

¶ Orange Juice Sales at Wasatch Grocery Chain

Identification of Significant Predictor Variables and Predictive Modelling of Customer Preference in Minute Maid Sales

Chris Gearheart and Chris Porter

2022-11-28

```
set.seed(1234)
df <- read.csv(url("http://data.mishra.us/files/project/OJ_data.csv"))
df[-1] <- lapply(df[-1], as.numeric)
df$Purchase <- as.factor(df$Purchase)
purchase_testtrain <- initial_split(df, prop = 0.75, strata = Purchase)
train <- training(purchase_testtrain)
test <- testing(purchase_testtrain)
```

Introduction

Wasatch Grocery Chain (WGC) is a regional grocery chain operating in the Intermountain West of the US. WGC sells two brands of orange juice in its stores, Citrus Hill (CH) and Minute Maid (MM) of which MM is the more profitable to the company. This report will identify **what customer factors** within available data **contribute to purchase of Minute Maid over Citrus Hill**, as well as **to what degree these factors influence customer choice**. In addition, a predictive model is created that will allow the Sales Department to identify other customers within our customer base that are more likely to purchase Minute Maid brand orange juice, thus driving profitability across the company.

Available Data The data set used in this report contains 13 possible predictor variables as well as 1 outcome variable, Purchase, which records whether or not a customer purchased MM. There are a total of 1070 observations in the data set. The data set was further partitioned into a **training** data set, containing 801 observations, and a validation **testing** data set containing 269 observations.

The code below imports the data set, converts the binary **Purchase** outcome into a factor, and pulls out 25% of the observations as a hold-out set or test set against which our final model can be tested. Doing so helps

us avoid the mistake of training a model that performs well against the sample data, but fails to generalize to a new data set from the same population.

Methods

Logistic Regression: WGC's management team wants to know **which variables contribute to an customer outcome of "Yes; Purchased Minute Maid."** Their goal matches the strengths of a logistic regression, which can explain the strength and direction of independent variables' effects on a binary classification outcome (often yes/no or is/is not). This algorithm will tell management which variables push customers towards or away from a Minute Maid purchase, plus which variables have no bearing on the outcome. Significant variables proven to have big enough effects can become levers for action or intervention for management.

Pre-processing Logistic regressions work when:

1. Qualitative variable have been turned into quantitative dummy variables.
2. No columns are uniformly filled with one unique value
3. There is no missing data.
4. There is no correlation between the variables.

Fortunately, the first three conditions were already true of our dataset.

1. Dummy variables are unnecessary because only 'Purchase' is a factor, and it's already expressed us

2. No columns are uniformly filled with one unique value - there is spread in each of the 13 independ
summary(train)

```
## Purchase PriceCH PriceMM DiscCH DiscMM
## 0:312 Min. :1.690 Min. :1.690 Min. :0.00000 Min. :0.0000
## 1:489 1st Qu.:1.790 1st Qu.:1.990 1st Qu.:0.00000 1st Qu.:0.0000
## Median :1.860 Median :2.090 Median :0.00000 Median :0.0000
## Mean :1.865 Mean :2.088 Mean :0.05192 Mean :0.1281
## 3rd Qu.:1.990 3rd Qu.:2.180 3rd Qu.:0.00000 3rd Qu.:0.2400
## Max. :2.090 Max. :2.290 Max. :0.50000 Max. :0.8000
## SpecialCH SpecialMM LoyalCH SalePriceMM
## Min. :0.0000 Min. :0.0000 Min. :0.000014 Min. :1.19
## 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.320000 1st Qu.:1.69
## Median :0.0000 Median :0.0000 Median :0.585435 Median :2.09
## Mean :0.1548 Mean :0.1685 Mean :0.555908 Mean :1.96
## 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.836160 3rd Qu.:2.18
## Max. :1.0000 Max. :1.0000 Max. :0.999947 Max. :2.29
## SalePriceCH PriceDiff PctDiscMM PctDiscCH
## Min. :1.390 Min. : -0.6700 Min. :0.00000 Min. :0.00000
## 1st Qu.:1.750 1st Qu.: 0.0000 1st Qu.:0.00000 1st Qu.:0.00000
## Median :1.860 Median : 0.2400 Median :0.00000 Median :0.00000
## Mean :1.813 Mean : 0.1464 Mean :0.06164 Mean :0.02739
## 3rd Qu.:1.890 3rd Qu.: 0.3200 3rd Qu.:0.11834 3rd Qu.:0.00000
## Max. :2.090 Max. : 0.6400 Max. :0.40201 Max. :0.25269
## ListPriceDiff
## Min. :0.0000
## 1st Qu.:0.1400
## Median :0.2400
```

