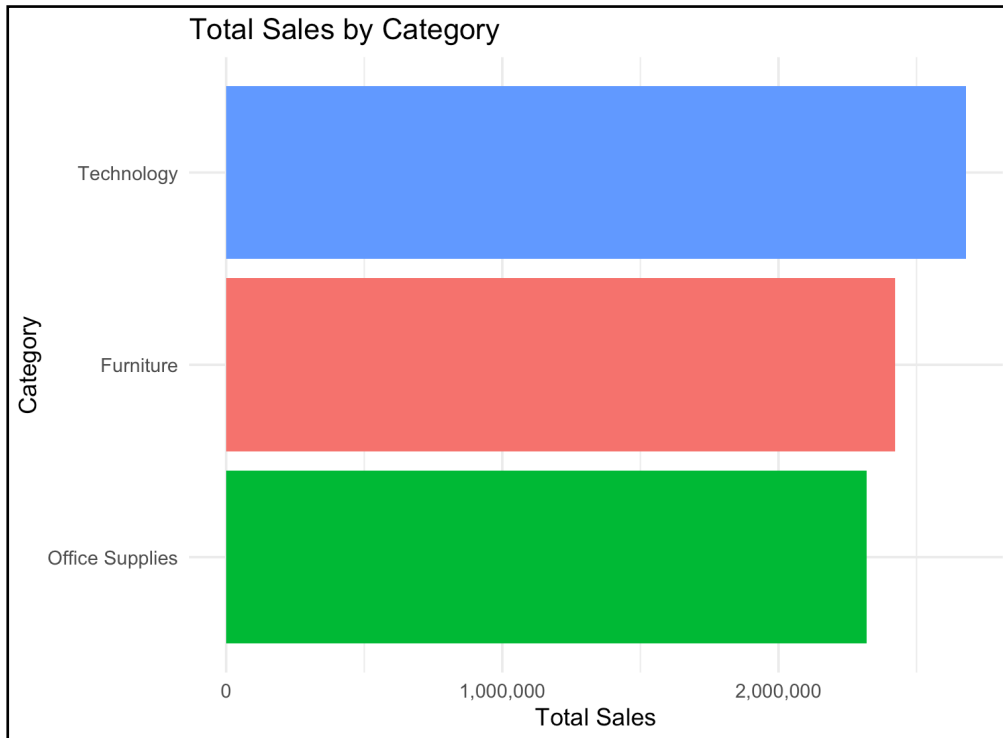


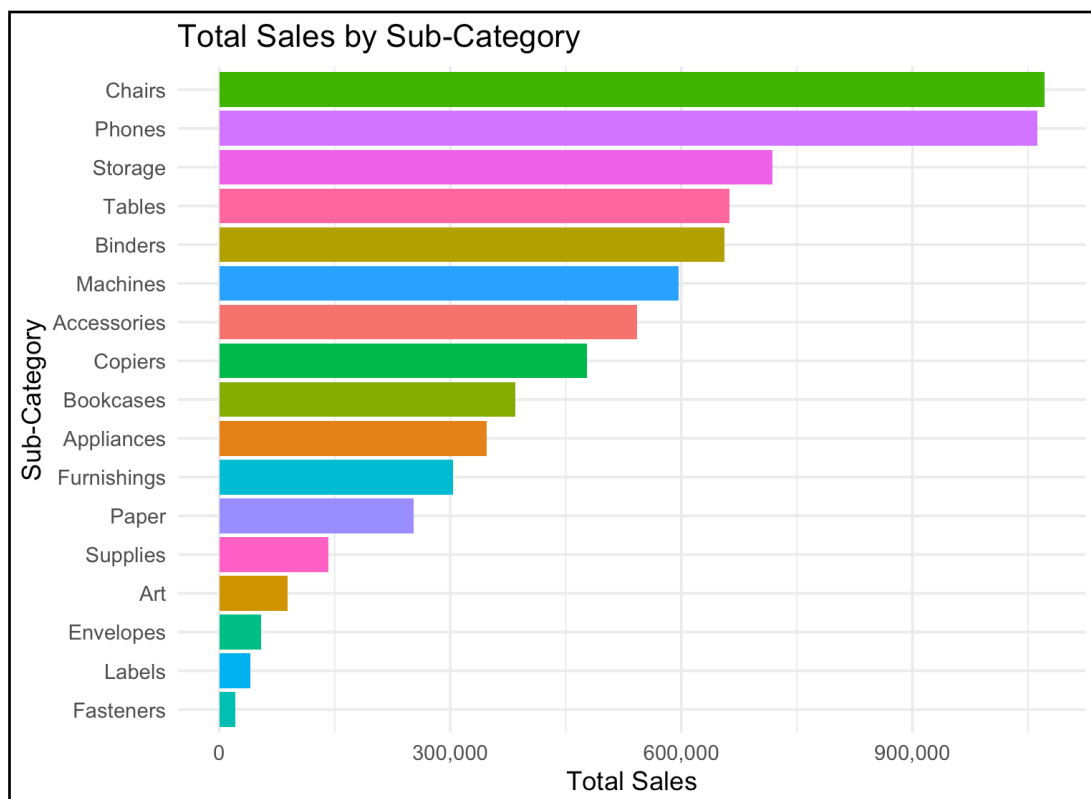
Prediction and Forecasting for Business Success

