

**Customer Reviews Analysis**:

Unlocking Insights into Consumer Satisfaction and Preferences

BY HAMZA YOUSAF

## **BUSINESS QUESTIONS & MOTIVATION BEHIND THEM**

- 1. The first question aims to determine if the customer is likely to recommend the product or not based on their review? (Positive recommendations from satisfied customers can lead to new customers trying out the product, while negative recommendations can have the opposite effect.)
- 2. The second question seeks to identify if the customer will give the product a rating of 5 or not?

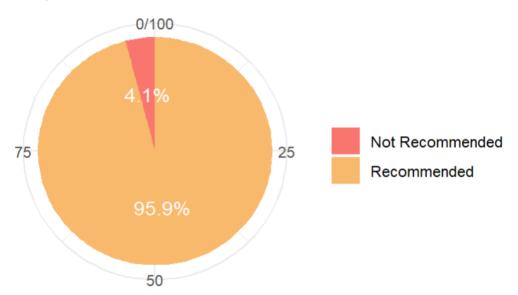
  (Marketers need to focus on providing excellent customer experiences and addressing any issues that arise promptly to increase the chances of customers giving higher ratings.)
- 3. The third question examines the customer's rating behavior, specifically if they will give a rating of 5, 4, or below 4? (Understanding the specific rating a customer is likely to give a product can help marketers identify areas for improvement.)

# DESCRIPTIVE ANALYSIS:

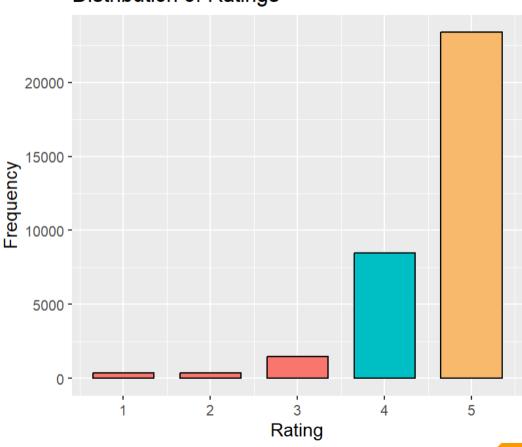
21 variables and 34,660 observations in this dataset.4 variables were selected for the study, and an additional5 product group variables were created from one of the original variables, 'categories'.

### **Imbalanced dataset**

### Proportion of Recommendations



### Distribution of Ratings



The dataset covers a time frame of 3.5 years, ranging from Oct 2014 to Apr2018.

# DESCRIPTIVE ANALYSIS:

### **VARIABLE GROUPS:**

- 1. Tablets & E-Readers
- 2. Amazon Devices
- 3. Accessories
- 4. Electronics
- 5. office & Business

Group	Percent_Recommended	Not_Recommended_Percentage	Average_Rating	Total_Reviews
OffnBuss	98.1738	1.826197	4.772355	3176
Elec	96.68956	3.310442	4.682338	5679
Accz	96.32243	3.677574	4.670519	6662
AmzDvcs	93.52708	6.472919	4.560106	757
Tblts	95.16973	4.830272	4.490526	17204

## Top 10 Products with Highest Not Recommended Percentage

Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, Fire HD 8 Tablet with Alexa, 8 HD Display, 16 GB, Tangerine - with Special Offers",

Amazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,, Amazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,,

Brand New Amazon Kindle Fire 16gb 7 Ips Display Tablet Wifi 16 Gb Blue...

New Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,, New Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,, Kindle Oasis E-reader with Leather Charging Cover -Black, 6 High-Resolution Display (300 ppi), Wi-Fi -Includes Special Offers,,

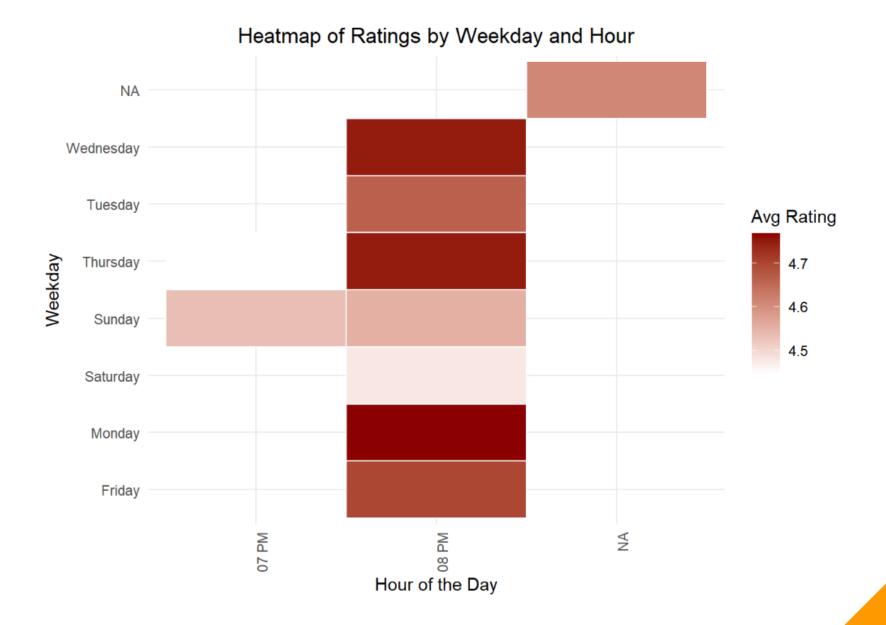
Echo (White),,, Echo (White),,,

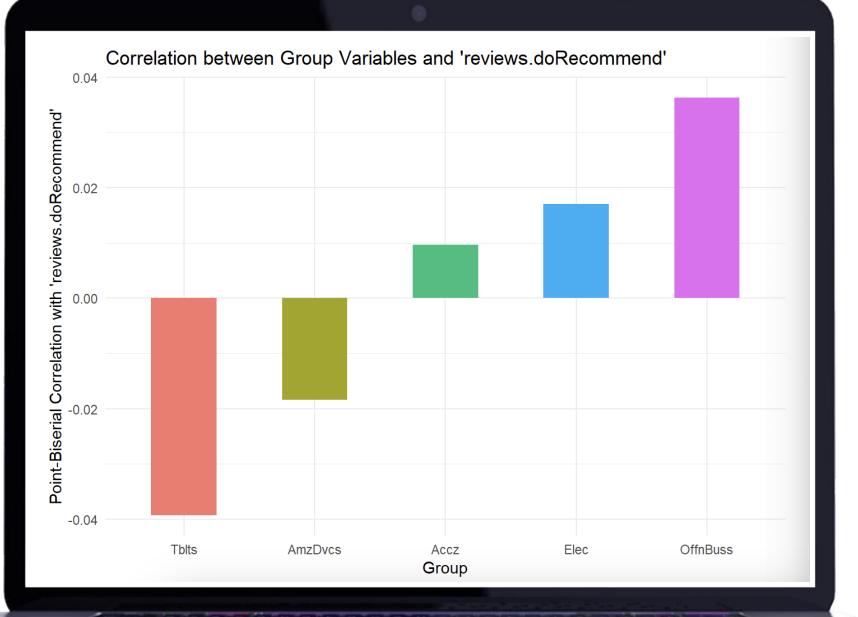
Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,,

Fire Tablet, 7
Display, Wi-Fi, 8 GB
- Includes Special
Offers, Magenta

# DATA VISUALIZATION:

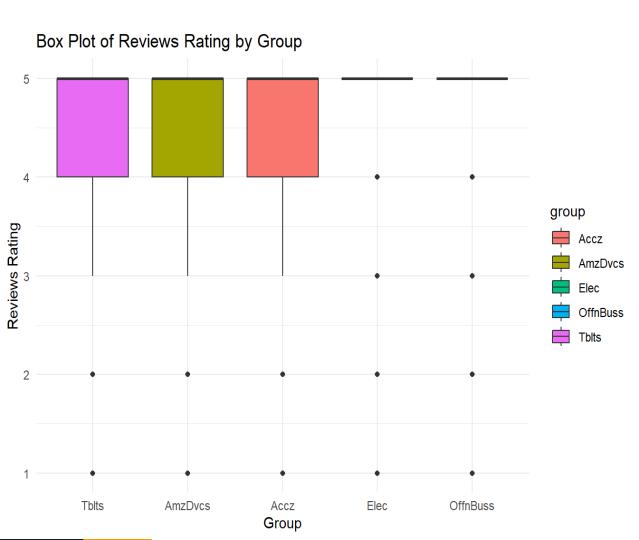
- heatmap average customer ratings by weekdays and hours most ratings - 8 PM
- Sundays 7 PM ratings range from 4.48 to 4.7
- Mondays and Wednesdays highest ratings
- Sundays slightly lower average of 4.55

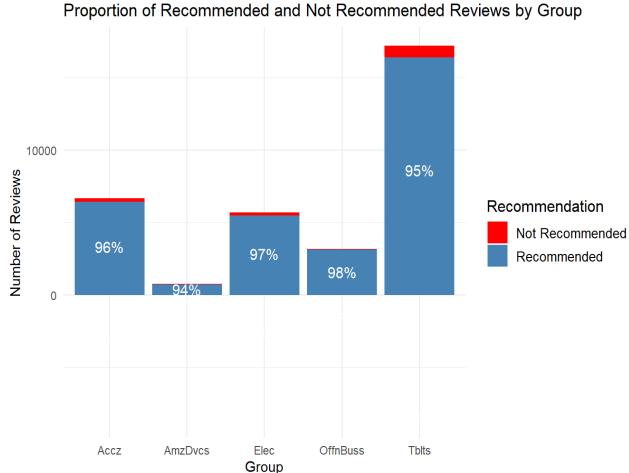




- **1.Tblts**: weak negative correlation (-0.039) with 'reviews.doRecommend'
- **2.AmzDvcs**: weak negative correlation (-0.018) with 'reviews.doRecommend'
- **3.Accz**: weak positive correlation (0.010) with 'reviews.doRecommend'
- **4.Elec**: weak positive correlation (0.017) with 'reviews.doRecommend'
- **5.OffnBuss**: higher positive correlation (0.036) with 'reviews.doRecommend'

# EVALUATING PRODUCTS REPUTATION





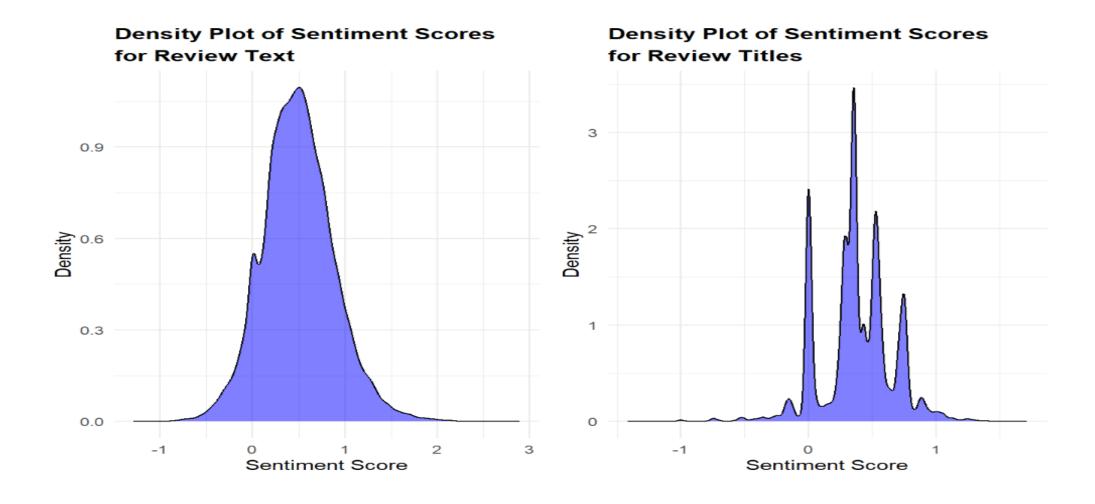
## TEXT DATA PREPARATION AND SENTIMENT SCORE EXTRACTION

Lowercasing

Removing special characters

Stopwords removal

Stemming/ Lemmatization Sentiment Score conversion

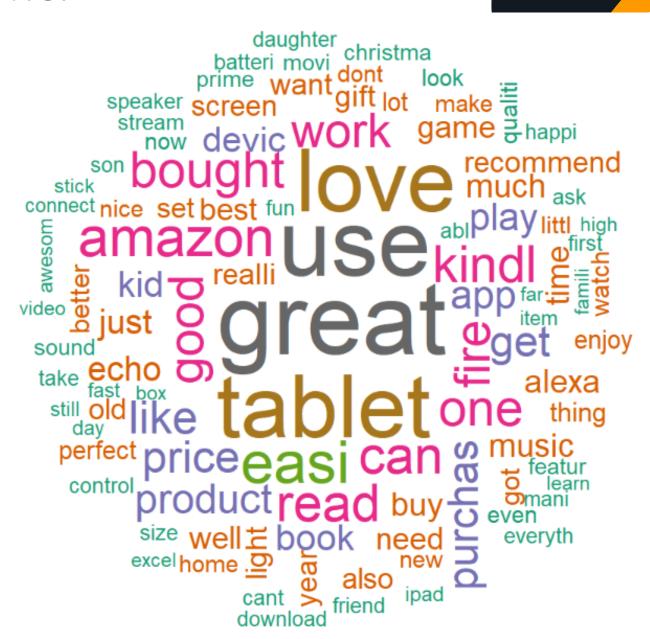


## **VISUALIZING CUSTOMER SENTIMENTS:**

Word Cloud (Reviews-based)



- 1. 'great' and 'use' are most frequently customers find the product to be of high quality
- 2. 'tablet' and 'love' are the second most common customers are expressing their appreciation for the tablets
- 3. 'amazon', 'bought', 'good', 'read', 'work', 'kindle', and 'one'. significant in frequency



MODEL EVALUATION — BUSINESS QUESTION 1

1.00 0.75 -Model 0.25 0.00 0.00 0.25 0.50 0.75 1.00 1 - Specificity

Predict whether customer is going to recommend the product or not?