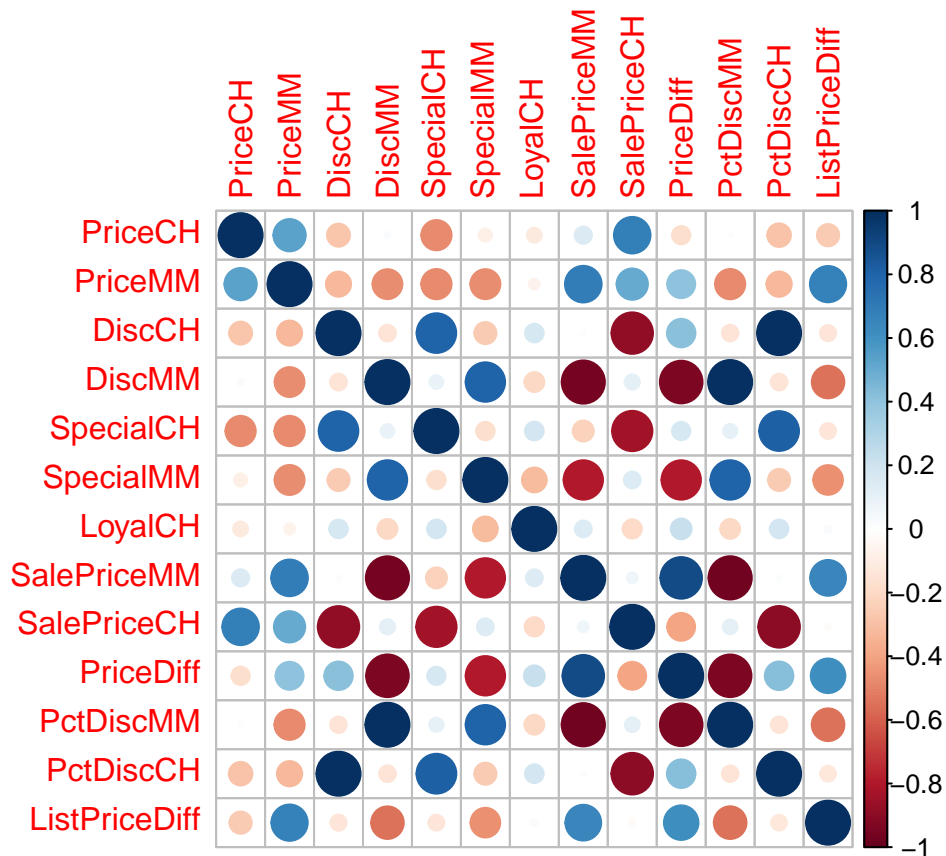
```
## Mean :0.2225
## 3rd Qu.:0.3000
## Max. :0.4400
```

```
# 3. There is no missing data - imputation is not necessary
sum(is.na(train))
```

```
## [1] 0
```

A correlogram confirms that there is high correlation between the thirteen variables. Some of them appear to be multicollinear, or not fully independent of one another.

```
corr <- cor(df[-1])
corr %>% cor() %>% corrplot()
```



Accordingly, our team decided to use the “Lasso” method of logistic regression that regresses all variables against all other variables, modifying each variable’s predictive weight based on its correlation to to other variables by strengthening, weakening, or even nullifying its effect.

Variable selection and model design The `cv.glmnet` function below does that work by finding THEBESTWAYTOSUCCINCTLYDESCRIBEIT, ultimately printing out coefficients for each variable that have been penalized or nullified if their relationship to other variables is multicollinear.

Additionally, in a microcosm of the training/test split we set up at the beginning of the project, this method cross-validates the results of the trained regression by testing it against seven different one-seventh chunks of the entire set.

The code below performs a logistic regression, but it uses Lasso (`alpha = 1`), giving us something to say about the magnitude and direction of variables, plus which variables’ influences were shrunk to zero when