Authored by SHIH-HSUAN, YAN

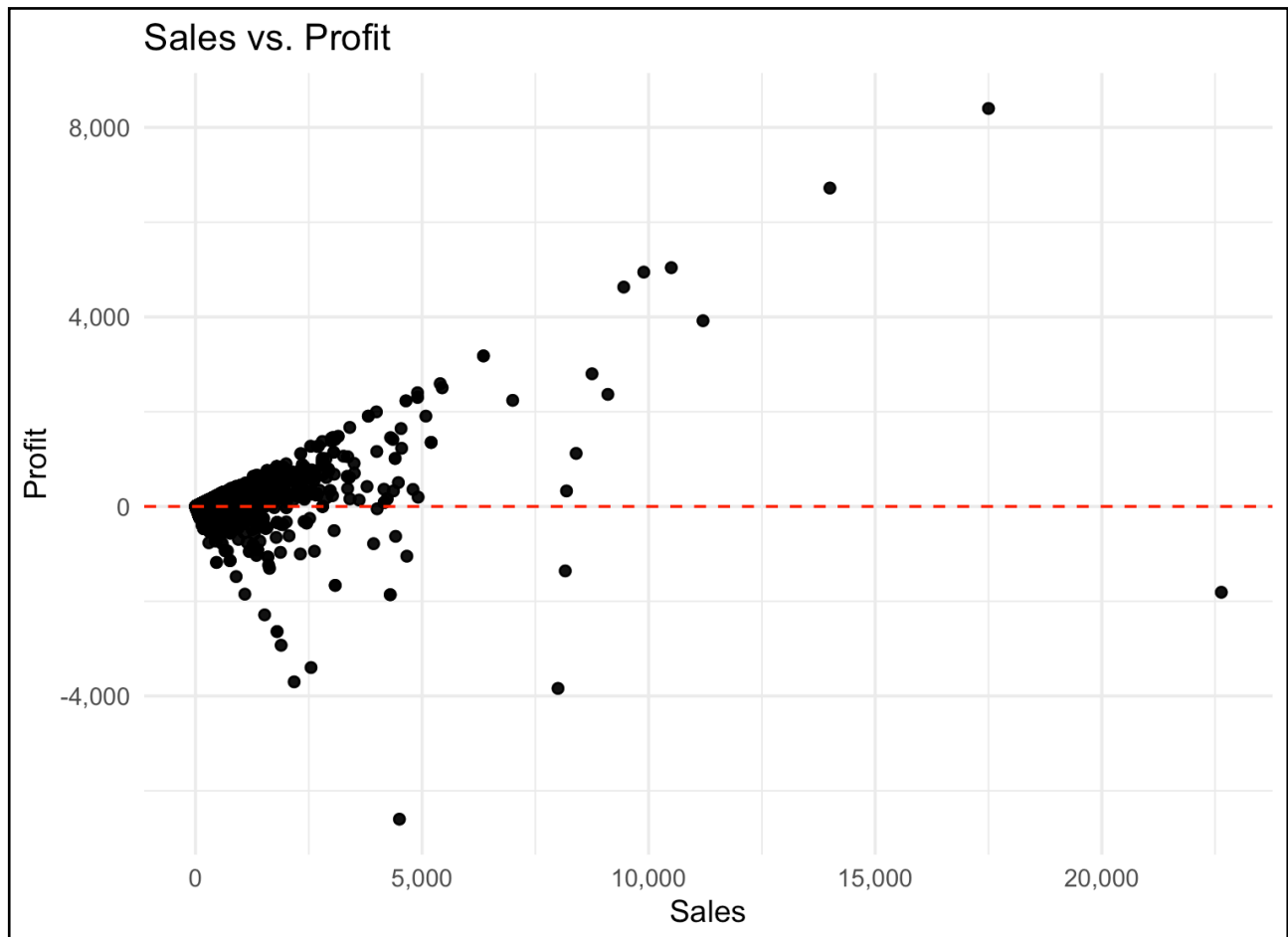
In this report, I will guide you through the dataset from a superstore which includes data from 2014 to 2024. There are a lot of columns in this data. However, given that I will bring you some interesting insights and also tell you about how the business can be successful, the main variables I will focus on are sales, profits, and discounts. Before I guide you to the main insight, let us take a look at how much sales come from different categories.