	Logistic	<b>Classification Tree</b>	Random Forest
Sensitivity	49.4%	61.5%	52.8%
specificity	99.6%	99.1%	99.5%
Accuracy	97.7%	97.6%	97.7%
AUROC	94.6%	91.5%	92.2%
Balanced Accuracy	74.5%	80.3%	76.1%

**ROC Curves** 

As our focus is on identifying the customer who will not recommend the product, so we would focus on sensitivity metric, and it is considerably higher for Classification Tree and sensitivity is also same across all the models.

Decision Tree

Logistic Regression
Random Forest

Classification Tree model demonstrates the best overall performance in terms of balanced accuracy, making it the preferred choice for this business question.

# MODEL EVALUATION – BUSINESS QUESTION 2

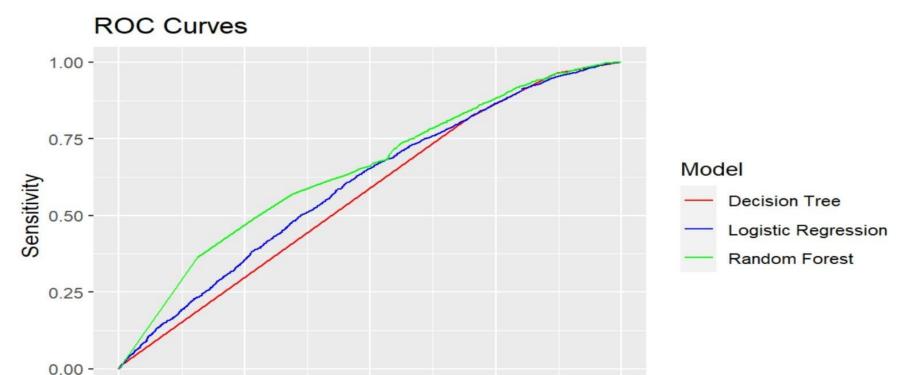
0.25

0.00

0.50

1 - Specificity

	Logisitic	Classification Tree	Random Forest
Sensitivity	55.2%	29.79%	10.20%
specificity	61.6%	82.45%	97.00%
Accuracy	59.1%	62.00%	63.40%
AUROC	62.0%	63.11%	63.92%
Balanced Accuracy	58.3%	56.12%	53.60%



0.75

1.00

Predict whether the customer is going to give rating of 5 or not?

As our focus is on identifying customers who will not give a rating of 5, we would pay to the attention more specificity metric, as it measures the model's ability correctly identify to customers who will not give a rating of 5. Random Forest model is the best choice, as it has the highest specificity of 97%

# MODEL EVALUATION - BUSINESS QUESTION 3

	Multinomial logistic	Classification Tree	Random Forest
Senstivity	57.1%	56.7%	57.0%
specificity	74.2%	74.1%	74.1%
Accuracy	54.1%	57.6%	54.4%
<b>Balanced Accuracy</b>	65.5%	65.4%	65.5%

Predict which rating customer is going to give 5, 4 or below 4 (3-1)?

Classification Tree model has the highest overall accuracy at 57.6%, which indicates it performs better overall in classifying the ratings correctly compared to the other models.

The <u>Classification Tree model appears to be the best choice for this business question</u>.



- collect more data on empty variables, such as reviews.didPurchase and reviews.id. knowing if the customer giving the review purchased the product can help to differentiate between genuine and fake reviews. reviews.id can provide a unique identifier for each review, which can be used to avoid duplicates and improve accuracy of analysis
- collect more data on empty variables, such as reviews.userCity and reviews.userProvince gain insights into geographic distribution of customers and tailor marketing strategies
- identified that Monday is most favorable day to request reviews from customers as they are more likely to leave positive feedback information valuable for businesses strategically plan their review requests and increase likelihood of receiving positive reviews businesses can leverage this insight to enhance their customer satisfaction and overall success
- recommendations, ratings, and use of tablets have a significant impact on the overall good reviews inform marketing and product development strategies
- collect more data that focuses on negative customer feedback, especially on variables such as reviews.doRecommend, reviews.ratings, snt\_scr\_text, snt\_scr\_title, Tblts, AmzDvcs, Accz, Elec, and OffnBuss, to enhance accuracy

# Thank you!

# APPENDIX!

DA	DATA DICTIONARY:				
No.	Variable Name	Data Type	Data type used for analysis	Variables choose for this study	Description
1	id	Alphanumeric	Chr		Unique identifier for the product.
2	name	Character	Chr		Name of the product.
3	asins	Character	Chr		Amazon Standard Identification Numbers for the product.
4	brand	Character	Factor		Brand of the product.
5	categories	Character	Chr	ü (made 5 product groups from this feature)	Categories the product belongs to, separated by commas.
6	keys	Character	Chr		Various identifiers and keys related to the product, separated by commas.
7	manufacturer	Character	Factor		Manufacturer of the product.
8	reviews.date	Date	Date		Date when the review was written.
9	reviews.dateAdded	Date	Date		Date when the review was added to the dataset.
10	reviews.dateSeen	Character	Chr		Dates when the review was seen online, separated by commas.
11	reviews.didPurchase	Logical	logi		Indicates if the reviewer purchased the product (TRUE) or not (FALSE).
12	reviews.doRecommend	Logical	logi	ü	Indicates if the reviewer recommends the product (TRUE) or not (FALSE).
13	reviews.id	Numeric	Int		Unique identifier for the review.
14	reviews.numHelpful	Numeric	num		Number of people who found the review helpful.
15	reviews.rating	Numeric	num	ü	Rating given by the reviewer on a scale of 1 to 5, with 1 being the lowest and 5 the highest.
16	reviews.sourceURLs	Character	Chr		URLs where the review was found, separated by commas.
17	reviews.text	Character	Chr	ü	Text of the review.

ü

Title of the review.

City of the reviewer.

Province of the reviewer.

Username of the reviewer.

reviews.title

reviews.userCity

reviews.userProvince

reviews.username

18

19

20

21

Chr

logi

logi Chr

Character

Character

Character

Character

## **MULTINOMIAL LOGISTIC REGRESSION:**

Metric	Class: 4	Class: 5	Class: below_4
Accuracy	54.14%	54.14%	54.14%
Sensitivity	63.90%	50.25%	<b>57.10</b> %
Specificity	51.80%	71.40%	99.50%
Precision	30.50%	78.60%	89.10%
Balanced Accuracy	57.60%	60.70%	78.30%

## **DECISION TREE:**

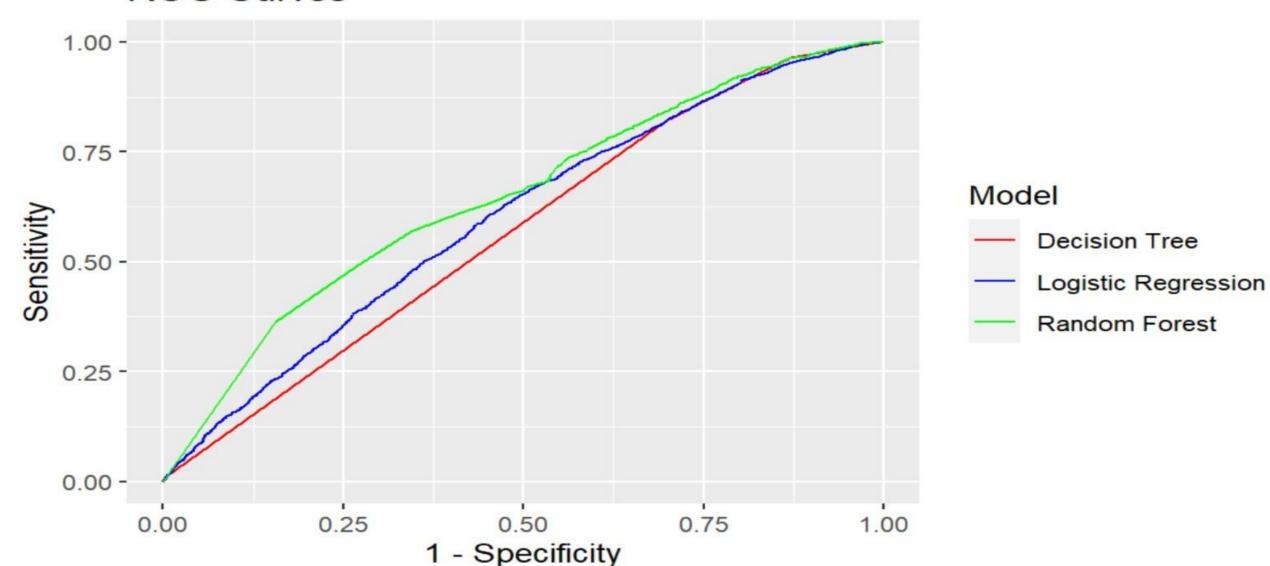
Metric	Class: 4	Class: 5	Class: below_4
Accuracy	57.6%	<b>57.6</b> %	57.6%
Sensitivity	54.2%	58.8%	57.1%
Specificity	60.0%	62.9%	99.5%
Precision	31.0%	77.1%	89.4%
Balanced Accuracy	60.8%	<b>57.1</b> %	78.2%

## **RANDOM FOREST:**

Metric	Class: 4	Class: 5	Class: below_4
Accuracy	54.4%	54.4%	54.4%
Sensitivity	63%	51.0%	57.1%
Specificity	52.6%	70.3%	99.4%
Precision	30.6%	78.4%	89.3%
Balanced Accuracy	57.8%	60.6%	78.2%

# ROC CURVES FOR Q2 FOR MODELS BUILT ON IMBALANCED DATA





### Confusion Matrix and Statistics

Actual
Predicted FALSE TRUE
FALSE 381 230
TRUE 2261 3941

Accuracy: 0.6344

95% CI: (0.6228, 0.6458)

No Information Rate: 0.6122

P-Value [Acc > NIR] : 8.669e-05

Kappa : 0.1037

Mcnemar's Test P-Value : < 2.2e-16

Sensitivity: 0.14421

Specificity: 0.94486

Pos Pred Value: 0.62357

Neg Pred Value: 0.63544

Prevalence: 0.38779

Detection Rate: 0.05592

Detection Prevalence: 0.08968

Balanced Accuracy : 0.54453

'Positive' Class : FALSE

Q2: LOGISTIC MODEL

**PREDICATIVE** 

**PERFORMANCE** 

**USING** 

**IMBALANCED** 

**DATA** 

### Confusion Matrix and Statistics

Actual
Predicted FALSE TRUE
FALSE 381 230
TRUE 2261 3941

Accuracy: 0.6344

95% CI: (0.6228, 0.6458)

