

Project 1: Predicting Catalog Demand

Step 1: Business and Data Understanding

Provide an explanation of the key decisions that need to be made. (500 word limit)

Key Decisions:

Answer these questions

1. What decisions needs to be made?

To determine if they should send the catalog to the 250 new customers, if the expected profit exceeded \$10,000, the company will send the catalog. Otherwise, they will not send it.

2. What data is needed to inform those decisions?

the data needed to predict sales and calculate expected profit are customer segment, average number of product purchased, Score yes, and margin and cost of catalog.

Step 2: Analysis, Modeling, and Validation

Provide a description of how you set up your linear regression model, what variables you used and why, and the results of the model. Visualizations are encouraged. (500 word limit)

Important: Use the *p1-customers.xlsx* to train your linear model.

At the minimum, answer these questions:

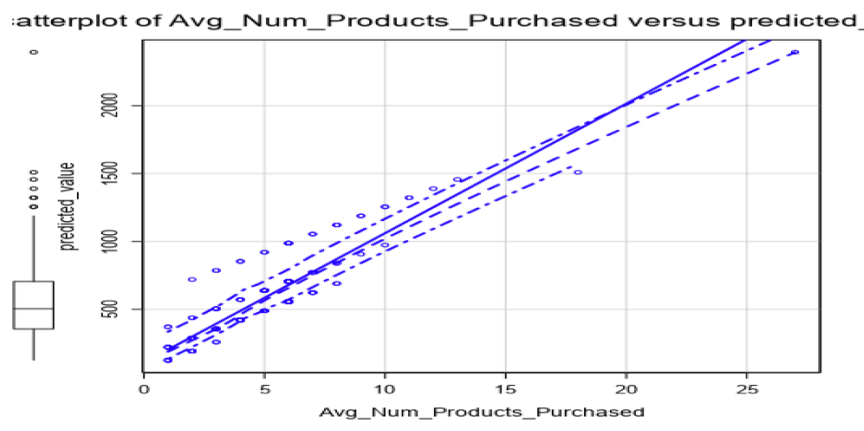
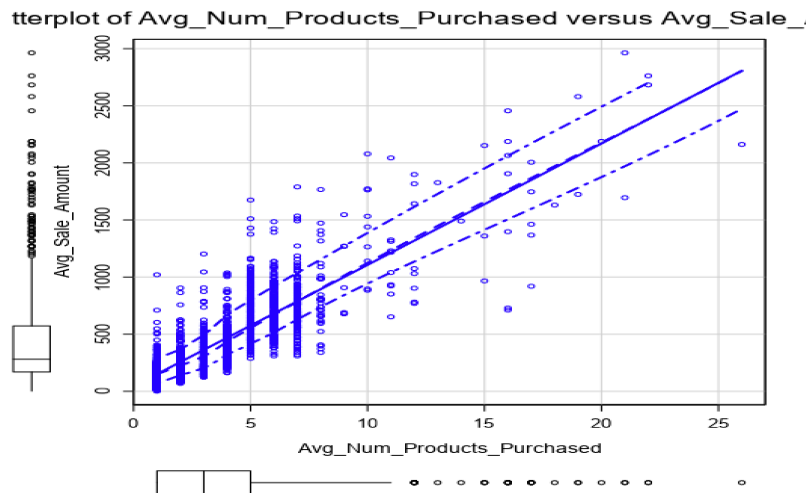
1. How and why did you select the predictor variables in your model? You must explain how your continuous predictor variables you've chosen have a linear relationship with the target variable. Please refer back to the "Multiple Linear Regression with Excel" lesson to help you explore your data and use scatterplots to search for linear relationships. You must include scatterplots in your answer.

After exploring the data, I did not select any personal information because it is not important in this model. State is always fixed CO, and Responded_To_Last_Catalog does not exist in the mailing list dataset.

Response: Avg_Sale_Amount

	Sum Sq	DF	F value	Pr(>F)
Customer_Segment	28715078.96	3	506.4	< 2.2e-16 ***
Avg_Num_Products_Purchased	36939582.5	1	1954.31	< 2.2e-16 ***
Residuals	44796869.07	2370		

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1



2. Explain why you believe your linear model is a good model. You must justify your reasoning using the statistical results that your regression model created. For each variable you selected, please justify how each variable is a good fit for your model by using the p-values and R-squared values that your model produced.

- As you can see, the P value of all variables is less than 0.05.
- R-squared is 0.8371 (strong relation).
- As a result, the model is considered a good one.

