

# Project Presentation

Max Chis





# Methodology

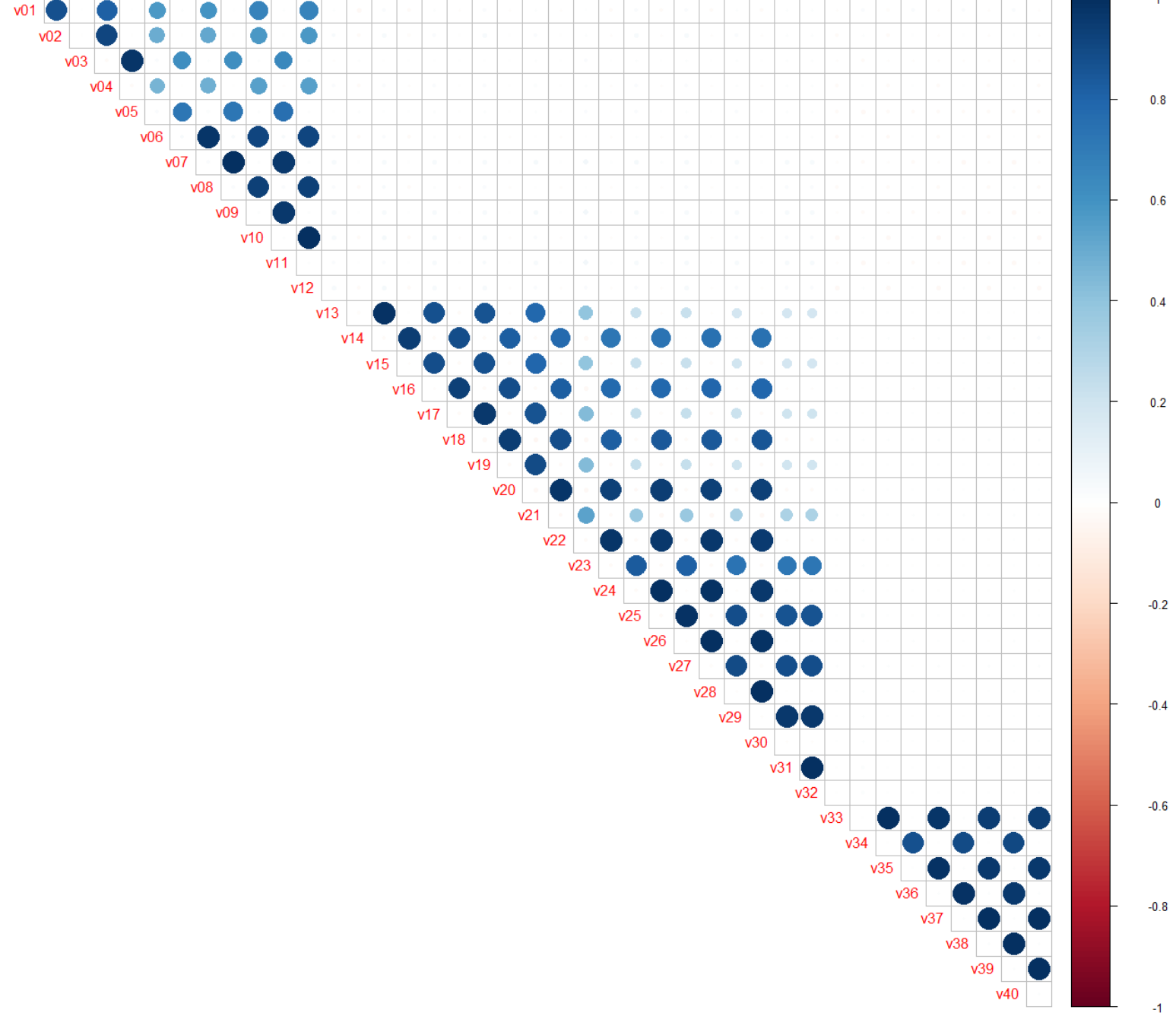
Data exploration with emphasis on variable correlation, relationship of variable values with event probability and continuous output

Regression and Classification Models trained, including:

- Linear and Quadratic Additive
- All pairwise interactions
- Interactions of splines of selected variables
- Elastic Net
- K Nearest Network
- Support Vector Machine
- Gradient Boosted Tree
- Random Forest
- Neural Networks

