

Introduction to Data Science

Supermarket Sales

Introduction

Objective

This dataset of sales observations aims to predict customer satisfaction levels for two groups: members and normal visitors. Satisfaction is categorized into two levels: low (0) and high (1). Key influencing factors include Quantity, Total, Gross Income, and Unit Price. The goal is to use Rating as the predicted variable, as customer satisfaction helps identify potential revenue generators.

Data Source

The dataset of 1000 supermarket, Supermarket_sales.xlsx, was sourced from Kaggle to analyze customer satisfaction trends. Categorizing the rating variable provides clear insights into overall satisfaction, aiding in sustainable growth analysis and developing effective marketing strategies.

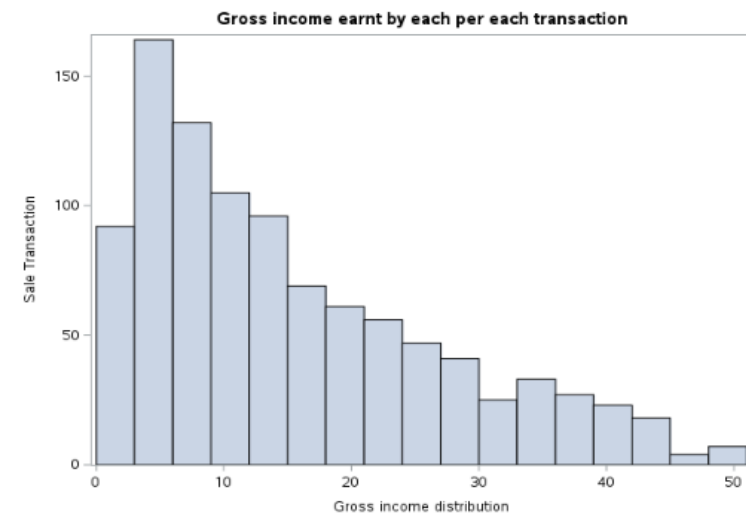
Sn o	Variable	Description
1	Customer_Type	Member, Normal
2	Gender	Male, Female
3	Product_line	Product_line, Health and beauty, Electronic accessories, Home and lifestyle, Sports and travel, Food and beverages, Fashion accessories
4	Unit_Price	Price of a single unit purchased
5	Quantity	Number of items purchased
6	Tax	Tax charged on the purchase, calculated as $(\text{Unit Price} * \text{Quantity}) * 5\%$
7	Total	Total revenue for a transaction, calculated as $(\text{Unit Price} * \text{Quantity}) + \text{Tax}$
8	Date	Date of purchase
9	Time	Time of purchase
10	Payment	Mode of Payment
11	Cogs	Cost of goods purchased
12	Gross_income	Total revenue minus COGS (Total - Cogs)

13	Rating	0 (Unsatisfied), 1 (Satisfied)
----	--------	--------------------------------

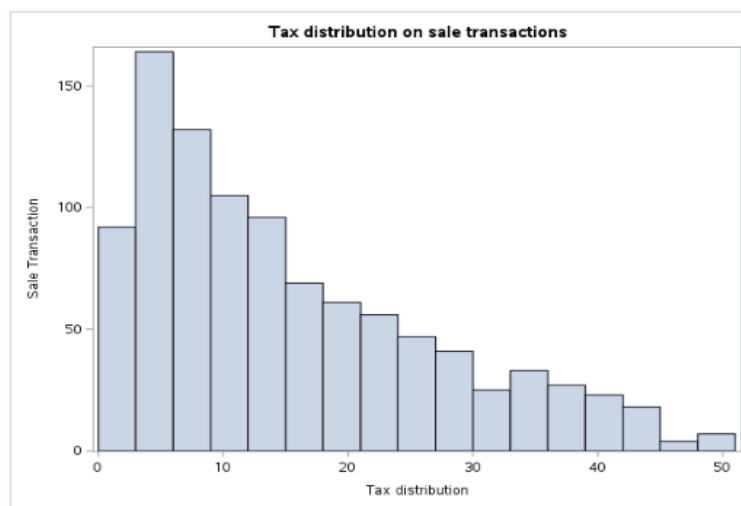
Data Partition

Customers with satisfaction scores between 0-7 are categorized as low-level, while scores above 7 are high-level. The data is split 80-20 for training and testing, ensuring a balanced dataset without the need for over-sampling or under-sampling.

