers, while the Home Office segment has the highest mean order frequency.
- Each segment shows a wide range of order frequency, with some customers placing orders very frequently and others infrequently.
- Understanding these variations can help tailor marketing strategies and customer engagement efforts for each segment.

## Time Series forecasting using XGB

In this section, we explore the application of machine learning, specifically the XGBoost algorithm, for time series forecasting within the context of our Superstore Sales dataset analysis.

**Methodology**

**Data Preparation**

Ensure that the dataset includes a timestamp (or date) column and the target variable we aim to forecast, such as 'Sales.'

**Data Splitting**

Split the dataset into training and testing sets, typically using the first 80% of the data for training and the remaining 20% for testing.

**Model Training**

Initialize an XGBoost regressor, configured for time series forecasting. Use the training dataset to fit (train) the XGBoost model. Specify hyperparameters like the number of trees (n_estimators), learning rate (eta), maximum depth (max_depth), etc.

**Results and Findings**

**Model Evaluation**

Evaluate the forecasting accuracy of the XGBoost model using appropriate metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), or Mean Absolute Error (MAE).

**Conclusion**

Incorporating XGBoost into our time series analysis has allowed us to leverage advanced machine learning techniques for forecasting future sales. The XGBoost model, trained on historical sales data, provides valuable insights and predictions for future sales trends.

# Conclusion

- **Seasonal Sales Variation:** Sales exhibit a distinct seasonal pattern, with the highest sales occurring in Q3 and Q4, while the lowest sales are observed in Q1 and Q2. This trend highlights the importance of aligning inventory and marketing efforts with seasonal demand.

- **Preferred Ship Mode:** The "Standard Class" ship mode emerges as the most preferred choice among customers. Additionally, there is a significant relationship between the chosen ship mode and the time it takes to ship orders, emphasizing the importance of optimizing shipping strategies.

- **Customer Segment Distribution:** The majority of customers fall into the "Consumer" segment, highlighting the significance of tailoring marketing and customer engagement strategies to this segment's preferences and needs.

- **Regional Disparities:** The "West" region stands out with the highest sales count, followed by "East" and "Central." Understanding regional sales dynamics can inform expansion and marketing strategies.

- **Category Contribution:** Technology products contribute the most to sales, accounting for 49.3% of total sales, followed by Furniture at 37.9%, and Office Supplies at 12.9%. This insight can guide inventory management and marketing priorities.