all variables were regressed against each other (Price MM, Disc CH, SalePriceCH, and PctDiscCH). **Doing that gives us an AUC of 0.9**, and that's after inline k-fold validation of 7 when training the model.

```
predictors <- train[,c(2:13)]
purchase_only <- train$Purchase

str(predictors)

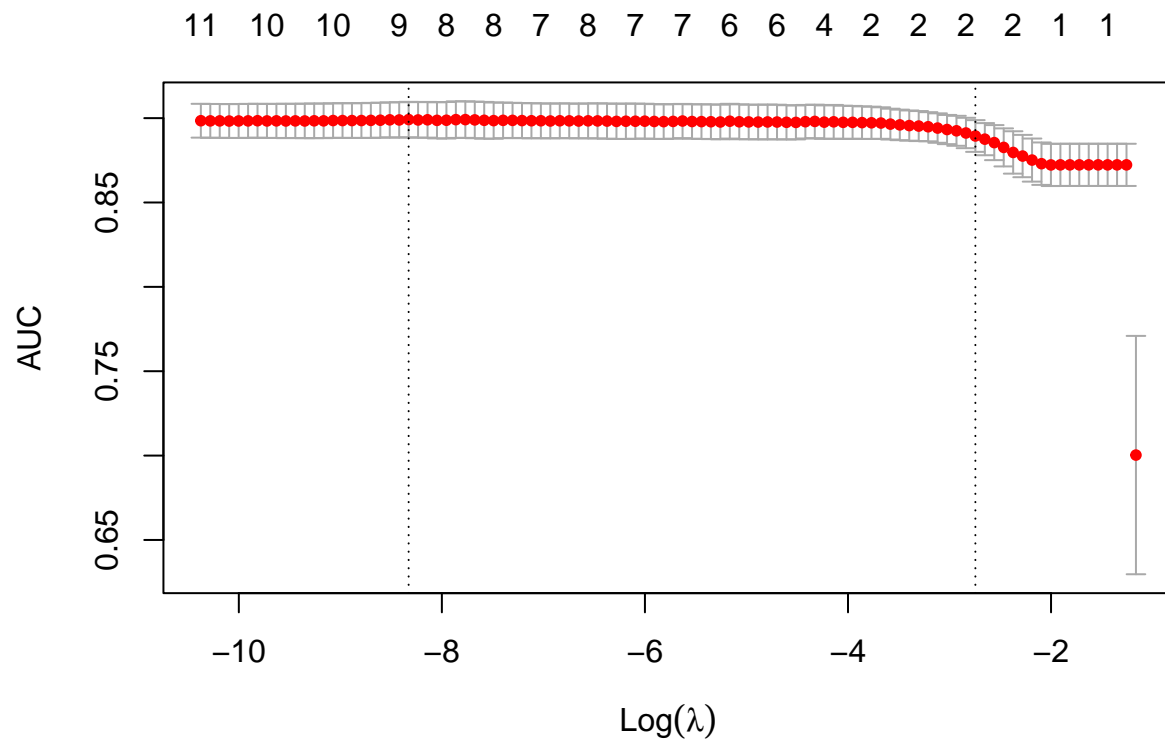
## 'data.frame': 801 obs. of 12 variables:
## $ PriceCH : num 1.75 1.86 1.69 1.99 1.99 1.75 1.86 1.99 1.86 1.75 ...
## $ PriceMM : num 1.99 2.18 1.69 2.09 2.09 1.99 2.18 2.29 2.13 1.99 ...
## $ DiscCH : num 0 0 0 0.1 0.1 0 0 0 0 0 ...
## $ DiscMM : num 0.3 0 0 0 0 0.3 0 0 0.24 0.4 ...
## $ SpecialCH : num 0 0 1 0 0 0 0 0 0 0 ...
## $ SpecialMM : num 1 0 0 0 0 1 0 0 0 0 ...
## $ LoyalCH : num 0.5 0.32 0.68 0.944 0.7 ...
## $ SalePriceMM: num 1.69 2.18 1.69 2.09 2.09 1.69 2.18 2.29 1.89 1.59 ...
## $ SalePriceCH: num 1.75 1.86 1.69 1.89 1.89 1.75 1.86 1.99 1.86 1.75 ...
## $ PriceDiff : num -0.06 0.32 0 0.2 0.2 -0.06 0.32 0.3 0.03 -0.16 ...
## $ PctDiscMM : num 0.151 0 0 0 0 ...
## $ PctDiscCH : num 0 0 0 0.0503 0.0503 ...

predMod <- glm(train$Purchase ~ ., data = train, family = binomial(link='logit'))

predictors <- data.matrix(predictors)

set.seed(1234)
cv.binomial <- cv.glmnet(x = predictors, y = train$Purchase,
                        alpha = 1, family = "binomial",
                        nfolds = 7, standardize = TRUE, type.measure = "auc")

plot(cv.binomial)
```



```
(best.lambda <- cv.binomial$lambda.min)
```

```
## [1] 0.0002418264
```

```
y4<- coef(cv.binomial, s="lambda.min", exact=FALSE)
print(y4)
```

```
## 13 x 1 sparse Matrix of class "dgCMatrix"
```

```
##              s1
## (Intercept) -4.1361101
## PriceCH     -0.3496704
## PriceMM     .
## DiscCH      .
## DiscMM     -13.2136017
## SpecialCH   0.1216139
## SpecialMM  -0.2974192
## LoyalCH     6.6471993
## SalePriceMM 0.6000532
## SalePriceCH .
## PriceDiff   3.5224475
## PctDiscMM   31.5219099
## PctDiscCH   .
```

Since the Lasso function of the regression has shrunk the effects of PriceMM, DiscCH, SalesPriceCH and PctDiscCH to zero or “.” in light of multicollinearity, management can be confident that those variables are not meaningful levers for action.