Record

Report

1

Report for Linear Model Linear_Regression_5

2

Basic Summary

3

Call:
lm(formula = Avg_Sale_Amount ~ Customer_Segment + Avg_Num_Products_Purchased, data = the.data)

4

Residuals:

5

	Min	1Q	Median	3Q	Max
	-663.8	-67.3	-1.9	70.7	971.7

6

Coefficients:

7

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	303.46	10.576	28.69	< 2.2e-16 ***
Customer_SegmentLoyalty Club Only	-149.36	8.973	-16.65	< 2.2e-16 ***
Customer_SegmentLoyalty Club and Credit Card	281.84	11.910	23.66	< 2.2e-16 ***
Customer_SegmentStore Mailing List	-245.42	9.768	-25.13	< 2.2e-16 ***
Avg_Num_Products_Purchased	66.98	1.515	44.21	< 2.2e-16 ***

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

8

Residual standard error: 137.48 on 2370 degrees of freedom

Multiple R-squared: 0.8369, Adjusted R-Squared: 0.8366

F-statistic: 3040 on 4 and 2370 degrees of freedom (DF), p-value < 2.2e-16

9

Type II ANOVA Analysis

10

Response: Avg_Sale_Amount

	Sum Sq	DF	F value	Pr(>F)
Customer_Segment	28715078.96	3	506.4	< 2.2e-16 ***
Avg_Num_Products_Purchased	36939582.5	1	1954.31	< 2.2e-16 ***
Residuals	44796869.07	2370		

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

3. What is the best linear regression equation based on the available data? Each coefficient should have no more than 2 digits after the decimal (ex: 1.28)

$$\text{Avg_Sale_Amount} = 303.46 - 149.36 \times (\text{Loyalty Club Only}) + 281.84 \times (\text{Loyalty Club and Credit Card}) - 245.42 \times (\text{Store Mailing List}) + 0 \times (\text{Credit Card Only}) + 66.98 \times (\text{Avg_Num_Products_Purchased})$$

Important: The regression equation should be in the form:

$$Y = \text{Intercept} + b_1 * \text{Variable}_1 + b_2 * \text{Variable}_2 + b_3 * \text{Variable}_3 + \dots$$

For example:
$$Y = 482.24 + 28.83 * \text{Loan_Status} - 159 * \text{Income} + 49 (\text{If Type: Credit Card}) - 90 (\text{If Type: Mortgage}) + 0 (\text{If Type: Cash})$$

Note that we **must** include the 0 coefficient for the type Cash.

Note: For students using software other than Alteryx, if you decide to use Customer Segment as one of your predictor variables, please set the base case to Credit Card Only.

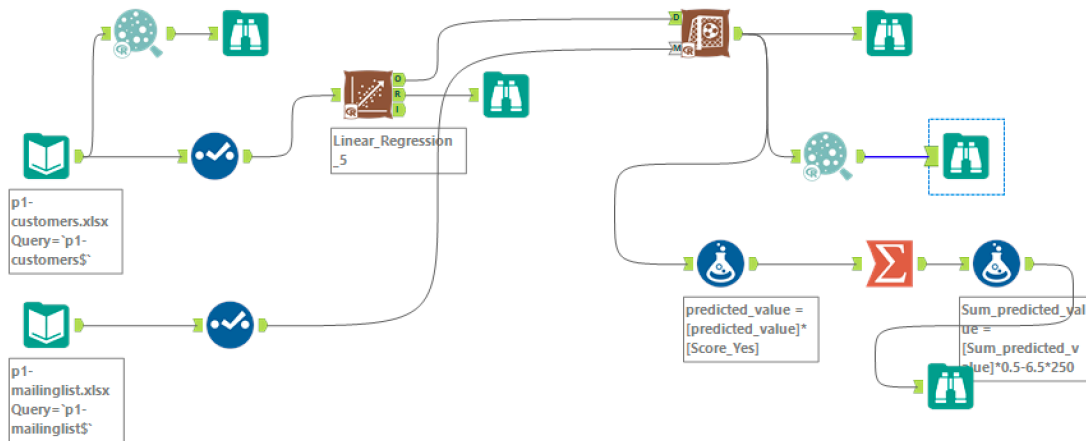
Step 3: Presentation/Visualization

Use your model results to provide a recommendation. (500 word limit)

At the minimum, answer these questions:

1. What is your recommendation? Should the company send the catalog to these 250 customers?
 - The expected profit is \$21,987.
 - The company should send the catalog to the 250 new customers because the expected profit is higher than \$10,000.
2. How did you come up with your recommendation? (Please explain your process so reviewers can give you feedback on your process)

Using linear regression model, the expected revenue from each customer is determined by multiplying expected sale amount with Score_Yes value, then multiply the sum by gross margin of 50% - \$6.50 (catalog price) and multiply it by 250 (new customers).



3. What is the expected profit from the new catalog (assuming the catalog is sent to these 250 customers)?

The predicted Revenue = SUM (Predicted Avg_Sale_Amount * Score_Yes) = \$47,224.

The expected profit = $\$47,224 \times 0.5 - \$6.5 \times 250 = \$21,987$.