

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 need to be made?

Do we send out the catalogues to the new customers or not?

2. What data is needed to inform those decisions?

What is the expected profit from the 250 new customers?

Step 2: Analysis, Modeling, and Validation

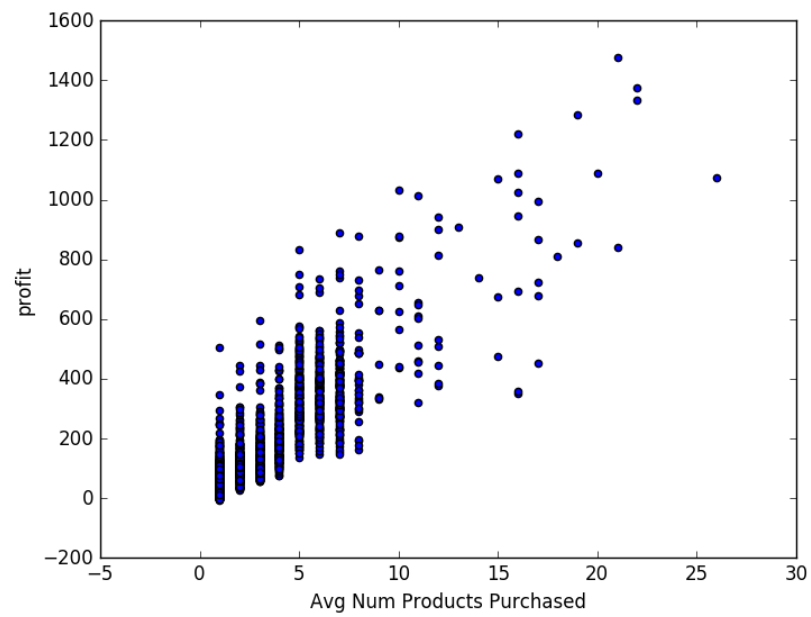
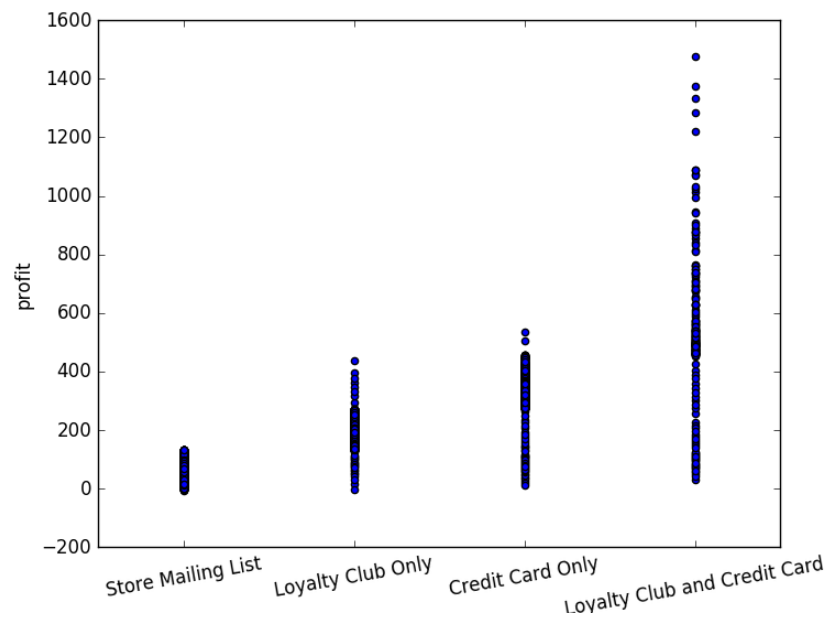
Provide a description of how you set up your linear regression model, what variables you used and why, and the results of the model. Visualizations are encouraged. (500 word limit)

