

# Data Report – Predicting Catalog Demand

by Travis Gillespie

## Table of Contents

- [Introduction](#)
- [Business and Data Understanding](#)
  - [Key Decisions](#)
- [Analysis, Modeling, and Validation](#)
  - [Question 1](#)
  - [Question 2](#)
  - [Question 3](#)
- [Presentation/Visualization](#)
  - [Question 1](#)
  - [Question 2](#)
  - [Question 3](#)

## Introduction

You recently started working for a company that manufactures and sells high-end home goods. Last year the company sent out its first print catalog, and is preparing to send out this year's catalog in the coming months. The company has 250 new customers from their mailing list that they want to send the catalog to.

Your manager has been asked to determine how much profit the company can expect from sending a catalog to these customers. You, the business analyst, are assigned to help your manager run the numbers. While fairly knowledgeable about data analysis, your manager is not very familiar with predictive models.

You've been asked to predict the expected profit from these 250 new customers. Management does not want to send the catalog out to these new customers unless the expected profit contribution exceeds \$10,000.

## Business and Data Understanding

### Key Decisions

## Question 1

*What decisions needs to be made?*

Whether or not to send this year's company catalog to new customers; dependant on the profit exceeding \$10,000.

## Question 2

*What data is needed to inform those decisions?*

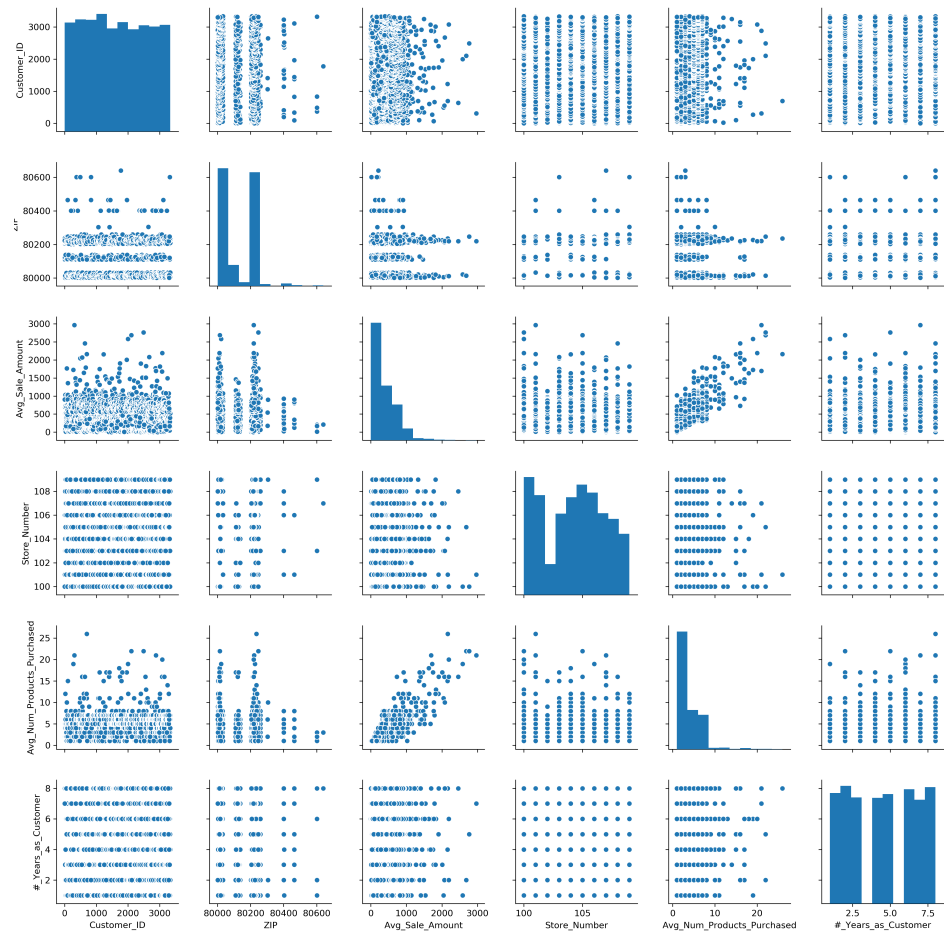
- The expected revenue from 250 new customers.
- The probability a custome will buy the catalog.
- Number of catalogs purchased.
- Categorical variables converted to dummy variables