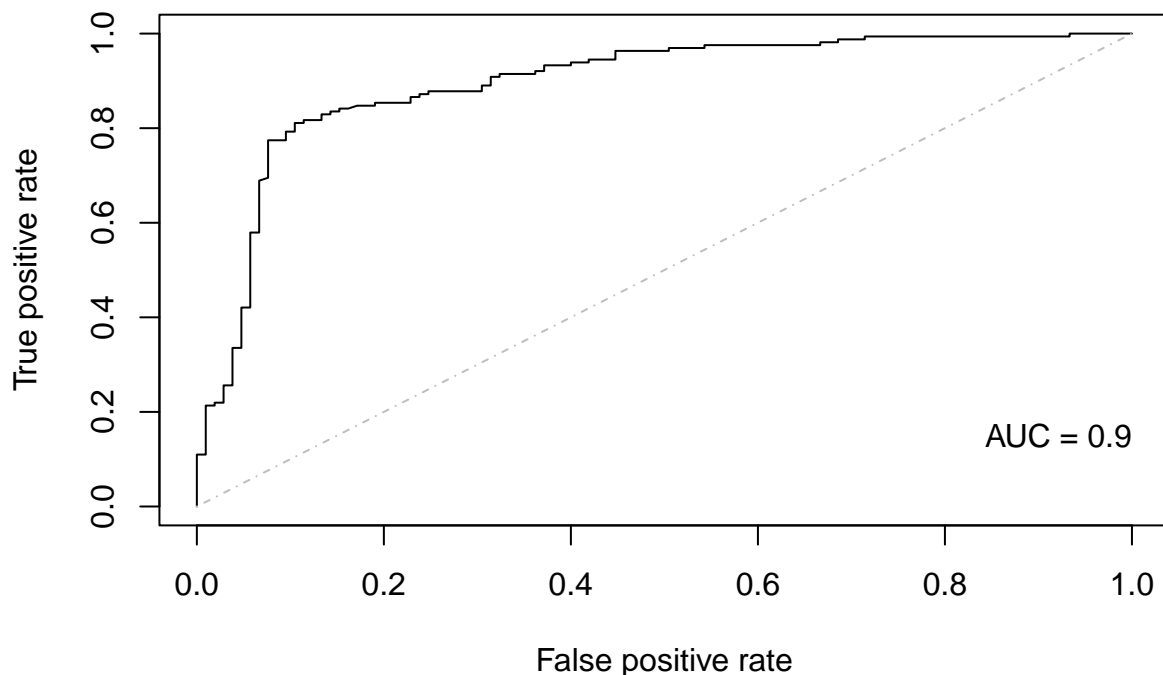
Performance against test data The predictions of this logistic regression performed well against the ground truth outcomes in the test set held in reserve at the beginning of our analysis.

Our regression turned variables into percentage likelihoods, but it is up to the analyst to decide what percentage triggers a label of “Yes; Purchased MinuteMaid”, a decision called the “classification threshold.” The area-under-the-curve (AUC) metric is a sign of an model’s general performance across those different thresholds — a higher AUC means a model is good at balancing the risk of true positives to true negatives.

The area-under-curve for this model is 0.90.

```
test_predictors <- test[,c(2:13)]
test_predictors <- data.matrix(test_predictors)
pred <- predict(cv.binomial, newx = test_predictors,
               type = "response", s = "lambda.min")
pred <- prediction(pred, test$Purchase)
perf <- performance(pred, "tpr", "fpr")
auc_ROCR <- performance(pred, measure = "auc")

plot(perf, colorize=FALSE, col="black")
lines(c(0,1),c(0,1),col = "gray", lty = 4 )
text(1,0.15,labels=paste("AUC = ",round(auc_ROCR@y.values[[1]],
                                     digits=2),sep=""),adj=1)
```



Gradient Boosted Decision Trees: Management also wants to be able to predict the likelihood that any given future customer will buy Minute Maid. Knowing how many customers are likely to purchase Minute Maid can help in (1) forecasting cash flow and supply chain demand and (2) targeting marketing to customers who are in the ideal position to buy and ignoring those who are not.