# Variable Correlation

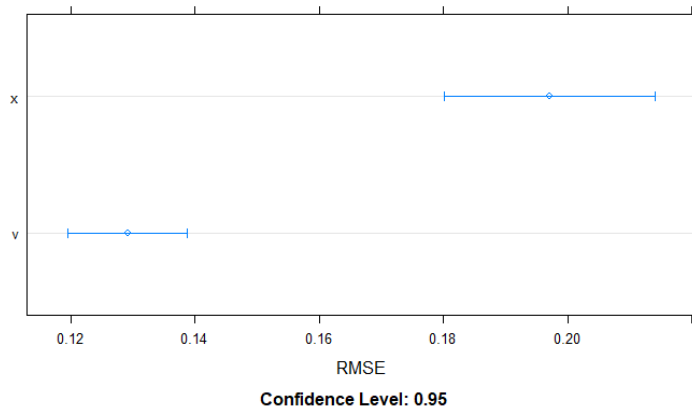
- V variables more likely to be strongly correlated, compared to X variables
- Highest X variable correlation was  $0.15\rho$
- Highest V variable correlation was  $1\rho$



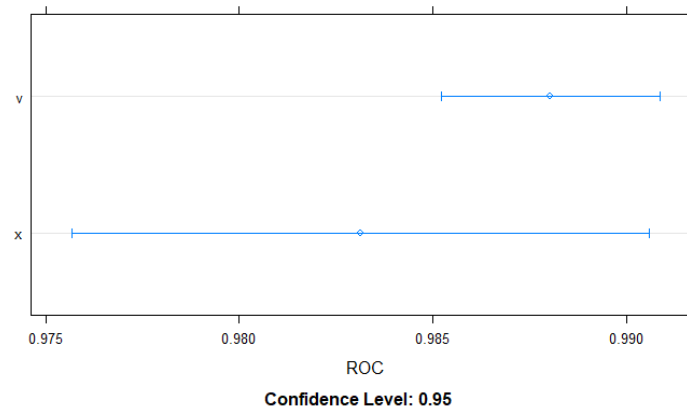
# Comparing Performance Between V and X Models

- Top-performing V model outperformed top-performing X model in Regression
- Reverse is true for classification in terms of Accuracy
- Comparable performance for classification in terms of ROC

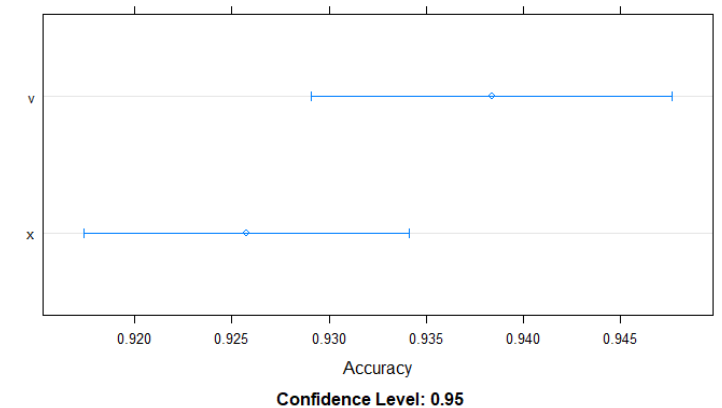
RMSE of Top-Performing Regression Models