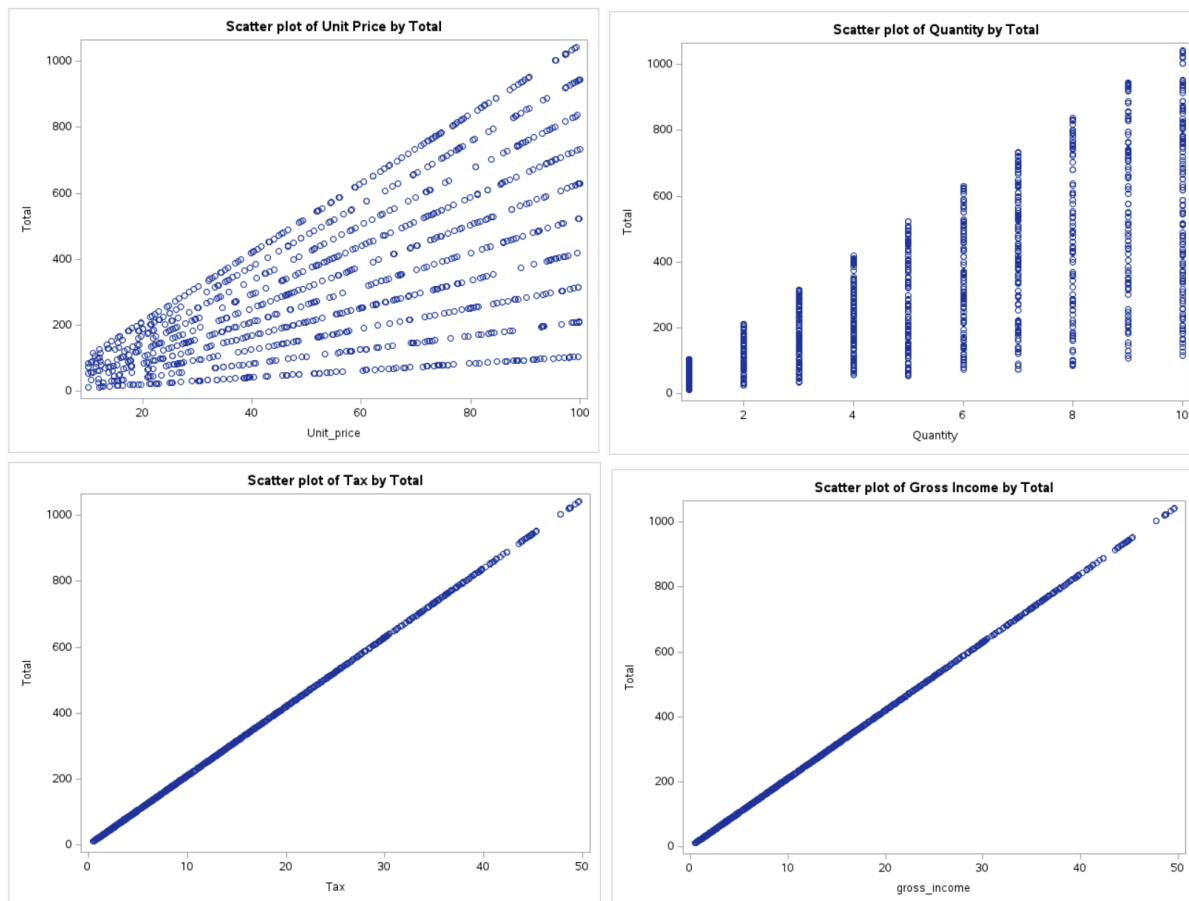
Graphs and Summary Statistics



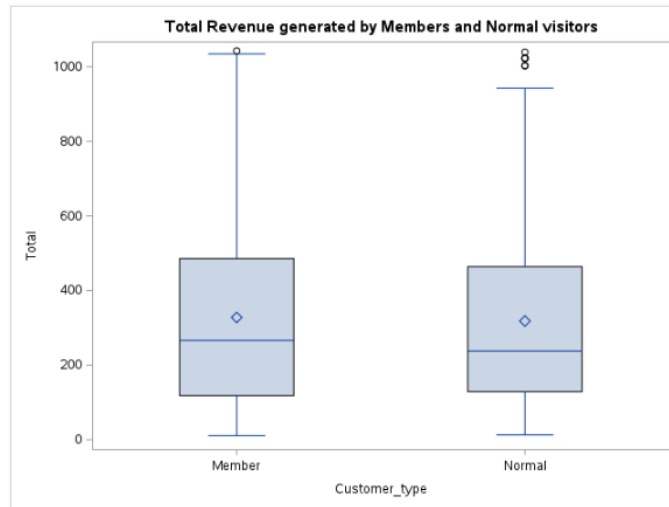
The histogram shows the distribution of sales transactions by gross income. It is right-skewed, indicating that high-gross-income transactions are rare. Most transactions fall within the lowest category, generating a gross income between 0 and 10.