Decision tree modelling models the data and assigns a probabilistic decision path to assign classification, in this case either to a likely Minute Maid purchase or not. However, the way decision trees are assembled can lead to overfitting to the data if the tree is too deep or has too many branches, in addition they are prone to fall prey to data sampling errors, creating trees that reflect the train sample better than they do the ground truth. To overcome this, Gradient Boosted Trees (GBT) are a machine learning algorithm that overcomes the propensity of decision tree algorithms to overfit the data and susceptibility to data sampling errors. GBT overcomes this by building a more accurate complex model iteratively by combining many smaller less predictive models. Each successive round of learning seeks to explain the remaining error left by the previously assembled tree.

```
set.seed(1234)
recipe_oj <- recipe(Purchase ~ ., train)

model_oj_bt <- boost_tree(trees = tune(), tree_depth = tune(), learn_rate = tune()) %>%
  set_engine('xgboost', verbosity = 0) %>%
  set_mode('classification')

hyperparameter_grid <- grid_regular(trees(), tree_depth(), learn_rate(), levels = 5)

purchase_folds <- vfold_cv(train, v=4) # 4-fold Cross validation

oj_workflow <- workflow() %>% add_model(model_oj_bt) %>% add_recipe(recipe_oj) #Set Workflow

# Tune Hyper-parameters
oj_tune <- oj_workflow %>% tune_grid(resamples = purchase_folds,
                                     grid = hyperparameter_grid,
                                     metrics = metric_set(accuracy))

best_bt_model <- oj_tune %>% select_best('accuracy') #Select best Hyper-parameters from grid

#best_bt_model
```

The hyperparameters for number of trees, tree depth, and learn rate for the boosted tree model were tuned using a grid with 5 levels and 4-fold cross validation. Hyperparameter performance was evaluated by overall model accuracy of prediction. The final hyperparameters for the model are number of trees (1000), tree depth (1), and learn rate (0.1).

```
oj_final_workflow <- oj_workflow %>% finalize_workflow(best_bt_model) # Create Final Workflow based upon best model

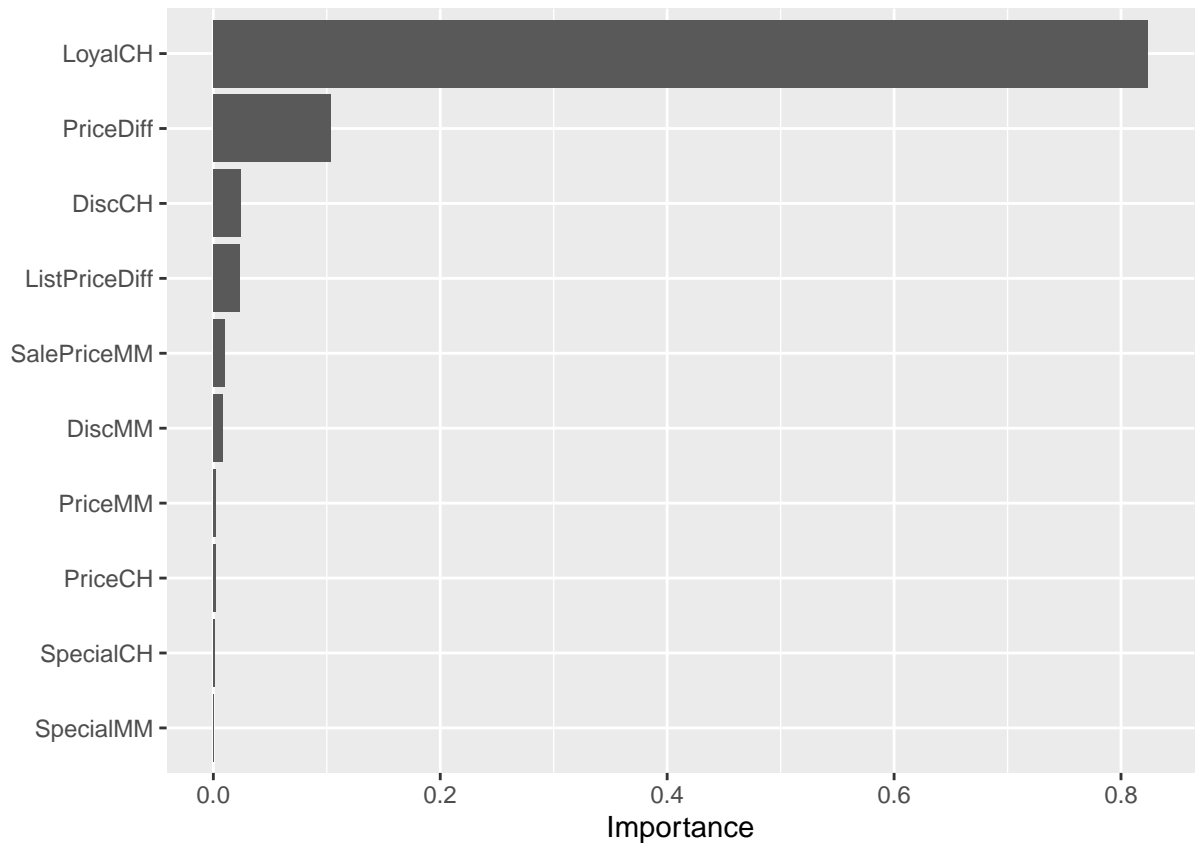
final_fit <- oj_final_workflow %>% last_fit(split = purchase_testtrain) # Final Fit Model

final_fit %>% collect_metrics()

## # A tibble: 2 x 4
##   .metric .estimator .estimate .config
##   <chr>    <chr>      <dbl> <chr>
## 1 accuracy binary      0.799 Preprocessor1_Model11
## 2 roc_auc  binary      0.892 Preprocessor1_Model11
```