ROC of Top-Performing Classification Models



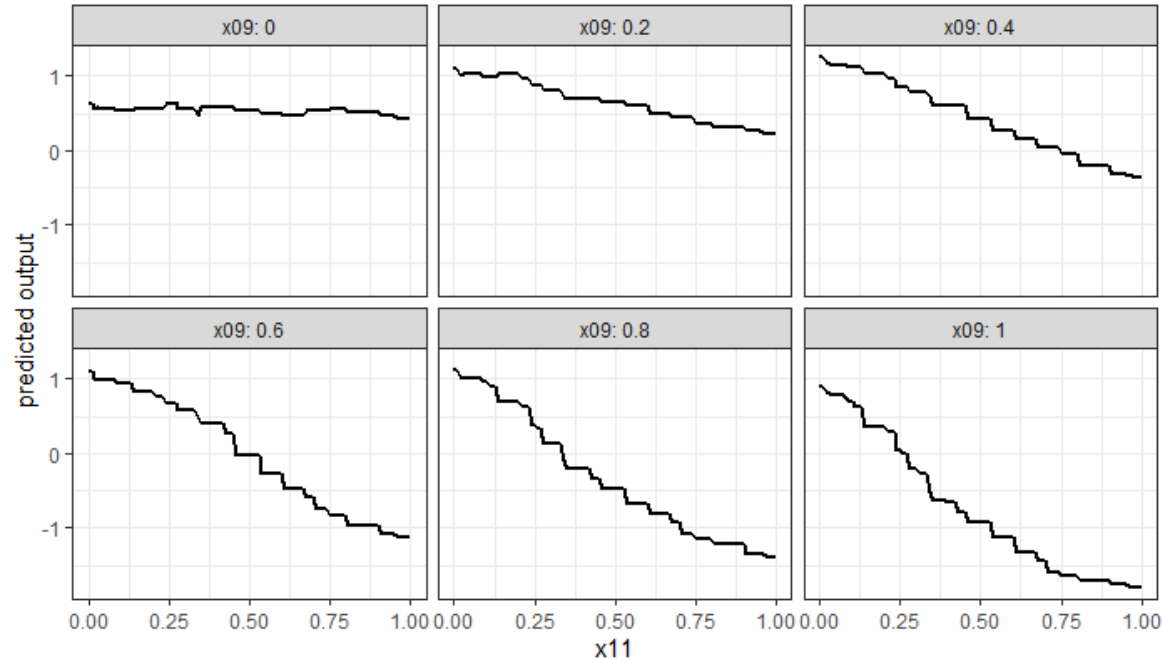
Accuracy of Top-Performing Classification Models



# Best-Performing Model – X Variables

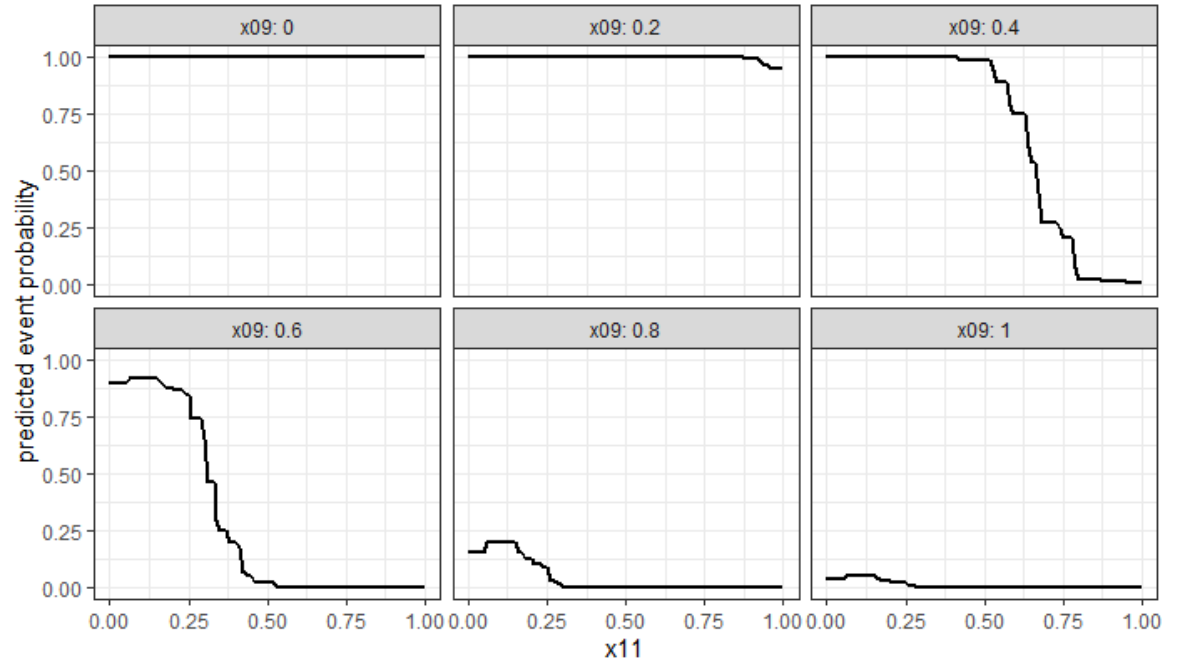
Continuous Output (Regression)

x11 Value vs Predicted Mean Trend, x09 facets, GB Tree



Event Probability (Classification)

x11 Value vs Predicted Event Probability, x09 facets, GB Tree (ROC)