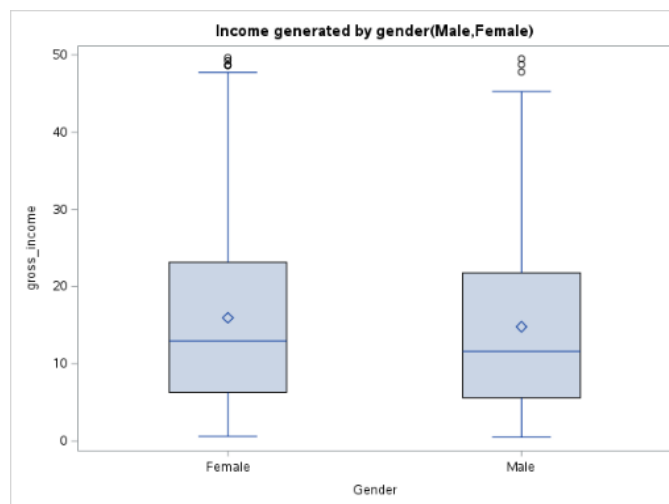
The tax distribution histogram mirrors the income distribution, showing a right skew. This suggests that most transactions have low tax amounts, with higher tax transactions being rare. The majority of transactions fall within the 0 to 10 tax range.



The scatter plots reveal a clear linear relationship between Total and other variables, such as Unit Price, Tax, Gross Income, and Quantity. As these variables increase, Total also rises. This suggests that higher values for Unit Price, Tax, Quantity, and Gross Income are positively correlated with an increased Total. Given the strong correlation observed between Total and the other variables in the plots, we chose Total as the primary variable for our models. Since the calculation of Total inherently includes contributions from Unit Price, Quantity, Tax, and Gross Income, incorporating Total alone is sufficient for our analysis.



Both member and normal customers show similar mean and median revenues, with members generating slightly higher revenue. Both categories have right-skewed distributions, with members displaying a slightly higher revenue trend. The left whiskers are nearly identical, and the interquartile ranges (IQRs) indicate similar revenue distributions. Notably, normal visitors have a higher deviation due to more outliers compared to members.



The plot shows that female customers generate slightly more revenue than male customers. Both categories have right-skewed distributions, with the whiskers extending upward. Additionally, the IQR, median, mean, and outliers are similar for both male and female customers.

Principal Component Analysis

Eigenvalues of the Correlation Matrix				
	Eigenvalue	Difference	Proportion	Cumulative
1	4.92116498	3.86736475	0.6151	0.6151
2	1.05380024	0.06972484	0.1317	0.7469
3	0.98407540	0.03337384	0.1230	0.8699
4	0.95070156	0.86044373	0.1188	0.9887
5	0.09025782	0.09025782	0.0113	1.0000
6	0.00000000	0.00000000	0.0000	1.0000
7	0.00000000	0.00000000	0.0000	1.0000
8	0.00000000		0.0000	1.0000

Eigenvectors									
		Prin1	Prin2	Prin3	Prin4	Prin5	Prin6	Prin7	Prin8
Unit_price	Unit_price	0.291799	0.209941	-.712162	-.055160	0.600468	0.000000	0.000000	0.000000
Quantity	Quantity	0.324628	-.182152	0.640577	0.047066	0.669987	0.000000	0.000000	0.000000
Tax	Tax	0.449800	0.004934	0.001368	0.004786	-.218244	-.288675	-.408248	0.707107
Total	Total	0.449800	0.004934	0.001368	0.004786	-.218244	-.288675	-.408248	-.707107
Date	Date	-.013470	0.681666	0.143481	0.717318	0.004280	0.000000	0.000000	0.000000
Time	Time	-.002695	0.676742	0.248782	-.692897	-.003891	0.000000	0.000000	0.000000
cogs	cogs	0.449800	0.004934	0.001368	0.004786	-.218244	0.866025	0.000000	0.000000
gross_income	gross_income	0.449800	0.004934	0.001368	0.004786	-.218244	-.288675	0.616497	0.000000

The PCA analysis suggests that correlation is a better approach, as the variables are measured in different units. Regarding the eigenvalues, the first principal component explains 61% of the data variability, while the second accounts for 13%. Together, the first three principal components explain 86% of the total variability.

The most important features in each principal component are as follows:

- **Prin1:** Tax, Total, COGS, and Gross Income
- **Prin2:** Date and Time
- **Prin3:** Quantity
- **Prin4:** Date and Time
- **Prin5:** Quantity
- **Prin6:** COGS
- **Prin7:** Gross Income
- **Prin8:** Total and Tax

Model Selection

We have evaluated and compared the performance of several predictive models, including Logistic Regression, CART, Neural Networks, and Discriminant Analysis.

The dataset is split into an 80% training set and a 20% validation set, comprising 800 records for model building. The dependent variable is the rating, with a focus on predicting **Rating 1** as the event of interest.

The SURVEYSELECT Procedure	
Selection Method	Simple Random Sampling
Input Data Set	SALES
Random Number Seed	12345
Sampling Rate	0.8
Sample Size	800
Selection Probability	0.8
Sampling Weight	0
Output Data Set	SALES_PART

Model 1: Logistic Regression

Dummy variables are created for categorical variables such as Customer_Type, Gender, Product_Line, and Payment to fit the logistic regression model appropriately.