The finalized model gave an AUC of 0.89, which is comparable, but slightly underperforms the logistic regression model previously discussed.

```
oj_final_workflow %>% fit(data = train) %>% extract_fit_parsnip() %>% vip(geom = 'col') #Plot most important features
```



```
vi_values <- oj_final_workflow %>% fit(data = train) %>% extract_fit_parsnip() %>% vi()
```

```
vi_values
```

```
## # A tibble: 11 x 2
##   Variable      Importance
##   <chr>         <dbl>
## 1 LoyalCH      0.824
## 2 PriceDiff    0.104
## 3 DiscCH       0.0239
## 4 ListPriceDiff 0.0232
## 5 SalePriceMM  0.0104
## 6 DiscMM       0.00852
## 7 PriceMM      0.00237
## 8 PriceCH      0.00237
## 9 SpecialCH    0.00133
## 10 SpecialMM   0.000512
## 11 SalePriceCH 0.000356
```

One drawback to using a black-box machine learning algorithm like Gradient Boosted Trees, is that understanding the insights the model provides are not immediately available, and the use of explanatory analysis is required to further understand what actions management can take to increase sales of Minute Maid. One such tool is the use of variable importance to understand which variables the model sees as most important in determining a customer outcome of **“Yes; Purchased Minute Maid”**.

The most important variable according to the Boosted Tree model is Customer Brand Loyalty to Citrus Hill(LoyalCH) with 82.35% importance, followed by Price Difference(PriceDiff) with 10.35% importance. All

other independent variables displayed importance of <3%.

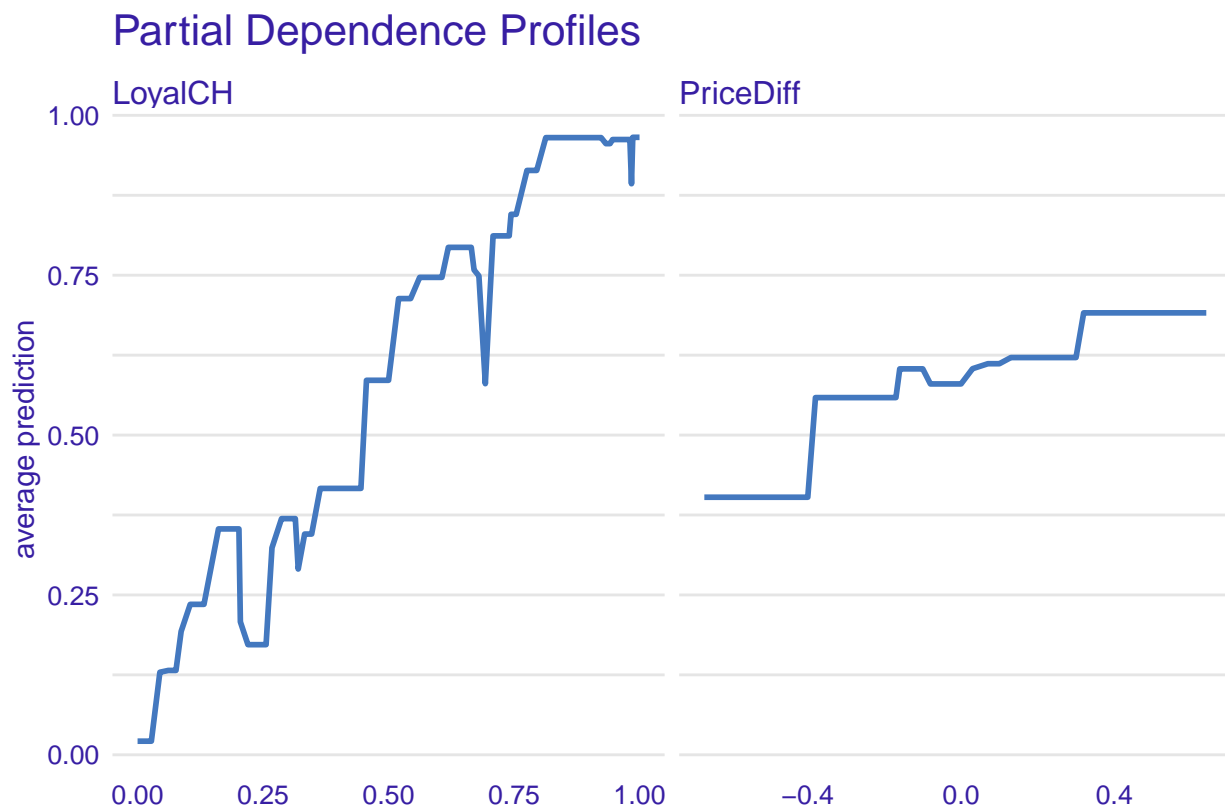
```
model_fitted <- oj_final_workflow %>% fit(data = train)

explainer_rf <- explain_tidymodels(model_fitted,
                                   data = train[,-1],
                                   y = train$Purchase,
                                   type = "pdp", verbose = FALSE)

pdp <- model_profile(explainer_rf,
                    variables = c("LoyalCH", "PriceDiff"),
                    N=NULL)
```

