

# Solo 3: Customer Segmentation and Target Marketing

Ji Jung

## 1 Introduction

XYZ's database of customers is leveraged for it contains many variables relating to sales and campaign results. This data will be used to first assess which variables have the greatest importance in determining both the chance of response to a mailing campaign as well as the likely spend of each customer. This reduced dataset will be then used to fit three classification models to predict chance of response, and three regression models to predict likely spend. Optimal models are selected and subsequently used to estimate the expected net revenue from conducting a new targeted mailing campaign, accounting for the cost of mailing each customer.

## 2 Data

The dataset used for this assessment is based on XYZ's database of customers. It includes 30,779 records and 554 features of customer data. The features capture customer sales, results from previous mailing campaigns, and Experian properties which provide additional insights about each customer. The features are broken up into 345 character, 48 integer and 161 numeric variable types.

Due to the scale of data, a subjective assessment of variable relevance is conducted prior to employing any pre-processing or modelling routines. This is conducted by assessing the descriptions for variables contained within the provided data dictionary and grading each by its perceived ability to predict both the chance of response and likely spend. This resulted in 227 variables being graded as having a 'low' relevance, which were subsequently excluded from the dataset. Note that a summary data file of each variable grade is available on request. The remaining variables can be broken up into 119 character, 48 integer and 160 numeric variable types.

From an initial look at the data, it is worth noting that the compiled R data frame fails to distinguish between numeric and factor variables. As such, prior to performing any data exploration, all character class variables were converted to factor type and retained all other variables as numeric type. All 'ANY\_MAIL\_x' and 'RESPONSEx' variables were manually converted to factor type for variables contained no character-based observations yet were observed to be categorical in nature.

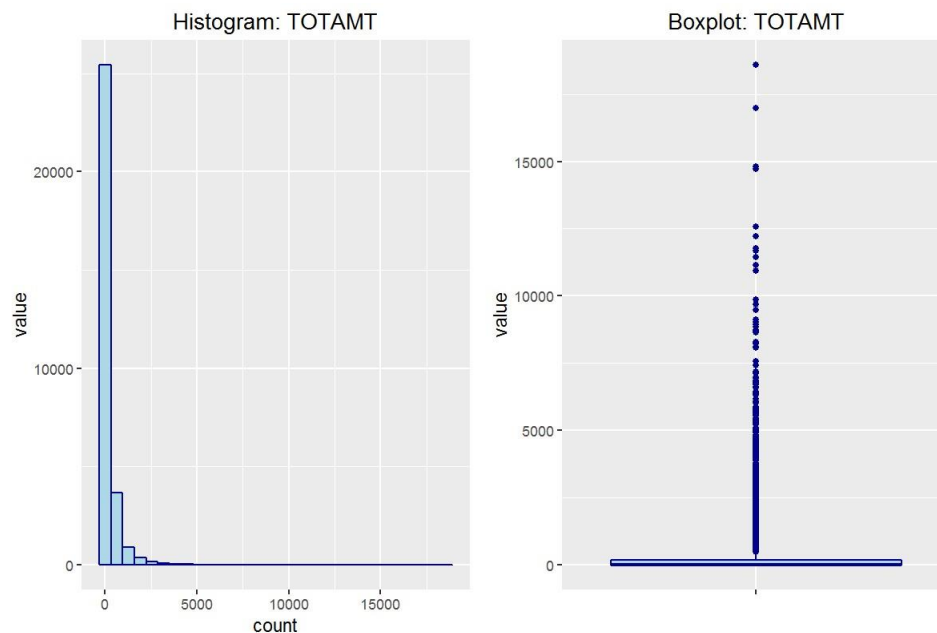
### 3 Data Exploration

Several exploration routines were conducted. These routines allowed to gain an understanding of potential data limitations, including identifying variables which have missing observations, outlier observations, or those variables which may benefit from transformation.

#### 3.1 Univariate Data Analysis

As part of the univariate data analysis, summary statistics for all the 160 retained numeric variables were calculated and observed. Many numeric variables do not suffer from missing values. Many variables have a minimum value of zero, suggesting zero-inflated data. Histogram and box plots were also generated and reviewed for a large subset of numeric variables, with total amount spent (TOTAMT) selected for further discussion below. Note that zero value observations were removed prior to generating each plot.