The LOGISTIC Procedure			Odds Ratio Estimates			
Model Information			Effect	Point Estimate	95% Wald Confidence Limits	
Data Set	WORK.SALES_TRAIN		Is_Member	0.986	0.745	1.304
Response Variable	Rating	Rating	Is_Female	1.066	0.804	1.413
Number of Response Levels	2		Is_Product_line_HB	1.122	0.686	1.834
Model	binary logit		Is_Product_line_EA	0.779	0.487	1.246
Optimization Technique	Fisher's scoring		Is_Product_line_HL	0.840	0.522	1.350
			Is_Product_line_ST	0.789	0.490	1.273
			Is_Product_line_FB	1.156	0.722	1.852
			Total	1.000	0.999	1.001
			Date	0.999	0.993	1.004
			Time	1.000	1.000	1.000
			Is_Payment_Ewallet	1.072	0.766	1.499
			Is_Payment_CC	1.137	0.804	1.609

Number of Observations Read	800
Number of Observations Used	800

Response Profile		
Ordered Value	Rating	Total Frequency
1	0	412
2	1	388

Probability modeled is Rating='1'.

Analysis of Maximum Likelihood Estimates					
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	24.3634	61.5348	0.1568	0.6922
Is_Member	1	-0.0143	0.1427	0.0100	0.9203
Is_Female	1	0.0638	0.1438	0.1957	0.6582
Is_Product_line_HB	1	0.1147	0.2511	0.2086	0.6478
Is_Product_line_EA	1	-0.2499	0.2397	1.0869	0.2972
Is_Product_line_HL	1	-0.1745	0.2423	0.5185	0.4715
Is_Product_line_ST	1	-0.2364	0.2438	0.9404	0.3322
Is_Product_line_FB	1	0.1453	0.2402	0.3657	0.5453
Total	1	0.000058	0.000285	0.0409	0.8398
Date	1	-0.00111	0.00285	0.1520	0.6966
Time	1	-8.34E-6	6.201E-6	1.8085	0.1787
Is_Payment_Ewallet	1	0.0692	0.1714	0.1628	0.6866
Is_Payment_CC	1	0.1287	0.1770	0.5283	0.4673
cogs	0	0	.	.	.

Testing Global Null Hypothesis: BETA=0			
Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	7.6143	12	0.8145
Score	7.5871	12	0.8165
Wald	7.5298	12	0.8207

Association of Predicted Probabilities and Observed Responses			
Percent Concordant	55.7	Somers' D	0.115
Percent Discordant	44.3	Gamma	0.115
Percent Tied	0.0	Tau-a	0.057
Pairs	159856	c	0.557

Fit Statistics for SCORE Data										
Data Set	Total Frequency	Log Likelihood	Error Rate	AIC	AICC	BIC	SC	R-Square	Max-Rescaled R-Square	AUC
WORK.SALES_TRAIN	800	-550.4	0.4625	1126.701	1127.164	1187.601	1187.601	0.009473	0.012634	0.557289
WORK.SALES_VALID	200	-138.5	0.4750	303.0507	305.0077	345.9288	345.9288	-0.00387	-0.00517	0.516732

The FREQ Procedure					The FREQ Procedure				
Frequency Row Pct	Table of Rating by I_Rating				Frequency Row Pct	Table of Rating by I_Rating			
	I_Rating(Into: Rating)					I_Rating(Into: Rating)			
	Rating(Rating)	0	1	Total		Rating(Rating)	0	1	Total
	0	269 65.29	143 34.71	412		0	61 57.01	46 42.99	107
	1	227 58.51	161 41.49	388		1	49 52.69	44 47.31	93
	Total	496	304	800		Total	110	90	200

Stepwise Selection Methods

Forward: Only the intercept is included, as none of the variables tested for inclusion meet the statistical significance criterion at the 0.05 level.

Backward: Similar results, with the elimination of the Total, Gross Income, and COGS variables.

Stepwise: Only the intercept is included, as no variables tested for inclusion meet the statistical significance criterion at the 0.05 level.

Logistic Regression Final Model

The final model was built using the variables suggested by the backward selection method, ensuring that all dummy variables of a variable are retained. All measures in the final model are identical to those in the initial model.

	Training				Validation			
	AUC	Error Rate	Sensitivity	Specificity	AUC	Error Rate	Sensitivity	Specificity
Logistic Regression Model	0.56	46%	42%	65%	0.52	48%	47%	57%
Logistic Regression Final Model	0.56	46%	42%	65%	0.52	48%	47%	57%

Evaluation of Logistic regression model

- Likelihood Ratio is 0.81 (p-value > 0.05), indicating the model is not statistically significant.
- Variable Significance, Is_Member, Total, Date, Time p-values less than 0.05, suggesting they are statistically significant predictors in this model.
- AUC (Area Under the Curve) is 0.56, which is below the desired threshold of 1. This indicates below-average performance, though within an acceptable range for initial analysis.

Performance Metrics

- Error Rate: High error rates for both model fit and accuracy indicate poor classification performance.
- Sensitivity and Specificity:
 - Training Set: Sensitivity: 42%, Specificity: 65%
 - Validation Set: Sensitivity improved to 47%, while specificity decreased to 57%.

