



## **FORECASTING AND PREDICTING THE FUTURE USING DATA**

### **A3: BUSINESS INSIGHT REPORT**

Professor: Thomas Kurnicki

Author: Gabriel Martinelli

## *Abstract*

This report analyses a dataset emulating the sales of an e-commerce store managed by a popular American influencer over the course of a year. The dataset includes various variables, ranging from the buyer's age and gender to shipping expenses charges to the customer and the ratings received by the site.

The analysis develops along three main lines: predictive, qualitative, and forecasting. Predictive models were built to identify which drivers led consumers to leave a 5-star review or lower. However, it emerged that none of the available variables proved significant in predicting the business outcome.

This result highlights that the collected variables are not at all adequate to predict what rating the consumer might give. Consequently, if the available variables are inconclusive, it becomes necessary to investigate and understand which might be relevant.

The qualitative analysis therefore aimed to define which variables future models should focus on. From an analysis of a random sample of 50 reviews rated 1 or 2 stars, and a sample of 50 5-star reviews, it emerged that quality, design, and delivery times are the factors that most affect customer satisfaction.

To complete the dataset exploration, forecasting models were built to estimate future growth rates of overall sales and sales by category. Once again, the available data proved insignificant. No seasonality was detected, and the confidence intervals were wide and centered around zero. Despite the limited relevance of the data, the analysis still generated valuable business insights to be considered in future planning. Data always carries a message, and even when results differ from expectations, an in-depth analysis will always prove useful.

## *Table of Contents*

<i>Introduction</i> .....	4
<i>Predictive Analysis</i> .....	5
<i>Predictive Analysis Insights</i> .....	8
<i>Forecast Analysis</i> .....	9
<i>Forecast Analysis Insights</i> .....	11
<i>Conclusion</i> .....	12
<i>Appendix</i> .....	13
<i>References</i> .....	19

## ***Introduction***

This paper presents an analysis of a dataset covering the global merchandise sales of a well-known American influencer from November 2023 to November 2024.

The main focus of the analysis is to predict, based on order characteristics, when a buyer is likely to leave a 5-star review, with the goal of identifying the most impactful aspects and enhancing customer satisfaction.

At the same time, forecasting models were developed to estimate future sales growth, both overall and for each product category.

The analysis faced major setbacks related to the significance of the available data, which proved unable to capture either the drivers of customer satisfaction or the volatility of sales. Nevertheless, the study still provides insights into how the challenges encountered by the predictive and forecasting models can be addressed by fundamentally redefining the type of data that should be collected.

## Predictive Analysis

Before starting any analysis, cleaning the data is essential. For the dataset examined, the cleaning process focused on removing empty rows or those lacking an order date, eliminating duplicates based on the order ID, and converting categorical variables (Product Category, Buyer Gender, and International Shipping) into numerical variables. Converting these variables was necessary to normalize them alongside the others, making them comparable in subsequent analyses.

Once the cleaning was completed, the business outcome was defined: a successful outcome corresponded to an order that received a 5-star rating. Several predictive models were then built to determine whether, based on order characteristics, it was possible to predict the ratings consumers would give.

The first of these models was a Gini decision tree (figure 1), designed to uncover splits that

separate 5-star ratings from others based on variables such as shipping charges, sales price, and buyer demographics. However, to visualize the branches, the `cp` parameter had to be set extremely low, a clear sign that these variables could not meaningfully discriminate the target.

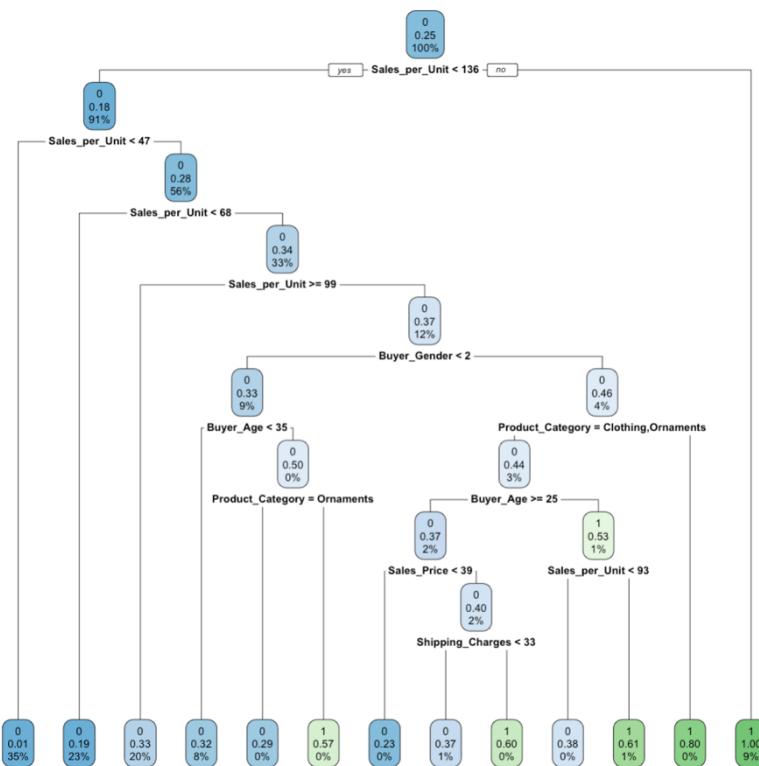


Figure 1: Gini Decision Tree, source: own elaboration

This inability was confirmed by the logistic regression (figure 2) fitted afterwards, which revealed that none of the variables considered in the model were significant. It is enough to note that the lowest p-value was 0.360, well above the typical threshold for significance.