## Analysis, Modeling, and Validation

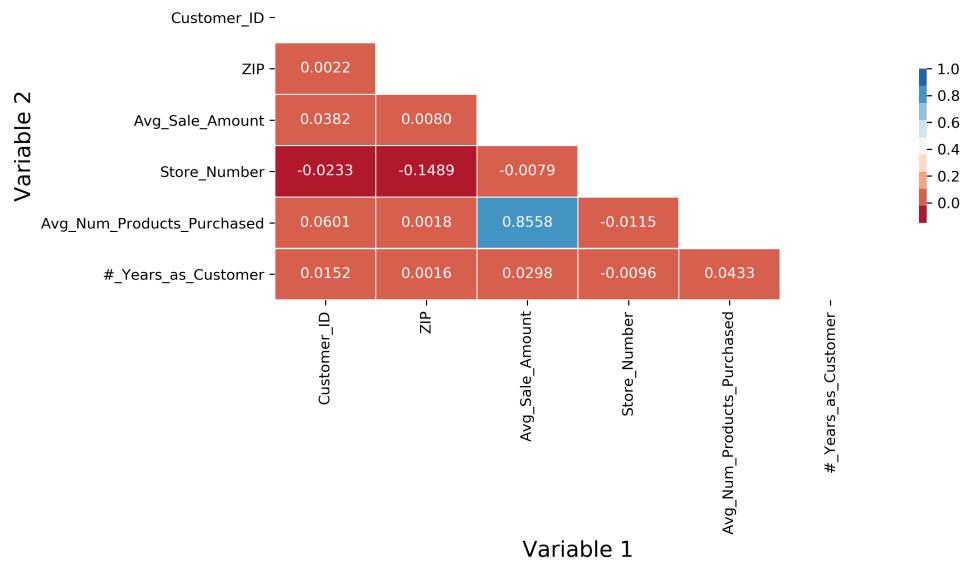
### Question 1

*How and why did you select the predictor variables in your model?*

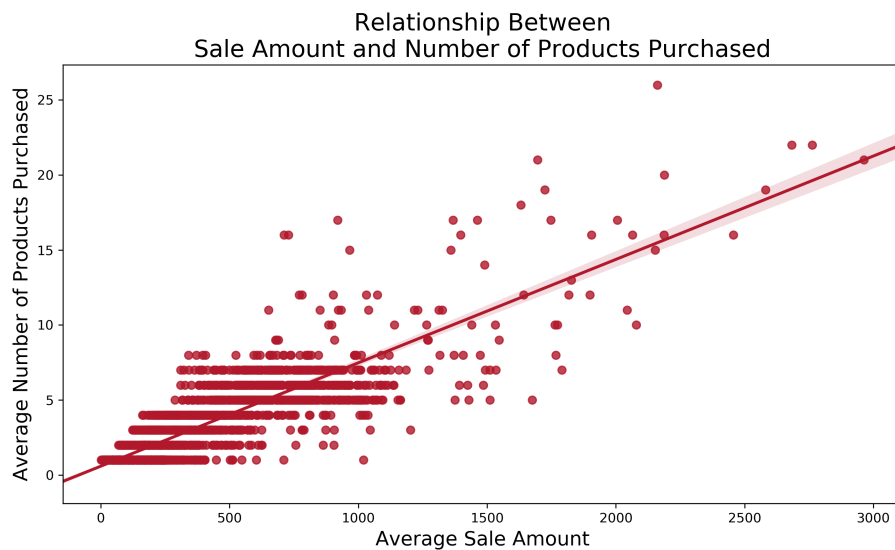
Using linear regression models I was able to assess which predictor variables have the strogest correlation.



Heatmap



The pair plot and Pearson Correlation matrix (above) suggests *Avg\_Sale\_Amount* and *Avg\_Num\_Products\_Purchased* have a strong positive correlation of approximately 0.8558.



The scatter plot provides a cleaner display of the relationship between *Avg\_Sale\_Amount* and *Avg\_Num\_Products\_Purchased* as a positive correlation.

## Question 2

*Explain why you believe your linear model is a good model.*

### OLS Regression Results

<b>Dep. Variable:</b>	Avg_Sale_Amount	<b>R-squared:</b>	0.837
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.837
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	3040.
<b>Date:</b>	Thu, 03 Jan 2019	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	21:22:38	<b>Log-Likelihood:</b>	-15061.
<b>No. Observations:</b>	2375	<b>AIC:</b>	3.013e+04
<b>Df Residuals:</b>	2370	<b>BIC:</b>	3.016e+04
<b>Df Model:</b>	4		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	303.4635	10.576	28.694	0.000	282.725	324.202
Avg_Num_Products_Purchased	66.9762	1.515	44.208	0.000	64.005	69.947
Customer_Segment_Loyalty_Club_Only	-149.3557	8.973	-16.645	0.000	-166.951	-131.760
Customer_Segment_Loyalty_Club_and_Credit_Card	281.8388	11.910	23.664	0.000	258.484	305.194
Customer_Segment_Store_Mailing_List	-245.4177	9.768	-25.125	0.000	-264.572	-226.263