No Information Rate : 0.6122

P-Value [Acc > NIR] : 8.669e-05

Kappa : 0.1037

Mcnemar's Test P-Value : < 2.2e-16

Sensitivity: 0.14421 -

Specificity: 0.94486

Pos Pred Value: 0.62357

Neg Pred Value: 0.63544

Prevalence: 0.38779

Detection Rate: 0.05592

Detection Prevalence: 0.08968

Balanced Accuracy: 0.54453

'Positive' class : FALSE

Q2:LOGISTIC
MODEL
PREDICATIVE
PERFORMANCE
USING

**IMBALANCED** 

**DATA** 

Reference
Predicted FALSE TRUE
FALSE 1457 1601
TRUE 1185 2570
Confusion Matrix and Statistics

Reference Predicted FALSE TRUE FALSE 1457 1601 TRUE 1185 2570

Accuracy: 0.5911

95% CI: (0.5793, 0.6028)

No Information Rate: 0.6122 P-Value [Acc > NIR]: 0.9998

Kappa: 0.1629

Mcnemar's Test P-Value: 3.767e-15

Sensitivity: 0.5515

Specificity: 0.6162

Pos Pred Value: 0.4765 Neg Pred Value: 0.6844

Prevalence: 0.3878

Detection Rate: 0.2139

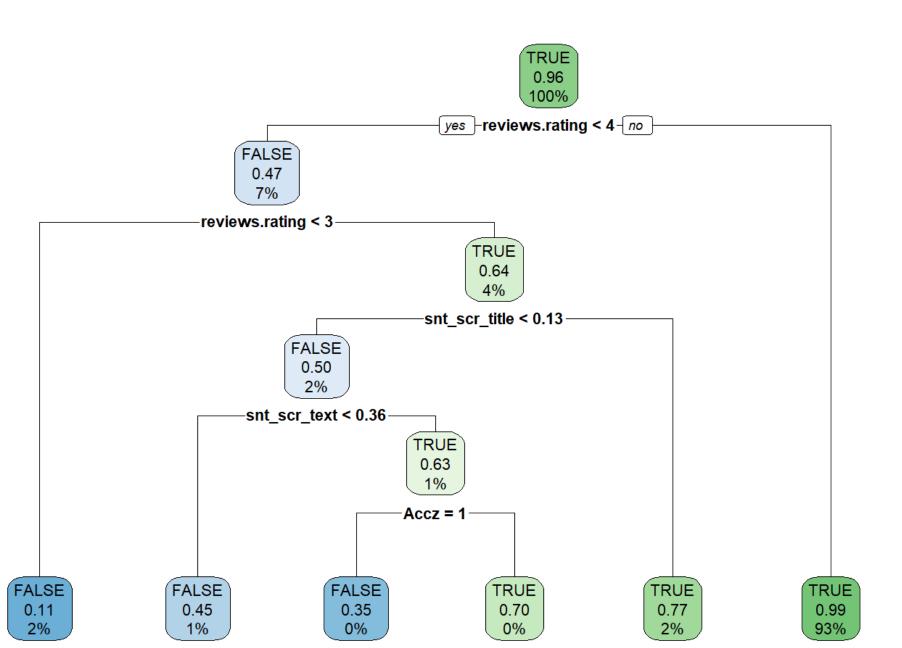
Detection Prevalence: 0.4488

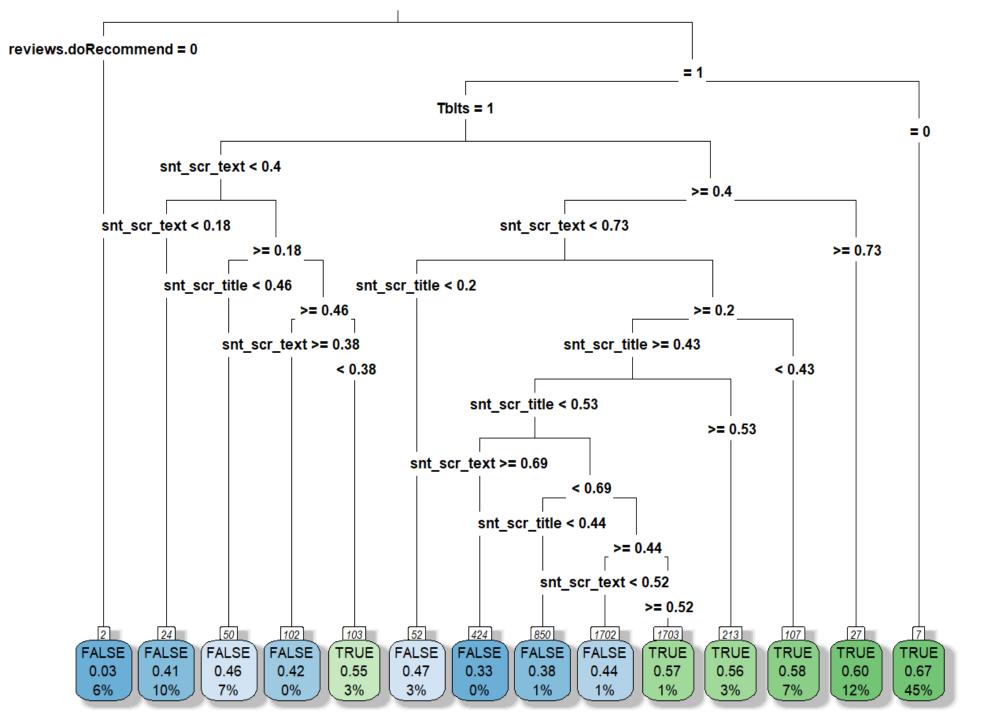
Balanced Accuracy : 0.5838

'Positive' Class : FALSE

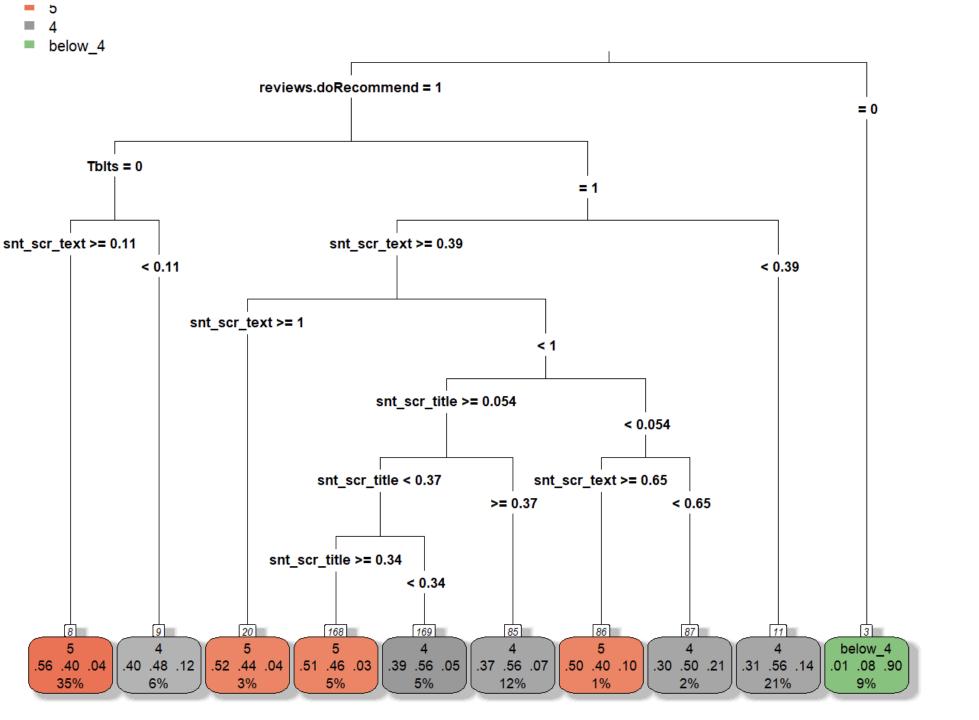
Q2:LOGISTIC
MODEL
PREDICATIVE
PERFORMANCE
USING
BALANCED
DATA

# RUNNING DECISION TREES ON UNBALANCE TRAINING DATA ON BUSINESS QUESTION 1:-











SIN		

3	MISSING VALUI			
	Variable Name	No.		
0				

8

15

**Description** 

**New Data** type

**Missing Values in** each column

%age of data missing from total obv of entire var

which method to use to fill in the missing values. deleting

Date when the review was written.

**Date** 

Date

39 10,621

0.1%

observations **Feature not** important/ not

added to the dataset. Indicates if the reviewer recommends the product (TRUE) or not (FALSE).

Number of people who

found the review helpful.

Rating given by the reviewer on a scale of 1 to

5, with 1 being the lowest

and 5 the highest.