Conclusion

The model demonstrates a high error rate in both its fit and accuracy, indicating poor classification performance. The Area Under the Curve (AUC) remains consistent at 0.56 across both models, reflecting below-average predictive capability.

Notably, the AIC and BIC values are significantly lower in the model accuracy assessment compared to the initial fit, suggesting slight improvements. Sensitivity and specificity metrics further highlight the model's limitations: on the training set, sensitivity is 42% and specificity is 65%, while validation results show a modest increase in sensitivity to 47% and a decrease in specificity to 57%.

There is no evidence of overfitting, but overall, the model's performance falls below expectations.

Model 2: CART

CART (Classification and Regression Trees) model, which is well-suited for our objective of handling a classification event. Given this classification focus, the model was developed using two key measures, the Gini Index and Entropy.

GINI

Model Information	
Split Criterion Used	Gini
Pruning Method	Cost-Complexity
Subtree Evaluation Criterion	Cost-Complexity
Number of Branches	2
Maximum Tree Depth Requested	10
Maximum Tree Depth Achieved	10
Tree Depth	0
Number of Leaves Before Pruning	142
Number of Leaves After Pruning	1
Model Event Level	1

Number of Observations Read	1000
Number of Observations Used	1000
Number of Training Observations Used	800
Number of Validation Observations Used	200

ENTROPY

Model Information	
Split Criterion Used	Entropy
Pruning Method	Cost-Complexity
Subtree Evaluation Criterion	Cost-Complexity
Number of Branches	2
Maximum Tree Depth Requested	10
Maximum Tree Depth Achieved	10
Tree Depth	10
Number of Leaves Before Pruning	122
Number of Leaves After Pruning	85
Model Event Level	1

Number of Observations Read	1000
Number of Observations Used	1000
Number of Training Observations Used	800
Number of Validation Observations Used	200

Model	Training				Validation			
	AUC	Error Rates	Sensitivity	Specificity	AUC	Error Rates	Sensitivity	Specificity
Gini	0.5	48.50%	0%	100%	0.5	48.50%	0%	100%
Entropy	0.89	21.25%	85%	73%	0.55	45.50%	62%	47.66%

The Entropy model produced a larger and more complex tree, resulting in 85 leaves after pruning. This complexity allowed it to capture more intricate patterns within the data. In contrast, the Gini model generated an overly simplistic tree with only one leaf, failing to capture any meaningful patterns.

The HPSPLIT Procedure

Confusion Matrices				
	Actual	Predicted		Error Rate
		0	1	
Training	0	300	112	0.2718
	1	58	330	0.1495
Validation	0	51	56	0.5234
	1	35	58	0.3763

Fit Statistics for Selected Tree									
	N Leaves	ASE	Mis-class	Sensitivity	Specificity	Entropy	Gini	RSS	AUC
Training	85	0.1334	0.2125	0.8505	0.7282	0.5667	0.2669	213.5	0.8906
Validation	85	0.3350	0.4550	0.6237	0.4766	0.7280	0.3535	134.0	0.5503

Performance Metrics

- Error Rate: High error rates for model fit is 0.2125, whereas for the model accuracy is increased to 0.4550 indicating poor classification performance.
- Sensitivity and Specificity:
 - Training Set: Sensitivity: 85.05%, Specificity: 72.82%
 - Validation Set: Sensitivity: 62.37%, Specificity: 47.66%

Model Fit Indicators

- AUC on the training data is 0.8906 and, on the test, data is 0.5503. The significant drop in AUC from training to test data suggests overfitting.

Conclusion

The sensitivity and specificity are high on the training set, but drop significantly on the validation set, suggesting poor model performance on test data. The AUC is 0.8906 on the training data but drops to 0.5503 on the test data, showing overfitting.

Model 3: Neural Networks

Models	Hidden Layers	Neuron in each Layer	Training			Valdiation		
			Error Rate	Sensitivity	Specificity	Error Rate	Sensitivity	Specificity
1	1	8	45%	47.68%	61.89%	43%	52.69%	60.75%
2	1	20	45%	49.74%	59.22%	44%	52.69%	58.88%
3	1	25	44%	44.59%	66.99%	48%	44.09%	58.88%
4	1	30	45%	45.10%	64.08%	43%	51.61%	61.68%
5	1	35	45%	42.78%	65.53%	5%	38.71%	58.88%
6	1	40	44%	46.13%	64.32%	48%	43.01%	59.81%
7	3	8	41%	36.60%	79.61%	46%	30.11%	74.77%

Performance after Adding 3 hidden layer with 8 neurons in each

Train: Misclassification Rate	Valid: Misclassification Rate
0.4125	0.4600

Performance Metrics

Error Rate: The misclassification rate has reduced for training (0.41) and slight change in validation sets (0.46) as compared to previous model.