- **Product Performance:** The best-selling product is the "Canon Image Class 2200 Advanced Copier," while the least-sold product is the "Eureka Disposable Bag." Identifying product performance can aid in inventory management and marketing campaigns.
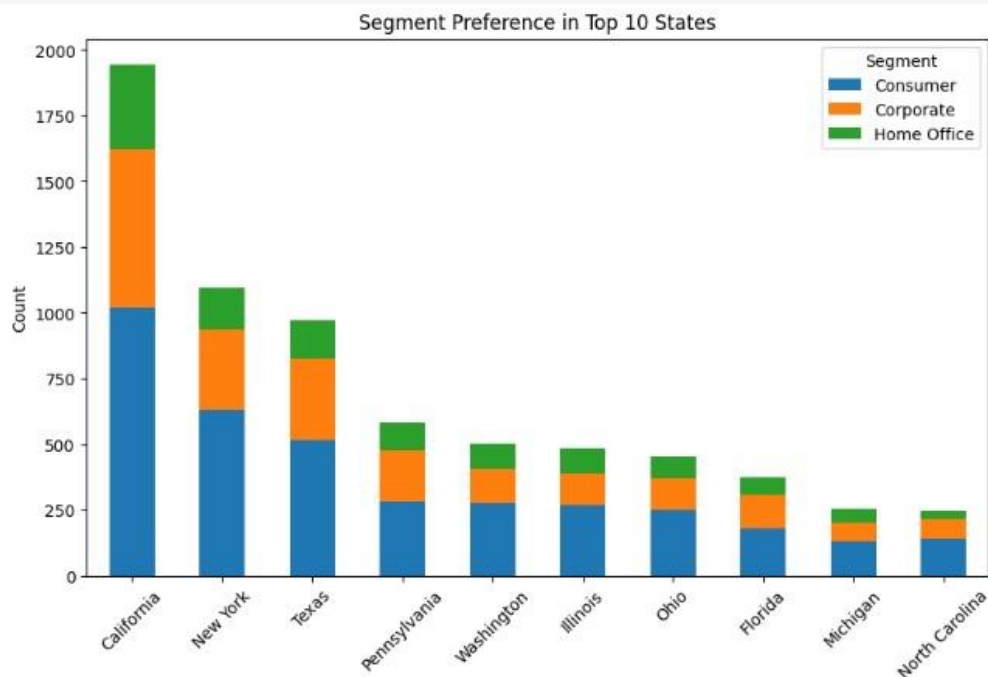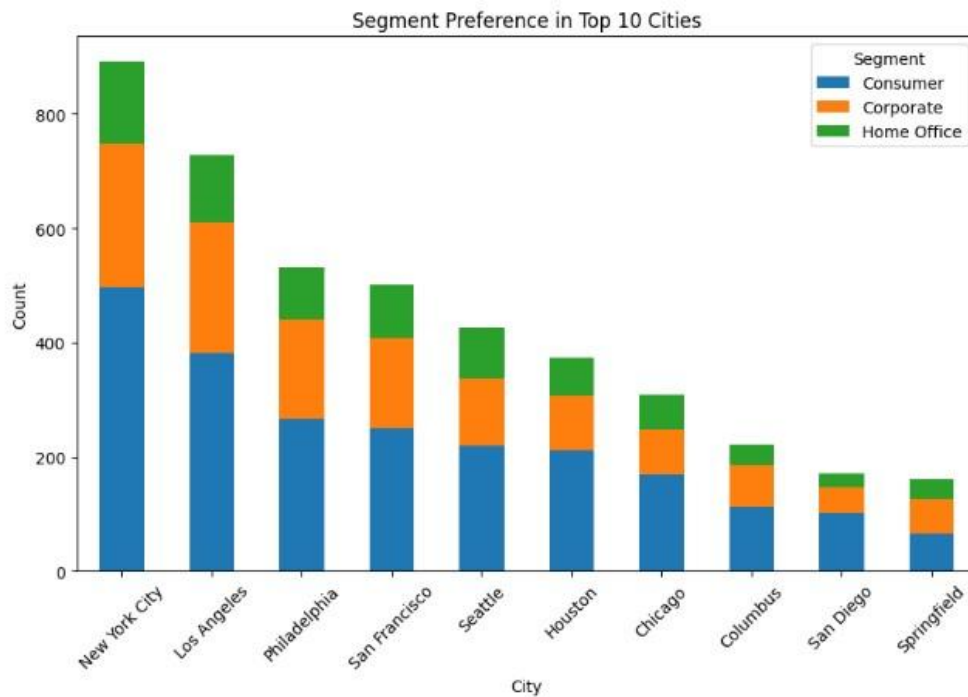
- **Customer Segmentation:** Customers are categorized into segments based on order frequency and sales behavior. Segments include "High Frequency, High Sales," "Medium Frequency, Medium Sales," and "Low Frequency, Low Sales." Tailoring strategies to these segments can enhance customer engagement.
- **Churn Rate and CLV:** The churn rate is relatively low at 1.6, indicating good customer retention rates. The Customer Lifetime Value (CLV) estimates that, on average, a customer is expected to generate around $1,739.64 in revenue over a one-year period. These metrics underscore the importance of customer loyalty and retention efforts.
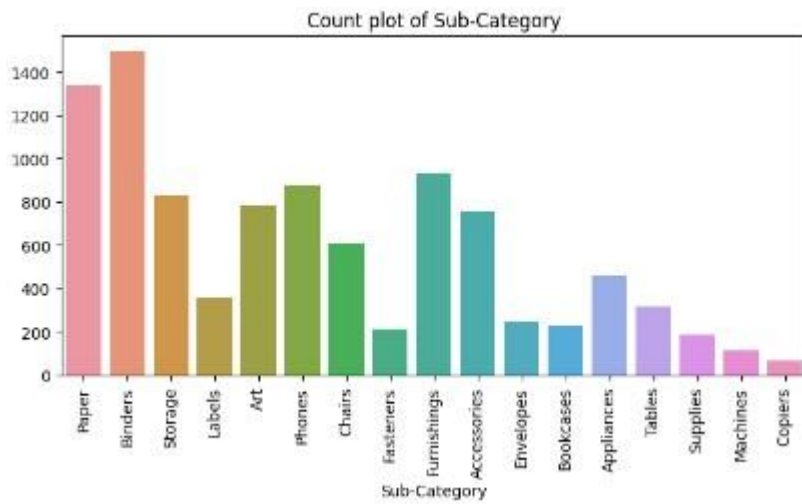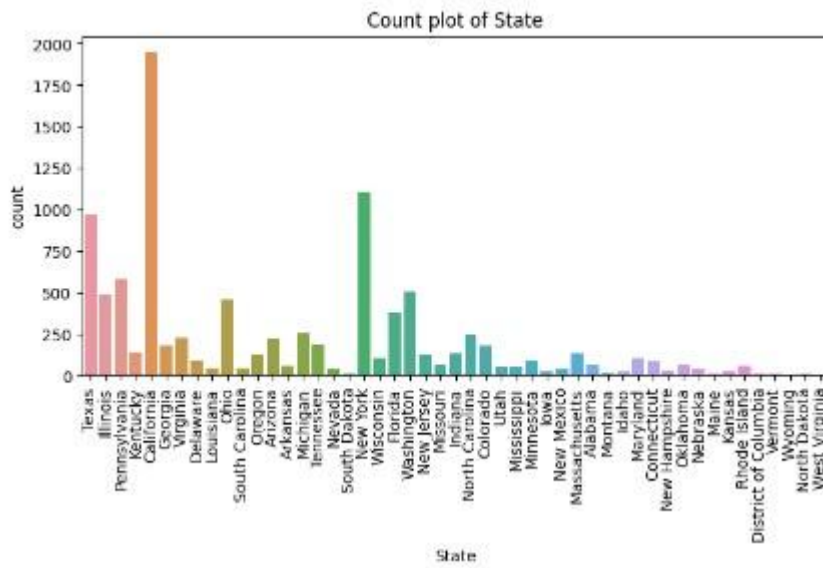
# Future Work