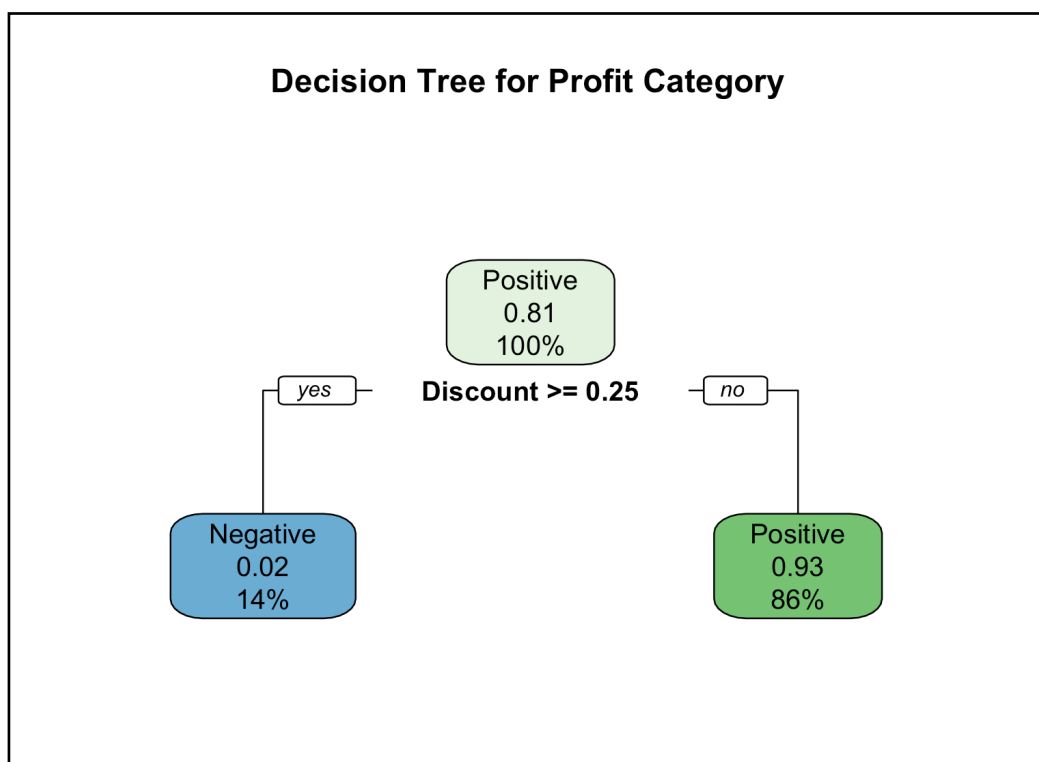
You can see that the category that makes the most sales is technology, which is pretty important in the 21st century. After technology, the sales are followed by furniture as the second category, and office supplies as the least when comparing to the previous two. These categories provide a comprehensive view of the product range and customer preferences, which is crucial for strategic planning.



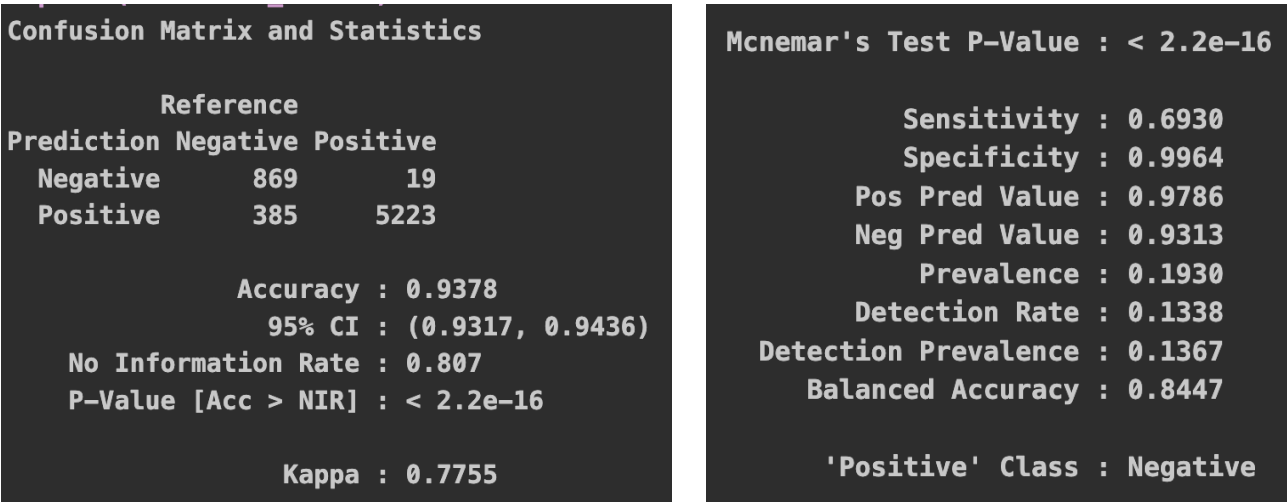
Now, let us look a bit deeper from category to sub-category. You can see that although technology has the highest sales, the highest sales for sub-category is chairs, which is in the category of furniture. However, does higher sales indicate business success? This question is essential as it highlights the need to analyze profitability, not just revenue, to ensure sustainable growth.