Sensitivity and Specificity:

Training Set: Sensitivity: 36.60%,
Specificity: 79.61%

Validation Set: Sensitivity:
30.11%, Specificity: 74.77%

In Model 7, compared to Model 1, the sensitivity decreased in both the training and validation sets. However, the specificity improved in both sets. Therefore, adding 3 hidden layers with the same number of neurons enhanced specificity but reduced sensitivity, resulting in a moderate

Thus, we are choosing Model 7 as final model in neural networks.

Model 4: Discriminant Analysis

Class Level Information					
Rating	Variable Name	Frequency	Weight	Proportion	Prior Probability
0	0	412	412.0000	0.515000	0.515000
1	1	388	388.0000	0.485000	0.485000

412 observations belong to class 0, ratings unsatisfactory. 388 observations belong to class 1, where ratings are satisfactory. Corresponds to 52% and 48%.

Discriminant Analysis	Training			Valdiation		
	Error Rate	Sensitivity	Specificity	Error Rate	Sensitivity	Specificity
Without Prior	48%	14.69%	86.41%	49%	8.60%	90.65%
With Prior(60-40)	40%	100%	0%	40%	100%	0%

Frequency
Row Pct

Table of Rating by I_Rating			
Rating(Rating)	I_Rating(Into: Rating)		
	0	1	Total
0	328 79.61	84 20.39	412
1	246 63.40	142 36.60	388
Total	574	226	800

Frequency
Row Pct

The FREQ Procedure

Table of Rating by I_Rating			
Rating(Rating)	I_Rating(Into: Rating)		
	0	1	Total
0	80 74.77	27 25.23	107
1	65 69.89	28 30.11	93
Total	145	55	200

Number of Observations and Percent Classified into Rating			
From Rating	0	1	Total
0	356 86.41	56 13.59	412 100.00
1	331 85.31	57 14.69	388 100.00
Total	687 85.88	113 14.13	800 100.00
Priors	0.515	0.485	

Number of Observations and Percent Classified into Rating			
From Rating	0	1	Total
0	97 90.65	10 9.35	107 100.00
1	85 91.40	8 8.60	93 100.00
Total	182 91.00	18 9.00	200 100.00
Priors	0.515	0.485	

Error Count Estimates for Rating			
	0	1	Total
Rate	0.1359	0.8531	0.4838
Priors	0.5150	0.4850	

Error Count Estimates for Rating			
	0	1	Total
Rate	0.0935	0.9140	0.4914
Priors	0.5150	0.4850	

Conclusion

The model performed poorly in terms of sensitivity but showed strong results for specificity. After applying the prior proportion, the model's performance declined, as it struggled to accurately identify specificity.

Model Comparison and Evaluation

Model	Training				Validation			
	AUC	Error Rate	Sensitivity	Specificity	AUC	Error Rate	Sensitivity	Specificity
Logistic Regression Model	0.56	46%	42%	65%	0.52	48%	47%	57%
CART Model(Entropy)	0.89	21%	85%	73%	0.55	46%	62%	48%
Neural Network (3 hidden layer with 8 Neurons each)		41%	37%	80%		46%	30%	75%
Discriminant Analysis		48%	15%	86%		49%	9%	91%

The CART model performs well on the training data, achieving 85% sensitivity, 73% specificity, and the lowest error rate of 21%. It also performed reasonably on the validation data, with 62% sensitivity, 48% specificity, and a 46% error rate.

While the model demonstrates a strong ability to fit the training data, making it a suitable choice for predicting customer satisfaction, the weaker performance on the validation set and the higher error rate suggest signs of overfitting.

Conclusion

Practical Significance and Implications

- Targeted Marketing: Insights on spending behaviors and satisfaction levels of any customer that can be used to create targeted promotions for specific customers including those who are members and those who are not.

- Operational Efficiency: Analysis of date and time trends might be useful for the proper staffing of employees and the proper supply chain management during peak hours.

Inferences

- Key factors influencing customer satisfaction through data-driven analysis - Date and Time.
- Models can predict satisfaction trends across similar customer demographics.

Shortcomings

- Potential biases in the data collection process could impact the generalizability of results.
- The dataset size is relatively small, which might limit the robustness of the findings.

Suggestions for Improvement

- Increasing the dataset size and incorporating the latest data, as we currently have data from 2019, will contribute to greater diversity and help capture a broader range of customer behaviors.
- Incorporate additional variables, such as customer reviews, to improve the accuracy and relevance of predictions.