In addition to understanding which variables are important for management to focus on, it is also important to understand how those variables interact with the prediction for Minute Maid purchases by the customer. It is useful to know that Brand Loyalty is important, but even more useful to know how to use that lever to identify potential crossover customers. Partial Dependence Profiling (PDP) allows some insight what is happening inside a blackbox model such as GBTs. The above plot shows the partial independent portion of a variable's influence on the dependent outcome variable. Comparable to information that can be obtained from linear or logistic regression.

```
plot(pdp, title='Partial Dependence Profiles', subtitle='')
```



Both variables display a positive relationship with the purchase of Minute Maid. Meaning, that the more Brand Loyalty a customer displays towards Citrus Hill and the larger the price difference between MM and CH (in Citrus Hill's favor) the more likely the customer was to purchase Minute Maid. This would seem to be counter-intuitive and so it was verified by looking at the original data, where this observation was supported (see below). This would seem to indicate that there is a unique positioning opportunity for Minute

Maid in WGC stores.

```
ggplot2::ggplot(df, aes(LoyalCH, PriceDiff, color=Purchase)) + geom_hline(yintercept = 0) + geom_point()
```



Conclusions and Recommendations At the beginning of this project we met with stakeholders in the Branding and Sales departments and identified key deliverables to ensure that this project provided actionable information and value to the company. Based upon our work we suggest the following interpretations and courses of action moving forward.

Brand

1. What predictor variables influence the purchase of MM?
2. Are all the variables in the dataset effective or are some more effective than others?
3. How confident are you in your recommendations?
4. Based on your analysis what are the specific recommendations you have for the brand manager?

Sales

1. Can you build a predictive model that can inform him the probability of customers buying MM?
2. How good is the model in its predictions?
3. How confident are you in your recommendations?