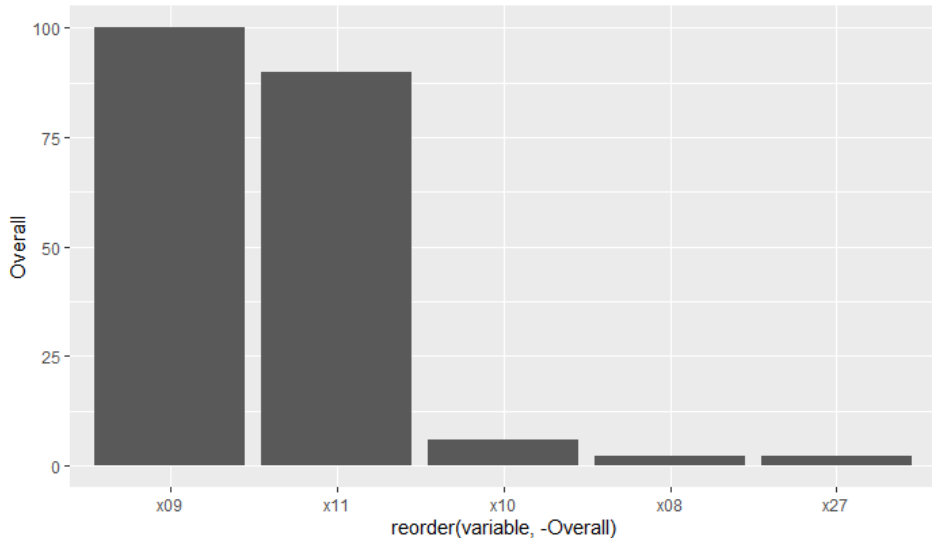
- Model fitted using gradient-boosted tree.
- Maximizing x09 and x11 minimizes continuous output and error probability



# Variable Importance – X Variables

- X09 and x11 consistently found to be most important variables
- All other variables have marginal importance to model.
- Importance of X09 and X11 consistently found across different models.

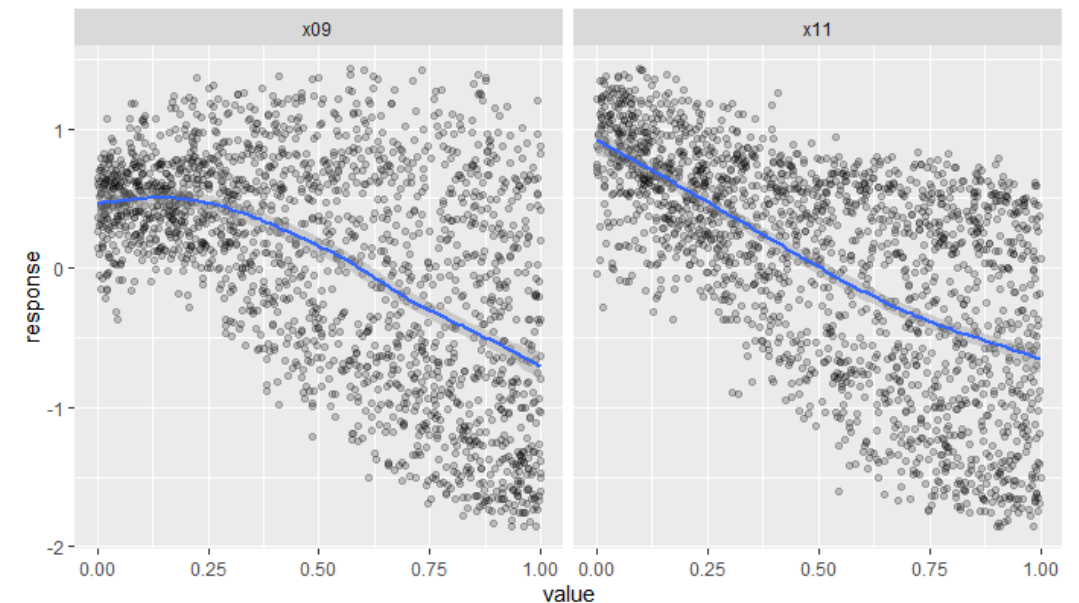
Top 5 most important variables in Gradient Boosted Tree Model