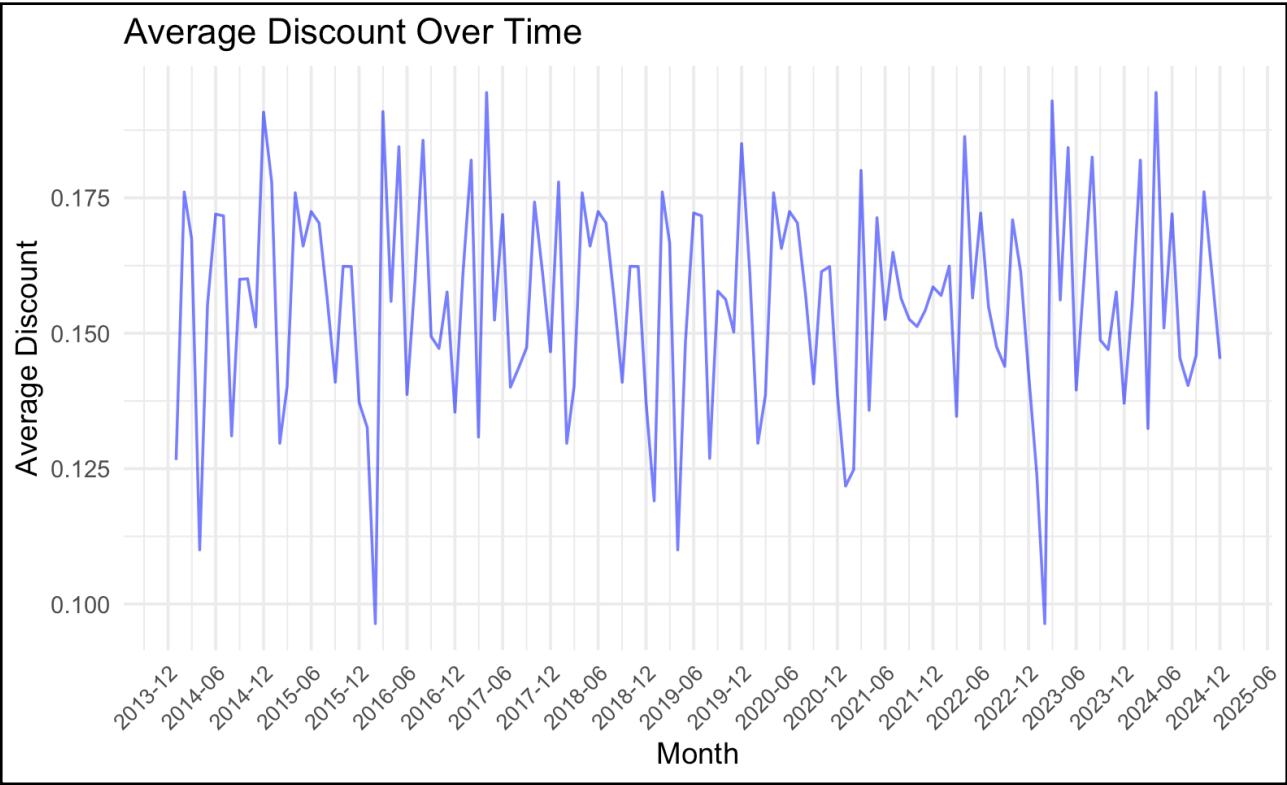
Shown in this graph, just because a product has higher sales does not mean that it provides higher profits. Therefore, I will show you what factor affects the profit, which might bring a massive business insight. Understanding these factors can help the business make informed decisions about pricing and discount strategies.



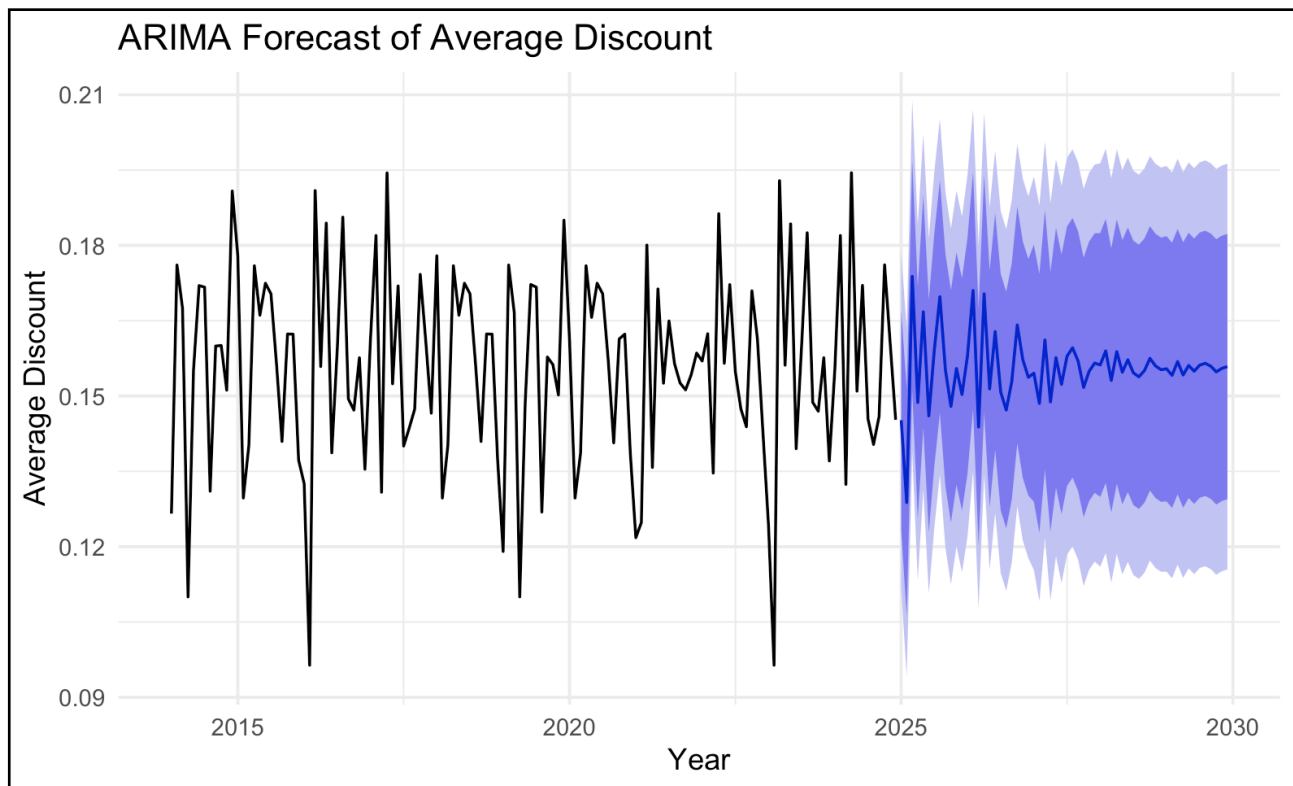
In this decision tree plot, you can see that even though I used multiple columns to predict whether the profit is positive or negative, the only factor that matters is if the discount is greater than 25 percent. This means that if a product's discount is greater than 25 percent, it has a very high chance of having a negative profit, which should be avoided. This insight can guide the business in setting optimal discount rates to maximize profit margins.