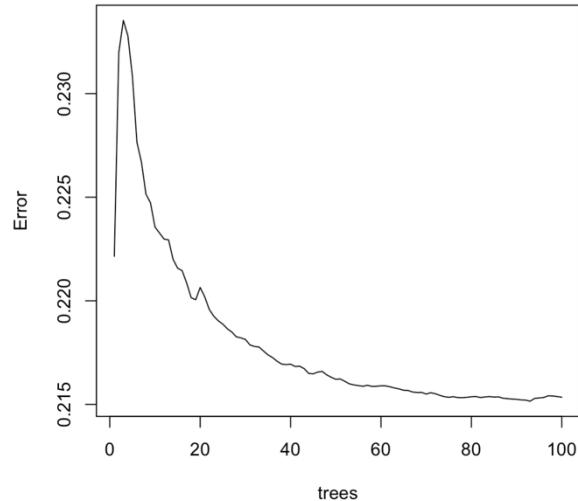
```

Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.1267988 0.2566997 -4.390 1.14e-05 ***
Product_Category 0.0504113 0.0701931 0.718 0.473
Buyer_Gender 0.0568188 0.0620955 0.915 0.360
Buyer_Age 0.0006603 0.0054504 0.121 0.904
International_Shipping -0.0081101 0.1337385 -0.061 0.952
Sales_Price 0.0009332 0.0013697 0.681 0.496
Shipping_Charges 0.0015383 0.0024672 0.623 0.533
Quantity 0.0174753 0.0260965 0.670 0.503
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Figure 2: Logistic Regression, source: own elaboration

To investigate more deeply the interactions between variables and the target, a random forest



(Figure 3) was then constructed. This approach undeniably reduced the error compared to the single decision tree, but it also quickly plateaued (around 70-80 trees) and settled at approximately 0.215, which remains relatively high.

Figure 3: Random Forest, source: own elaboration

The confusion matrix (Figure 4) further confirmed the models' inability to successfully predict

#### Confusion Matrix and Statistics

		Reference	
Prediction	0	1	
0	1025	408	Accuracy : 0.7025
1	32	14	95% CI : (0.6785, 0.7257)

No Information Rate : 0.7147  
P-Value [Acc > NIR] : 0.8565

the target, showing an accuracy of 70.25%, which is actually lower than simply always predicting the dominant class (71.47%).

Figure 4: Confusion Matrix, source: own elaboration

Given the models' failure to predict customer ratings, it became necessary to understand which variables should instead be considered. For this purpose, a qualitative analysis was performed on two samples of 50 customer reviews each: the first containing 1-2 star reviews, and the second containing only 5-star reviews. This step was crucial to identify which variables had been overlooked during data collection. The blend of predictive and qualitative analysis thus provided a deeper understanding of what drives customer satisfaction, and what does not, and consequently, where future models should focus.

### ***Predictive Analysis Insights***

The combination of the predictive analysis and qualitative analysis certainly provided important insights.

The models clearly demonstrated how the variables collected in the dataset, such as buyer age, shipping charges, product category, and sales price, were irrelevant at predicting 5-star reviews. The decision tree struggled to create meaningful branches, the logistic regression revealed that none of the variables had a significant p-value, and while the random forest did reduce error, it did so insufficiently and quickly reached a plateau.

However, these initially discouraging results prompted a deeper qualitative analysis. By examining two samples, one containing negative reviews and the other positive, it became possible to identify which variables truly mattered for customer satisfaction.

In the sample containing the negative reviews, 32% of customers complained about delivery, 16% about product quality, and another 16% about design. On the other hand, in the sample containing the positive reviews, 32% were satisfied with quality, 26% with delivery, and 16% with design.

This suggests that future data collection should focus on variables like perceived quality and delivery performance to better train models to predict the target outcome.

This first phase of analysis teaches us that sometimes understanding what does not drive outcomes is just as valuable. Indeed, knowing what fails to predict satisfaction provides a clear guide for where to redirect further analysis.

## Forecast Analysis

The forecasting analysis focused on estimating future sales trends by applying ARIMA models to the available data.

The first step was to aggregate monthly sales to create time series, which were then examined through ACF and PACF plots (figures 5 and 6).

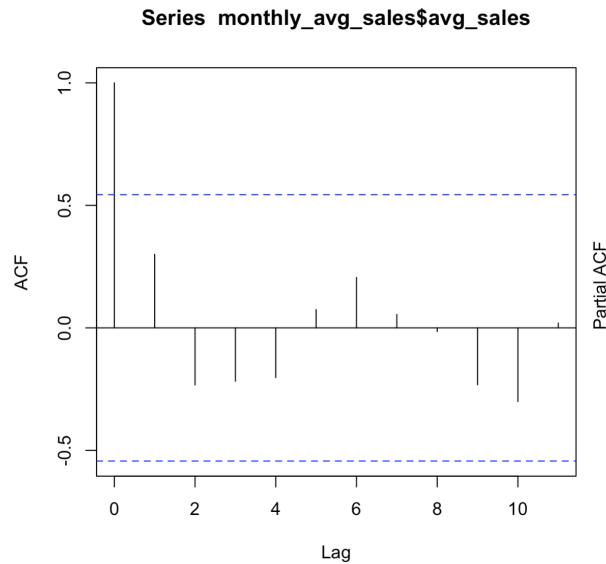


Figure 5: ACF Plot, source: own elaboration

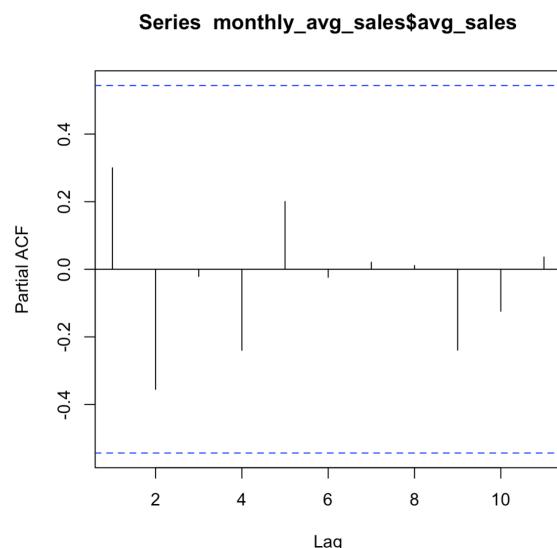


Figure 6: PACF Plot, source: own elaboration

These plots showed weak autocorrelation and an almost complete absence of seasonal patterns, indicating that sales fluctuated randomly around their mean.

As a result, the `auto.arima` function employed selected an ARIMA(0,0,0), effectively treating the series as white noise. The resulting forecast (Figure 7) thus settled around the historical mean, with extremely wide confidence intervals.

**Forecasts from ARIMA(0,0,0) with non-zero mean**

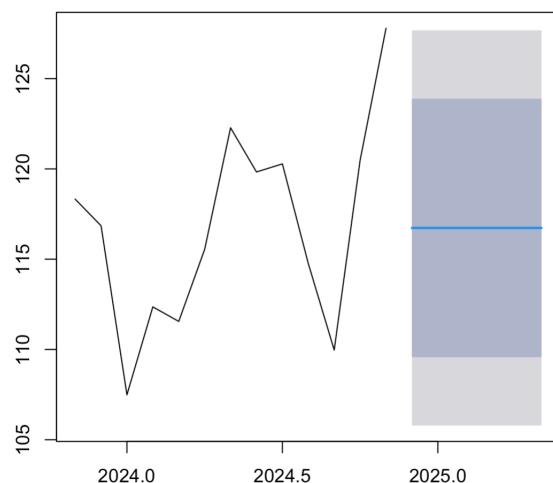
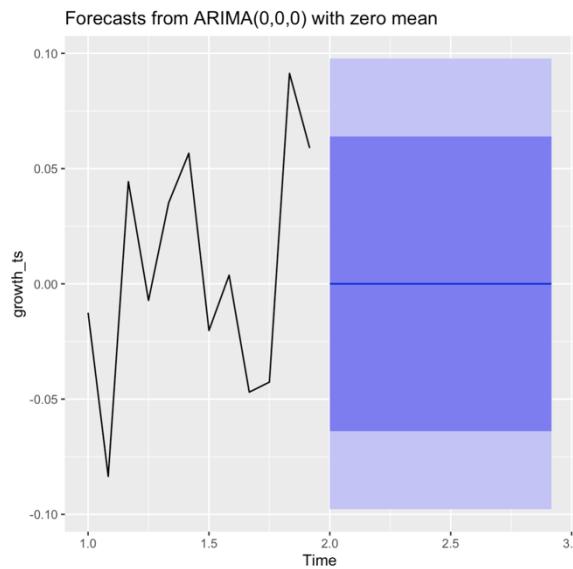


