# MA478 TEE

Cooper Klein

May 2024
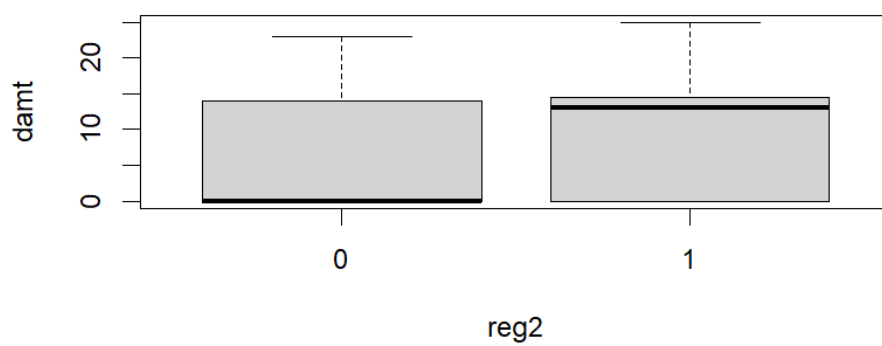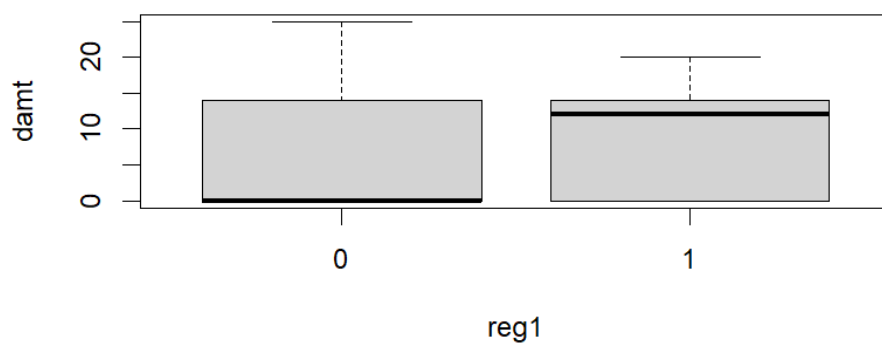
# 1 Data
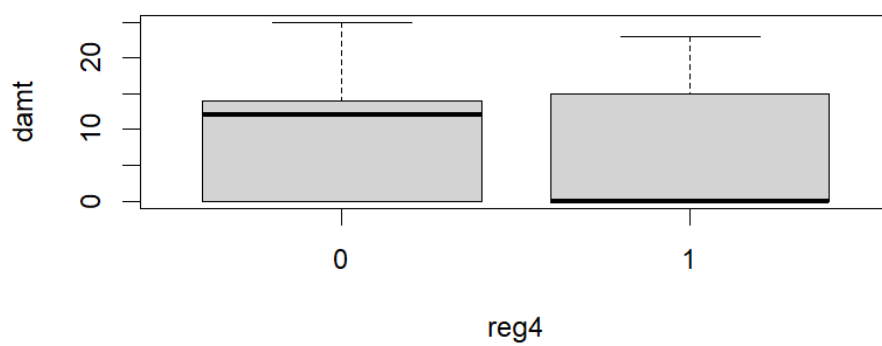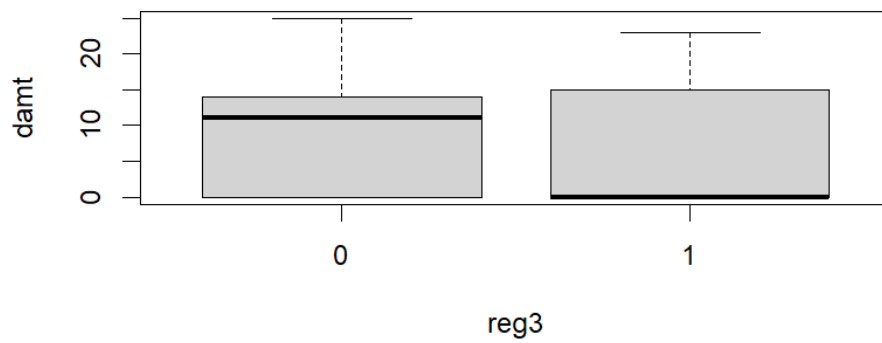
## 1.1 Numeric Predictors

| Description | Non-Donor Mean | Donor Mean |
|---|---|---|
| Number of Children | 2.33 | 0.830 |
| Household Income (7 categories) | 3.91 | 3.98 |
| Wealth rating (scale of 0-9) | 6.47 | 7.63 |
| Average Home Value in Neighborhood | 176 | 194 |
| Average Family Income in Neighborhood | 53.6 | 60.7 |
| Percent low-income in Neighborhood | 15.6 | 11.9 |
| Number of lifetime promotions | 57.5 | 65.7 |
| Dollar amount of lifetime gifts | 107 | 127 |
| Dollar amount of largest gift | 22.4 | 24.0 |
| Number of months between 1st and 2nd gift | 6.81 | 5.79 |

Table 1: Comparison of numeric statistics for Donors and Non-Donors

## 1.2 Region Comparison

# 2 Models

## 2.1 Model 1: Full Binary Classification and Linear Regression

Our first model predicts whether an individual will be a donor with a binary classification, and then predicts the amount donated with a linear regression. This model uses all available predictors for both the logistic and linear regression models.