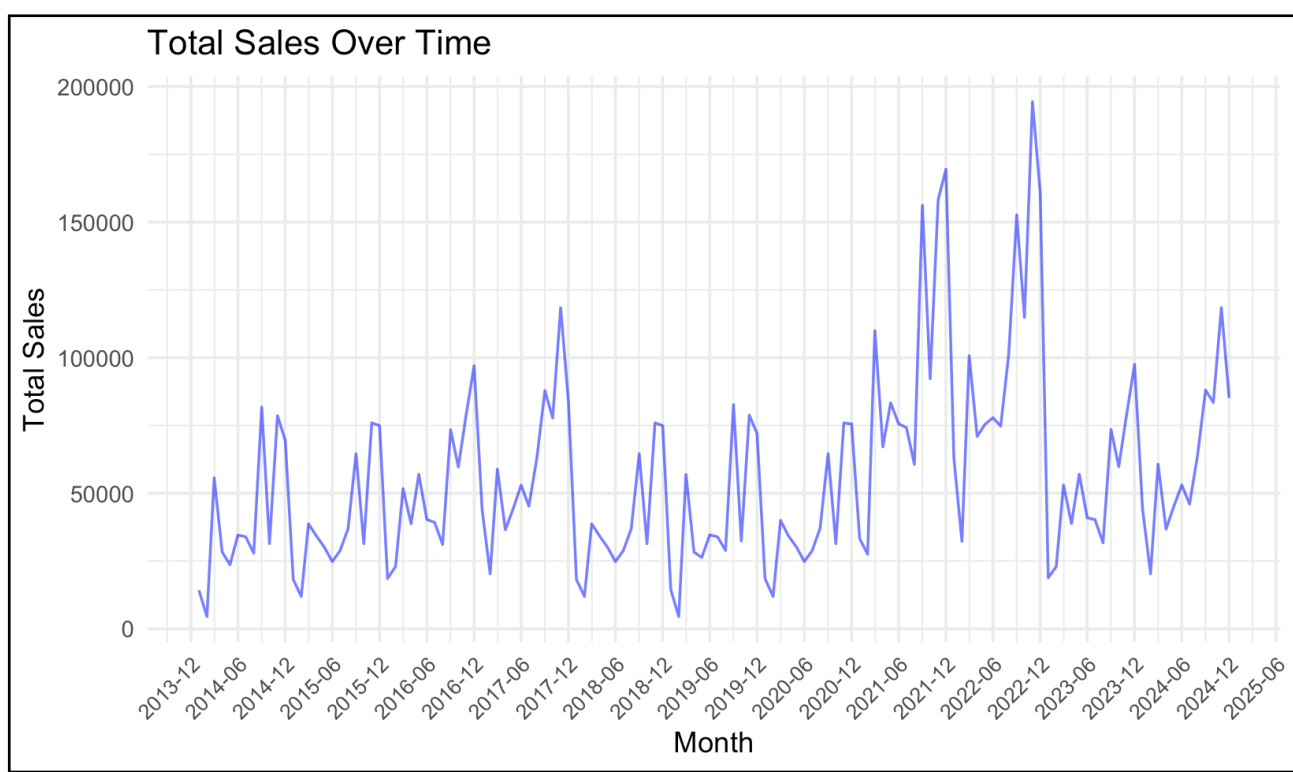
To prove that this analysis is reliable, the confusion matrix shows that the model has a 93.78% accuracy, which is pretty high. The model's sensitivity is 69.30%, meaning it correctly identifies 69.30% of the negative profit cases. The specificity is very high at 99.64%, indicating it almost always correctly identifies positive profit cases. The positive predictive value is 97.86%, meaning when the model predicts a negative profit, it is correct 97.86% of the time. The negative predictive value is 93.13%, meaning when the model predicts a positive profit, it is correct 93.13% of the time. The balanced accuracy, which is the average of sensitivity and specificity, is 84.47%. This shows the model is reliable in distinguishing between positive and negative profits. The high specificity and positive predictive value indicate that the model is particularly effective in identifying profitable products, which is crucial for strategic decision-making.



This graph shows the average discount over time by month from 2014 to 2024. It shows that the average discount varies from month to month, indicating some degree of fluctuation in discount rates offered by the superstore. However, this graph seems it is stationary which can apply for ARIMA to forecast the average discount for the next 5 years. The ARIMA model, which stands for AutoRegressive Integrated Moving Average, requires stationarity to make accurate forecasts. This characteristic makes ARIMA a powerful tool for predicting stable trends in time series data.