Figure 7: Sales Forecast, source: own elaboration

Quoting Professor Thomas Kurnicki from one of his lectures in the “Forecasting and Predicting the Future using Data” course at HULT International Business School: “*If your boss asks you to forecast sales, you should reply that, after statistical adjustments, you can forecast the growth rates of sales.*”

Consequently, the next phase of the analysis focused on forecasting sales growth rates (Figure 8).



However, once again, auto.arima picked the (0,0,0) specification, forecasting future growth at zero, with wide intervals that indicated an equal likelihood of positive or negative shifts.

Figure 8: Sales’ Growth Rate Forecast, source: own elaboration

The analysis then turned to forecasting both sales and their growth rates for each individual product category. Yet the pattern remained unchanged. None of the categories exhibited significant autocorrelation or seasonality (Figures 9,10 and 11).

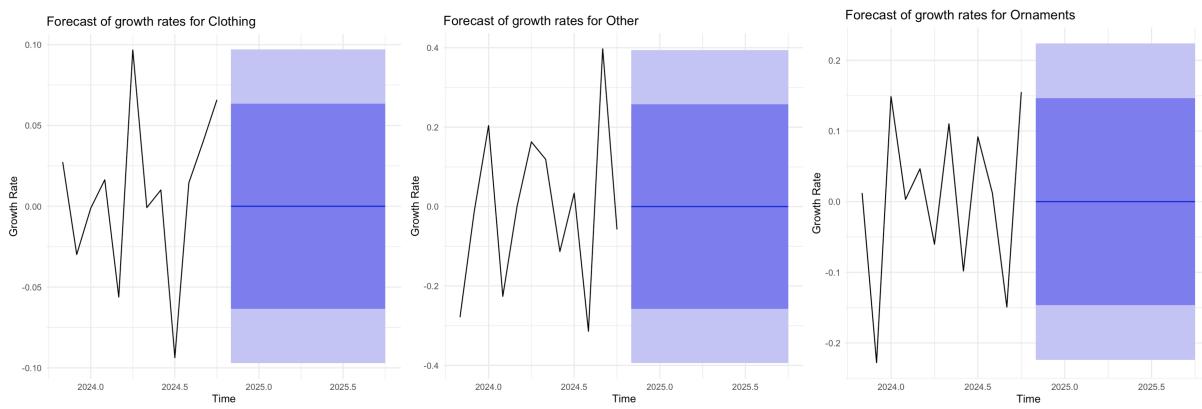


Figure 9, 10 and 11: Clothing, Other, Ornaments Growth Rates Forecast, source: own elaboration

Overall, the forecasting analysis highlighted the volatility and unpredictability of sales. Much like the predictive analysis, this suggested the need to reconsider both the approach and the type of data being collected to enable more meaningful forecasting in the future.

### ***Forecast Analysis Insights***

The forecasting analysis highlighted several important aspects regarding the sales dynamics captured by the dataset.

The forecasts consistently settled around the historical mean with wide confidence intervals, indicating pronounced volatility where future performance is equally likely to fluctuate above or below this average. The near-total absence of trends, seasonality, and autocorrelation, both in overall sales and in the sales of each individual category, suggests not only that the time series is too short but also that key variables may have been overlooked.

The same interpretation applies to the forecasts of growth rates, leaving the e-commerce business exposed to uncertainty. However, despite these discouraging findings, two key insights stand out.

The first is that to gain a clearer picture of the e-commerce trajectory, more time is needed, implying a longer observation period to build a sufficiently rich time series. The second is that, given this is an e-commerce operation selling an influencer's merchandise, it is misleading to estimate its performance without taking into account the influencer's social media activity.

These forecasts, when analysed alone, lose much of their meaning if not compared with the social campaigns being run. In fact, linking sales trends to the number of promotional posts could reveal patterns that would otherwise remain invisible when only considering the transactional data contained in this dataset.

## ***Conclusion***

In this paper, the combination of predictive, qualitative, and forecasting analyses was used to examine the annual sales performance of an e-commerce business selling merchandise for an American influencer.

All predictive models consistently demonstrated that the variables present in the dataset were unable to predict customer satisfaction, represented by the ratings customers gave after placing their orders. Indeed, the p-values for these variables were all insignificant, and the confusion matrix showed that the model could not even outperform naive guesses.

The forecasting analysis further revealed the absence of trends, seasonality, or autocorrelation. Forecast models repeatedly settled around historical means with wide confidence intervals, highlighting the volatility of the business and their inability to produce accurate estimates given the available data.

However, the qualitative analysis of customer reviews pointed to critical aspects that the structured data missed: quality, design, and delivery emerged as the real drivers of customer satisfaction. These factors should therefore be prioritized in future data collection and incorporated into upcoming predictive and forecasting models, as they have the potential to uncover hidden patterns.

Although the existing data proved incapable of predicting the business outcome, it still revealed a key message. It is essential to reshape the data collection strategy by integrating insights on sentiment, product quality and delivery satisfaction, and the influencer's social activity. Doing so will enable the building of more accurate models in the future.

## Appendix