$$\begin{aligned}
\text{logit}(p_i) = {} & -0.25 + 0.56 Region1_i + 1.21 Region2_i + 0.02 Region3_i - 0.01 Region4_i \\
& - 1.96 Child_i + 1.11 Homeowner_i + 0.82 WRAT_i - 0.03 Female_i \\
& + 0.37 INCM_i + 0.01 INCA_i - 0.19 PLOW_i - 0.44 TLAG_i \\
& + 0.032 AGIF_i - 0.46 TLAG_i - 0.23 TDON_i - 0.03 LGIF_i + 0.14 TGIF_i
\end{aligned} \tag{1}$$

where:

- $z_i = \begin{cases} 1 & \text{if Donor} \\ 0 & \text{if non-donor} \end{cases}$

- $z_i \sim Ber(p_i)$

The linear regression was fit as follows:

$$\begin{aligned}
\text{DAMT} = {} & 14.2 - 0.04 Region1_i - 0.07 Region2_i + 0.32 Region3_i + 0.64 Region4_i \\
& - 0.60 Child_i + 0.24 Homeowner_i - 0.01 WRAT_i - 0.03 Female_i \\
& + 0.37 INCM_i + 0.01 INCA_i - 0.19 PLOW_i - 0.44 TLAG_i + 0.03 AGIF_i \\
& - 0.003 TLAG_i + 0.07 TDON_i - 0.03 LGIF_i + 0.14 TGIF_i + \epsilon
\end{aligned} \tag{2}$$

## 2.2 Model 2: Subset Binary Classification and Linear Regression

Our second model similarly predicts whether an individual will be a donor with a binary classification, and then predicts the amount donated with a linear regression. This model uses a smaller subset of available predictors, that are believed to hold more significance from our data exploration.

The logistic regression was fit as follows:

$$\begin{aligned}
\text{logit}(p_i) = {} & -0.22 + 0.56 Region1_i + 1.18 Region2_i + 0.03 Region3_i + 0.01 Region4_i \\
& - 1.91 Child_i + 1.08 Homeowner_i + 0.79 WRAT_i - 0.02 Female_i \\
& + 0.35 INCM_i + 0.06 INCA_i - 0.19 PLOW_i - 0.44 TLAG_i + 0.032 AGIF_i
\end{aligned} \tag{3}$$

where:

- $z_i = \begin{cases} 1 & \text{if Donor} \\ 0 & \text{if non-donor} \end{cases}$

- $z_i \sim Ber(p_i)$

The linear regression was fit as follows:

$$\begin{aligned}
\text{DAMT} = {} & 14.28 - 0.04 Region1_i - 0.09 Region2_i + 0.34 Region3_i \\
& + 0.66 Region4_i - 0.53 Child_i + 0.26 Homeowner_i + \epsilon
\end{aligned} \tag{4}$$

## 2.3 Model 3: Zero-Inflated Poisson

$$\log \lambda_i = 2.66 - 0.01 Region1_i - 0.02 Region2_i + 0.07 Region3_i + 0.01 Region4_i - 0.02 Child_i + 0.002 TFIG_i \tag{5}$$

$$\text{logit}(\phi_i) = -0.29 - 1.13 Region1_i - 2.11 Region2_i - 0.04 Region3_i + 0.07 Region4_i + 1.15 Child_i - 0.01 TGIF \tag{6}$$

where:

- $y_i \sim \begin{cases} 0 & \text{with prob} \phi_i \\ \text{Pois}(\lambda_i) w/probability(1 - \phi_i) \end{cases}$

- $Z_i \sim \text{Bern}(\phi_i)$

- $Y_i \sim \text{Po}(\lambda_i)$

# 3 Results

We will evaluate our 3 models on their AIC and their MSE results on our withheld test set. Additionally, we will present the profit curve for our best model.

| Model | AIC | MSE |
|---|---|---|
| Model 1 (binary classification) | 2804.3 | 1.76 |
| Model 1 (linear regression) | 6646.78 | 1.76 |
| Model 2 (binary classification) | 2916.71 | 2.12 |
| Model 2 (linear regression) | 7880.144 | 2.12 |
| Model 3 (zero-inflated Poisson) | 13119.15 | 1.77 |

Table 2: AIC Performance of Models

From Table 2, we will recommend our model combining logistic regression and linear regression with all available predictors. We deduce this is our preferable model because it yields the lowest MSE and lowest AIC. We believe, however, that our zero-inflated poisson model offers the most potential, we just did not have enough time to fully develop the model with the optimal subset of predictors.
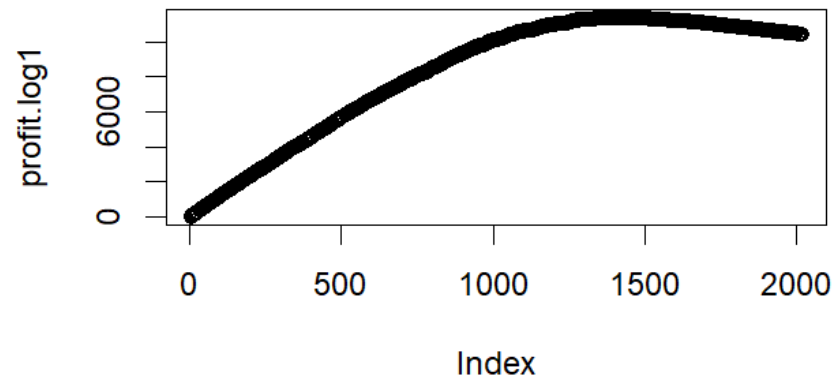
F

Figure 1: Expected Profit Curve of Model 1