- **Advanced Forecasting Techniques:** Explore and implement additional machine learning techniques beyond XGBoost for sales forecasting. Experiment with models like ARIMA, LSTM, or Prophet to compare and improve forecasting accuracy.
- **Hyperparameter Tuning:** Conduct thorough hyperparameter tuning for the forecasting models used. Optimizing model parameters can significantly enhance predictive performance, leading to more accurate sales forecasts.
- **Profit Margin Analysis**: Perform a detailed analysis of profit margins across product categories, regions, or customer segments. Understanding profit margins can help identify areas for cost reduction or revenue optimization.
- **Discount Strategy Analysis:** Dive deeper into the impact of discounts on sales and customer behavior. Analyze which discount strategies are most effective in driving sales and profitability.
- **Customer Behavior Analysis**: Extend the customer analysis to gain insights into customer preferences, shopping patterns, and product affinities. This can guide personalized marketing and product recommendations.
- **Supply Chain Optimization:** Investigate the supply chain efficiency, including inventory management and order fulfillment processes. Identify areas where improvements can lead to cost savings and faster delivery times.
- **Market Expansion Strategies:** Explore opportunities for expanding into new markets or regions based on the analysis of regional sales. Assess potential growth areas and develop market entry strategies.
- **Customer Retention Strategies:** Develop and implement customer retention strategies based on churn rate and Customer Lifetime Value (CLV) insights. Focus on retaining high-value customers and increasing their loyalty.