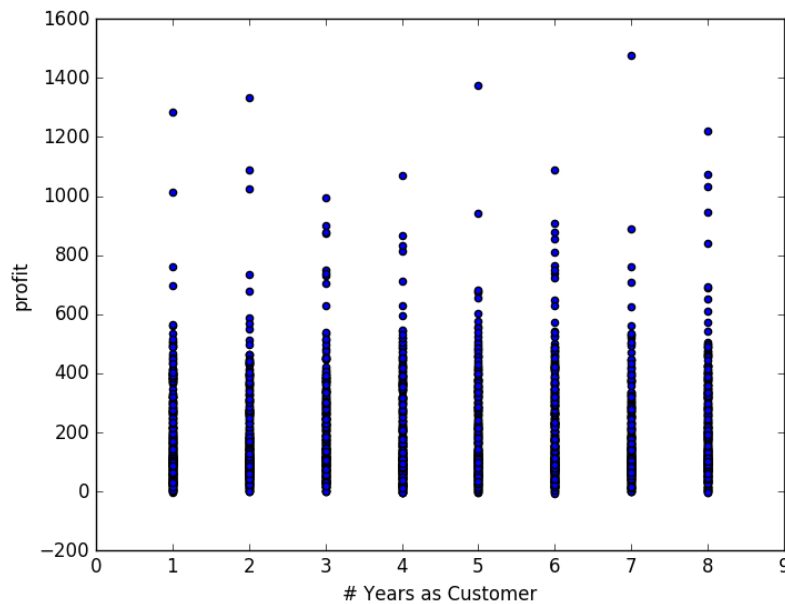
Important: Use the *p1-customers.xlsx* to train your linear model.

At the minimum, answer these questions:

1. How and why did you select the [predictor variables \(see supplementary text\)](#) in your model? You must explain how your continuous predictor variables you've chosen have a linear relationship with the target variable. Please refer to this [lesson](#) to help you explore your data and use scatterplots to search for linear relationships. You must include scatterplots in your answer.

I selected customer segment and avg num products purchased. These appear to have a linear relationship with profit from the scatterplots. I also looked at avg num years as a customer, and this did not appear related to profit. The p-value from a linear regression showed years as customer had a p-value of 0.104 (> 0.05), which means we accept the null hypothesis, which is that the predictor variable has no relationship to the target.





- 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.

The model has a decent adjusted r-squared value (0.837) and the p-values of the coefficients are 0.0, meaning they have a meaningful relationship to the target variable. The regression summary from statsmodels (Python) follows:

```

OLS Regression Results
=====
Dep. Variable:      profit    R-squared:      0.837
Model:              OLS      Adj. R-squared:    0.837
Method:             Least Squares    F-statistic:    3040.
Date:               Fri, 16 Dec 2016    Prob (F-statistic): 0.00
Time:               12:10:48    Log-Likelihood: -13415.
No. Observations:   2375    AIC:              2.684e+04
Df Residuals:       2370    BIC:              2.687e+04
Df Model:           4
Covariance Type:    nonrobust
=====
                    coef    std err          t      P>|t|      [95.0% Conf. Int.]
-----
const              104.8919      3.056     34.320      0.000      98.899      110.885
Credit Card Only    40.3398      3.163     12.754      0.000      34.137      46.542
Loyalty Club Only   -34.3380      2.598    -13.216      0.000     -39.433     -29.243
Loyalty Club and Credit Card 181.2592      4.926     36.794      0.000     171.599     190.920
Store Mailing List  -82.3691      2.648    -31.101      0.000     -87.563     -77.175
Avg Num Products Purchased 33.4881      0.758     44.208      0.000      32.003      34.974
=====
Omnibus:           359.638    Durbin-Watson:      2.045
Prob(Omnibus):     0.000    Jarque-Bera (JB):    4770.580
Skew:              0.232    Prob(JB):             0.00
Kurtosis:          9.928    Cond. No.             3.16e+15
=====

```

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)

$$\text{Profit} = \text{CC_only} * 40.34 + \text{Loyal_club_only} * -34.33 + \text{loyalty_club_and_cc} * 181.25 + \text{store_mailing_list} * -82.37 + \text{Avg_num_purch_prods} * 33.49 + 104.89$$

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?

I recommend sending the catalogue, because the model predicts about \$23K -- this is greater than the \$10K threshold that was set to determine go/no go.

2. How did you come up with your recommendation? (Please explain your process so reviewers can give you feedback on your process)

I made the model using dummied variables for the Customer Segment data, the Avg num prods purchased, and an intercept term. I fit the model to the past customer data, and used it to predict the total profit from the new customers. I then multiplied each prediction by the Score_Yes for each customer, and summed up the results to get \$23K.

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

\$23K