```

1 library(readr)
2 library(dplyr)
3 library(lubridate)
4 library(stringr)
5 library(rpart)
6 library(rpart.plot)
7 library(randomForest)
8 library(caret)
9 library(dplyr)
10 library(caret)
11 library(Forecast)
12 library(ggplot2)
13 library(stringr)
14 library(tseries)
15 library(rugarch)
16
17 #Importing the dataset
18 merch_sales <- read_csv("/Users/gab/Desktop/FORECASTING AND PREDICTING THE FUTURE USING DATA/INDIVIDUAL ASSIGNMENT/merch_sales.csv")
19
20 #####CLEANING/WRANGLING#####
21 #####
22 #####
23
24 #removing rows that have no order date
25 rows_before <- nrow(merch_sales)
26 merch_sales_clean <- merch_sales %>%
27   filter(!is.na("Order Date"))
28 rows_after <- nrow(merch_sales_clean)
29 merch_sales <- merch_sales_clean
30
31 #checking NA
32 anyNA(merch_sales)
33 #changing Order Date to date obj
34 merch_sales$"Order Date" <- as.Date(merch_sales$"Order Date")
35 #removing duplicates rows
36 merch_sales <- merch_sales[!duplicated(merch_sales$"Order ID"), ]
37
38 #changing chr columns to num in order to normalise and compare them later
39 merch_sales <- merch_sales %>%
40   mutate(`Product Category` = case_when(
41     `Product Category` == "Clothing" ~ 1,
42     `Product Category` == "Ornaments" ~ 2,
43     `Product Category` == "Other" ~ 3
44   ))
45
46 merch_sales <- merch_sales %>%
47   mutate(`Buyer Gender` = case_when(
48     `Buyer Gender` == "Male" ~ 1,
49     `Buyer Gender` == "Female" ~ 2,
50   ))
51
52 merch_sales <- merch_sales %>%
53   mutate(`International Shipping` = case_when(
54     `International Shipping` == "No" ~ 0,
55     `International Shipping` == "Yes" ~ 1,
56   ))
57
58 #removing the spaces and changing them with "_"
59 names(merch_sales) <- str_replace_all(names(merch_sales), " ", "_")
60
61 #normalise the variables to compare them
62
63 rescale <- function(x){
64   min_max <- (x-min(x))/(max(x)-min(x))
65   return(min_max)
66 }
67
68 merch_sales$Product_Category_norm <- rescale(x=merch_sales$Product_Category)
69 merch_sales$Buyer_Gender_norm <- rescale(x=merch_sales$Buyer_Gender)
70 merch_sales$Buyer_Age_norm <- rescale(x=merch_sales$Buyer_Age)
71 merch_sales$Sales_Price_norm <- rescale(x=merch_sales$Sales_Price)
72 merch_sales$Shipping_Charges_norm <- rescale(x=merch_sales$Shipping_Charges)
73 merch_sales$Sales_per_Unit_norm <- rescale(x=merch_sales$Sales_per_Unit)
74 merch_sales$Quantity_norm <- rescale(x=merch_sales$Quantity)
75 merch_sales$Total_Sales_norm <- rescale(x=merch_sales$Total_Sales)
76
77 #####PREDICTIVE ANALYSIS ON RATINGS#####
78 #####
79 #####
80
81 #define business outcome
82 business_suc <- merch_sales[which(merch_sales$Rating == 5), ]
83 business_fail <- merch_sales[which(merch_sales$Rating < 5), ]
84
85 merch_sales$businessoutcome <- c()
86 merch_sales[which(merch_sales$Rating == 5), c('businessoutcome')] <- 1
87 merch_sales[which(merch_sales$Rating < 5), c('businessoutcome')] <- 0
88
89 #define random sample
90 indx <- sample(x=1:nrow(merch_sales), size=0.8*nrow(merch_sales))
91
92 merch_sales_train <- merch_sales[ indx, ]
93 merch_sales_test <- merch_sales[-indx, ]
94
95 #build the decision tree
96 my_tree <- rpart(businessoutcome~Product_Category + Buyer_Gender + Buyer_Age + International_Shipping +
97                   Sales_Price + Shipping_Charges + Quantity,
98                   data=merch_sales_train, method="class", cp=0.0007252)

```