Density of data points at different values, classified by event: x09 and x11

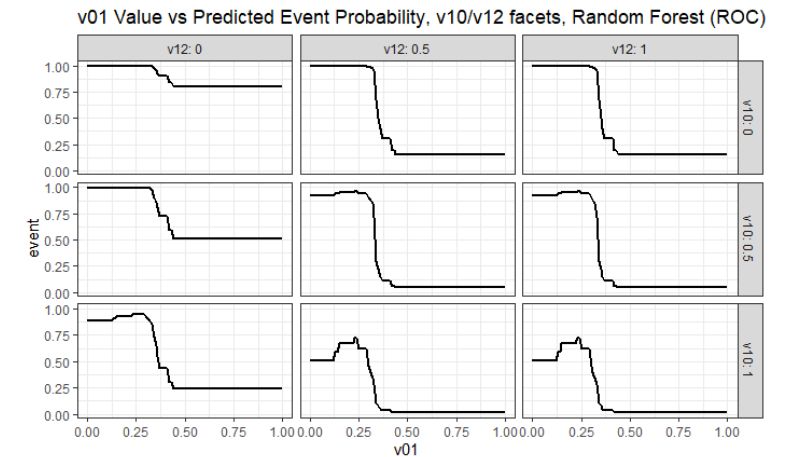
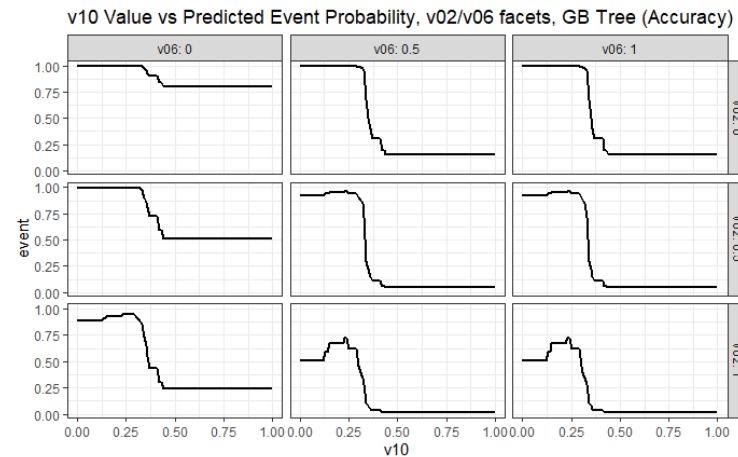
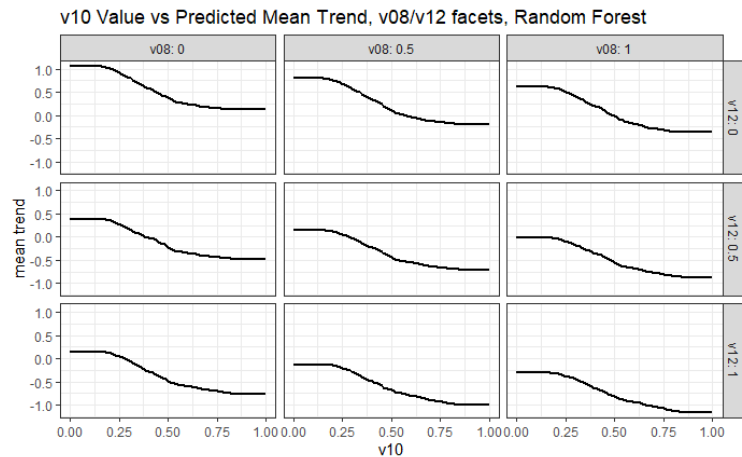