31% 1.7%

included in model building deleting

reviews.doRecommend reviews.numHelpful 14

reviews.rating

logical (T/F) numeric

numeric

594 529

33

11,816

1.5%

0.10%

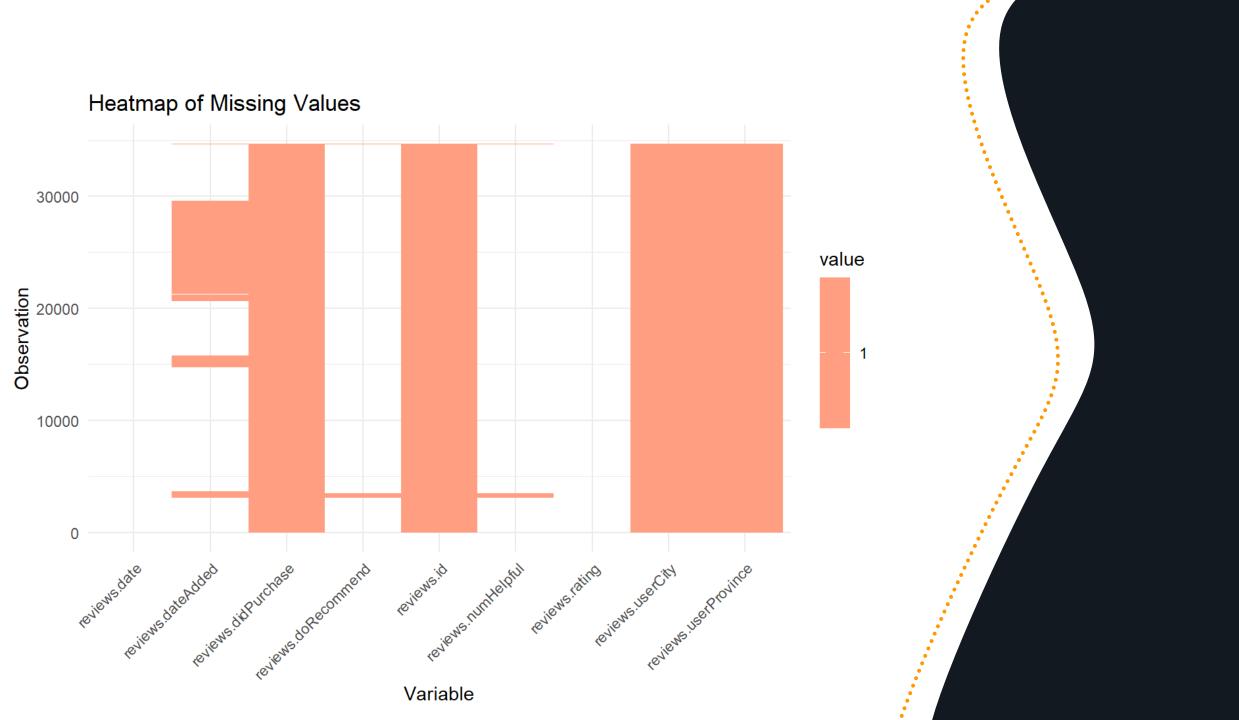
observations deleting observations

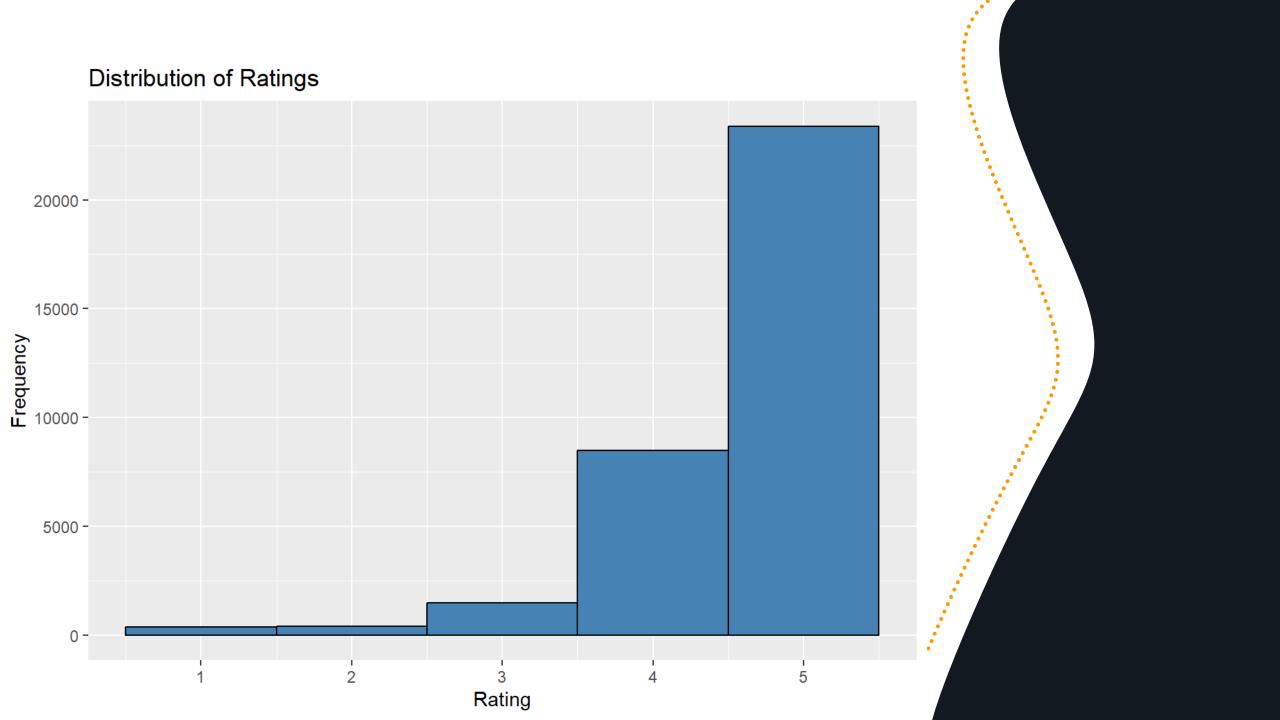
deleting

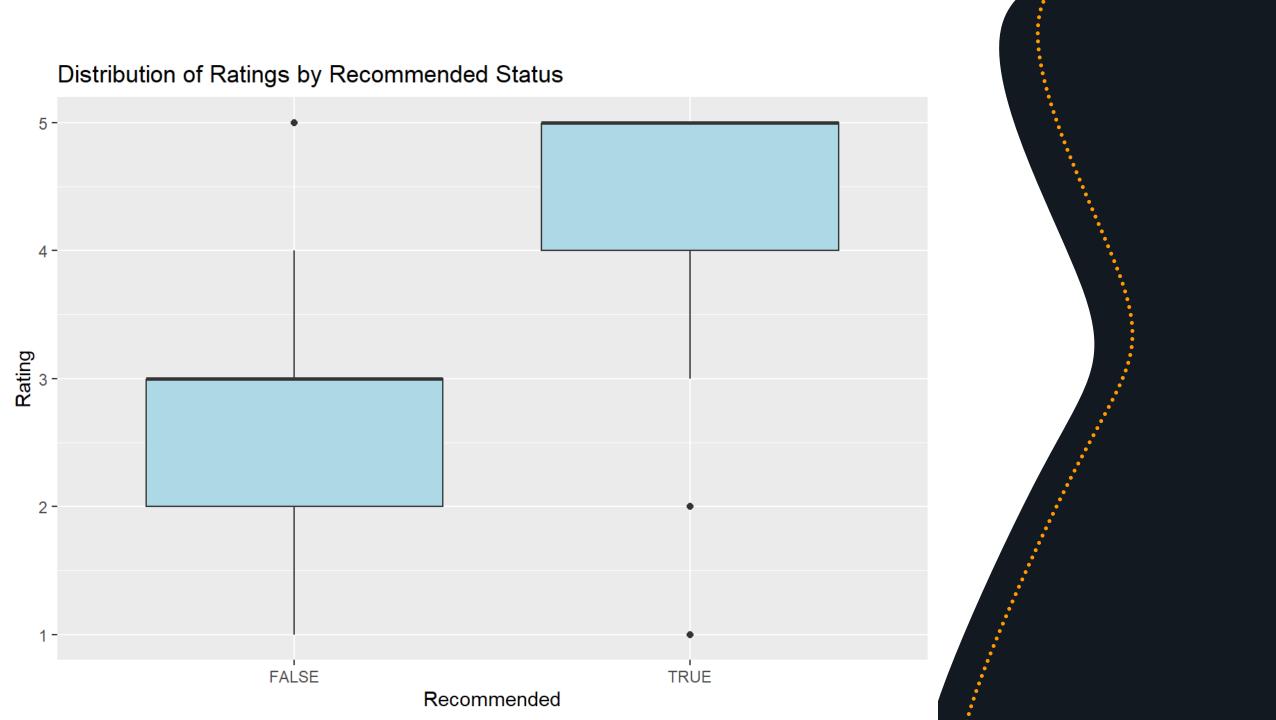
observations

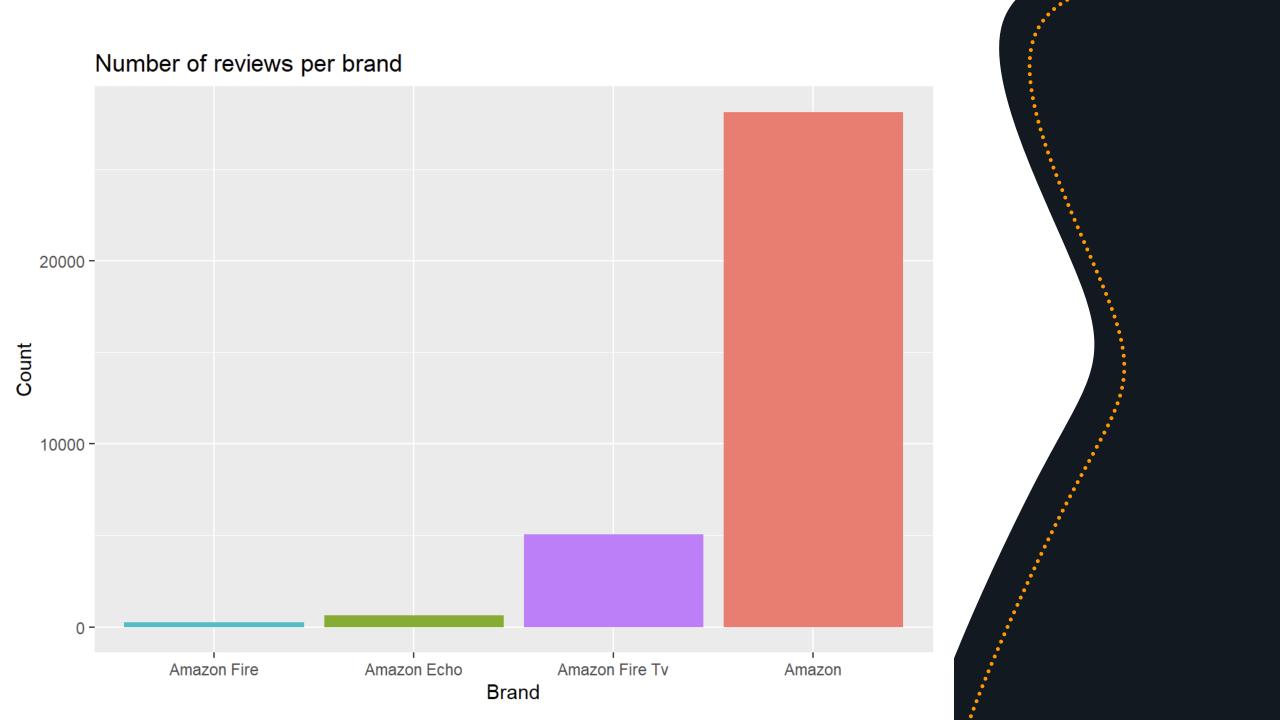
Date when the review was 9 reviews.dateAdded 12