<b>Omnibus:</b>	359.638	<b>Durbin-Watson:</b>	2.045
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	4770.580
<b>Skew:</b>	0.232	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	9.928	<b>Cond. No.</b>	25.0

First categorical variables were converted to dummy variables rather than assign random numbers to each category for the model. Thus preventing an erroneous relationship between the target variable and the category variable(s) due to arbitrarily assigned value(s).

From there a simple regression analysis was conducted using the dummy variables. The Adjusted R-squared value is 0.837, indicating a strong positive correlation between my predictor variables (listed below):

- Avg\_Num\_Products\_Purchased
- Customer\_Segment\_Loyalty\_Club\_Only
- Customer\_Segment\_Loyalty\_Club\_and\_Credit\_Card
- Customer\_Segment\_Store\_Mailing\_List

Note: The p-value is truncated after the third decimal place. Although, I cannot determine if p-value is exactly equal zero, the data suggests this a very small p-value (e.g.  $p < 0.05$ , and  $p < 0.001$ ) which is statistically significant.

### Question 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).*

$$\begin{aligned} \text{PredictedAverageSaleAmount} = & 303.46 + \\ & (66.98 \times \text{AvgNumProductsPurchased}) + \\ & (-149.36 \times \text{CustomerSegmentLoyaltyClubOnly}) + \\ & (281.84 \times \text{CustomerSegmentLoyaltyClubAndCreditCard}) + \\ & (-245.42 \times \text{CustomerSegmentStoreMailingList}) \end{aligned}$$

## Presentation/Visualization

### Question 1

*What is your recommendation? Should the company send the catalog to these 250 customers?*

Yes. The company should send the catalog out to the 250 new customers in their mailing list.

### Question 2

*How did you come up with your recommendation?*

First I used my multiple linear regression model to calculate the *Predicted\_Average\_Sale\_Amount*, then I calculated *Predicted\_Revenue*, and finally calculated *Predicted\_Profit*. I placed these values for these variables in corresponding columns in the dataset labeled [df\\_mailingList\\_dummies](#)

(./assets/data/df\_mailingList\_dummies.csv). To find the overall value for these variables I calculated the sum of their corresponding column within the dataset. The overall values were finally written to a csv titled [df\\_overallValues](#) (./assets/data/df\_overallValues.csv).

Further information on how each variable was calculated is listed below, while the actual calculations reside in the [Data Wrangling](#) (./Data%20Wrangling.ipynb#three) file.

- *Predicted\_Average\_Sale\_Amount* is calculated by following the *PredictedAverageSaleAmount* formula above. Substitute the formula's variables with corresponding column values for each of the 250 customers in the mailing list dataset. Finally sum the *Predicted\_Average\_Sale\_Amount* values. Example formula below:

$$\begin{aligned} \text{PredictedAverageSaleAmount} = & 303.46 + \\ & (66.98 \times \text{AvgNumProductsPurchased}) \\ & + \\ & (-149.36 \times \text{CustomerSegmentLoyaltyClubOnly}) + \\ & (281.84 \times \text{CustomerSegmentLoyaltyClubAndCreditCard}) + \\ & (-245.42 \times \text{CustomerSegmentStoreMailingList}) \end{aligned}$$

- *Predicted\_Revenue* is calculated by finding the product of *Average\_Sale\_Amount* and *Score\_Yes* (the probability a customer will respond and make a purchase), then taking the sum of all those values. Example formula below:

$$\text{Predicted\_Revenue} = \text{Predicted\_Average\_Sale\_Amount} * \text{Score\_Yes}$$

- *Predicted\_Profit* is calculated by subtracting the catalog cost (given \$6.50) from the product of *Predicted\_Revenue* and average gross margin (which is a given value of 50%). Example formula below:

$$\text{Predicted\_Profit} = (0.5 * \text{Predicted\_Revenue}) - 6.5$$

### Question 3

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

As shown in the image below, the *Overall Predicted Profit* is approximately \$21,987.96. This is more than double the \$10,000 breaking point.

	Variable Name	Values
0	Overall Predicted Average Sale Amount	\$138,295.16
1	Overall Predicted Revenue	\$47,225.91
2	Overall Predicted Profit	\$21,987.96

In [ ]: