

RemixAutoML Library Introduction

Adrian Antico

2019-04-27

Contact Info

LinkedIn: <https://www.linkedin.com/in/adrian-antico/>

Remix Institute: <https://www.remyxcourses.com> or <https://www.remixinstitute.ai>

Vignette Intent

This vignette is designed to give you the highlights of the set of automated machine learning functions available in the RemixAutoML package. To see the functions in action, visit the Remyx Courses website for the free course at <https://www.remyxcourses.com> and walk through them (and check out the other courses too!).

Package Goals

The **RemixAutoML** package (*Remix Automated Machine Learning*) is designed to automate and optimize the quality of machine learning, the pace of development, along with the handling of big data and the processing time of data management. The library has been a development task at [Remix Institute](https://www.remyxcourses.com) over the course of the past year to consolidate all of our winning methods for successfully completing machine learning and data science consulting projects. These were actual projects at Fortune 500 companies, Fortune 100 companies, tech startups, and other consulting clients. We are avid R users and feel that the R community could benefit from its release.

Package Design Philosophy

Core packages utilized

The two core packages RemixAutoML relies on are H2O and data.table. There are other packages used, for example, the forecast package, but H2O and data.table are used the most. I use data.table for data wrangling of all internal functions due to its ability to handle big data with minimal memory and the speed at which their functions process data. I chose to use H2O and their machine learning algorithms because of their high quality results, flexibility of use, ease of operationalization, and ability to manage big data. I use these functions routinely for machine learning projects and they continue to outperform every other package / software I test them against. Many of the other R packages for modeling or data manipulation have slow run times and fail once I get going with bigger data. As a simple example, sometimes I do a one-hot encoding for testing out keras models and I can just run my DummifyDT function to create those. I've tried out a few others with fast runtimes but they fail on bigger data, mess up the column ordering, and don't offer the flexibility of creating one-hot encoding features versus standard dichotomized features.

Machine learning methodology

The package is designed to give your models to the best chance at being the most accurate for your machine learning tasks. There are functions in here to help you get the most out of your data for the models to take advantage of. There are three types of features you need to manage for machine learning: numeric,

factor, and text data. The functions in here will help you squeeze as much information out of your data set as possible. Then you feed those data sets into the best machine learning algorithms where you can automatically grid-tune and let them compete against each other.

Handling numerical features

The typical challenge with numerical data is modeling nonlinear relationships. The models used handle that aspect so you don't have to worry about it. Interactions of numerical data are also handled internally as tree-based models and deep learning do well at capturing those relationships. The areas of numerical data typically missed relates to the effect of numerical features across time. The existence of a specific value can have impacts over time. Trends in the features and target are also predictive of future events. The length of history you account for also matters. The *GDL* functions will let you generate lags from the target variable and other numeric features, along with rolling statistics on the target variable and other numerical features. You can also specify to return information about the numerical time between records and generate their lags and rolling statistics on the target and other numerical features.

Handling categorical features

In other modeling frameworks, such as Python, you need to convert your categorical features into dummy variables. This means you need to take care of that coding task along with managing that in a production environment. With this package, categorical features are handled internally within the modeling function **AutoH2OModeler()**. Turning your categorical features into dummy variables is problematic for high cardinality factor variables. There are other approaches and the AutoH2OModeler will actually test out the other methods to see which ones offer the best performance. So you don't have to deal with the coding and management of factor variables and you get better performance.

Handling text data

We are living in a world now where text data is becoming more readily available. Now, you can simply run the **AutoWord2VecModeler()** function to build models for all your text columns, save the models, and recreate them on the fly with the **AutoH2OScoring()** function in your production setting.

Machine learning algorithms

I've seen several packages and automated machine learning software that brags about all the models they have available. However, when you've been doing this work as long as I have, you eventually realize that most of the models out there aren't competitive. In fact, H2O was smart not to go after all the possible algorithms out there and instead focused on enhancing the models with a track record of being the best for regression, quantile regression, classification, and multinomial classification. The grid-tuning capabilities are excellent for squeezing that last drop of gain out of your machine learning project. Most of the gains with come from the feature engineering but you won't be able to squeeze more gains out of other package models. On top of that, you can manage the amount of memory and threads you are willing to use for your machine learning tasks and when you work in a shared environment, that can come in handy. Besides H2O, we have a few automated catboost models available as well. One for classification tasks and one for regression tasks (which includes quantile regression). Catboost has the added benefit of being able to train using a GPU and the performance gains are quite significant. The accuracies are also the most competitive I've seen thus far.

Interpreting your models

Once you have your models developed, you or a boss may want to see what features are most important and their relationship to the target variable. The functions available are also run internally to build out partial

dependence calibration plots. These show you the modeled predicted relationship along with the empirical relationship in the same graph to show the end user how accurate the relationship is what the relationship is (even for categorical variables). Variable importance tables are also available to view along with evaluation calibration plots and boxplots.

Operationalizing models with ease

The scoring process in machine learning is typically straightforward. We have an all-purpose scoring function to score all your supervised machine learning models, text models, and clustering models, regardless of type (mojo or standard model files).

How to install H2O

Follow this link to install H2O if it isn't on your machine already. [Install H2O](#)

There are seven categories of functions (currently) in this library I'll go over:

- Automated Supervised Learning
- Automated Unsupervised Learning
- Automated Model Evaluation
- Automated Feature Interpretation
- Automated Cost Sensitive Optimization
- Automated Feature Engineering
- A Few Miscellaneous Functions

Automated Supervised Learning Functions

Functions include:

- `AutoCatBoostClassifier()`
- `AutoCatBoostRegression()`
- `AutoH2OModeler()`
- `AutoH2OScoring()`
- `AutoTS()`
- `AutoRecommender()`
- `AutoRecommenderScoring()`
- `AutoNLS()`

AutoH2OModeler()

The supervised learning functions handle multiple tasks internally. The **AutoH2OModeler** function can build any number of H2O models, automatically compare hyper-parameter tuned versions to baseline versions, selecting a winner, saving the model evaluation and feature interpretation metrics / graphs, along with storing models and their metadata to refer to them later in a production setting.

The models available include:

- Gradient Boosting Machines
- LightGBM (Linux only)
- Distributed Random Forest

- XGBoost (Linux only)
- Deeplearning
- AutoML (for Windows users XGBoost and LightGBM are not available)

For Windows users (Mac?), XGBoost is not available and therefore neither is LightGBM (XGBoost and LightGBM are not utilized in AutoML model selection with Windows).

AutoH2OScoring()

This function is the complement of the **AutoH2OModeler**, **AutoKMeans**, and **AutoWord2VecModeler** functions. Specify which rows of your model metadata collection file to run and **AutoH2OScoring** will return a list of predicted values, where each element of the list is a set of predicted values from the model it ran. For the **AutoH2OModeler** you will generate a file called `grid_tuned_paths.Rdata` which contains the path to your models (among other items) that you can pass along to the **AutoH2OScoring** function to automatically score your models. For the **AutoKMeans** you will generate a file called `KMeansModelFile.Rdata` which contains the paths to the models for scoring your GLRM and KMeans models. For the **AutoWord2VecModeler** you will generate a file called `StoreFile.Rdata` which contains the paths to your word2vec models for scoring.

In total, the AutoH2OScoring function can score

- Regression models
- Quantile regression Models
- Binary classification Models
- Multinomial classification Models
- Multioutcome multinomial classification models
- Generalized low rank dimensionality reduction models
- KMeans clustering models
- Word2vec models

AutoTS()

Another automated supervised learning function we have is an automated time series modeling function (AutoTS) that optimally builds out seven types of time series forecasting models, compares them on holdout data, picks a winner, rebuilds the winner on full data, and generates the forecasts for the number of desired periods. The intent is to make these processes fast, easy, and of high quality. Every model makes use of the optimal settings of their parameters to give them the best chance of being the best. Each model uses a Box-Cox transformation on the target variable and all predictions are back-transformed. It also compares model-based frequency determination versus user-supplied (for the TimeUnit argument) along with the option to have imputation and outlier replacement conducted.

The competing models include:

- DSHW (Double Seasonal Holt Winters)
- ARFIMA (Autoregressive Fractional Integrated Moving Average)
- ARIMA (Autoregressive Integrated Moving Average)
- ETS (Exponential Smoothing and Holt Winters)
- TBATS (Exponential Smoothing State Space Model with Box-Cox Transformation, ARMA Errors, Trend and Seasonal Components)
- TSLM (Time Series Linear Model)
- NN (Autoregressive Neural Network)

AutoRecommender()

This function builds out several variations of collaborative filtering models on a binary ratings matrix. To automatically build the binary ratings matrix, see **RecomDataCreate**.

The competing models include:

- RandomItems
- PopularItems
- UserBasedCF
- ItemBasedCF
- AssociationRules

AutoRecommenderScoring()

This function will automatically score your winning model. Simply feed in your data and the winning model returned from the **AutoRecommender** function and this function will generate a table of several recommended products (by rank) for each entity.

AutoNLS()

This automated supervised learning function builds nonlinear regression models for a more niche set of tasks. It's set up to generate interpolation predictions, such as smoothing cost curves for optimization tasks. It returns the interpolated data, the winning model name, the model object, and the evaluation metrics table.

The competing models include:

- Asymptotic
- Asymptotic through origin
- Asymptotic with offset
- Bi-exponential
- Four parameter logistic
- Three parameter logistic
- Gompertz
- Michal Menton
- Weibull
- Polynomial regression or monotonic regression

Example of AutoNLS()

Find more demos at <https://www.remixcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
library(RemixAutoML)

# Create Growth Data
data <-
  data.table::data.table(Target = seq(1, 500, 1),
                        Variable = rep(1, 500))
for (i in as.integer(1:500)) {
  if (i == 1) {
    var <- data[i, "Target"][[1]]
```

```

      data.table::set(data,
                      i = i,
                      j = 2L,
                      value = var * (1 + runif(1) / 100))
    } else {
      var <- data[i - 1, "Variable"][[1]]
      data.table::set(data,
                      i = i,
                      j = 2L,
                      value = var * (1 + runif(1) / 100))
    }
  }
}

# Add jitter to Target
data[, Target := jitter(Target,
                        factor = 0.50)]

# To keep original values
data1 <- data.table::copy(data)

# Build models
data11 <- AutoNLS(
  data = data,
  y = "Target",
  x = "Variable",
  monotonic = TRUE
)

# Join predictions to source data
data2 <- merge(
  data1,
  data11$PredictionData,
  by = "Variable",
  all = FALSE
)

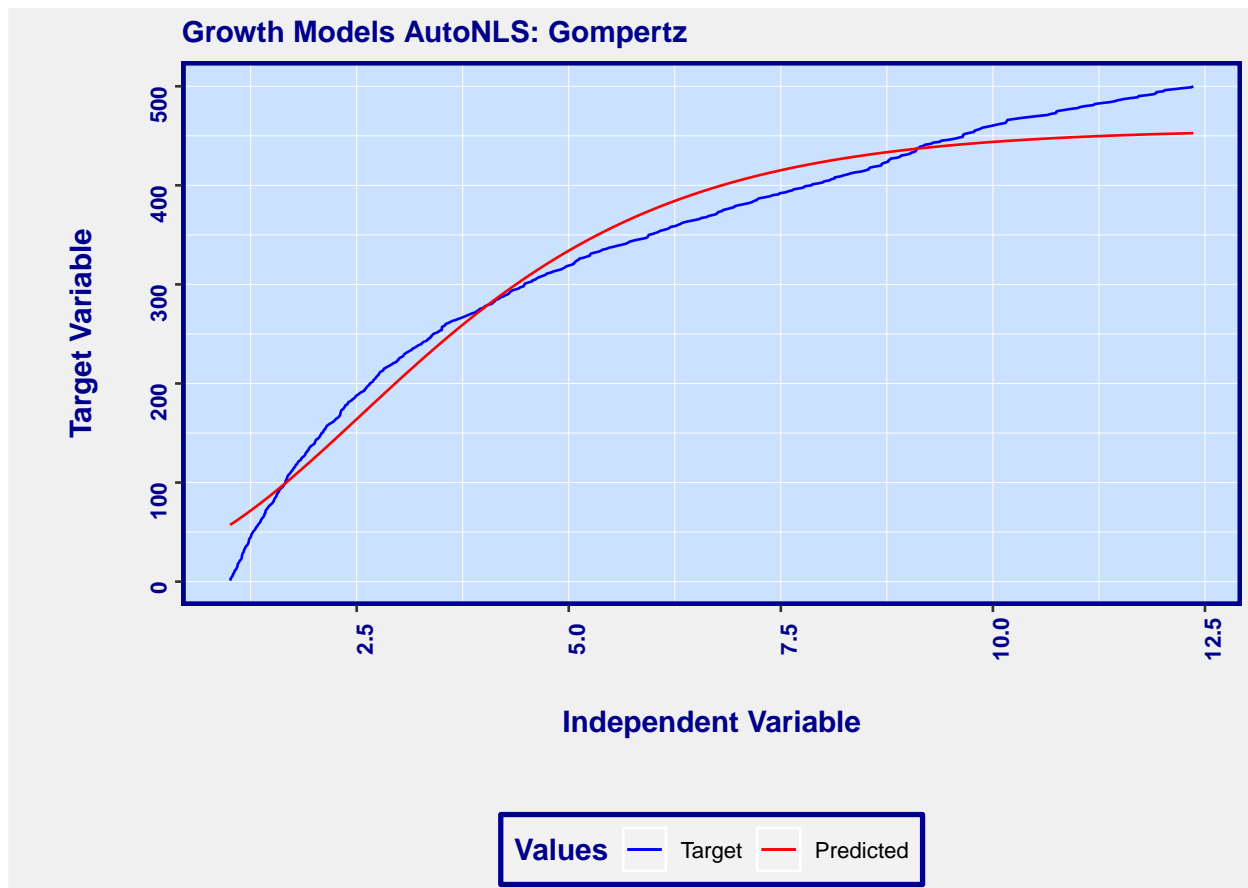
# Plot output
ggplot2::ggplot(data2, ggplot2::aes(x = Variable)) +
  ggplot2::geom_line(ggplot2::aes(y = data2[["Target.x"]],
                                   color = "Target")) +
  ggplot2::geom_line(ggplot2::aes(y = data2[["Target.y"]],
                                   color = "Predicted")) +
  RemixAutoML::ChartTheme(Size = 12) +
  ggplot2::ggtitle(paste0("Growth Models AutoNLS: ",
                          data11$ModelName)) +
  ggplot2::ylab("Target Variable") +
  ggplot2::xlab("Independent Variable") +
  ggplot2::scale_colour_manual("Values",
                              breaks = c("Target",
                                           "Predicted"),
                              values = c("red",
                                           "blue"))

```

```

# Print model makeup and evaluation metrics
summary(data11$ModelObject)
#>
#> Formula: Target ~ SSGompertz(Variable, Asym, b2, b3)
#>
#> Parameters:
#>      Estimate Std. Error t value Pr(>|t|)
#> Asym 4.571e+02  2.908e+00  157.16  <2e-16 ***
#> b2    3.348e+00  6.956e-02   48.14  <2e-16 ***
#> b3    6.227e-01  5.829e-03  106.83  <2e-16 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 22.2 on 497 degrees of freedom
#>
#> Number of iterations to convergence: 0
#> Achieved convergence tolerance: 1.416e-06
data11$EvaluationMetrics
#>      ModelName MeanAbsError
#> 1:      Gompertz      18.88935
#> 2: Michal_Menton      20.19742
#> 3:      Logistic      24.25077
#> 4:         Poly      81.25082

```



Find more demos at <https://www.remyxcourses.com/course?courseid=intro-to-remixautoml-in-r>

8


```

SkipModels = NULL,
StepWise   = TRUE,
TSClean    = TRUE)

#> [1] "DSHW FITTING"
#> [1] "ARFIMA FITTING"
#> [1] "ARIMA FITTING"
#> [1] "ETS FITTING"
#> [1] "TBATS FITTING"
#> [1] "TSLM FITTING"
#> [1] "NNet FITTING"
#> [1] "NNet Iteration: 1"
#> [1] "NNet Iteration: 2"
#> [1] "NNet Iteration: 3"
#> [1] "NNet Iteration: 4"
#> [1] "NNet Iteration: 5"
#> [1] "NNet 2 Iteration: 1"
#> [1] "NNet 2 Iteration: 2"
#> [1] "NNet 2 Iteration: 3"
#> [1] "NNet 2 Iteration: 4"
#> [1] "NNet 2 Iteration: 5"
#> [1] "NNet 3 Iteration: 1"
#> [1] "NNet 3 Iteration: 2"
#> [1] "NNet 3 Iteration: 3"
#> [1] "NNet 3 Iteration: 4"
#> [1] "NNet 3 Iteration: 5"
#> [1] "NNet 4 Iteration: 1"
#> [1] "NNet 4 Iteration: 2"
#> [1] "NNet 4 Iteration: 3"
#> [1] "NNet 4 Iteration: 4"
#> [1] "NNet 4 Iteration: 5"
#> [1] "FIND WINNER"
#> [1] "GENERATE FORECASTS"
#> [1] "FULL DATA ARIMA FITTING"

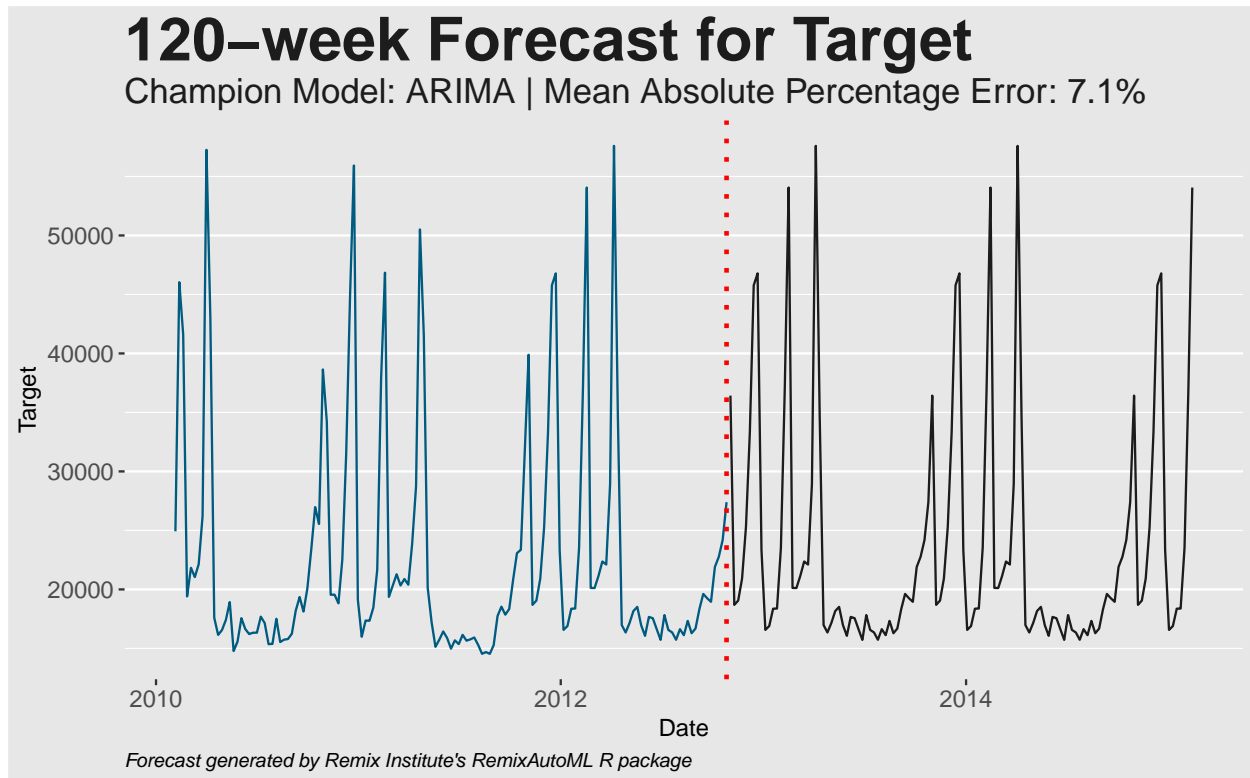
# Print the evaluation metric and model makeup
knitr::kable(output$EvaluationMetrics)

```

ModelName	MeanResid	MeanPercError	MAPE	ID
ARIMA	609.55	0.04692	0.07130	1
TBATS_TSC	1574.01	0.09270	0.11628	2
TBATS	688.19	0.03370	0.12458	3
TBATS_ModelFreq	688.19	0.03370	0.12458	4
TBATS_ModelFreqTSC	688.19	0.03370	0.12458	5
ARFIMA_TSC	431.74	0.02090	0.12866	6
NN	2167.72	0.11605	0.13049	7
ARIMA_ModelFreq	-1759.54	-0.08468	0.13118	8
ARIMA_ModelFreqTSC	-1759.54	-0.08468	0.13118	9
ARIMA_TSC	2121.11	0.14593	0.14820	10
ARFIMA	-2074.51	-0.09640	0.14925	11
ARFIMA_ModelFreq	-2074.51	-0.09640	0.14925	12
ARFIMA_ModelFreqTSC	-2074.51	-0.09640	0.14925	13
TSLM_ModelFreqTSC	-1306.83	-0.05191	0.17996	14
TSLM_TSC	2526.50	0.17461	0.18285	15

ModelName	MeanResid	MeanPercError	MAPE	ID
ETS	3271.82	0.19669	0.20523	16
ETS_ModelFreq	3271.88	0.19670	0.20523	17
ETS_ModelFreqTSC	3271.88	0.19670	0.20523	18
DSHW_ModelFreq	3634.13	0.22601	0.23281	19
DSHW_ModelFreqTSC	3634.13	0.22601	0.23281	20
NN_ModelFreq	3832.10	0.23880	0.24068	21
NN_ModelFreqTSC	3948.52	0.24826	0.24960	22
NN_TSC	19899.98	19898.97667	19898.97667	23

```
summary(output$TimeSeriesModel)
#> Series: dataTSTrain
#> ARIMA(0,0,1)(0,1,0)[52]
#> Box Cox transformation: lambda= TRUE
#>
#> Coefficients:
#>          ma1
#>      0.6695
#> s.e.  0.0719
#>
#> sigma^2 estimated as 52502546: log likelihood=-937.74
#> AIC=1879.49 AICc=1879.62 BIC=1884.51
#>
#> Training set error measures:
#>           ME      RMSE      MAE      MPE      MAPE      MASE
#> Training set -8.656056 5748.353 2431.659 -1.819206 9.799424 0.5835215
#>           ACF1
#> Training set 0.05599928
```



Automated Unsupervised Learning Functions

The suite of functions in this category currently handle optimized row-clustering and anomaly detection. For the row-clustering, we utilize H2O's Generalized Low Rank Model and their KMeans algorithm, with hyper-parameter tuning for both. The function automatically adds the clusters to your data and can save the models for scoring new data with the **AutoH2OScoring** function. We have a few others currently in development and will release those when they are complete. The anomaly detection functions we have currently are for time series applications. We have a control chart methodology version that lets you build upper and lower confidence bounds by up to two grouping variables along with a time series modeling version. The clustering function and the control chart method function update your data set that you feed in with new columns that store the clusterID or anomaly information. The time series function updates your data, supplies you with the final time series model built, and a data.table that only contains anomalies.

Functions include:

- GenTSAnomVars()
- ResidualOutliers()
- AutoKMeans()

Demo of ResidualOutliers()

Find more demos at <https://www.remixcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
# Run on (Target - Predicted)
library(RemixAutoML)
data <- data.table::data.table(DateTime = as.Date(Sys.time()),
```

```

                                Target = as.numeric(stats::filter(rnorm(1000,
                                                                    mean = 50,
                                                                    sd = 20),
                                                                    filter=rep(1,10),
                                                                    circular=TRUE)))
data[, temp := seq(1:1000)][, DateTime := DateTime - temp][, temp := NULL]
data <- data[order(DateTime)]
data[, Predicted := as.numeric(stats::filter(rnorm(1000,
                                                    mean = 50,
                                                    sd = 20),
                                                    filter=rep(1,10),
                                                    circular=TRUE)))]

# Run function and collect results
stuff <- ResidualOutliers(data = data,
                          DateColName = "DateTime",
                          TargetColName = "Target",
                          PredictedColName = "Predicted",
                          TimeUnit = "day",
                          maxN = 5,
                          tstat = 2)

data      <- stuff$FullData
model     <- stuff$ARIMA_MODEL
outliers  <- data[type != "<NA>"]

# Create Plots
p1 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Preds),
                    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "AO", "DateTime"],
                    ggplot2::aes(xintercept = outliers[
                      type == "AO"][["DateTime"]]),
                    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Additive Outliers")

p2 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Residuals),
                    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "IO", "DateTime"],
                    ggplot2::aes(xintercept = outliers[
                      type == "IO"][["DateTime"]]),
                    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Innovational Outliers")

p3 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Preds),
                    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "LS", "DateTime"],
                    ggplot2::aes(xintercept = outliers[
                      type == "LS"][["DateTime"]]),

```

```

        linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Level Shift")

p4 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Residuals),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "TC", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "TC"][["DateTime"]]),
    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Transient Change")

# Print plots
RemixAutoML::multiplot(plotlist = list(p1,p2,p3,p4), cols = 2)

# Run on Target data
data <- data.table::data.table(DateTime = as.Date(Sys.time()),
  Target = as.numeric(stats::filter(rnorm(1000,
    mean = 50,
    sd = 20),
    filter=rep(1,10),
    circular=TRUE)))

data[, temp := seq(1:1000)][, DateTime := DateTime - temp][, temp := NULL]
data <- data[order(DateTime)]
data[, Predicted := as.numeric(stats::filter(rnorm(1000,
  mean = 50,
  sd = 20),
  filter=rep(1,10),
  circular=TRUE)))]

# Run function and collect results
stuff <- ResidualOutliers(data = data,
  DateColName = "DateTime",
  TargetColName = "Target",
  PredictedColName = NULL,
  TimeUnit = "day",
  maxN = 5,
  tstat = 2)

data <- stuff$FullData
model <- stuff$ARIMA_MODEL
outliers <- data[type != "<NA>"]

# Create Plots
p11 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Preds),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "AO", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "AO"][["DateTime"]]),
    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Additive Outliers")

```

```

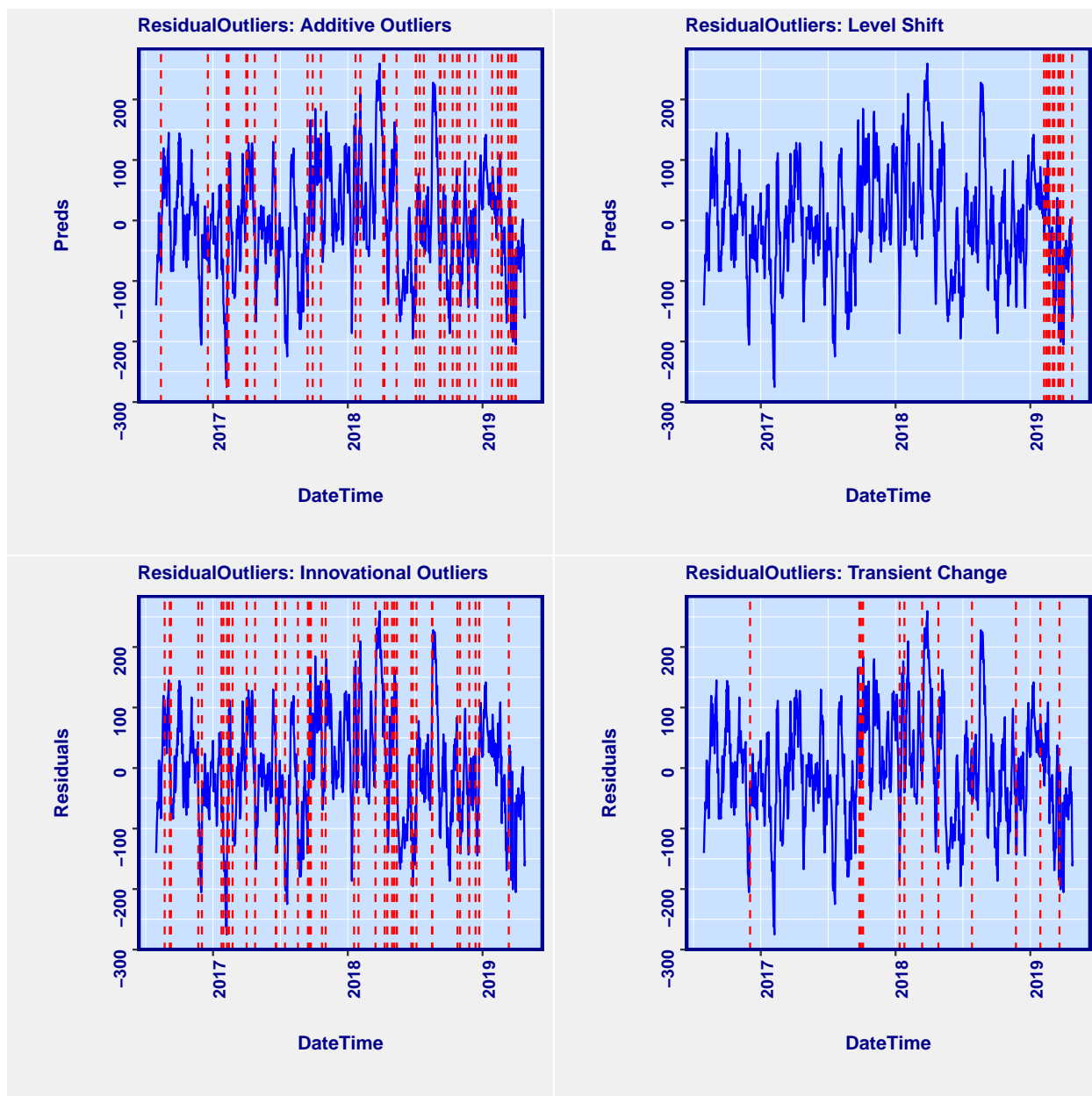
p22 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Residuals),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "IO", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "IO"][["DateTime"]]),
    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Innovational Outliers")

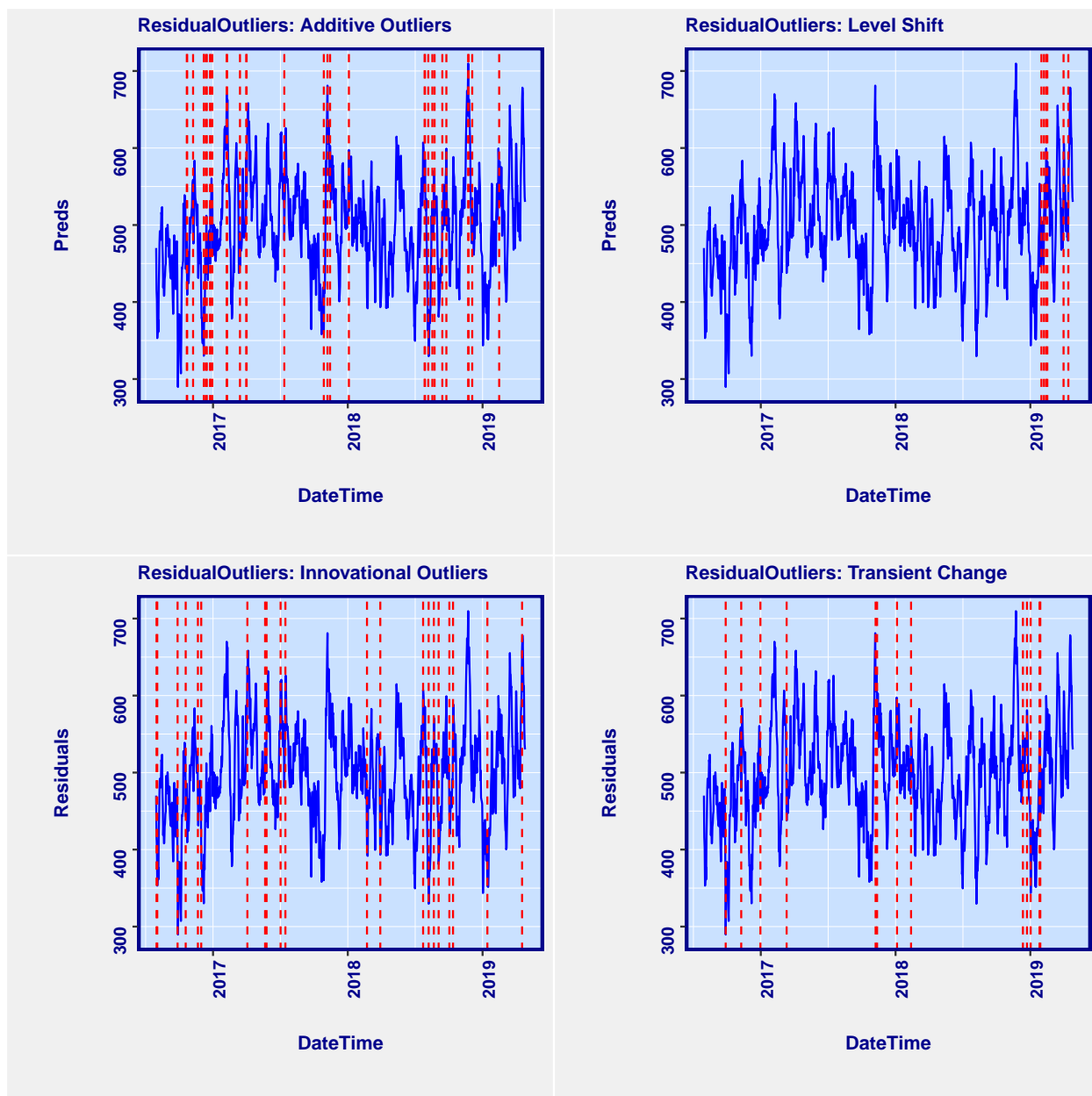
p33 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Preds),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "LS", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "LS"][["DateTime"]]),
    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Level Shift")

p44 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Residuals),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "TC", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "TC"][["DateTime"]]),
    linetype = 8, colour = "red") +
  ggplot2::ggtitle("ResidualOutliers: Transient Change")

# Print plots
RemixAutoML::multiplot(plotlist = list(p11,p22,p33,p44), cols = 2)

```





Automated Model Evaluation, Feature Interpretation, and Cost Sensitive Optimization Functions

The model evaluation graphs are calibration plots or calibration boxplots. The calibration plots are used for regression (expected value and quantile regression), classification, and multinomial modeling problems. The calibration boxplots are used for regression (expected value and quantile regression). These graphs display both the actual target values and the predicted values, grouped by the number of bins that you specify. The calibration boxplots are useful to understand not only the model bias but also the model variance, across the range of predicted values.

Functions include:

- EvalPlot()

Demo of EvalPlot()

Find more demos at <https://www.remixxcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
library(RemixAutoML)

# Data generator function
dataGen <- function(Correlation = 0.95) {
  Validation <- data.table::data.table(target = runif(1000))
  Validation[, x1 := qnorm(target)]
  Validation[, x2 := runif(1000)]
  Validation[, predict := pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2))]
  return(Validation)
}

# Store data sets
data1 <- dataGen(Correlation = 0.50)
data2 <- dataGen(Correlation = 0.75)
data3 <- dataGen(Correlation = 0.90)
data4 <- dataGen(Correlation = 0.99)

# Generate EvalPlots (calibration)
p1 <- RemixAutoML::EvalPlot(data = data1,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "calibration",
                             PercentileBucket = 0.05,
                             aggrfun = function(x) mean(x,
                                                         na.rm = TRUE))
p1 <- p1 + ggplot2::ggtitle("Calibration Evaluation p1: Corr = 0.50") +
  RemixAutoML::ChartTheme(Size = 10)

p2 <- RemixAutoML::EvalPlot(data = data2,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "calibration",
                             PercentileBucket = 0.05,
                             aggrfun = function(x) mean(x,
                                                         na.rm = TRUE))
p2 <- p2 + ggplot2::ggtitle("Calibration Evaluation p2: Corr = 0.75") +
  RemixAutoML::ChartTheme(Size = 10)

p3 <- RemixAutoML::EvalPlot(data = data3,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "calibration",
                             PercentileBucket = 0.05,
                             aggrfun = function(x) mean(x,
                                                         na.rm = TRUE))
p3 <- p3 + ggplot2::ggtitle("Calibration Evaluation p3: Corr = 0.90") +
```

```

RemixAutoML::ChartTheme(Size = 10)

p4 <- RemixAutoML::EvalPlot(data = data4,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "calibration",
                             PercentileBucket = 0.05,
                             aggrfun = function(x) mean(x,
                                                         na.rm = TRUE))
p4 <- p4 + ggplot2::ggtitle("Calibration Evaluation p4: Corr = 0.99") +
  RemixAutoML::ChartTheme(Size = 10)
RemixAutoML::multiplot(plotlist = list(p1,p2,p3,p4), cols = 2)

# Generate EvalPlots (boxplots)
p1 <- RemixAutoML::EvalPlot(data = data1,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p1 <- p1 + ggplot2::ggtitle("Calibration Evaluation p1: Corr = 0.50") +
  RemixAutoML::ChartTheme(Size = 10)

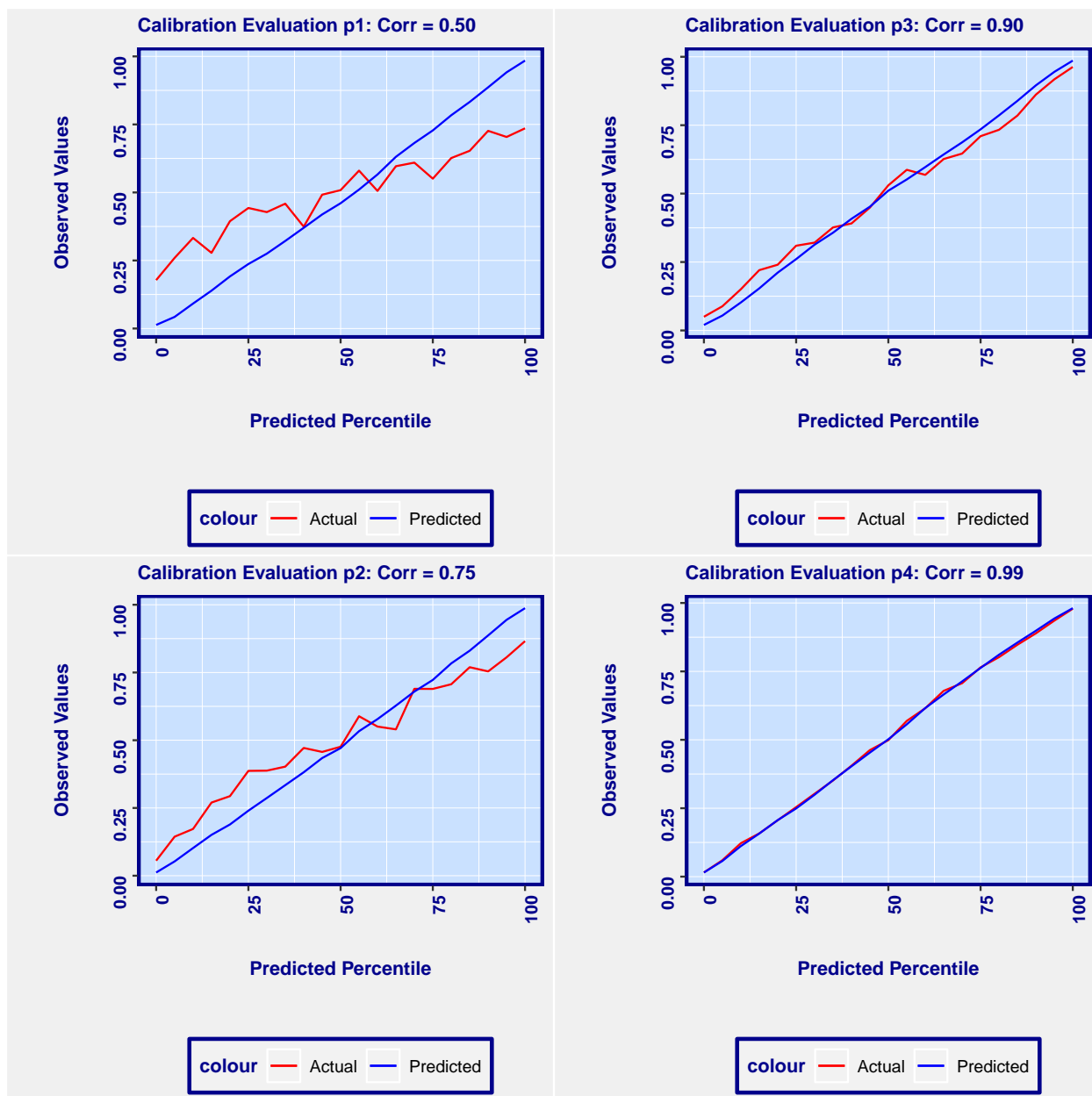
p2 <- RemixAutoML::EvalPlot(data = data2,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p2 <- p2 + ggplot2::ggtitle("Calibration Evaluation p2: Corr = 0.75") +
  RemixAutoML::ChartTheme(Size = 10)

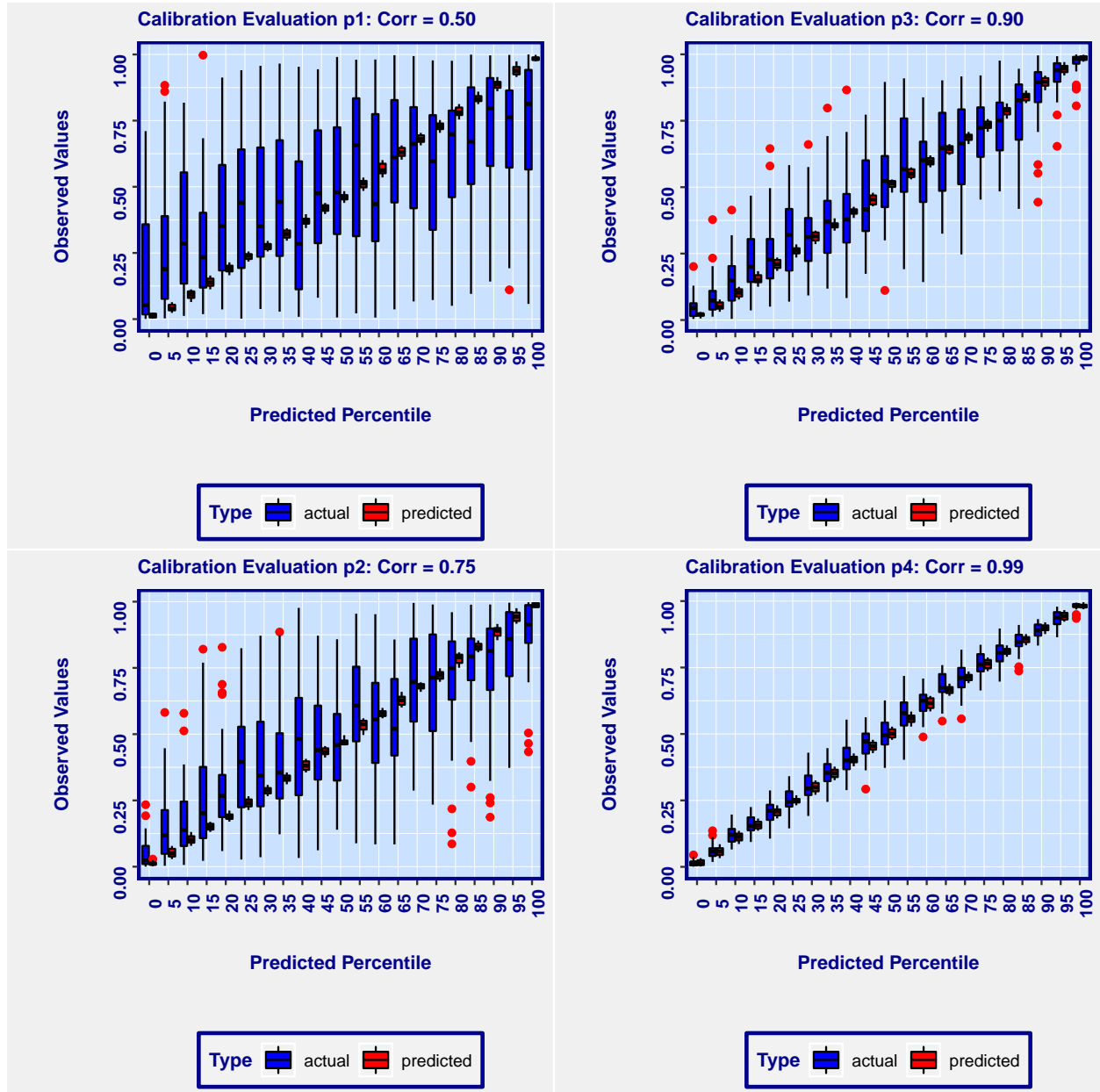
p3 <- RemixAutoML::EvalPlot(data = data3,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p3 <- p3 + ggplot2::ggtitle("Calibration Evaluation p3: Corr = 0.90") +
  RemixAutoML::ChartTheme(Size = 10)

p4 <- RemixAutoML::EvalPlot(data = data4,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p4 <- p4 + ggplot2::ggtitle("Calibration Evaluation p4: Corr = 0.99") +
  RemixAutoML::ChartTheme(Size = 10)