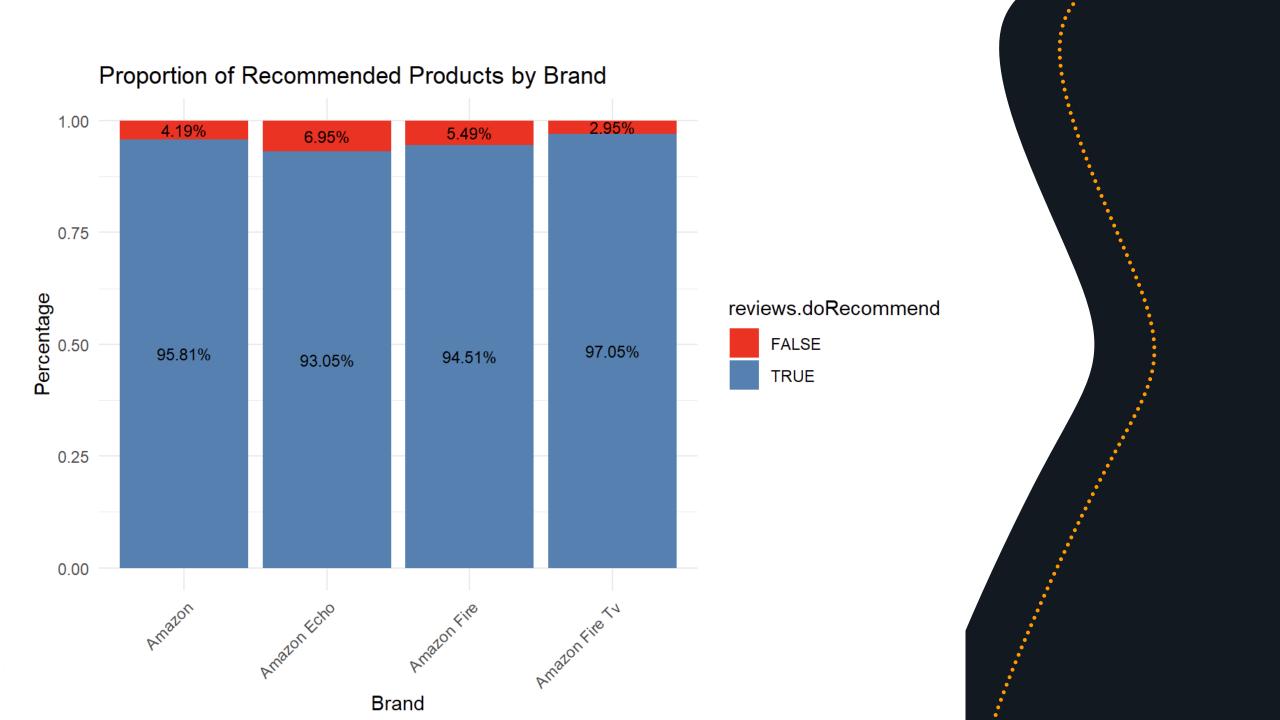
reviews.date

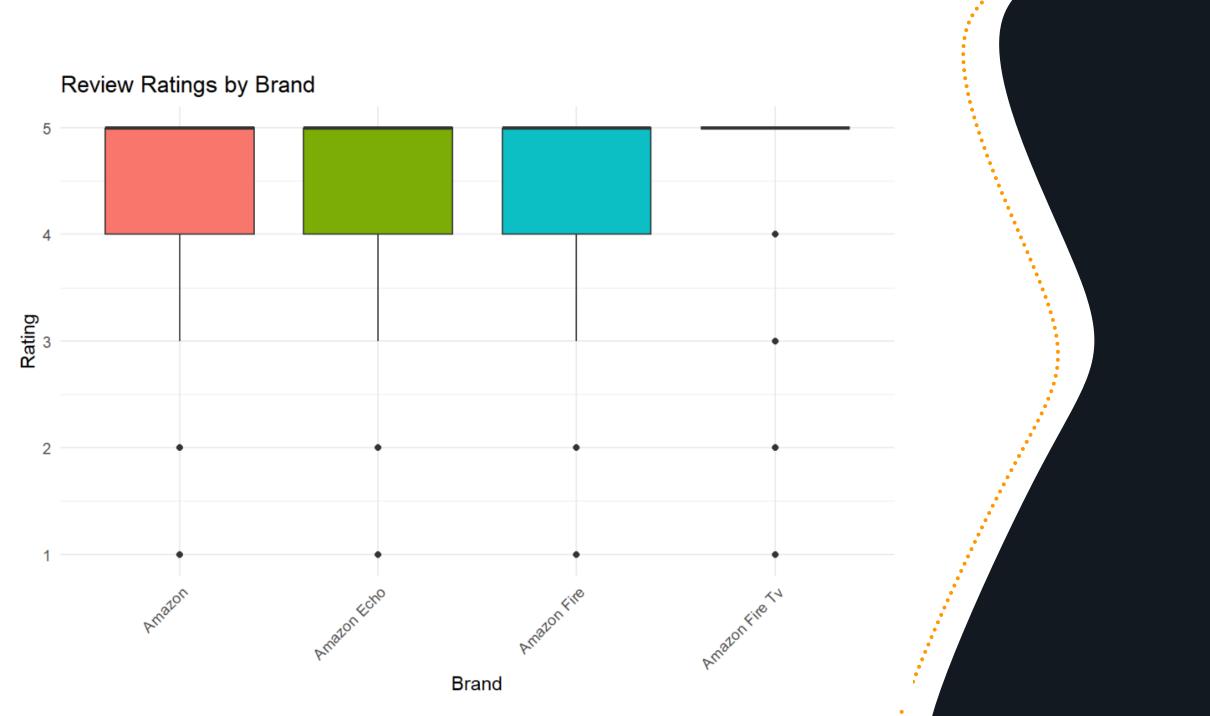


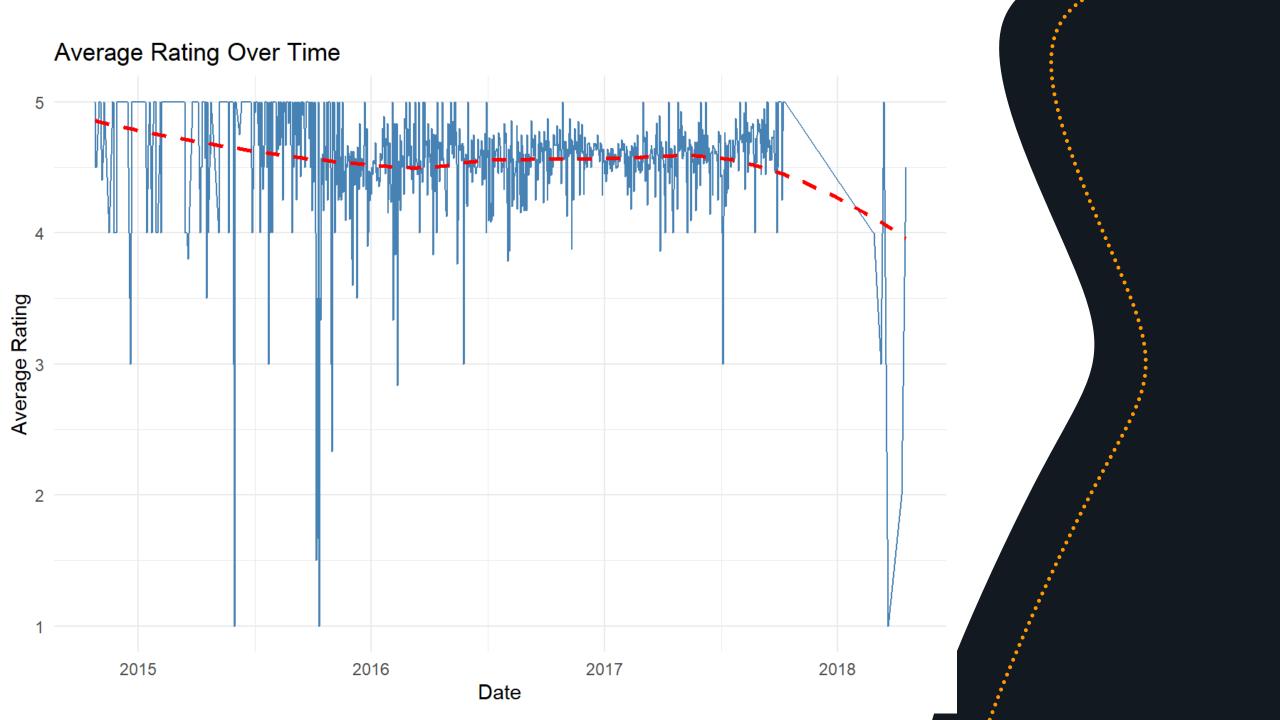




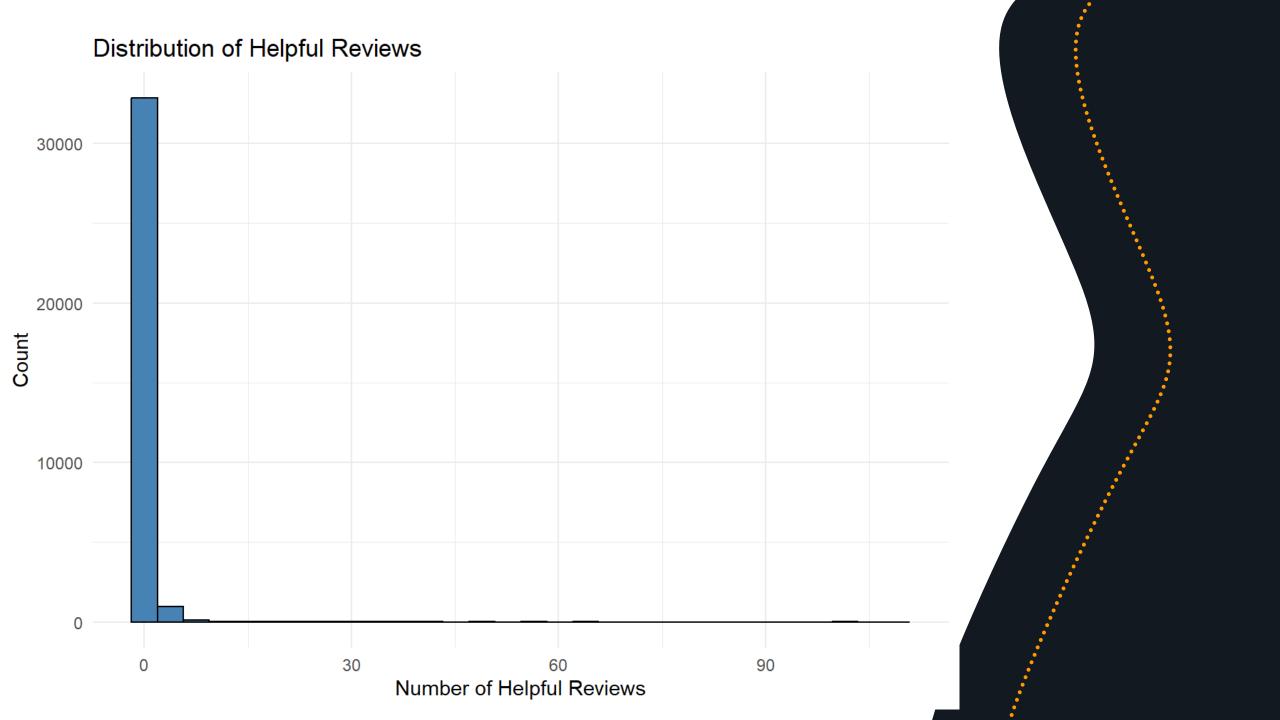


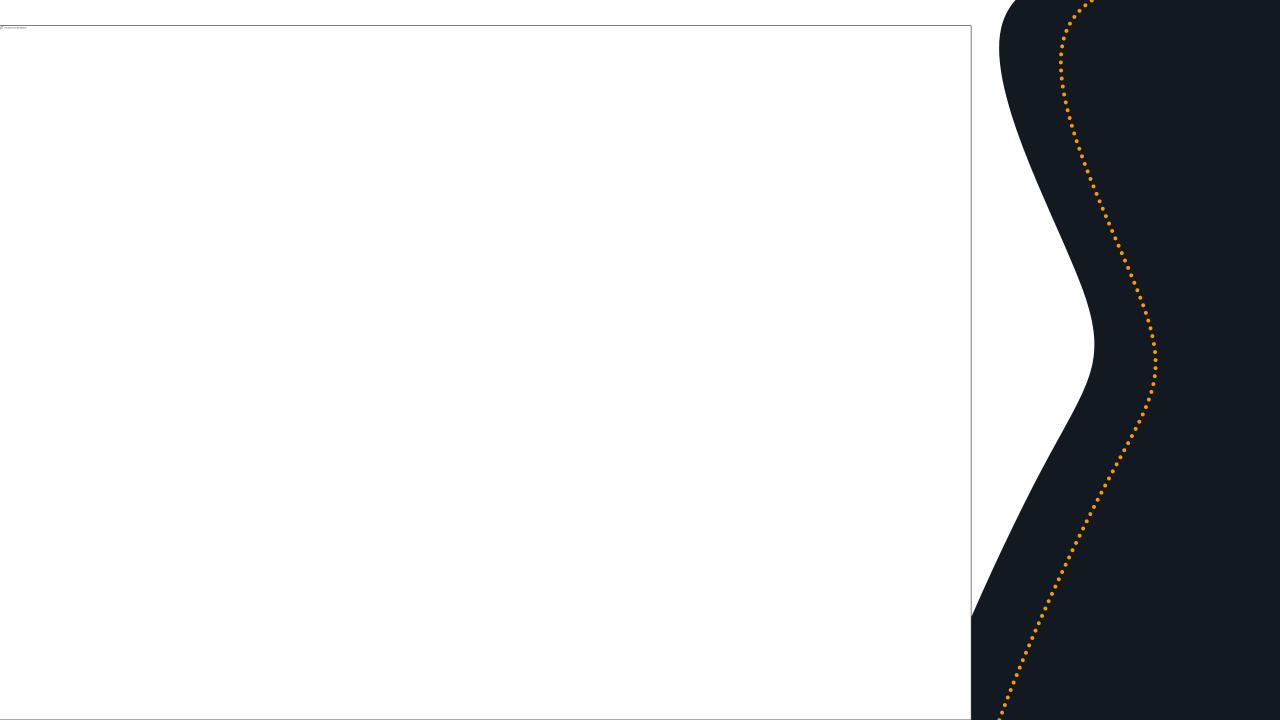




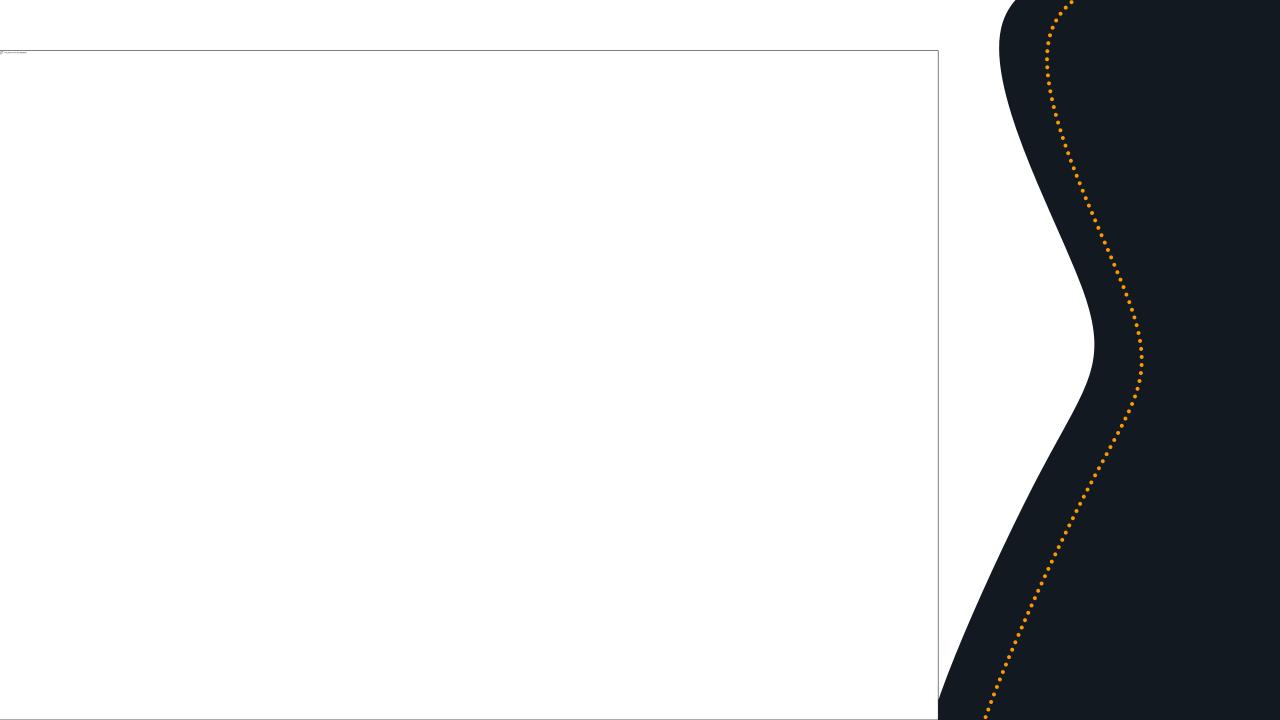






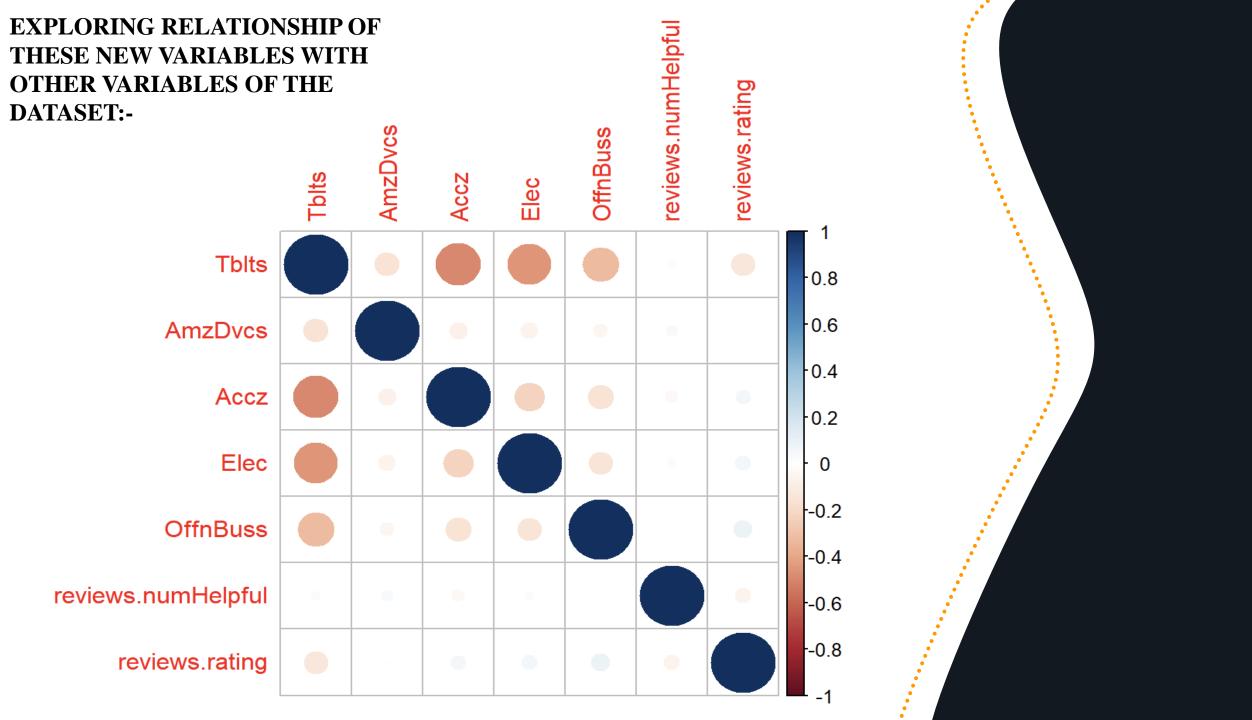


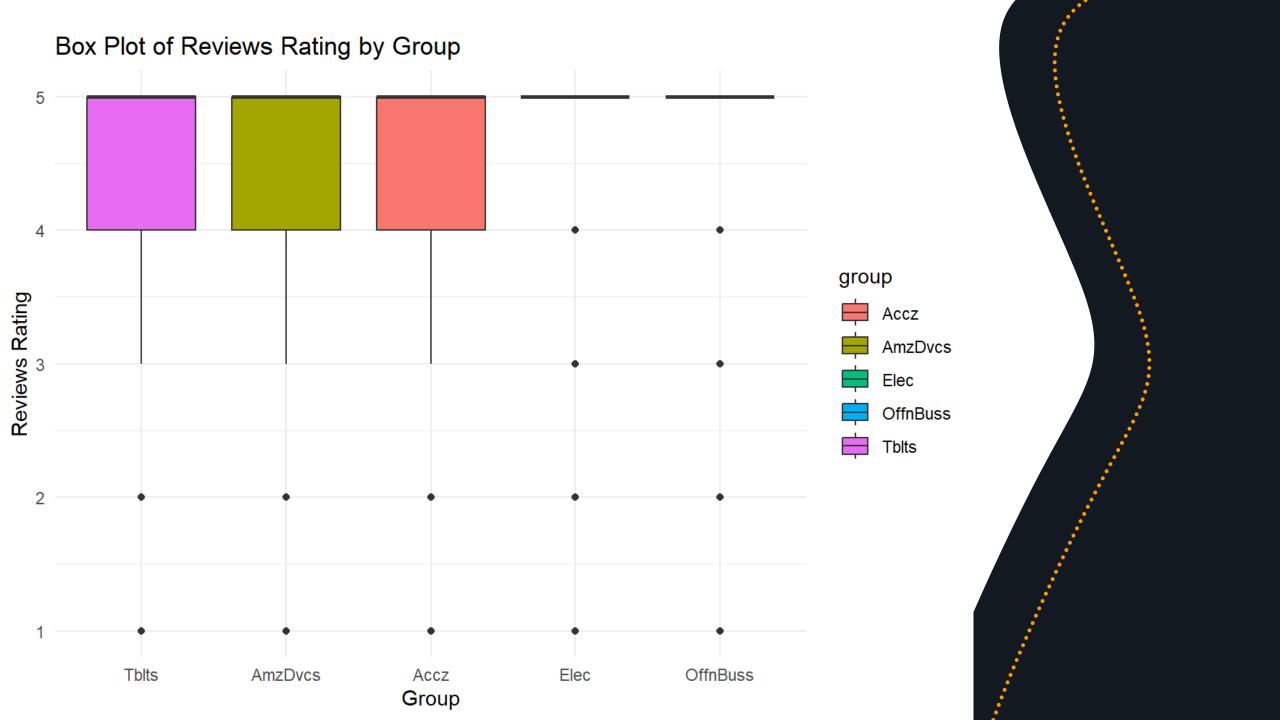


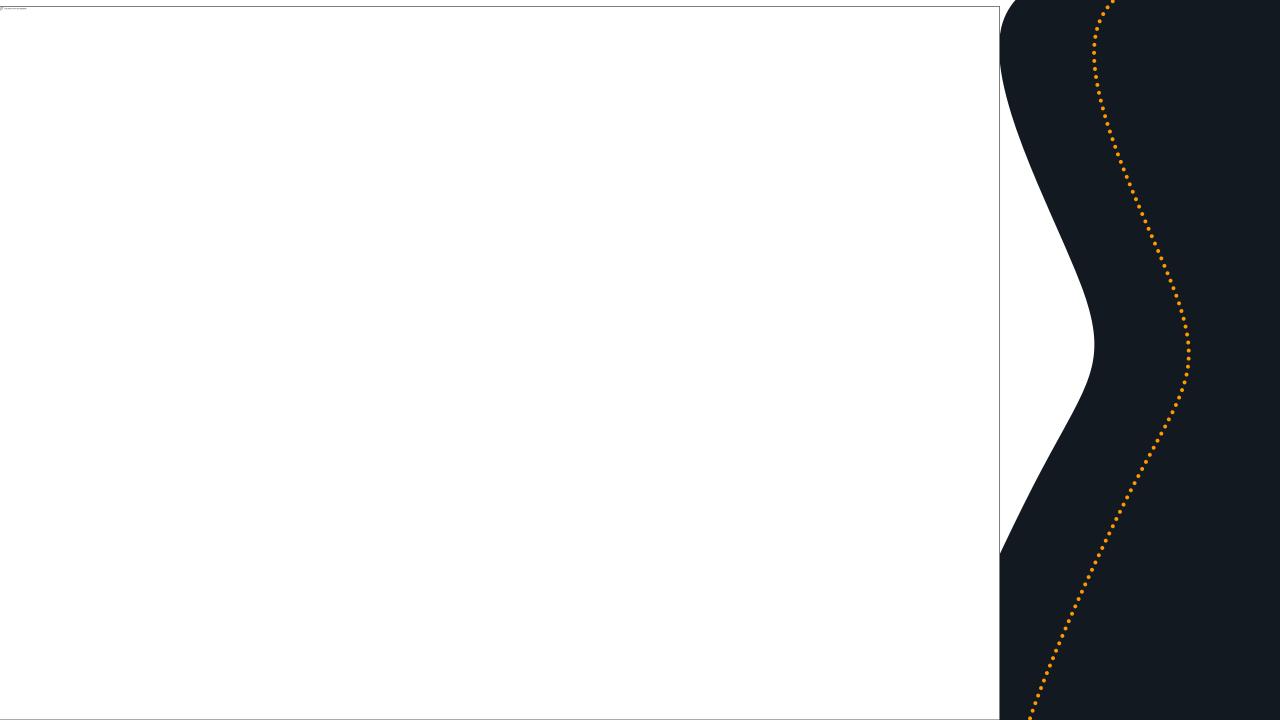


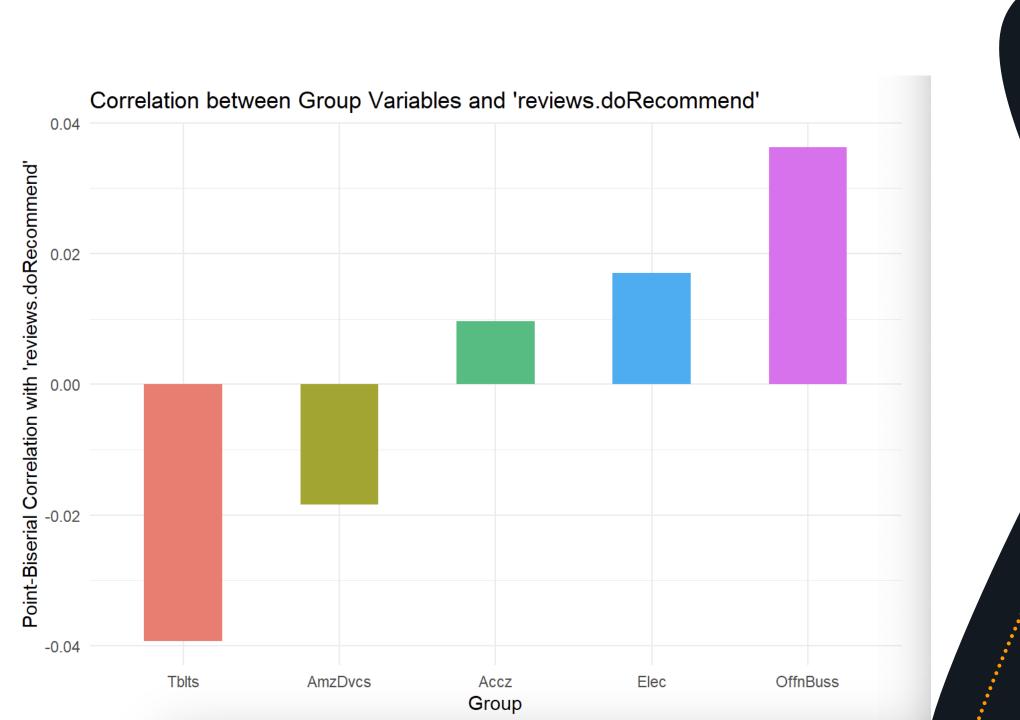
### **AVERAGE RATING BY DAY:**

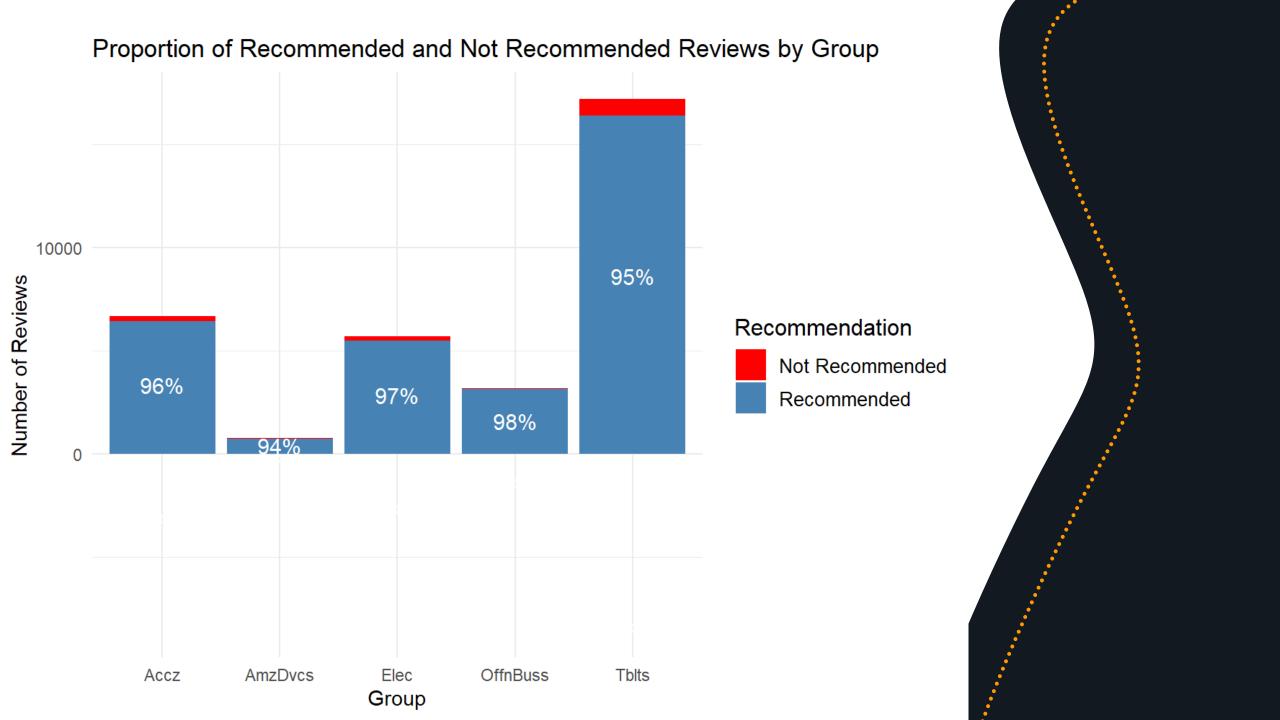
weekday	avg_rating
Monday	4.77
Wednesday	4.75
Friday	4.7