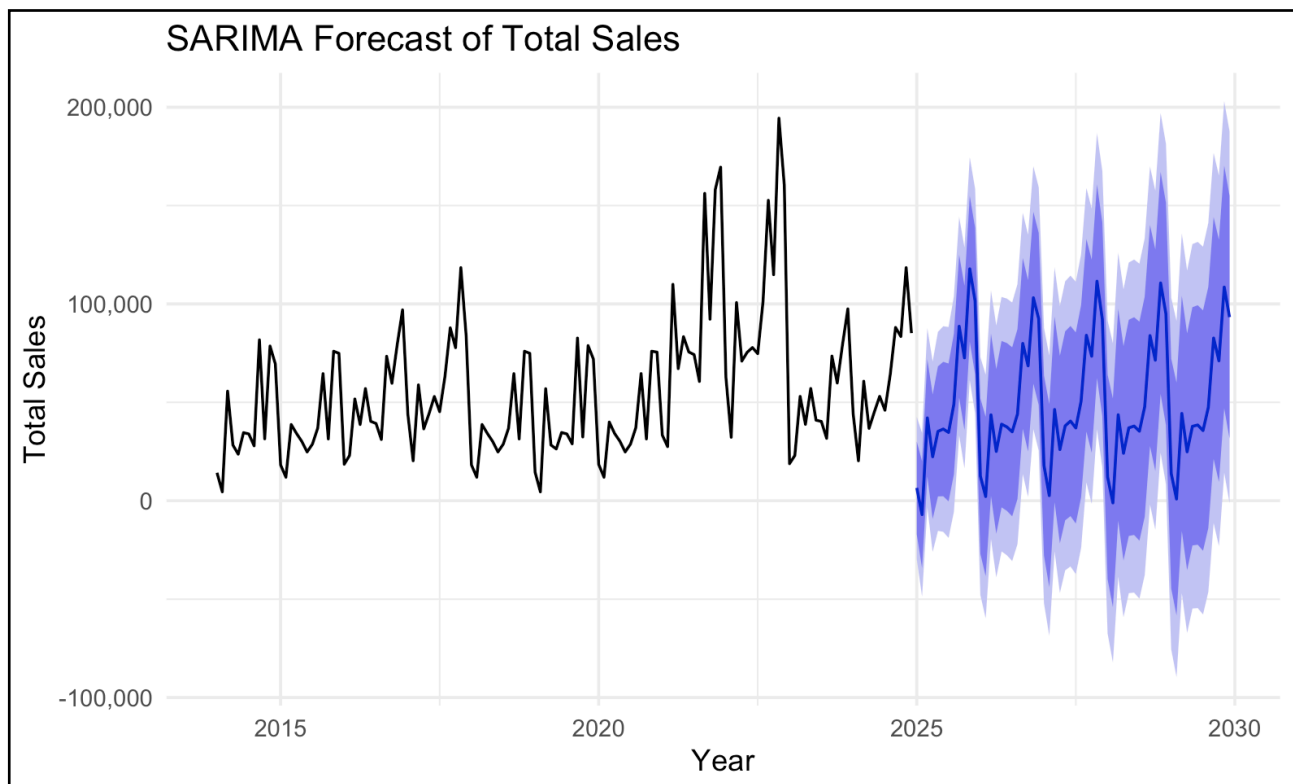


After forecasting the average discount, we observe that the forecast predicts a relatively stable trend in the average discount rates for the next 5 years, as shown by the blue shaded area representing the confidence intervals. This forecast can help the business plan its pricing and discount strategies more effectively, ensuring that discount rates are set at optimal levels to maintain profitability without overly reducing prices. By understanding the expected trends in discount rates, the business can better align its promotional activities and inventory management to the anticipated market conditions. This proactive approach ensures that the business remains competitive and responsive to market dynamics.



This graph shows the total sales over time by month from 2014 to 2024. We can see that total sales show a clear seasonal pattern, with peaks and troughs repeating each year. However, there is an anomaly from 2021 to 2022 where the pattern looks different, with increased volatility likely due to the impacts of COVID-19. Despite this period of greater volatility, the overall yearly pattern remains evident.

Given this clear seasonality, I chose to use the SARIMA (Seasonal AutoRegressive Integrated Moving Average) model to forecast total sales for the next 5 years. SARIMA is well-suited for time series data with seasonal patterns, making it ideal for this sales data that shows repeating annual trends. Analyzing these seasonal patterns helps in understanding consumer behavior and market trends over different periods.



After applying the SARIMA model, the forecast shows a continuation of the seasonal pattern, with expected peaks and troughs over the next 5 years. The blue shaded area represents the confidence intervals, indicating the range within which future sales are likely to fall. This forecast can help the business plan for seasonal fluctuations in sales, allowing for better inventory management, staffing, and promotional activities. By anticipating periods of high and low sales, the business can optimize its operations and marketing efforts to maximize revenue and efficiency. This foresight enables the business to allocate resources more effectively and enhance overall performance.

In conclusion, this analysis provides valuable insights into the sales, profit, and discount patterns of the superstore from 2014 to 2024. By examining the data, we identified key factors affecting profitability, such as the impact of high discounts on profit margins. The decision tree model highlighted the significance of maintaining discounts below 25 percent to avoid negative profits. Furthermore, the ARIMA and SARIMA models allowed us to forecast future trends in average discounts and total sales, respectively. These forecasts enable the business to plan strategically for the next 5 years, optimizing pricing strategies, inventory management, and promotional activities. Overall, this comprehensive analysis equips the business with the knowledge to make informed decisions, enhance profitability, and achieve sustainable growth.