Value vs. Response of x09, x11



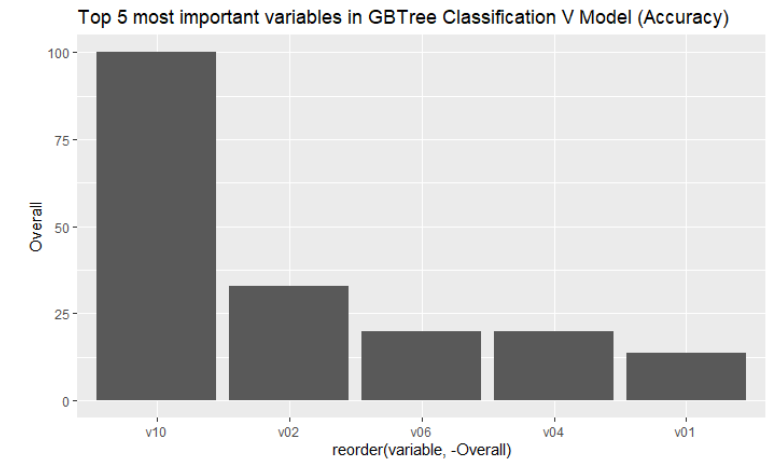
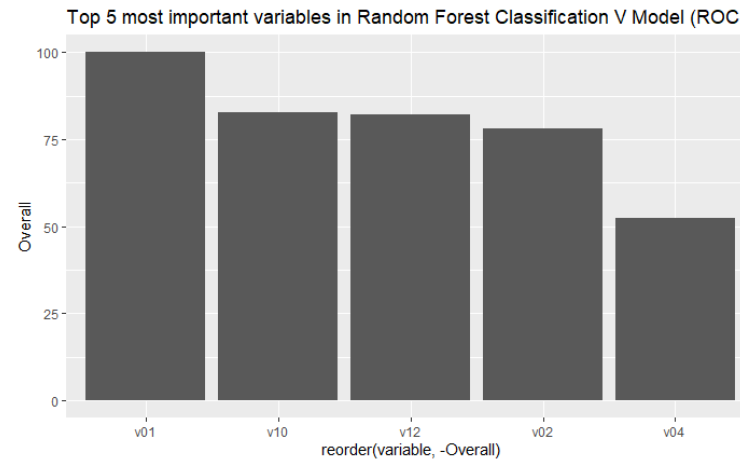
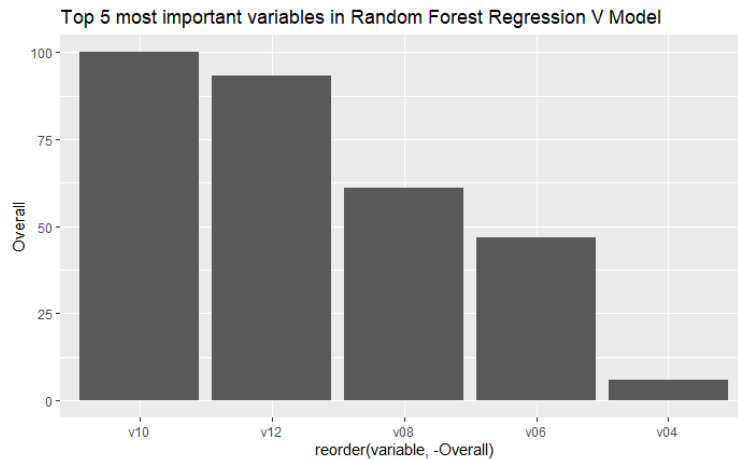
## Best Performing Model – V Variables

- Gradient-Boosted tree maximized Accuracy (in Classification) and RMSE (in Regression)
- Random Forest maximized ROC scores in Classification
- Maximizing v01,v02, v06, v08, v10, v12 minimizes continuous output and event probability.