References

<https://www.kaggle.com/code/aryantiwari123/supermarket-sales-prediction/inpu>

SAS Code



Group8_Project_Final_
Code.sas

```

1  /* Import */
2
3  proc import out=supermarket datafile="/home/u63814242/sasuser.v94/Data Science/Project/supermarket_sales.xlsx"
4  dbms=xlsx replace;
5  run;
6
7  /* Plots */
8
9  /*Histogram*/
10 proc sgplot data=supermarket;
11     histogram gross_income/scale=count;
12     title "Gross income earned by each per each transaction";
13     xaxis label="Gross income distribution";
14     yaxis label="Sale Transaction";
15 run;
16
17 proc sgplot data=supermarket;
18     histogram Tax/scale=count;
19     title "Tax distribution on sale transactions";
20     xaxis label="Tax distribution";
21     yaxis label="Sale Transaction";
22 run;
23
24 /*Side by Side boxplot*/
25 proc sgplot data=supermarket;
26     vbox Total /category=Customer_type;
27     title "Total Revenue generated by Members and Normal visitors";
28 run;
29
30 proc sgplot data=supermarket;
31     vbox gross_income /category=gender;
32     title "Income generated by gender(Male,Female)";
33 run;

```

```

35 /*Scatter Plot*/
36 proc sgplot data=supermarket;
37     scatter x=unit_price y=gross_income;
38 run;
39
40 proc sgplot data=supermarket;
41     scatter x=quantity y=gross_income;
42 run;
43
44 proc sgplot data=supermarket;
45     scatter x=tax y=gross_income;
46 run;
47
48 proc sgplot data=supermarket;
49     scatter x=total y=gross_income;
50 run;
51
52 /*numerical summary(descriptive statistics)*/
53 proc means data=supermarket n mean median skew maxdec=2;
54     var gross_income Tax ;
55 run;
56
57 /*numerical summary(descriptive statistics)*/
58 proc means data=supermarket n mean median skew maxdec=2;
59     var total cogs ;
60 run;
61
62 /*Principal component analysis*/
63 data supermarket1;
64     set supermarket;
65     drop Rating;
66 run;
67
68 proc princomp data=supermarket1 out=supermarket2;
69 run;
70
71 /* logistic regression */
72 data supermarket;
73     set supermarket;
74     if Customer_type='Member' then Is_Member=1; else Is_Member=0;
75     if Gender='Female' then Is_Female=1; else Is_Female=0;
76     if Payment='Ewallet' then Is_Payment_Ewallet=1; else Is_Payment_Ewallet=0;
77     if Payment='Credit card' then Is_Payment_CC=1; else Is_Payment_CC=0;
78     if Product_line='Health and beauty' then Is_Product_line_HB=1; else Is_Product_line_HB=0;
79     if Product_line='Electronic accessories' then Is_Product_line_EA=1; else Is_Product_line_EA=0;
80     if Product_line='Home and lifestyle' then Is_Product_line_HL=1; else Is_Product_line_HL=0;
81     if Product_line='Sports and travel' then Is_Product_line_ST=1; else Is_Product_line_ST=0;
82     if Product_line='Food and beverages' then Is_Product_line_FB=1; else Is_Product_line_FB=0;
83 run;
84
85 proc surveyselect data=supermarket samprate=.8 method=srs outall out=supermarket_Part seed=12345;
86 run;
87
88 data supermarket_train supermarket_valid;
89     set supermarket_Part;
90     if selected=1 then output supermarket_train; else output supermarket_valid;
91 run;
92

```

```

93 proc logistic data=supermarket_train outmodel=model_logistic;
94     model Rating (event="1")= Is_Member Is_Female Is_Product_line_HB Is_Product_line_EA
95                               Is_Product_line_HL Is_Product_line_ST Is_Product_line_FB
96                               Total Date Time Is_Payment_Ewallet
97                               Is_Payment_CC cogs;
98 run;
99
100 proc logistic inmodel=model_logistic;
101     score data=supermarket_train fitstat out=supermarket_trainOut;
102     score data=supermarket_valid fitstat out=supermarket_validOut;
103 run;
104
105 proc freq data=supermarket_trainOut;
106     table Rating*i_Rating /nopercnt nocol;
107 run;
108
109 proc freq data=supermarket_validOut;
110     table Rating*i_Rating/nopercnt nocol;
111 run;
112
113
114 /* Selection method */
115 proc logistic data=supermarket_train outmodel=model1;
116     model Rating (event="1")= Is_Member Is_Female Is_Product_line_HB Is_Product_line_EA
117                               Is_Product_line_HL Is_Product_line_ST Is_Product_line_FB
118                               Total Date Time Is_Payment_Ewallet
119                               Is_Payment_CC cogs/selection=forward;
120 run; /*only intercept*/
121
122
123 proc logistic data=supermarket_train outmodel=model1;
124     model Rating (event="1")= Is_Member Is_Female Is_Product_line_HB Is_Product_line_EA
125                               Is_Product_line_HL Is_Product_line_ST Is_Product_line_FB
126                               Total Date Time Is_Payment_Ewallet
127                               Is_Payment_CC cogs/selection=backward;
128 run; /*few variables eliminated*/
129
130 proc logistic data=supermarket_train outmodel=model1;
131     model Rating (event="1")= Is_Member Is_Female Is_Product_line_HB Is_Product_line_EA
132                               Is_Product_line_HL Is_Product_line_ST Is_Product_line_FB
133                               Total Date Time Is_Payment_Ewallet
134                               Is_Payment_CC cogs/selection=stepwise;
135 run; /*only intercept*/
136
137 /*Final logistic Model*/
138 proc logistic data=supermarket_train outmodel=model_logistic_final;
139     model Rating (event="1")= Is_Member Is_Female Is_Product_line_HB Is_Product_line_EA
140                               Is_Product_line_HL Is_Product_line_ST Is_Product_line_FB
141                               Total Date Time Is_Payment_Ewallet
142                               Is_Payment_CC;
143 run;
144
145 proc logistic inmodel=model_logistic_final;
146     score data=supermarket_train fitstat out=supermarket_trainOut_final;
147     score data=supermarket_valid fitstat out=supermarket_validOut_final;
148 run;
149
150 proc freq data=supermarket_trainOut_final;
151     table Rating*i_Rating /nopercnt nocol;
152 run;
153
154 proc freq data=supermarket_validOut_final;
155     table Rating*i_Rating/nopercnt nocol;
156 run;