Reference

Patel, M. (2024). Superstore Dataset 2014-2024. Kaggle.
<https://www.kaggle.com/datasets/mananapatel99/superstore-dataset-2014-2024>

OpenAI. (2024). ChatGPT-4 (May 2024 version) [Large language model].
<https://www.openai.com>

Appendix

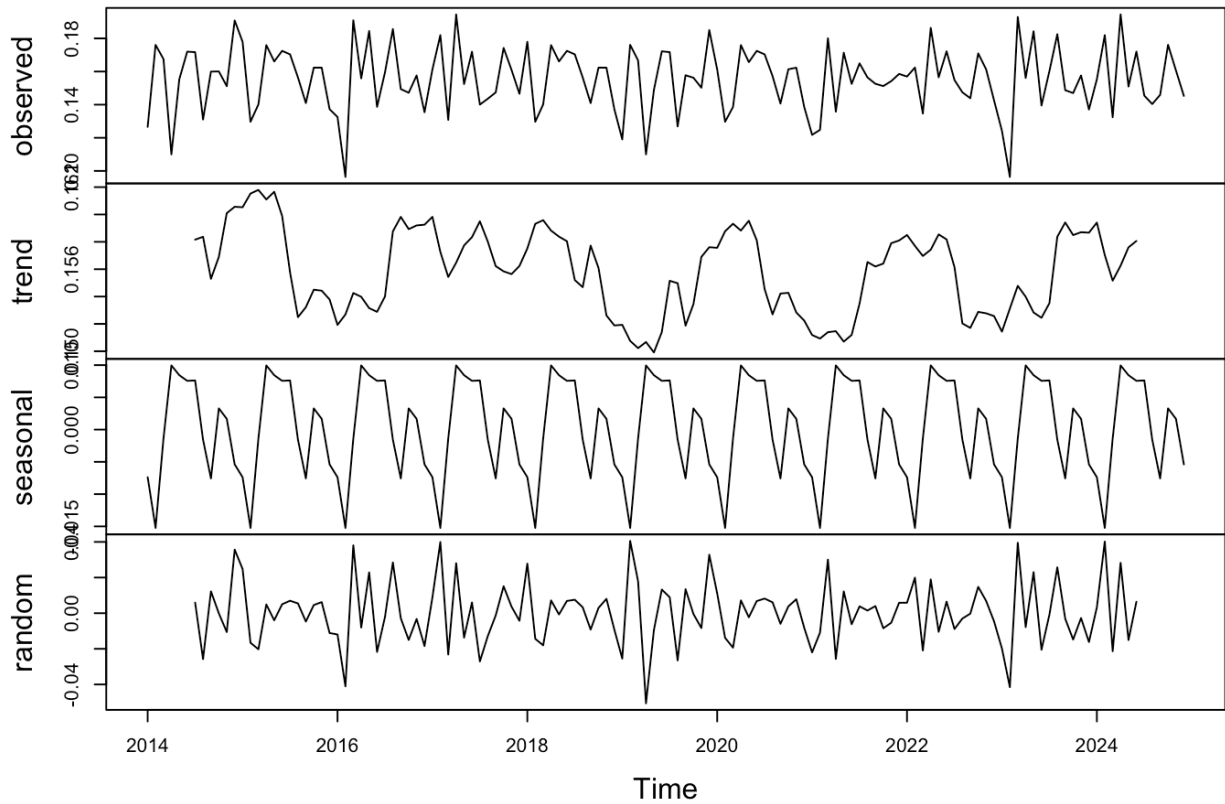
```
1 library(readr)
2 library(dplyr)
3 library(rpart)
4 library(rpart.plot)
5 library(caret)
6 library(corrplot)
7 library(ggplot2)
8 library(lubridate)
9 library(forecast)
10
11 # Read dataset
12 dataset <- read_csv("data/dataset.csv")
13
14 #####
15 # Overview Dataset
16 #####
17
18 # Remove NA values
19 dataset <- dataset %>% filter(!is.na(Category) & !is.na(`Sub-Category`))
20
21 # Bar plot of Sales by Category (transposed)
22 category_sales <- dataset %>%
23   group_by(Category) %>%
24   summarise(Total_Sales = sum(Sales, na.rm = TRUE))
25
26 ggplot(category_sales, aes(x = Total_Sales, y = reorder(Category, Total_Sales), fill = Category)) +
27   geom_bar(stat = "identity") +
28   theme_minimal() +
29   labs(title = "Total Sales by Category", x = "Total Sales", y = "Category") +
30   theme(legend.position = "none") +
31   scale_x_continuous(labels = scales::comma)
32
33 # Bar plot of Sales by Sub-Category (transposed)
34 subcategory_sales <- dataset %>%
35   group_by(`Sub-Category`) %>%
36   summarise(Total_Sales = sum(Sales, na.rm = TRUE))
37
38 ggplot(subcategory_sales, aes(x = Total_Sales, y = reorder(`Sub-Category`, Total_Sales), fill = `Sub-Category`)) +
39   geom_bar(stat = "identity") +
40   theme_minimal() +
41   labs(title = "Total Sales by Sub-Category", x = "Total Sales", y = "Sub-Category") +
42   theme(legend.position = "none") +
43   scale_x_continuous(labels = scales::comma)
44
45 # Scatter plot of Sales vs. Profit with y = 0 line
46 ggplot(dataset, aes(x = Sales, y = Profit)) +
47   geom_point(alpha = 0.6) +
48   geom_hline(yintercept = 0, color = "red", linetype = "dashed") +
49   theme_minimal() +
50   labs(title = "Sales vs. Profit", x = "Sales", y = "Profit") +
51   scale_x_continuous(labels = scales::comma) +
52   scale_y_continuous(labels = scales::comma)
53
54 #####
55 # Decision Tree Prediction
56 #####
57
58
59 # Create a binary target variable for Profit (Positive or Negative)
60 dataset <- dataset %>%
61   mutate(Profit_Category = ifelse(Profit > 0, "Positive", "Negative"))
62
63 # Convert categorical variables to factors
64 dataset$Category <- as.factor(dataset$Category)
65 dataset$`Sub-Category` <- as.factor(dataset$`Sub-Category`)
66 dataset$`Ship Mode` <- as.factor(dataset$`Ship Mode`)
67 dataset$Segment <- as.factor(dataset$Segment)
68 dataset$Region <- as.factor(dataset$Region)
69 dataset$Profit_Category <- as.factor(dataset$Profit_Category)
70
71 # Select features and target for the decision tree
72 features <- dataset %>% select(Sales, Quantity, Discount, Category, `Sub-Category`, `Ship Mode`, Segment, Region)
73 target <- dataset$Profit_Category
74
```

```

75 # Split the data into training and testing sets
76 set.seed(42)
77 train_index <- createDataPartition(target, p = 0.8, list = FALSE)
78 train_data <- dataset[train_index, ]
79 test_data <- dataset[-train_index, ]
80
81 # Train a decision tree model
82 decision_tree <- rpart(Profit_Category ~ Sales + Quantity + Discount + Category + `Sub-Category` + `Ship Mode` + Segment + Region,
83   data = train_data, method = "class")
84
85 # Predict on the test set
86 predictions <- predict(decision_tree, test_data, type = "class")
87
88 # Generate confusion matrix
89 confusion_matrix <- confusionMatrix(predictions, test_data$Profit_Category)