**Figure 3.1.1 Histogram and Boxplot: TOTAMT**



It can be noticed that the variable suffers from a heavy-positive skew, which is a common attribute over many numeric variables. The result is several observations which could be classed as outliers.

#### 3.2 Bivariate Data Analysis

As a prediction model will be built to determine both the chance of response to a mailing campaign (RESPONSE16) and the likely spend of each customer (TOTAMT16), identifying variables is crucial since they have explanatory power over these two variables. Hence, the Pearson correlation coefficient is calculated and reviewed for all

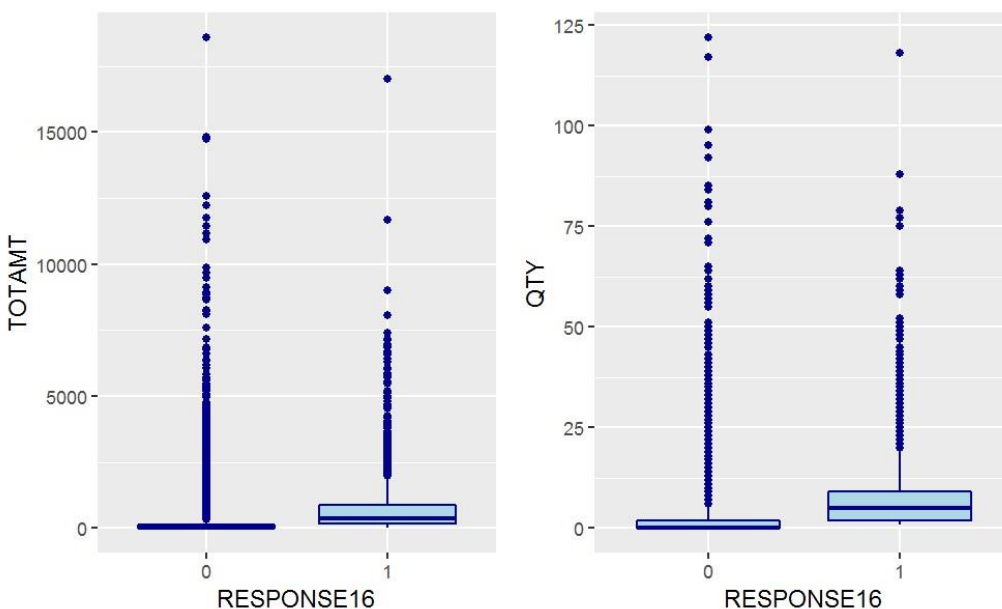
numeric variables against the numeric response variable, TOTAMT16. Correlations for the 10 most correlated numeric variables against TOTAMT16 are shown in the table below.

**Table 3.2.1 Correlations vs. TOTAMT16 (Top 10 Correlations)**

variable	correl coeff
YTD_SALES_2009	0.3723
YTD_TRANSACTIONS_2009	0.2938
LTD_SALES	0.2218
LTD_TRANSACTIONS	0.2067
PRE2009_TRANSACTIONS	0.1587
PRE2009_SALES	0.1445
TOTAL_MAIL_13	0.1182
TOTAL_MAIL_14	0.1178
TOTAL_MAIL_15	0.1178
SUM_MAIL_12	0.1158

None of the variables have reported a strong correlation with TOTAMT16, with the greatest absolute correlation being reported by total sales through to 2009 (YTD\_SALES\_2009) and total transactions through to 2009 (YTD\_TRANSACTIONS\_2009) at 0.37 and 0.29, respectively. Finally, bar plots are used to explore the relationship between the categorical response variable (RESPONSE16) and each numeric variable. Two of these plots have been selected for further discussion below.

**Figure 3.2.1 Boxplot: RESPONSE16 vs. TOTAMT / QTY**



There are recognizable differences in both the mean and distribution of several numeric variables depending on whether they are associated with a positive or negative response to an advertising mailing campaign.

## 4 Data Pre-processing

As part of the data pre-processing routine, imputing data was focused for missing observations (~11% of the dataset). This was initially attempted using the `rflmpute` function from the Random Forest package in R, but processing time eliminated this as a viable option. Instead, each variable were looped over and imputed observations for numeric variables with the variables' median value, and at the same time, imputed observations for character variables with the variables' most common value. Those variables which required imputation were copied and renamed to include the suffix '\_IMP' while the original (non-imputed) equivalent was removed from the dataset. Note that while 89 of the retained categorical variables were identified as having missing values, only one of the retained numeric variables, `ECHVPCT`, was identified as having missing values.