Appendix 1: Data Characteristics

```
summary(df)
```

```
## Purchase PriceCH PriceMM DiscCH DiscMM
## 0:417 Min. :1.690 Min. :1.690 Min. :0.00000 Min. :0.00000
```

```
## 1:653 1st Qu.:1.790 1st Qu.:1.990 1st Qu.:0.00000 1st Qu.:0.0000
## Median :1.860 Median :2.090 Median :0.00000 Median :0.0000
## Mean :1.867 Mean :2.085 Mean :0.05186 Mean :0.1234
## 3rd Qu.:1.990 3rd Qu.:2.180 3rd Qu.:0.00000 3rd Qu.:0.2300
## Max. :2.090 Max. :2.290 Max. :0.50000 Max. :0.8000
## SpecialCH SpecialMM LoyalCH SalePriceMM
## Min. :0.0000 Min. :0.0000 Min. :0.000011 Min. :1.190
## 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.325257 1st Qu.:1.690
## Median :0.0000 Median :0.0000 Median :0.600000 Median :2.090
## Mean :0.1477 Mean :0.1617 Mean :0.565782 Mean :1.962
## 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.850873 3rd Qu.:2.130
## Max. :1.0000 Max. :1.0000 Max. :0.999947 Max. :2.290
## SalePriceCH PriceDiff PctDiscMM PctDiscCH
## Min. :1.390 Min. : -0.6700 Min. :0.0000 Min. :0.00000
## 1st Qu.:1.750 1st Qu.: 0.0000 1st Qu.:0.0000 1st Qu.:0.00000
## Median :1.860 Median : 0.2300 Median :0.0000 Median :0.00000
## Mean :1.816 Mean : 0.1465 Mean :0.0593 Mean :0.02731
## 3rd Qu.:1.890 3rd Qu.: 0.3200 3rd Qu.:0.1127 3rd Qu.:0.00000
## Max. :2.090 Max. : 0.6400 Max. :0.4020 Max. :0.25269
## ListPriceDiff
## Min. :0.000
## 1st Qu.:0.140
## Median :0.240
## Mean :0.218
## 3rd Qu.:0.300
## Max. :0.440
```

```
summary(test)
```

```
## Purchase PriceCH PriceMM DiscCH DiscMM
## 0:105 Min. :1.690 Min. :1.690 Min. :0.00000 Min. :0.0000
## 1:164 1st Qu.:1.790 1st Qu.:1.990 1st Qu.:0.00000 1st Qu.:0.0000
## Median :1.860 Median :2.090 Median :0.00000 Median :0.0000
## Mean :1.874 Mean :2.079 Mean :0.05167 Mean :0.1094
## 3rd Qu.:1.990 3rd Qu.:2.180 3rd Qu.:0.00000 3rd Qu.:0.2000
## Max. :2.090 Max. :2.290 Max. :0.50000 Max. :0.8000
## SpecialCH SpecialMM LoyalCH SalePriceMM
## Min. :0.0000 Min. :0.0000 Min. :0.000011 Min. :1.190
## 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.384000 1st Qu.:1.780
## Median :0.0000 Median :0.0000 Median :0.635200 Median :2.090
## Mean :0.1264 Mean :0.1413 Mean :0.595184 Mean :1.969
## 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.875808 3rd Qu.:2.130
## Max. :1.0000 Max. :1.0000 Max. :0.999870 Max. :2.290
## SalePriceCH PriceDiff PctDiscMM PctDiscCH
## Min. :1.390 Min. : -0.6700 Min. :0.00000 Min. :0.00000
## 1st Qu.:1.750 1st Qu.: 0.0000 1st Qu.:0.00000 1st Qu.:0.00000
## Median :1.860 Median : 0.2300 Median :0.00000 Median :0.00000
## Mean :1.823 Mean : 0.1468 Mean :0.05231 Mean :0.02709
## 3rd Qu.:1.890 3rd Qu.: 0.3000 3rd Qu.:0.09569 3rd Qu.:0.00000
## Max. :2.090 Max. : 0.6400 Max. :0.40201 Max. :0.25269
## ListPriceDiff
## Min. :0.0000
## 1st Qu.:0.1000
## Median :0.2400
## Mean :0.2045
```

```
## 3rd Qu.:0.2900
## Max. :0.4400
```

```
summary(train)
```

```
## Purchase      PriceCH      PriceMM      DiscCH      DiscMM
## 0:312   Min. :1.690   Min. :1.690   Min. :0.00000   Min. :0.0000
## 1:489   1st Qu.:1.790   1st Qu.:1.990   1st Qu.:0.00000   1st Qu.:0.0000
##         Median :1.860   Median :2.090   Median :0.00000   Median :0.0000
##         Mean   :1.865   Mean   :2.088   Mean   :0.05192   Mean   :0.1281
##         3rd Qu.:1.990   3rd Qu.:2.180   3rd Qu.:0.00000   3rd Qu.:0.2400
##         Max.   :2.090   Max.   :2.290   Max.   :0.50000   Max.   :0.8000
## SpecialCH      SpecialMM      LoyalCH      SalePriceMM
## Min. :0.0000   Min. :0.0000   Min. :0.000014   Min. :1.19
## 1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.320000   1st Qu.:1.69
## Median :0.0000   Median :0.0000   Median :0.585435   Median :2.09
## Mean   :0.1548   Mean   :0.1685   Mean   :0.555908   Mean   :1.96
## 3rd Qu.:0.0000   3rd Qu.:0.0000   3rd Qu.:0.836160   3rd Qu.:2.18
## Max.   :1.0000   Max.   :1.0000   Max.   :0.999947   Max.   :2.29
## SalePriceCH      PriceDiff      PctDiscMM      PctDiscCH
## Min. :1.390   Min. : -0.6700   Min. :0.00000   Min. :0.00000
## 1st Qu.:1.750   1st Qu.: 0.0000   1st Qu.:0.00000   1st Qu.:0.00000
## Median :1.860   Median : 0.2400   Median :0.00000   Median :0.00000
## Mean   :1.813   Mean   : 0.1464   Mean   :0.06164   Mean   :0.02739
## 3rd Qu.:1.890   3rd Qu.: 0.3200   3rd Qu.:0.11834   3rd Qu.:0.00000
## Max.   :2.090   Max.   : 0.6400   Max.   :0.40201   Max.   :0.25269
## ListPriceDiff
## Min. :0.0000
## 1st Qu.:0.1400
## Median :0.2400
## Mean   :0.2225
## 3rd Qu.:0.3000
## Max.   :0.4400
```

```
corr <- cor(df[-1]) #correlogram of numeric variables, excluding outcome variable
testDf <- cor.mtest(df[-1], conf.level = 0.95) #compute significance of correlation
# Plot correlogram
```

```
corrplot(corr, p.mat = testDf$p, method = 'number', type = 'lower', insig='blank',
          addCoef.col = 'black', number.cex = 0.6, order = 'AOE', diag=FALSE, tl.srt = 45, tl.col = 'black')
```