```

99 rpart.plot(my_tree)
100
101 #building a Logistic regression
102 my_logit_norm <- glm(businessoutcome~Product_Category + Buyer_Gender + Buyer_Age + International_Shipping +
103                         Sales_Price + Shipping_Charges + Quantity,
104                         data=merch_sales_train, family = "binomial")
105 summary(my_logit_norm)
106
107 #predicting the test data
108 tree_pred <- predict(my_tree, merch_sales_test)
109
110 #building a confusion matrix
111 confusionMatrix(data = as.factor(as.numeric(tree_pred[,2]>0.5)),
112                   reference = as.factor(as.numeric(merch_sales_test$businessoutcome)))
113
114 #building a random forest
115 my_forest <- randomForest(
116   x = merch_sales_train[, c("Product_Category", "Buyer_Gender", "Buyer_Age", "International_Shipping",
117                           "Sales_Price", "Shipping_Charges", "Quantity")],
118   y = merch_sales_train$businessoutcome,
119   ntree = 100)
120 forest_pred <- predict(my_forest, merch_sales_test)
121 plot(my_forest)
122 varImpPlot(my_forest)
123
124 #checking the reviews, is there any relationship between 1 or 2 ratings and delivery issues?
125
126 filtered_rows <- merch_sales[merch_sales$Rating >= 1 & merch_sales$Rating <= 2, ]
127 sample_50 <- filtered_rows[sample(1:nrow(filtered_rows), 50), ]
128 rating_review_table <- sample_50[, c("Rating", "Review")]
129 View(rating_review_table)
130
131 #this table created with a sample of 50 orders that received a rating between 1 and 2 highlights that
132 #in the 32% of the orders there where, according to the reviews, problems with the quality, in the 16% of the orders
133 #problems with the delivery and in another 16% problems with an underwhelming design.
134 #these insights could be the groud for further reasoning.
135
136 #checking the reviews, is there any relationship between 5 ratings and quality, design e delivery?
137
138 filtered_rows_5stars <- merch_sales[merch_sales$Rating == 5 , ]
139 sample_50_5stars <- filtered_rows_5stars[sample(1:nrow(filtered_rows_5stars), 50), ]
140 rating_5_review_table <- sample_50_5stars[, c("Rating", "Review")]
141 View(rating_5_review_table)
142
143 #this table created with a sample of 50 orders that received a 5-star rating highlights that
144 #in the 32% of the orders the customers were satisfied about the delivery, in the 26% about the quality and in the
145 #16% about the design of the product.
146
147 #####

```

```

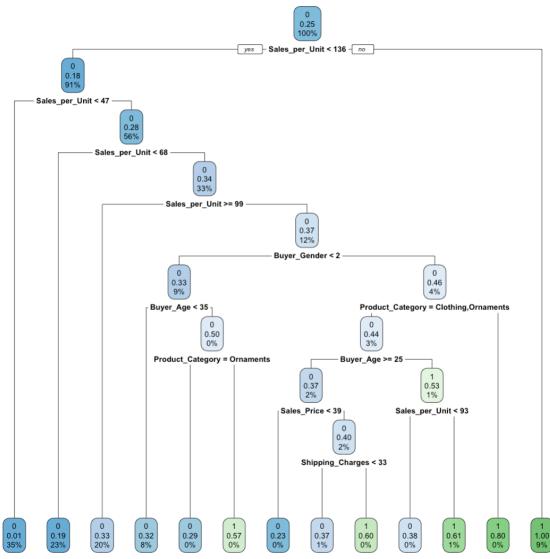
148 ######PREDICTIVE ANALYSIS ON SALES#####
149 #####DO NOT INCLUDE IN THE ASSIGNMENT#####
150 #####
151
152 summary(merch_sales$Total_Sales)
153 |
154 #define business outcome
155 business_suc_sales <- merch_sales[which(merch_sales$Total_Sales >= 137), ]
156 business_fail_sales <- merch_sales[which(merch_sales$Total_Sales < 137), ]
157
158 merch_sales$businessoutcome_sales <- c()
159 merch_sales[which(merch_sales$Total_Sales >= 137), c('businessoutcome_sales')] <- 1
160 merch_sales[which(merch_sales$Total_Sales < 137), c('businessoutcome_sales')] <- 0
161
162 #define random sample
163 indx <- sample(x=1:nrow(merch_sales), size=0.8*nrow(merch_sales))
164
165 total_merch_sales_train <- merch_sales[indx, ]
166 total_merch_sales_test <- merch_sales[-indx, ]
167
168 #build the decision tree
169 my_tree_sales <- rpart(businessoutcome_sales~Product_Category + Buyer_Gender + Buyer_Age + International_Shipping +
170                         Sales_Price + Shipping_Charges + Sales_per_Unit,
171                         data=total_merch_sales_train, method="class", cp=0.001)
172 rpart.plot(my_tree_sales)
173
174 #building a Logistic regression
175 my_logit_sales <- glm(businessoutcome~Product_Category + Buyer_Gender + Buyer_Age + International_Shipping +
176                         Sales_Price + Shipping_Charges + Sales_per_Unit,
177                         data=total_merch_sales_train, family = "binomial")
178 summary(my_logit_sales)
179
180 #predicting the test data
181 tree_sales_pred <- predict(my_tree_sales, total_merch_sales_test)
182
183 #building a confusion matrix
184 confusionMatrix(data = as.factor(as.numeric(tree_sales_pred[,2]>0.5)),
185                   reference = as.factor(as.numeric(total_merch_sales_test$businessoutcome_sales)))
186
187 #####
188 #####FORECAST ANALYSIS ON SALES#####
189 #####
190 #####
191
192 converted_order_dates <- as.Date(merch_sales$Order_Date, format="%Y-%m-%d")
193 sum(is.na(converted_order_dates))
194 min(converted_order_dates)
195 max(converted_order_dates)
196

```

```

197 #plot the average sales trend
198 merch_sales <- merch_sales %>%
199   mutate(Month = floor_date(Order_Date, "month"))
200
201 monthly_avg_sales <- merch_sales %>%
202   group_by(Month) %>%
203   summarise(avg_sales = mean(Total_Sales, na.rm=TRUE))
204
205 ggplot(monthly_avg_sales, aes(x = Month, y = avg_sales)) +
206   geom_line() +
207   labs(title = "Monthly average sales trend",
208       x = "Months",
209       y = "Average sales per month") +
210   theme_minimal()
211
212 adf.test(monthly_avg_sales$avg_sales)
213
214 acf(monthly_avg_sales$avg_sales)
215 pacf(monthly_avg_sales$avg_sales)
216
217 #ARIMA Model
218
219 sales_ts <- ts(monthly_avg_sales$avg_sales, start = c(2023, 11), frequency = 12)
220 arima_model <- auto.arima(sales_ts)
221 summary(arima_model)
222
223 forecast_arima <- forecast(arima_model, h = 6)
224 plot(forecast_arima)
225
226 #forecasting the growth rates
227 monthly_avg_sales_growth <- monthly_avg_sales %>%
228   arrange(Month) %>%
229   mutate(growth_rate = c(NA, diff(log(avg_sales))))
230
231 growth_ts <- ts(na.omit(monthly_avg_sales_growth$growth_rate), frequency=12)
232 growth_model <- auto.arima(growth_ts)
233 growth_forecast <- forecast(growth_model, h=12)
234 autoplot(growth_forecast)
235
236 #####
237 #####FORECAST EACH CATEGORY SEPARATELY#####
238 #####
239
240 merch_sales <- merch_sales %>%
241   mutate(Product_Category = case_when(
242     Product_Category == 1 ~ "Clothing",
243     Product_Category == 2 ~ "Ornaments",
244     Product_Category == 3 ~ "Other",
245     TRUE ~ as.character(Product_Category)
246   ))
247
248 categories <- unique(merch_sales$Product_Category)
249 for (cat in categories) {
250   cat("\n--- Forecast for category:", cat, "---\n")
251   category_data <- merch_sales %>%
252     filter(Product_Category == cat) %>%
253     mutate(Month = floor_date(Order_Date, "month")) %>%
254     group_by(Month) %>%
255     summarise(avg_sales = mean(Total_Sales, na.rm=TRUE)) %>%
256     arrange(Month)
257
258 ts_data <- ts(category_data$avg_sales, frequency=12, start=c(2023,11))
259
260 model <- auto.arima(ts_data)
261 forecast_h <- forecast(model, h=6)
262
263 print(
264   autoplot(forecast_h) +
265   ggtitle(paste("Forecast of average monthly sales for", cat)) +
266   ylab("Average Sales") + xlab("Time") +
267   theme_minimal()
268 )
269 }
270
271 #forecasting the growth rates for each category
272
273 for (cat in unique(merch_sales$Product_Category)) {
274   cat_data <- merch_sales %>%
275     filter(Product_Category == cat) %>%
276     mutate(Month = floor_date(Order_Date, "month")) %>%
277     group_by(Month) %>%
278     summarise(avg_sales = mean(Total_Sales, na.rm=TRUE)) %>%
279     arrange(Month) %>%
280     mutate(growth_rate = c(NA, diff(log(avg_sales))))
281
282 growth_ts_category <- ts(na.omit(cat_data$growth_rate), frequency = 12, start = c(2023, 11))
283
284 growth_model_category <- auto.arima(growth_ts_category)
285 growth_forecast_category <- forecast(growth_model_category, h=12)
286
287 print(
288   autoplot(growth_forecast_category) +
289   ggtitle(paste("Forecast of growth rates for", cat)) +
290   ylab("Growth Rate") + xlab("Time") +
291   theme_minimal()
292 )
293 }
294

```



```
> summary(my_logit_norm)
```

```
Call:
glm(formula = businessoutcome ~ Product_Category + Buyer_Gender +
  Buyer_Age + International_Shipping + Sales_Price + Shipping_Charges +
  Quantity, family = "binomial", data = merch_sales_train)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.1267988	0.2566997	-4.390	1.14e-05 ***
Product_Category	0.0504113	0.0701931	0.718	0.473
Buyer_Gender	0.0568188	0.0620955	0.915	0.360
Buyer_Age	0.0006603	0.0054504	0.121	0.904
International_Shipping	-0.0081101	0.1337385	-0.061	0.952
Sales_Price	0.0009332	0.0013697	0.681	0.496
Shipping_Charges	0.0015383	0.0024672	0.623	0.533
Quantity	0.0174753	0.0260965	0.670	0.503

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 7222.3 on 5914 degrees of freedom  
Residual deviance: 7218.9 on 5907 degrees of freedom  
AIC: 7234.9

Number of Fisher Scoring iterations: 4

### my\_forest

```
> confusionMatrix(data = as.factor(as.numeric(tree_pred[,2]>0.5)),
+   reference = as.factor(as.numeric(merch_sales_test$businessoutcome)))
Confusion Matrix and Statistics

Reference
Prediction    0      1
      0 1025  408
      1    32   14

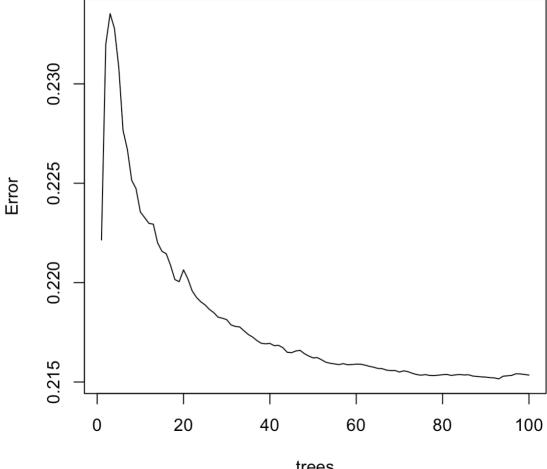
Accuracy : 0.7025
95% CI : (0.6785, 0.7257)
No Information Rate : 0.7147
P-Value [Acc > NIR] : 0.8565

Kappa : 0.004

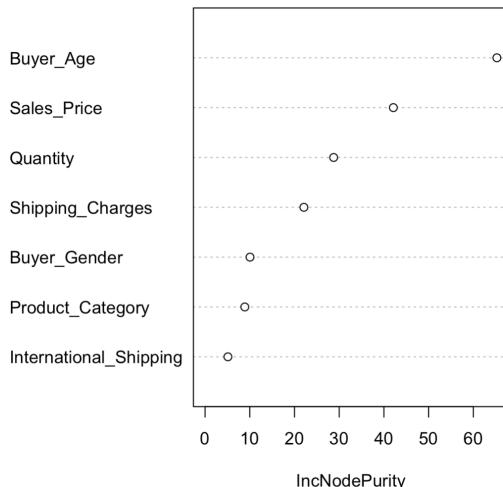
McNemar's Test P-Value : <2e-16

Sensitivity : 0.969933
Specificity : 0.03318
Pos Pred Value : 0.71528
Neg Pred Value : 0.30435
Prevalence : 0.71467
Detection Rate : 0.69304
Detection Prevalence : 0.96890
Balanced Accuracy : 0.50145

'Positive' Class : 0
```



### my\_forest



```
> summary(arima_model)
Series: sales_ts
ARIMA(0,0,0) with non-zero mean
```

Coefficients:

mean
116.7311

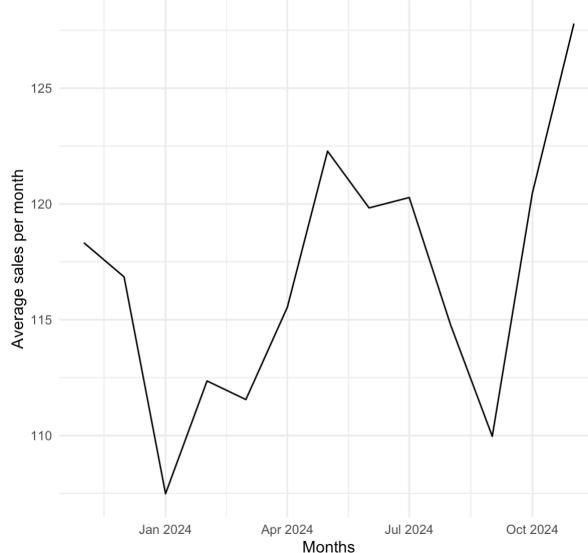
s.e. 1.4881

sigma^2 = 31.19: log likelihood = -40.29  
AIC=84.57 AICc=85.77 BIC=85.7

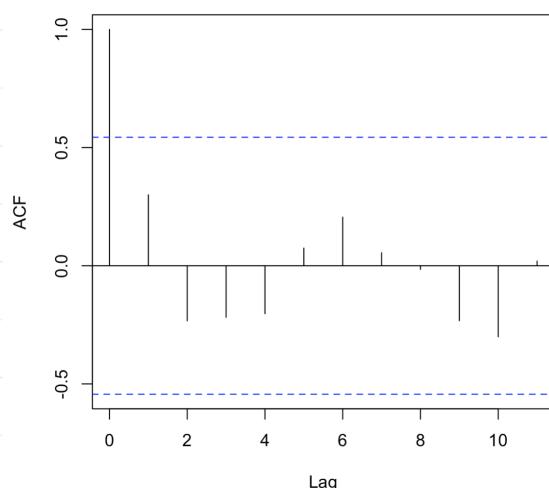
Training set error measures:

	ME	RMSE	MAE	MPE	MAPE	MASE	ACF1
Training set	-2.691864e-14	5.365523	4.417826	-0.2108027	3.793569	0.4668477	0.3000907