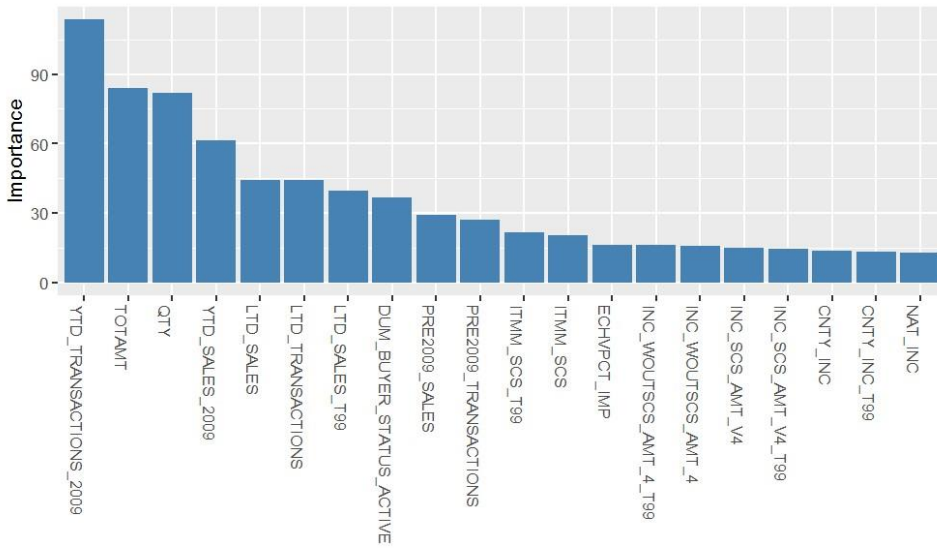
Outlier observations were processed for numeric variables. In many cases, the task of identifying outlier observations can be a subjective practice. Thus, a statistical approach was taken and utilized those observations which fell outside the 1st and 99th percentile range. Observations which met these criteria were replaced using the `squish` function as part of the `scales` package in R, effectively resulting in a newly created set of trimmed variables. Trimmed variables added to the dataset and can be recognized by the suffix '\_T99'.

Lastly, dummies were created for each of the retained factor variables. The first preference was to convert all factor variables to dummies yet many had a high-level count and several levels with relatively low occurrence. Dummies were created for only those factors with 10 or less levels. Dummy variables were named to include the prefix 'DUM\_', along with a suffix to represent the factor level. Note that  $k-1$  dummies were created, where  $k$  is the original number of levels for each variable. This resulted in the creation of 351 dummies.