## Variable Importance – V Variables

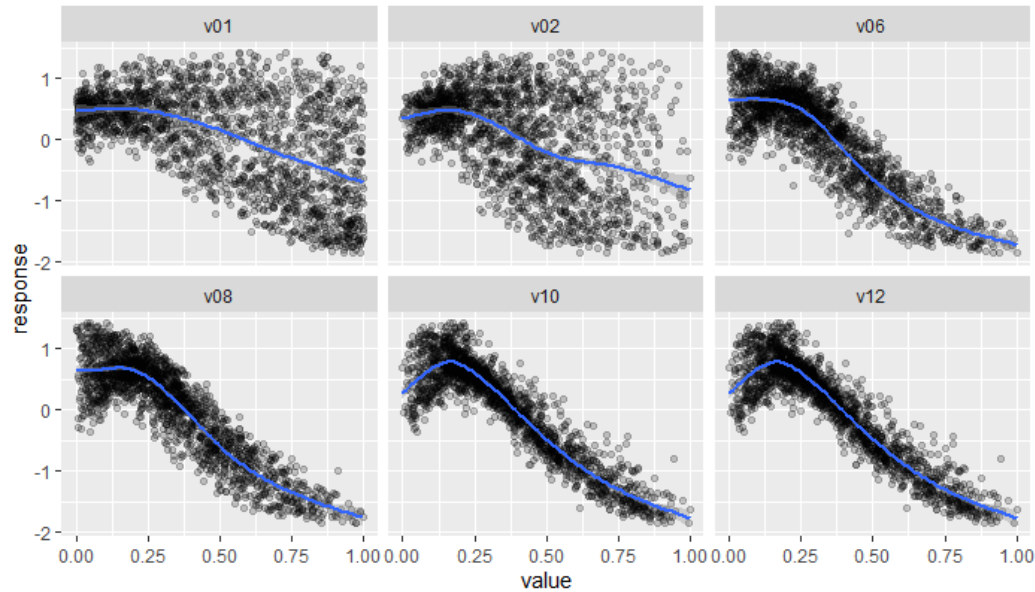
- V models disagreed more than X models on which variables were most important
- Consistent top-performers included v01, v02, v06, v08, v10, and v12



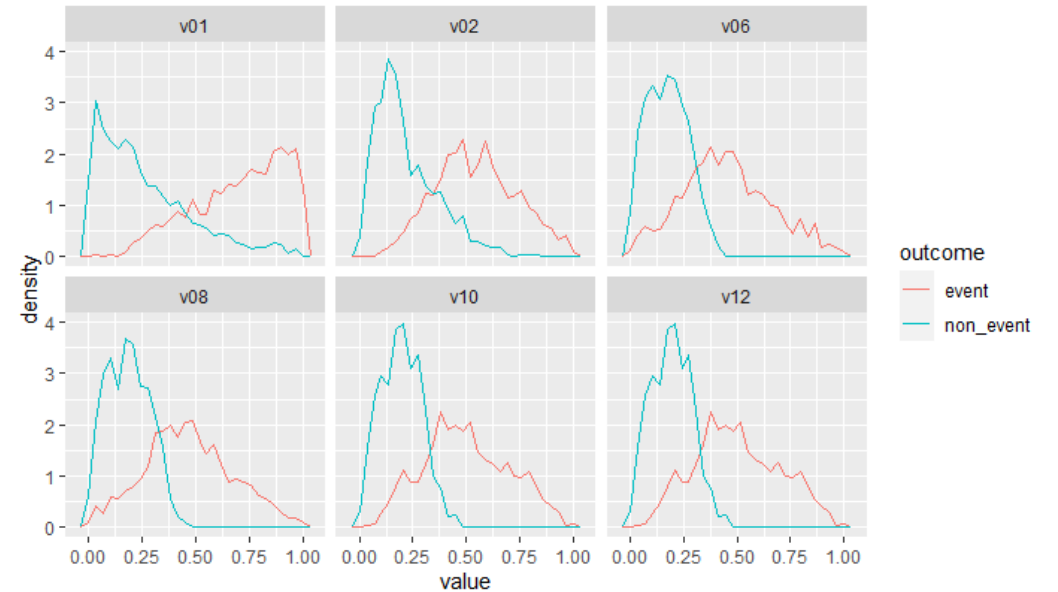


# Plotting Most Important V Variables

Value vs Response: Most Important V Variables



Density of V data points vs. value by event classification: Most important variables



# Performance on Holdout Set

X Holdout Set (GBTree) utilized

The holdout test set performance metrics associated with your model are displayed below.

Note: nothing is shown until you upload the CSV file.

.metric	.estimator	.estimate
rmse	standard	0.21
rsq	standard	0.93
mae	standard	0.16
accuracy	binary	0.91
mn_log_loss	binary	0.21
roc_auc	binary	0.98