Tuesday	4.67
NA	4.61
Sunday	4.55
Thursday	4.51
Saturday	4.48











### Top 10 Products with Highest Not Recommended Percentage

Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G -Black,,, Fire HD 8 Tablet with Alexa, 8 HD Display, 16 GB, Tangerine - with Special Offers",

New Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,, **New Amazon Kindle Fire Hd 9w Powerfast** Adapter Charger + Micro Usb Angle Cable,,, Amazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,, Amazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders...

Kindle Fire 16gb 7 **Ips Display Tablet** Wifi 16 Gb Blue,,,

**Brand New Amazon** 

Kindle Oasis E-reader with Leather Charging Cover -Black, 6 High-Resolution Display (300 ppi), Wi-Fi -Includes Special Offers,,

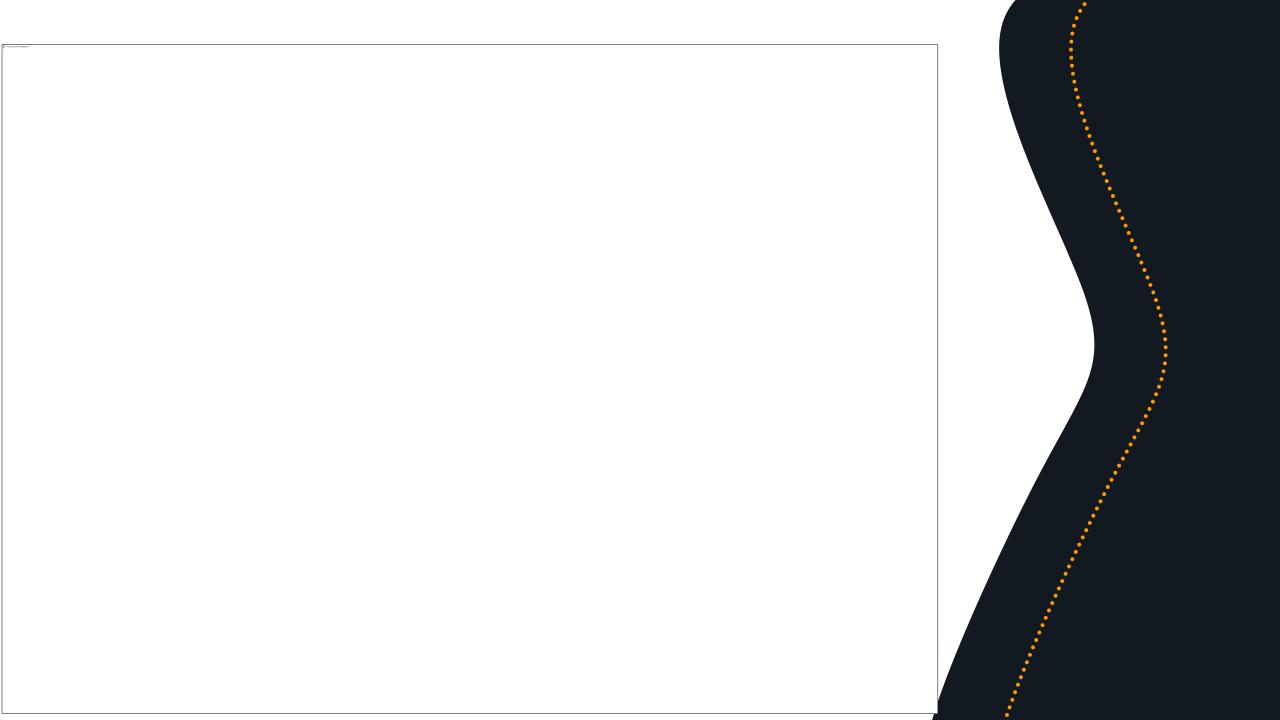
Echo (White),,, Echo (White),,,

Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,...

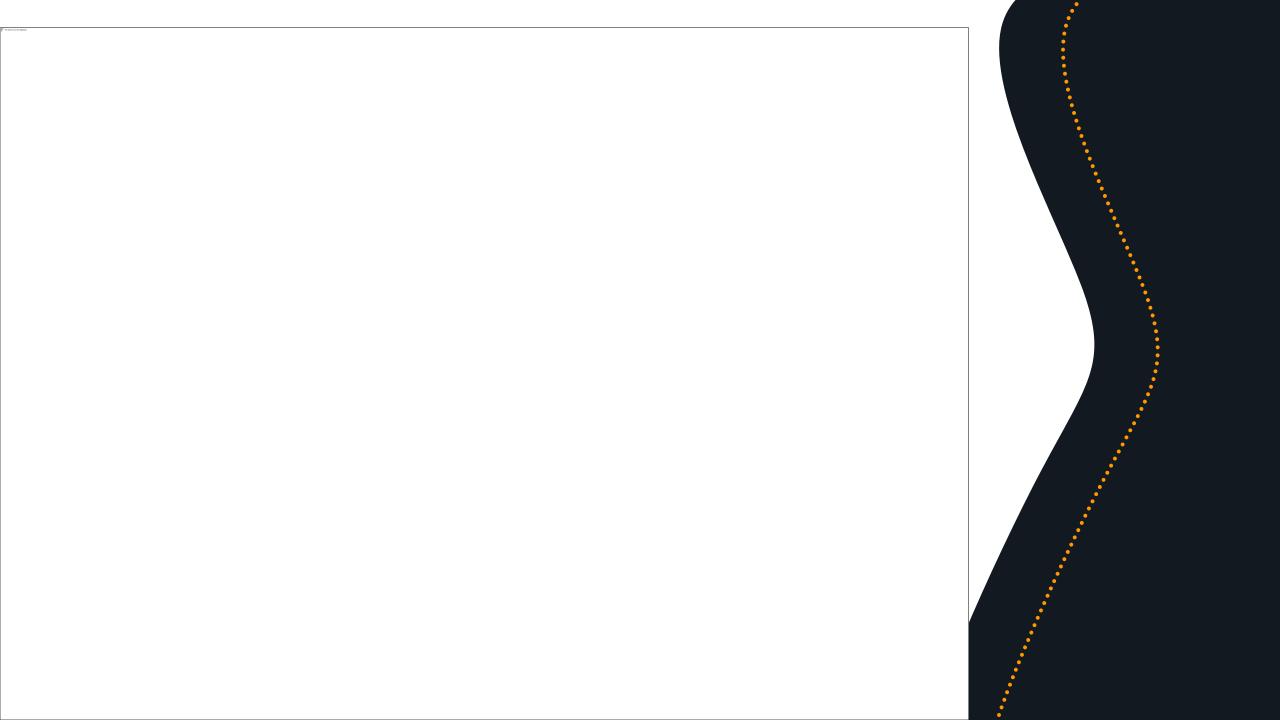
Fire Tablet, 7 Display, Wi-Fi, 8 GB - Includes Special Offers, Magenta

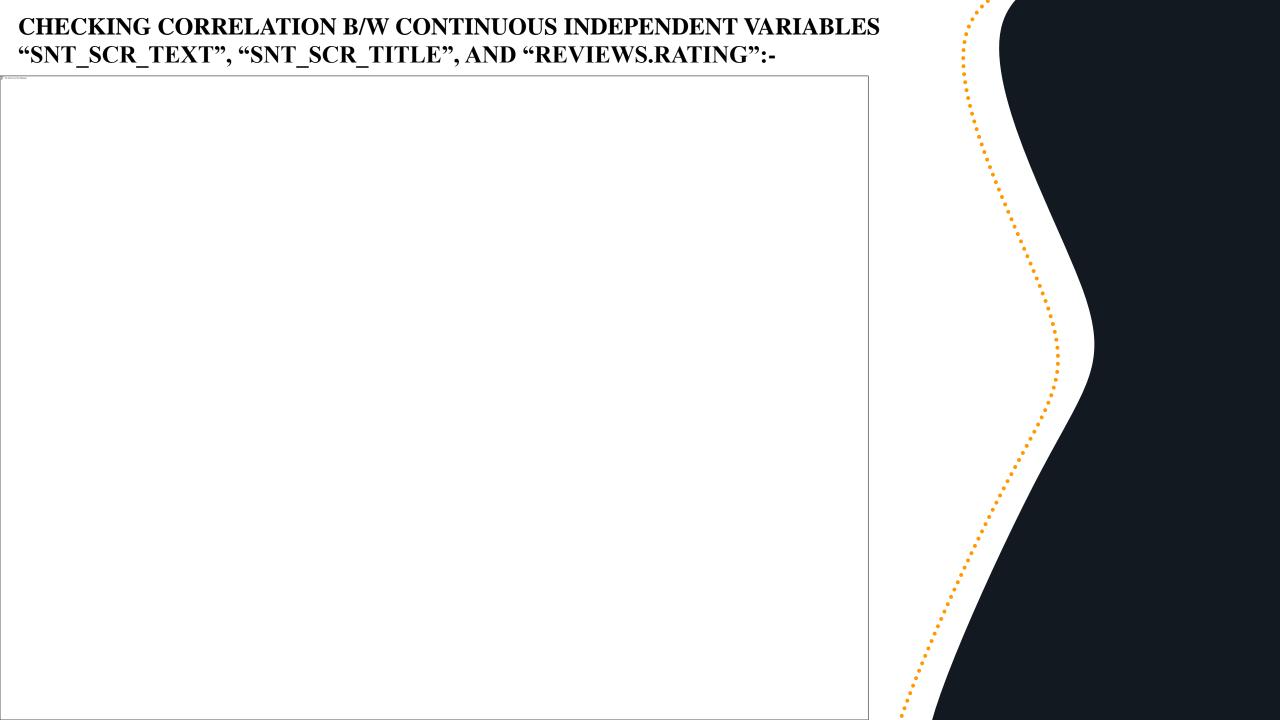
# PRODUCT GROUP SUMMARY: REVIEWS, RECOMMENDATIONS, AND AVERAGE RATINGS (PIVOT TABLE):-