174
175
176 # Decomposing the time series for Average Discount
177 discount_decomp <- decompose(ts_discount)
178 plot(discount_decomp)
179 title(main = "-Average Discount-")
180
181 # Decomposing the time series for Total Sales
182 sales_decomp <- decompose(ts_sales)
183 plot(sales_decomp)
184 title(main = "-Total Sales-")
185
186 # Convert Order Date to Date format
187 dataset_discount$`Order Date` <- as.Date(dataset_discount$`Order Date`, format = "%Y/%m/%d")
188
189 # Aggregate discount by month
190 discount_over_time <- dataset_discount %>%
191   mutate(Month = floor_date(`Order Date`, "month")) %>%
192   group_by(Month) %>%
193   summarise(Average_Discount = mean(Discount, na.rm = TRUE))
194
195 # Prepare the time series data for Discount
196 ts_discount <- ts(discount_over_time$Average_Discount, start = c(year(min(discount_over_time$Month)), month(min(discount_over_time$Month))), frequency = 12)
197
198 # Plot the original discount data
199 ggplot(discount_over_time, aes(x = Month, y = Average_Discount)) +
200   geom_line(color = "blue", alpha = 0.6) +
201   theme_minimal() +
202   labs(title = "Average Discount Over Time", x = "Month", y = "Average Discount") +
203   scale_x_date(date_labels = "%Y-%m", date_breaks = "6 months") +
204   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 8))
205
206 # Fit ARIMA model for Discount
207 arima_discount_model <- auto.arima(ts_discount)
208
209 # Forecast for the next 5 years (60 months) for Discount
210 forecast_horizon <- 5 * 12
211 arima_discount_forecast <- forecast(arima_discount_model, h = forecast_horizon)
212
213 # Plot the forecasted values for Discount
214 autoplot(arima_discount_forecast) +
215   theme_minimal() +
216   labs(title = "ARIMA Forecast of Average Discount", x = "Year", y = "Average Discount") +
217   scale_y_continuous(labels = scales::comma)
218
219 #####
220 # SARIMA Forecasting for Sales
221 #####
222
223 # Remove NA values for Sales
224 dataset_sales <- dataset %>% filter(!is.na(Sales))
225
226 # Convert Order Date to Date format
227 dataset_sales$`Order Date` <- as.Date(dataset_sales$`Order Date`, format = "%Y/%m/%d")
228
229 # Aggregate sales by month
230 sales_over_time <- dataset_sales %>%
231   mutate(Month = floor_date(`Order Date`, "month")) %>%
232   group_by(Month) %>%
233   summarise(Total_Sales = sum(Sales, na.rm = TRUE))
234
235 # Prepare the time series data for Sales
236 ts_sales <- ts(sales_over_time$Total_Sales, start = c(year(min(sales_over_time$Month)), month(min(sales_over_time$Month))), frequency = 12)
237
238 # Plot the original sales data
239 ggplot(sales_over_time, aes(x = Month, y = Total_Sales)) +
240   geom_line(color = "blue", alpha = 0.6) +
241   theme_minimal() +
242   labs(title = "Total Sales Over Time", x = "Month", y = "Total Sales") +
243   scale_x_date(date_labels = "%Y-%m", date_breaks = "6 months") +
244   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 8))
245
246 # Fit SARIMA model for Sales
247 sarima_sales_model <- auto.arima(ts_sales, seasonal = TRUE)
248
249 # Forecast for the next 5 years (60 months) for Sales
250 sarima_forecast_horizon <- 5 * 12
251 sarima_sales_forecast <- forecast(sarima_sales_model, h = sarima_forecast_horizon)
252
253 autoplot(sarima_sales_forecast) +
254   theme_minimal() +
255   labs(title = "SARIMA Forecast of Total Sales", x = "Year", y = "Total Sales") +
256   scale_y_continuous(labels = scales::comma)

```

Decomposition of additive time series
-Average Discount-



Decomposition of additive time series
-Total Sales-