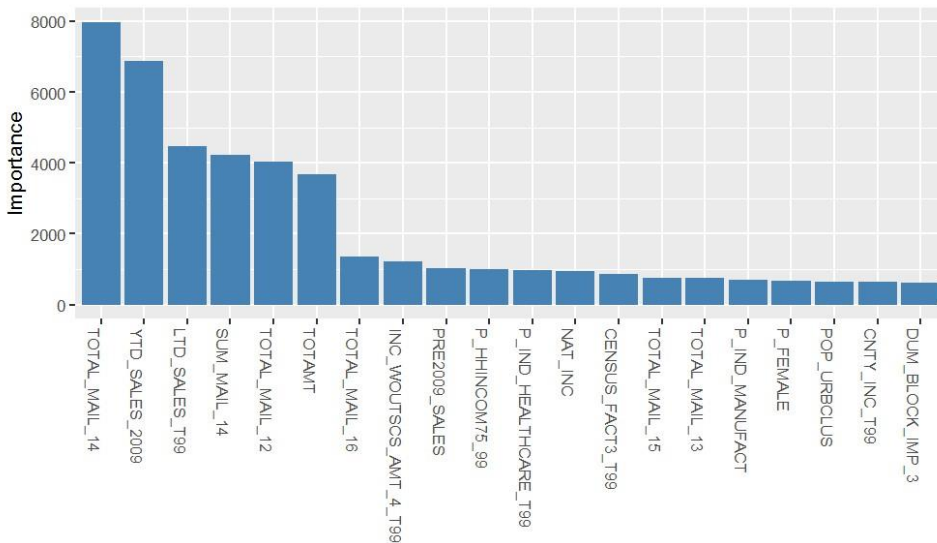
## 5 Variable Importance

The data processing routine produced a data frame of 586 numeric variables. With such a large dataset, it was clear that any subsequent model estimation would benefit from a further reduction in variable count. To achieve this, the `varImp` function as part of the `caret` package in R to calculate the variable importance according to both response variables. In both cases, variable importance was calculated by fitting a Random Forest model, with 'importance' measured by the mean decrease in node impurity. Bar plots of the 20 most important variables for both response variables are below.

**Figure 5.1 Variable Importance: RESPONSE16**



**Figure 5.2 Variable Importance: TOTAMT16**



There is a commonality between variable importance plots, with TOTAMT, YTD\_SALES\_2009 and PRE2009\_SALES all within the top-10 rank for both response variables. A quick drop-off is noticeable in variable importance beyond the first five variables within both plots. Based on these results, top-50 ranked variables were passed by importance through to the model estimation phase. This reduction provides a suitable trade-off both in terms of accuracy and performance.

## 6 Model Estimation

### 6.1 Classification Modelling: Chance of Response

Three classification-based models were fitted to predict the chance a customer will respond to a mailing campaign. These models include a Naive Bayes, Random Forest and a Lasso and Elastic-Net Regularized Generalized Linear Model (GLMnet) classifier. In each case, the train function was used as a part of the caret package with a 3-fold cross-validation sampling method, which was applied to a 70% subset of training data and tested against a 30% subset. Default parameters were used for each model.

The in and out-of-sample Receiver Operating Characteristic (ROC) Curves are shown for each model in Appendix A. The Naive Bayes classifier managed to deliver an in-sample Area Under the Curve (AUC) of 0.85, while its out-of-sample AUC was 0.82. Clearly, the classifier has avoided over-fitting the training data and has managed to maintain similar classification performance over both the training and test sets. The Random Forest classifier on the other hand, reported an in-sample AUC of 1.00 and an out-of-sample AUC of 0.85, which suggests that it greatly suffers from over-fitting. Finally, the GLMnet classifier generated an in-sample AUC of 0.87 and out-of-sample AUC of 0.86.

Below is the out-of-sample confusion matrix for each model.

**Table 6.1.1 Confusion Matrix: Classification Model Comparison**

	Naive Bayes		Random Forest			GLMnet		
	Pred: 0	Pred: 1	Pred: 0	Pred: 1		Pred: 0	Pred: 1	
Actual: 0	3503	246	Actual: 0	3992	452	Actual: 0	4004	416
Actual: 1	496	218	Actual: 1	7	12	Actual: 1	20	38

While the Random Forest and GLMnet classifiers were able to produce similar out-of-sample performance according to their ROC curves, the Random Forest classifier has done so by providing less true positive and true negative values. The performance metrics in the table below shows that the GLMnet classifier was able to obtain a superior true positive rate and true negative rate.

**Table 6.1.2 Performance Metrics: Classification Model Comparison**

	Naive Bayes	Random Forest	GLMnet
Accuracy	0.8337	0.8972	0.9026
95% CI	(0.8225, 0.8446)	(0.8879, 0.9059)	(0.8936, 0.9112)
Kappa	0.2793	0.0419	0.1284

Sensitivity	0.4698	0.0259	0.0837
Specificity	0.8760	0.9983	0.9950
Pos Pred Value	0.3053	0.6316	0.6552
Neg Pred Value	0.9344	0.8983	0.9059
Prevalence	0.1040	0.1040	0.1014
Detection Rate	0.0489	0.0027	0.0085
Detection Prevalence	0.1600	0.0043	0.0130
Balanced Accuracy	0.6729	0.5121	0.5394

From a view of the performance metrics, the GLMnet model has demonstrated a superior AUC, accuracy, sensitivity and specificity compared to the other models. The GLMnet classifier is used to predict the chance of response.

## 6.2 Regression Modelling: Amount Spent

Three regression-based models are fitted to predict the amount a customer will spend. These models include a Multiple Linear Regression (MLR), Random Forest and eXtreme Gradient Boost linear regression estimator. Note that for the MLR, a stepwise variable selection technique was used based on the Akaike Information Criterion (AIC). As with the previous classification models, a 3-fold cross-validation sampling method was employed and maintained the same 30/70 split between test and training data subsets.