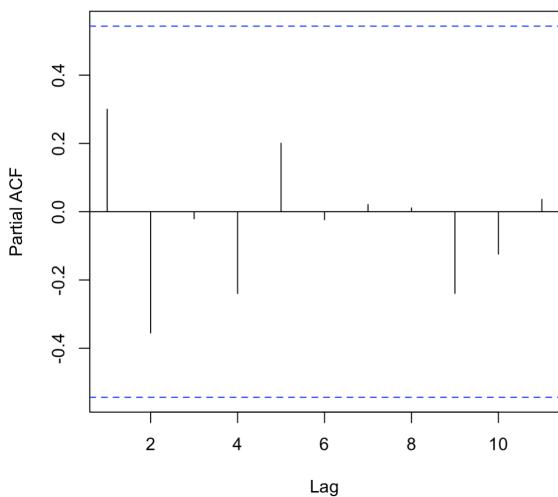
Monthly average sales trend



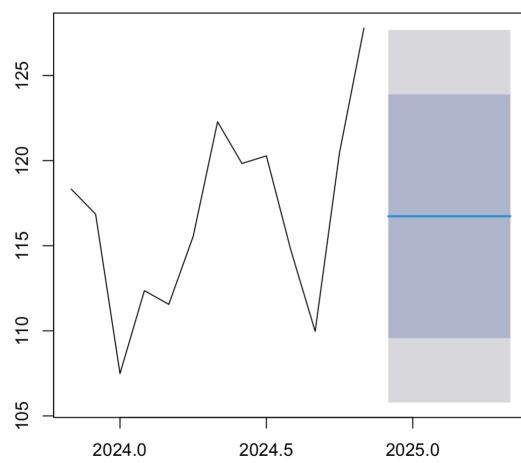
Series monthly\_avg\_sales\$avg\_sales



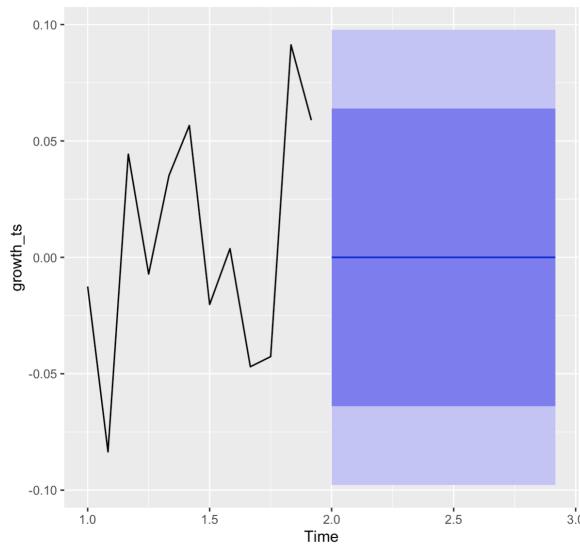
Series monthly\_avg\_sales\$avg\_sales



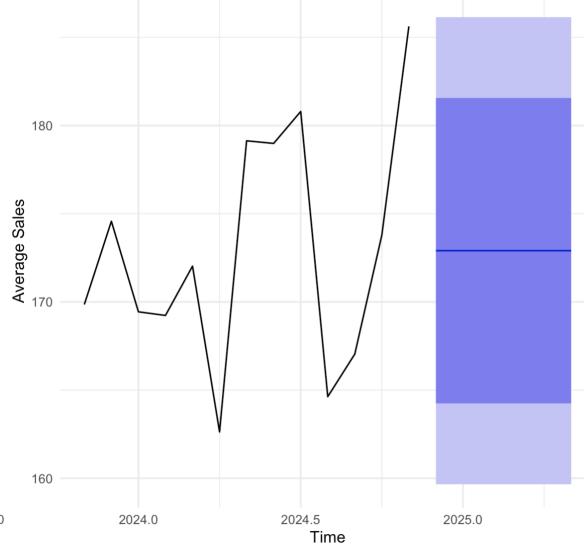
Forecasts from ARIMA(0,0,0) with non-zero mean

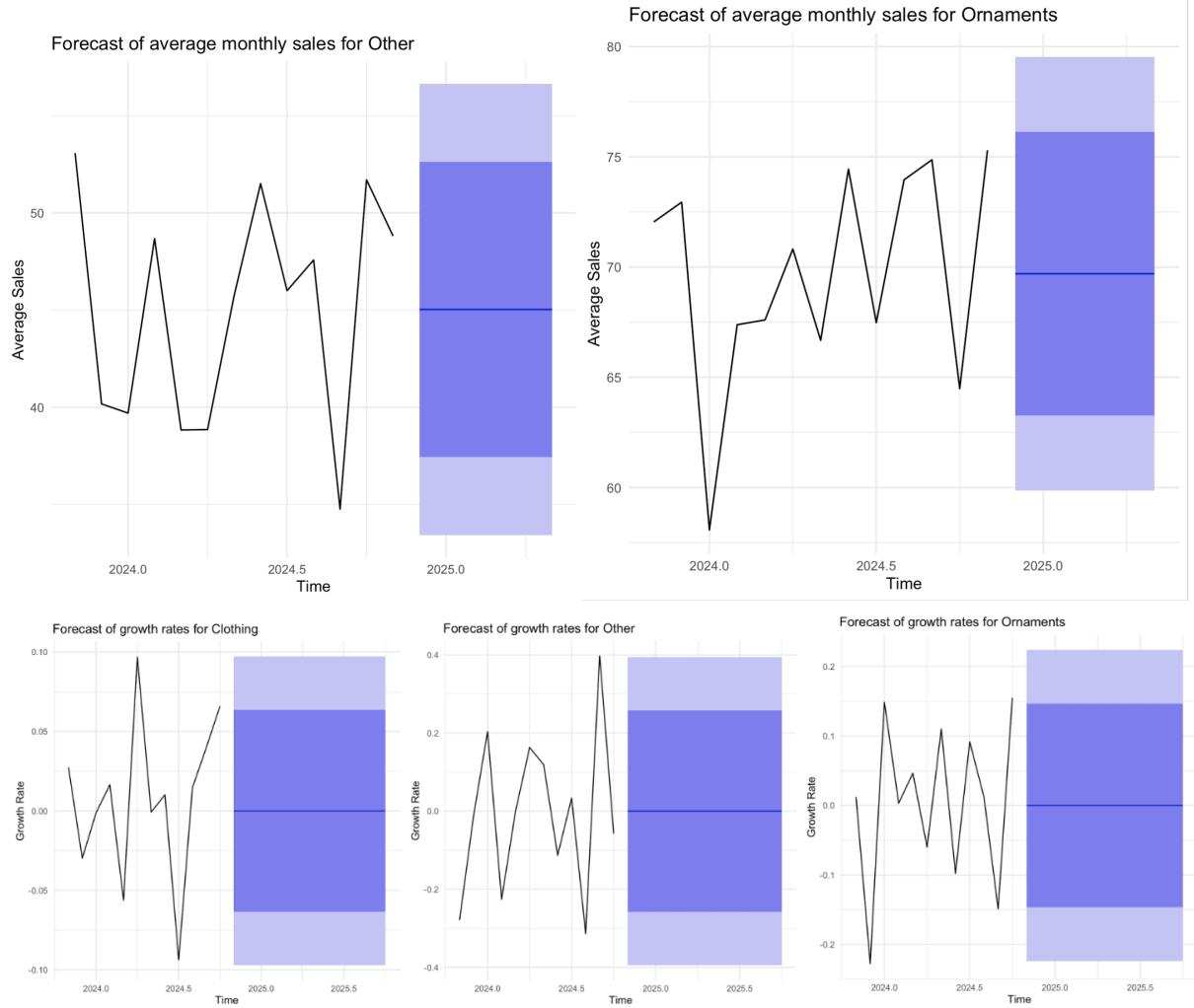


Forecasts from ARIMA(0,0,0) with zero mean



Forecast of average monthly sales for Clothing





### ***References***

Global Retail sales Data: Orders, reviews & Trends. (2024, December 10). Kaggle.  
<https://www.kaggle.com/datasets/adarsh0806/influencer-merchandise-sales>