Tblts	AmzDvcs	Accz	Elec	Group	Percent_Recommended	Not_Recommended_Percentage	Average_Rating	Total_Reviews
0	0	0	0	NA	97.78157	2.218430	4.730375	586
0	0	0	0	OffnBuss	98.17380	1.826197	4.772355	3176
0	0	0	1	Elec	96.68956	3.310442	4.682338	5679
0	0	1	0	Accz	96.32243	3.677574	4.670519	6662
0	1	0	0	AmzDvcs	93.52708	6.472919	4.560106	757
1	0	0	0	Tblts	95.16973	4.830272	4.490526	17204

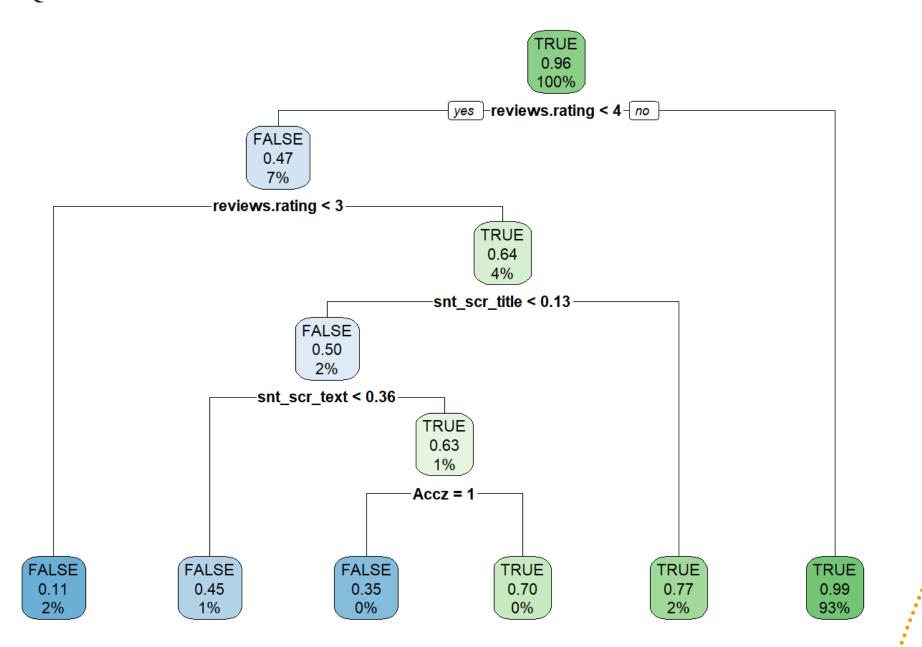


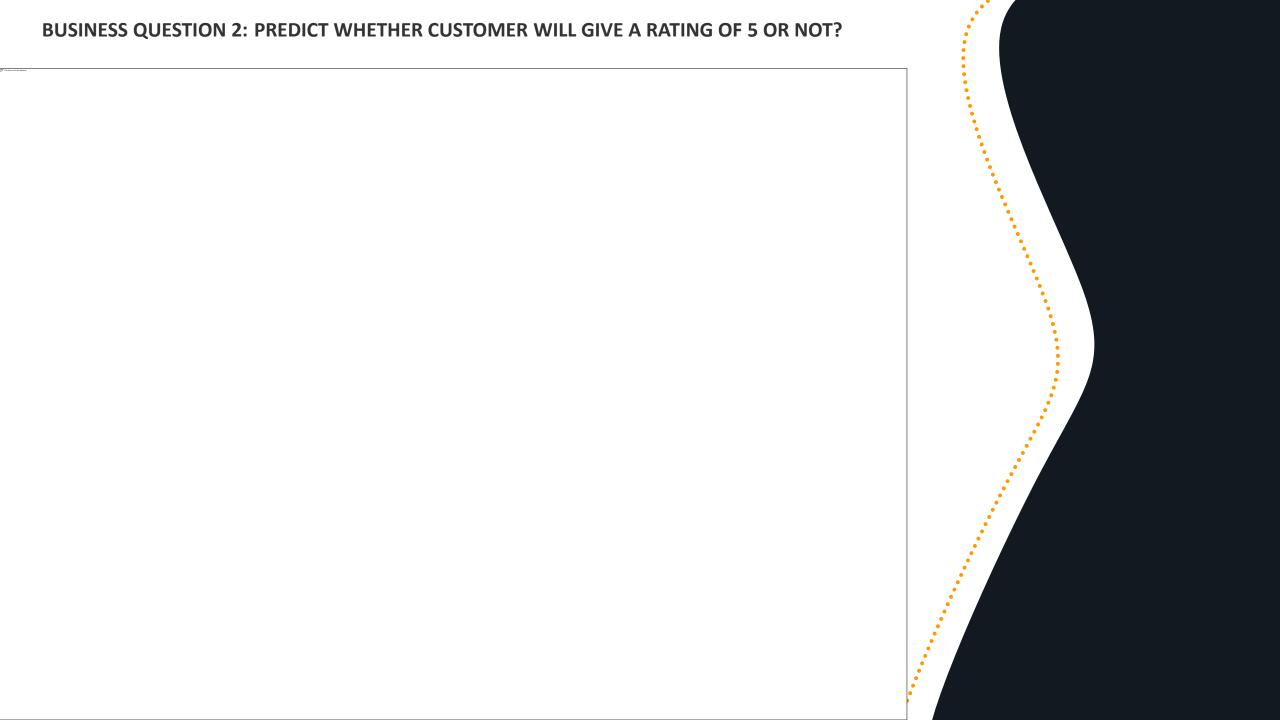


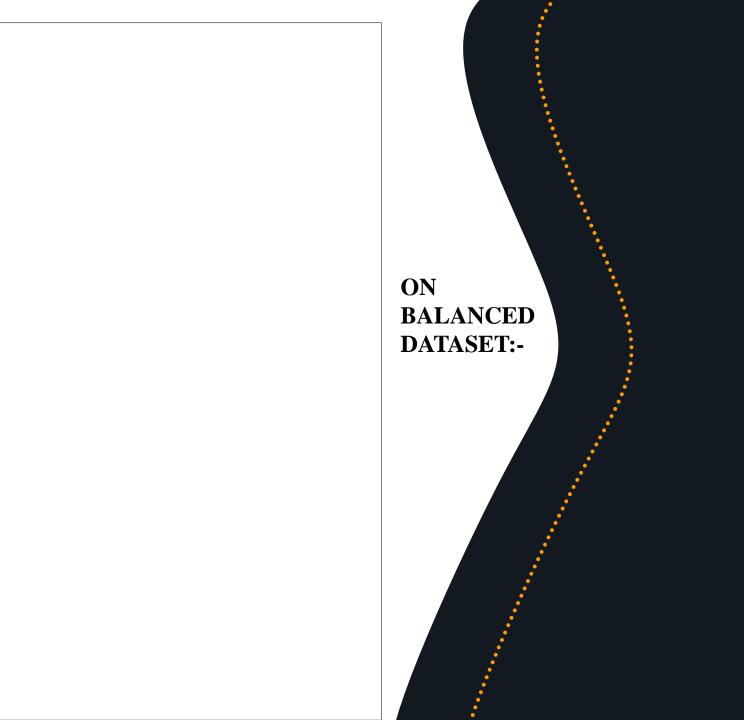




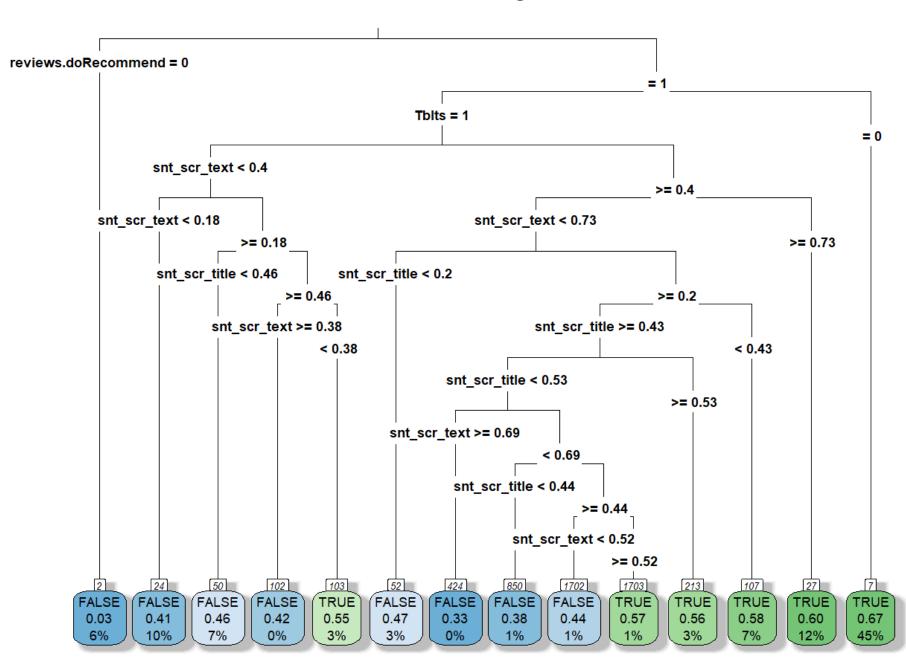
## RUNNING DECISION TREES ON UNBALANCE TRAINING DATA ON BUSINESS QUESTION 1:-



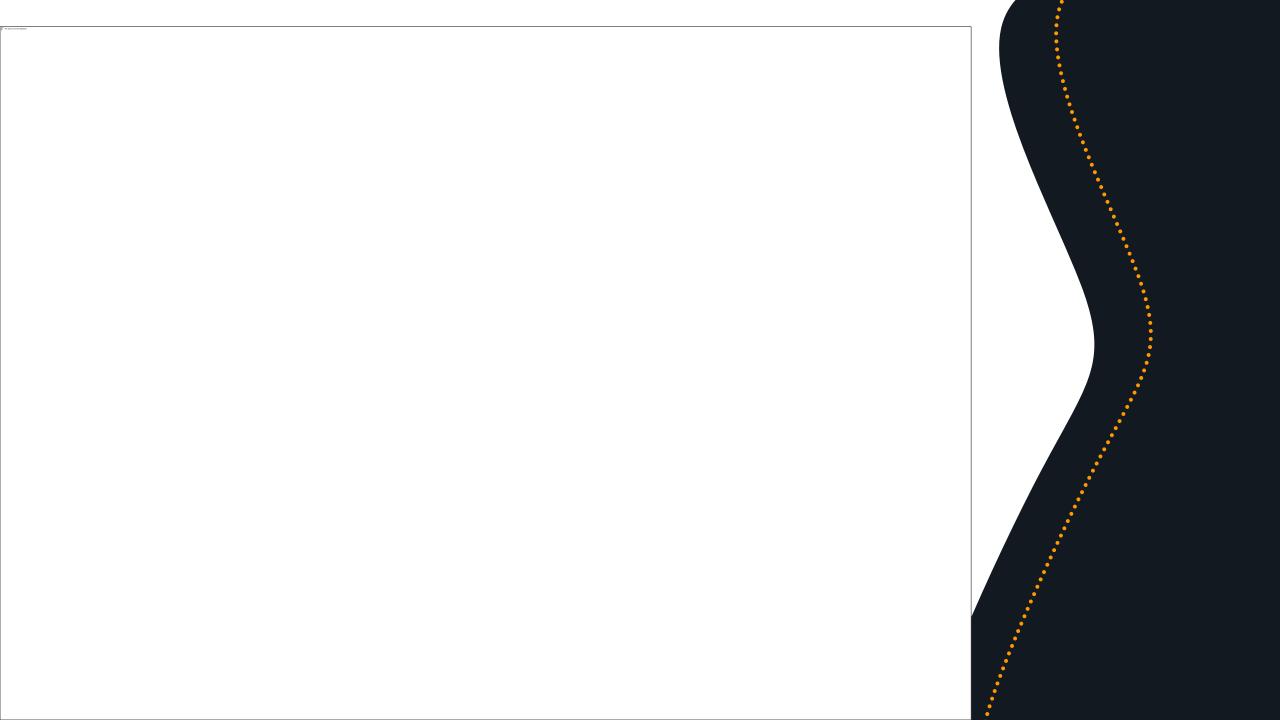




#### **RUNNING DECISION TREE ON BUSINESS QUESTION: 2**



**BUSINESS QUESTION 3:PREDICT WHICH RATING CUSTOMER WILL GIVE 5, 4** OR BELOW 4 (3-1)?



### **RUNNING DECISION-TRESS ON BUSINESS QUESTION 3:**