The in- and out-of-sample actuals versus predictions for each model are shown in Appendix A. Each model seems to struggle with both outlier observations and the zero-inflated predictor data. The model assessment can be exerted by observing the model fit statistics for each in the below table. Note that negative response values were taken as zero for the statistics shown in the table below. Negative spend amounts to be zero.

**Table 6.2.1 Performance Metrics: Regression Model Comparison**

	MLR	Random Forest	Grad Boost
Training set			
MAE	52.74	28.42	18.96
MSE	27243.39	9856.1	2005.58
RMSE	165.06	99.28	44.78
R <sup>2</sup>	0.18	0.7034	0.9396
Adj R <sup>2</sup>	0.1782	0.7019	0.9393
Test set			
	177.96	179.68	189.49

RMSE

R <sup>2</sup>	0.1366	0.1199	0.0211
Adj R <sup>2</sup>	0.1321	0.1098	0.0099

Each regression model has performed quite poorly over the test set of data. Combining the predictions with chance of response will aid in dealing with the zero-inflated response data. The MLR regression will be adopted to predict the amount spent as its training performance metrics were among the most favorable.

## 7 Customer Scoring

A customer score based on the combined predictions of the Random Forest classification and MLR model discussed above is constructed. This function is to represent the expected value from conducting a new advertising campaign, based on predictions against a subset of customers who have not yet been mailed. The customer score function is shown below.

$$CustomerScore = P(response) \cdot E(netrevenue) - CostofMail$$

From the function, 'P(response)' represents the probability of response as predicted by the chosen classification model, Random Forest. 'E(net revenue)' represents expected net revenue, which is assumed to be 10% of the amount spent as predicted by the chosen regression model, MLR. And finally, the 'cost of mail' represents the cost of mailing customers as part of a new advertising campaign which is assumed to be equal to \$1.00 per customer. The sum of customer scores represents the expected value from conducting a new advertising campaign.

The customer score above is used to propose four possible marketing strategies. The first strategy, ALL\_MAIL involves mailing all customers who have not yet been mailed, regardless of the probability of response or expected net revenue. This would obviously be a costly strategy, considering the cost of mailing all customers. The second strategy, ALLSCORE\_MAIL involves mailing only those customers who return a positive customer score according to the above function. For the third strategy, HIGHPROB\_MAIL, only those customers who are predicted to have a probability of response greater than or equal to 0.7 are mailed. Note that this strategy ignores the customer score and may capture customers who have a negative expected return when accounting for their predicted spend amount. Finally, the fourth strategy, HIGHVAL\_MAIL, involves mailing only those customers who have a predicted spend amount of greater than or equal to \$500 (\$50 net revenue). This strategy also ignores the customer score and may capture customers who have a negative expected return when accounting for their probability of response.

**Table 7.1 Customer Score Summary**

strategy	criteria	no. mailed	expected value	value per customer
----------	----------	------------	----------------	--------------------



ALL_MAIL	ANY_MAIL_16 = 0	15,857	-\$10,916.25	-\$0.69
ALLSCORE_MAIL	Customer Score >= 0	1,111	\$2,640.33	\$2.38
HIGHPROB_MAIL	P(response) >= 0.7	19	\$390.06	\$20.53
HIGHVAL_MAIL	E(net revenue) >= 50	20	\$692.07	\$34.60

It should be no surprise that the greatest expected value comes from the strategy which involves targeting all customers with a positive customer score. Viable strategies are found from mailing only those customers with a high probability of response or a high predicted spend amount. These two strategies can achieve a much greater expected value per customer. It may be that the most effective marketing strategy would be to target those customers as flagged by HIGHPROB\_MAIL and HIGHVAL\_MAIL in the first instance. And then, depending on the success of that campaign, proceed to target the remaining customers flagged by ALLSCORE\_MAIL. Note that a list of un-mailed customer account numbers according to the above strategies is available on request.