```

```

156
157 /*CART*/
158
159 proc hpsplit data=supermarket_part nodes=detail;
160     partition rolevar=selected(train='1' validate='0');
161     class Customer_type Gender Product_line Payment Rating;
162     model Rating(event="1")= Customer_type Gender Product_line Total Date Time Payment cogs;
163     grow gini;
164     prune cc;
165 run;
166
167 proc hpsplit data=supermarket_part nodes=detail;
168     partition rolevar=selected(train='1' validate='0');
169     class Customer_type Gender Product_line Payment Rating;
170     model Rating(event="1")= Customer_type Gender Product_line Total Date Time Payment cogs;
171     grow entropy;
172     prune cc;
173 run;
174
175 /* Neural Network */
176
177 proc hpneural data=supermarket_part;
178     partition rolevar=selected(train=1);
179     target Rating/level=nom;
180     input Customer_type Gender Product_line Payment/level=nom;
181     input Total Date Time cogs/level=int;
182     hidden 8;
183     train maxiter=1000 numtries=10;
184     id Rating selected;
185     score out=supermarket_out1;
186 run;
187
188 proc freq data=supermarket_out1;
189     table Rating*i_Rating/nopercent nocol;
190     where selected=1;
191 run;
192
193 proc freq data=supermarket_out1;
194     table Rating*i_Rating/nopercent nocol;
195     where selected=0;
196 run;
197
198 proc hpneural data=supermarket_part;
199     partition rolevar=selected(train=1);
200     target Rating/level=nom;
201     input Customer_type Gender Product_line Payment/level=nom;
202     input Total Date Time cogs/level=int;
203     hidden 8;
204     hidden 8;
205     hidden 8;
206     train maxiter=1000 numtries=10;
207     id Rating selected;
208     score out=supermarket_out2;
209 run;
210
211 proc freq data=supermarket_out2;
212     table Rating*i_Rating/nopercent nocol;
213     where selected=1;
214 run;
215
216 proc freq data=supermarket_out2;
217     table Rating*i_Rating/nopercent nocol;
218     where selected=0;
219 run;
220
221 proc hpneural data=supermarket_part;
222     partition rolevar=selected(train=1);
223     target Rating/level=nom;
224     input Customer_type Gender Product_line Payment/level=nom;
225     input Total Date Time cogs/level=int;
226     hidden 20;
227     train maxiter=1000 numtries=10;
228     id Rating selected;
229     score out=supermarket_out3;
230 run;

```



```

232 proc freq data=supermarket_out3;
233     table Rating*i_Rating/nopercent nocol;
234     where selected=1;
235 run;
236
237 proc freq data=supermarket_out3;
238     table Rating*i_Rating/nopercent nocol;
239     where selected=0;
240 run;
241 proc hpneural data=supermarket_part;
242     partition rolevar=selected(train=1);
243     target Rating/level=nom;
244     input Customer_type Gender Product_line Payment/level=nom;
245     input Total Date Time cogs/level=int;
246     hidden 40;
247     train maxiter=1000 numtries=10;
248     id Rating selected;
249     score out=supermarket_out4;
250 run;
251
252 proc freq data=supermarket_out4;
253     table Rating*i_Rating/nopercent nocol;
254     where selected=1;
255 run;
256
257 proc freq data=supermarket_out4;
258     table Rating*i_Rating/nopercent nocol;
259     where selected=0;
260 run;
261
262 /* Discriminant Analysis */
263
264 data supermarket_train supermarket_valid;
265     set supermarket_part;
266     if selected=1 then output supermarket_train; else output supermarket_valid;
267     drop Selected Customer_type Gender Product_line Total Date Time Payment cogs;
268 run;
269
270 proc discrim data=supermarket_train testdata=supermarket_valid;
271     class Rating;
272     priors proportional;
273 run;
274
275 proc discrim data=supermarket_train testdata=supermarket_valid;
276     class Rating;
277     priors '1'=.6 '0'=.4;
278 run;

```