- **Dynamic Pricing:** Investigate dynamic pricing strategies that adjust prices based on factors like demand, seasonality, or competitor pricing. Dynamic pricing can optimize revenue and profit margins.
- **Data Integration:** Consider integrating external data sources, such as economic indicators or market trends, to enhance forecasting accuracy and gain a broader understanding of market dynamics.
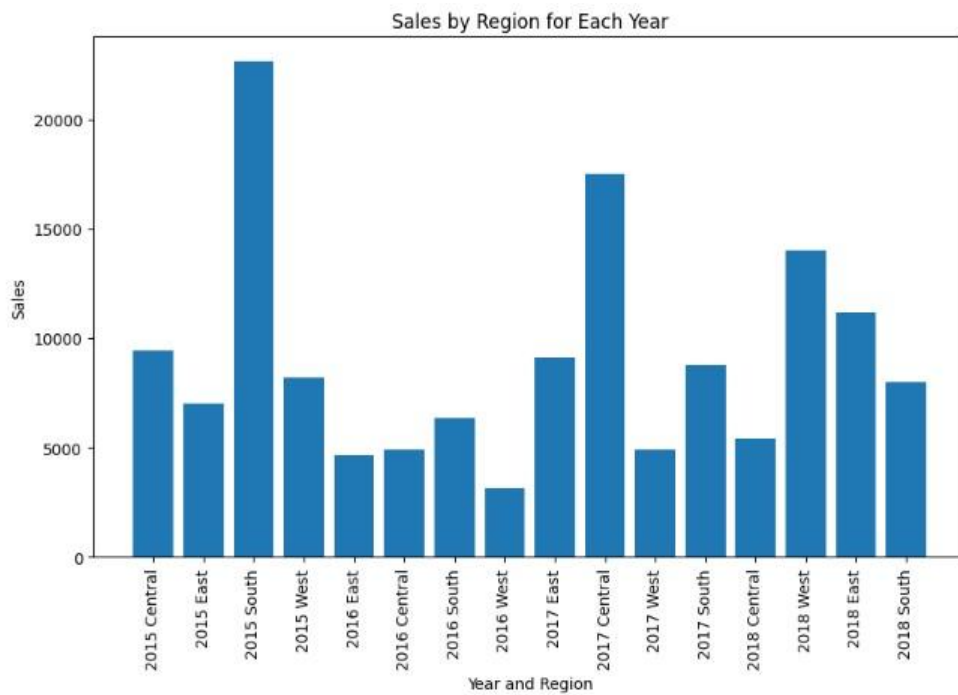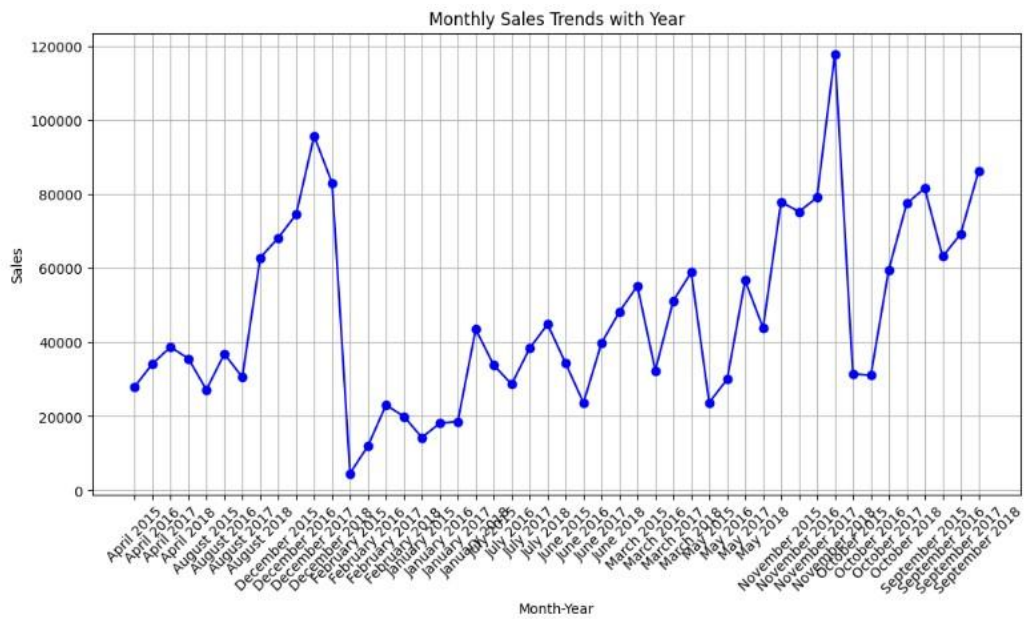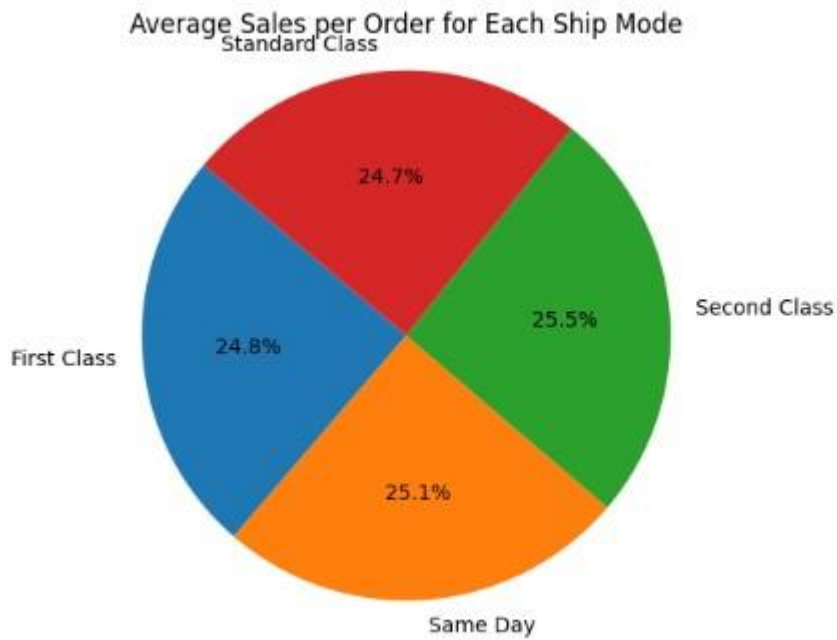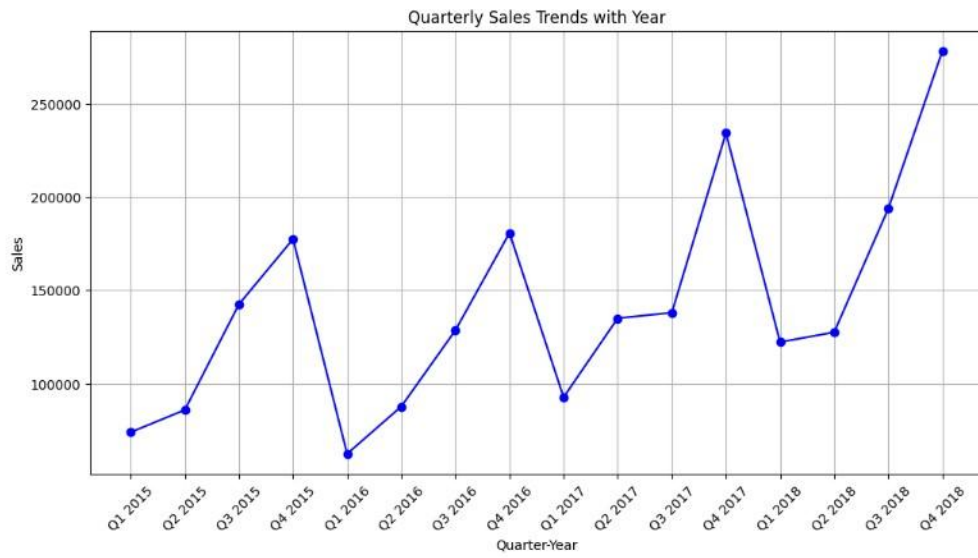
# Additional Graphs

Segment Counts Within Each Category



Average Time to Ship by State

Count plot of State


Count plot of Sub-Category


Yearly Sales Trends

Monthly Sales Trends with Year



Sales by Region for Each Year

## Quarterly Sales Trends with Year



## Average Sales per Order for Each Ship Mode

**Ship Mode Counts Within Each Category**



**Segment Preference in Top 10 Cities**



NOTE: The link to the code is my github profile: github.com/Sood-Akriti