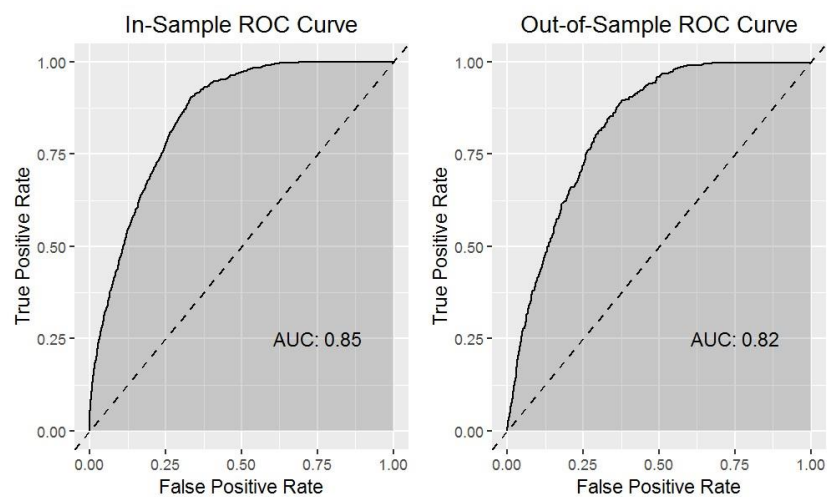
Interesting observations can be noted when reviewing the customers which were flagged by ALLSCORE\_MAIL. Most of these customers were flagged as 'active' customers, who are homeowners, with incomes between \$50-\$150k, and are educated with a bachelor's degree or higher. Other interesting observations are found when comparing the predicted customer response against those who have previously been mailed. Firstly, of the 14,922 customers who were previously mailed, 1,440 customers did in-fact respond (~10% response rate). The chosen classifier suggests that only 71 of the 15,857 un-mailed customers have a probability of response greater than 0.5 (~0.4% response rate). In addition, of the customers who were previously mailed, the average spend amount of those customers was \$342. The chosen regression model also suggests that the average spend amount for those customers with a probability of response greater than or equal to 0.5, is only \$204. It may be that both the chosen classification and regression models are quite conservative in their predictions, or that those customers who were already mailed carried a higher probability of response and higher predicted spend.

## 8 Conclusion

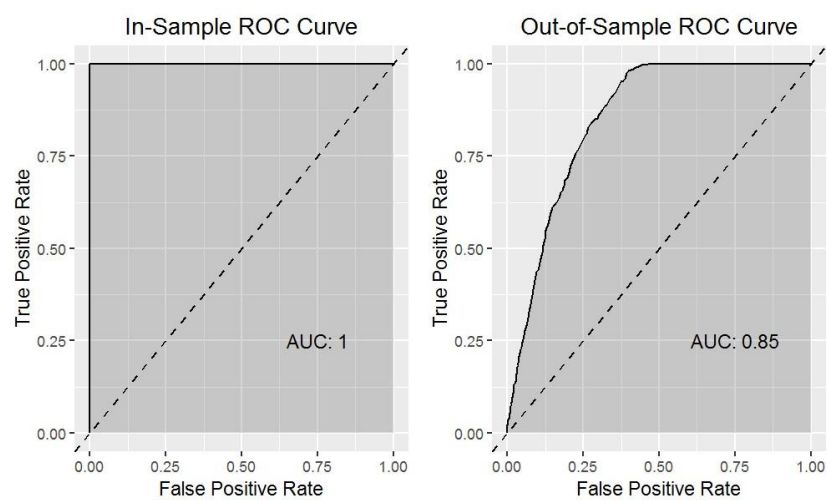
Three classification models were fitted to predict chance of response, and three regression models to predict likely spend. From the fitted models, a GLMnet based model is thought to be superior in predicting chance of response, and the MLR model to be superior in predicting spend. Optimal models were selected and subsequently used to estimate the expected net revenue from conducting a new targeted mailing campaign, accounting for the cost of mailing each customer. This score was then used to propose four possible marketing strategies, ranging from mailing all customers who had not yet been mailed to mailing only those customers who have a predicted spend amount greater than or equal to \$500. Results suggest a viable strategy to mail customers with a high probability of response and/or high expected spend in the first instance, and to follow this by mailing the remaining customers with a positive customer score.

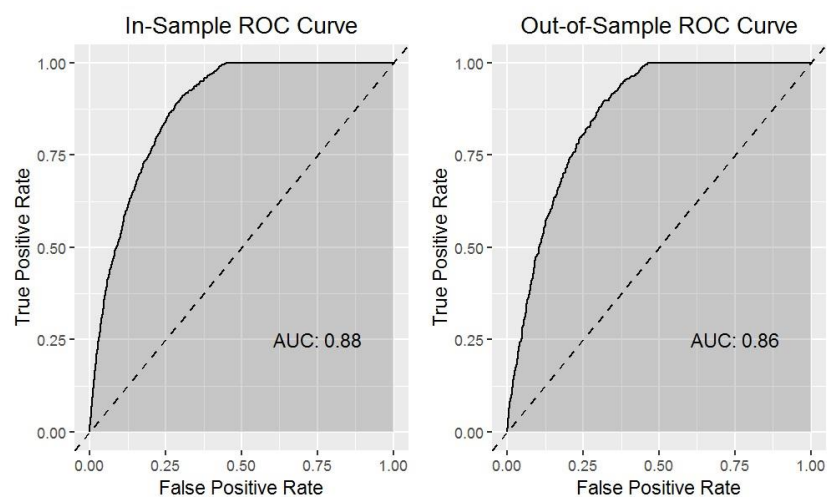
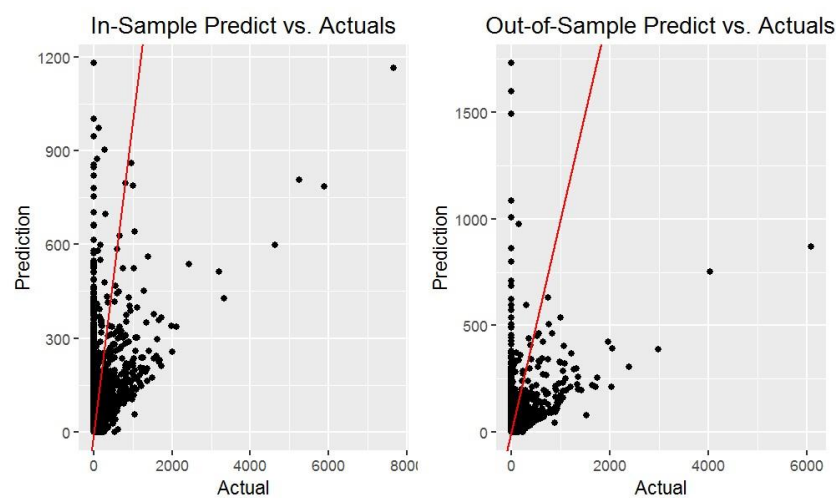
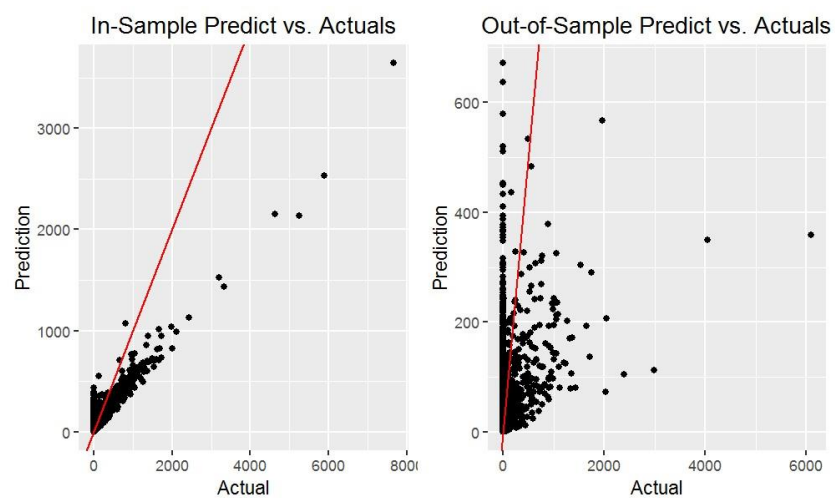
## 9 Appendix

**Figure A.1 ROC Curve: Naive Bayes**



**Figure A.2 ROC Curve: Random Forest**



**Figure A.3 ROC Curve: GLMnet****Figure A.4 Actuals vs. Predictions: Multiple Linear Regression****Figure A.5 Actuals vs. Predictions: Random Forest**

**Figure A.6 Actuals vs. Predictions: eXtreme Gradient Boost**