RemixAutoML::multiplot(plotlist = list(p1,p2,p3,p4), cols = 2)

```





The feature interpretation function graphs are very similar in nature to the model evaluation graphs. They display partial dependence calibration line plots, partial dependence calibration boxplots, and partial dependence calibration bar plots (for factor variables with the ability to limit the number of factors shown with the remainder grouped into “other”). The line graph version is for numerical features and have the ability to aggregate by quantile for quantile regression.

Functions include:

- `ParDepCalPlots()`

The cost sensitive optimization functions provide the user the ability to generate utility-optimized thresholds for classification tasks. There are two of these functions: one for generating a single threshold based on the values supplied to your cost confusion matrix outcomes and the second one provides two thresholds, where

your final predicted classification could be (0|1) and “do something else”. With the latter function, you would also need to supply a cost to the “do something else” option.

Demo of ParDepCalPlots() for calibration, boxplots, and barplots

Find more demos at <https://www.remixcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
library(RemixAutoML)

# Data generator function
dataGen <- function(Correlation = 0.95) {
  Validation <- data.table::data.table(target = runif(1000))
  Validation[, x1 := qnorm(target)]
  Validation[, x2 := runif(1000)]
  Validation[, predict := pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2))]
  Validation[, Feature1 := (pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2)))^1.25]
  Validation[, Feature2 := (pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2)))^0.25]
  Validation[, Feature3 := (pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2)))^(-1)]
  Validation[, Feature4 := pnorm(Correlation * x1 +
                                sqrt(1 - Correlation ^2) * qnorm(x2))]
  Validation[, Feature4 := ifelse(Feature4 < 0.5, "A",
                                ifelse(Feature4 < 1, "B",
                                ifelse(Feature4 < 1.5, "C", "D")))]

  return(Validation)
}

# Store data sets
data1 <- dataGen(Correlation = 0.95)

# Generate EvalPlots (calibration)
p1 <- RemixAutoML::ParDepCalPlots(data = data1,
  PredictionColName = "predict",
  TargetColName = "target",
  IndepVar = "Feature1",
  GraphType = "calibration",
  PercentileBucket = 0.05,
  Function = function(x) mean(x,
                                na.rm = TRUE),
  FactLevels = 10)
p1 <- p1 + ggplot2::ggtitle("Partial Dependence Calibration p1") +
  RemixAutoML::ChartTheme(Size = 10)

p2 <- RemixAutoML::ParDepCalPlots(data = data1,
  PredictionColName = "predict",
  TargetColName = "target",
  IndepVar = "Feature2",
  GraphType = "calibration",
  PercentileBucket = 0.05,
  Function = function(x) mean(x,
```

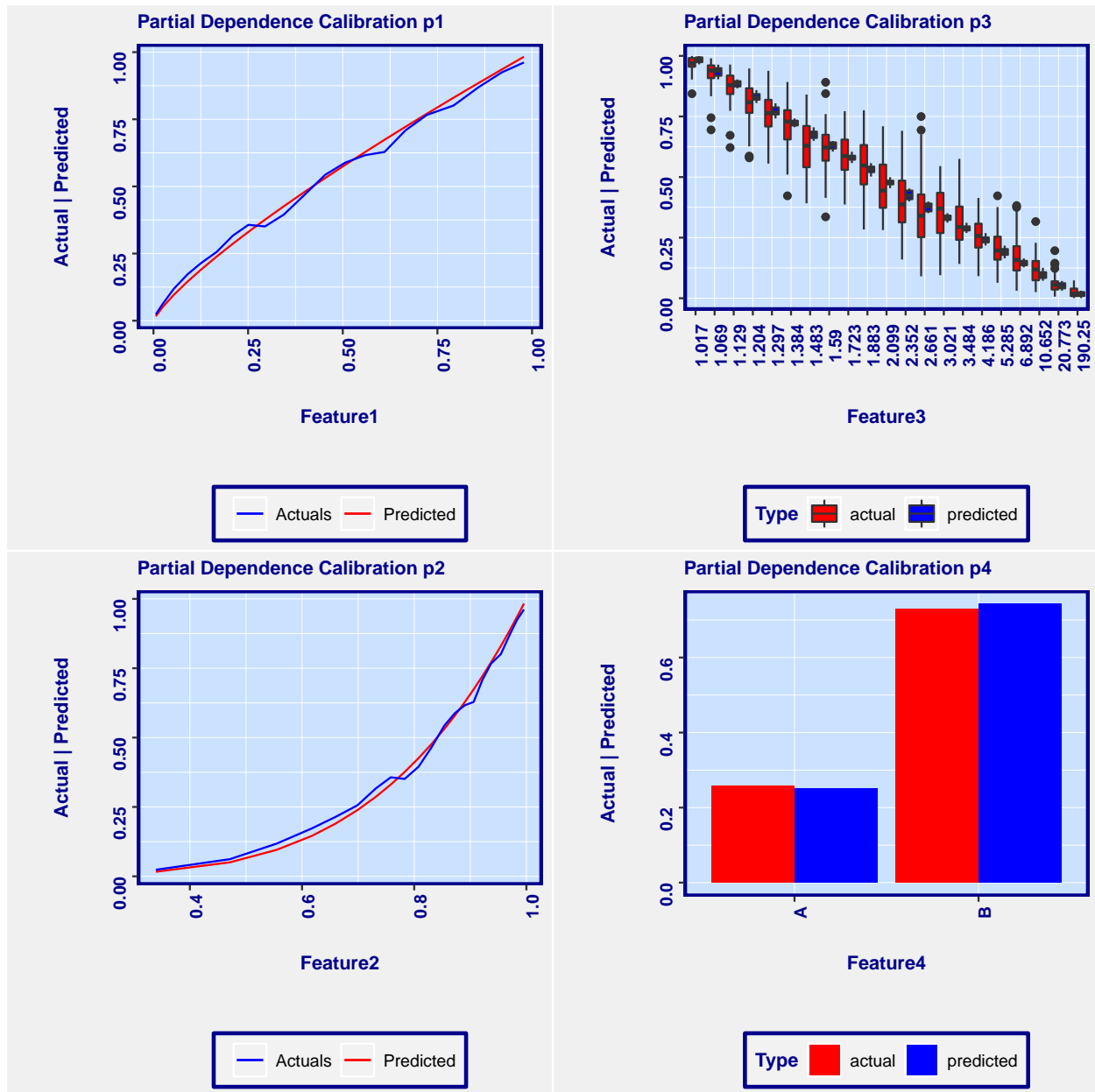
```

na.rm = TRUE),
    FactLevels = 10)
p2 <- p2 + ggplot2::ggtitle("Partial Dependence Calibration p2") +
  RemixAutoML::ChartTheme(Size = 10)

p3 <- RemixAutoML::ParDepCalPlots(data = data1,
  PredictionColName = "predict",
  TargetColName = "target",
  IndepVar = "Feature3",
  GraphType = "boxplot",
  PercentileBucket = 0.05,
  Function = function(x) mean(x,
    na.rm = TRUE),
  FactLevels = 10)
p3 <- p3 + ggplot2::ggtitle("Partial Dependence Calibration p3") +
  RemixAutoML::ChartTheme(Size = 10)

p4 <- RemixAutoML::ParDepCalPlots(data = data1,
  PredictionColName = "predict",
  TargetColName = "target",
  IndepVar = "Feature4",
  GraphType = "calibration",
  PercentileBucket = 0.05,
  Function = function(x) mean(x,
    na.rm = TRUE),
  FactLevels = 10)
p4 <- p4 + ggplot2::ggtitle("Partial Dependence Calibration p4") +
  RemixAutoML::ChartTheme(Size = 10)
RemixAutoML::multiplot(plotlist = list(p1,p2,p3,p4), cols = 2)

```



Functions include:

- threshOptim()
- RedYellowGreen()

RedYellowGreen Output (simulated data)

Find more demos at <https://www.remixcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
library(RemixAutoML)
Correl <- 0.70
aa <- data.table::data.table(target = runif(1000))
```

```

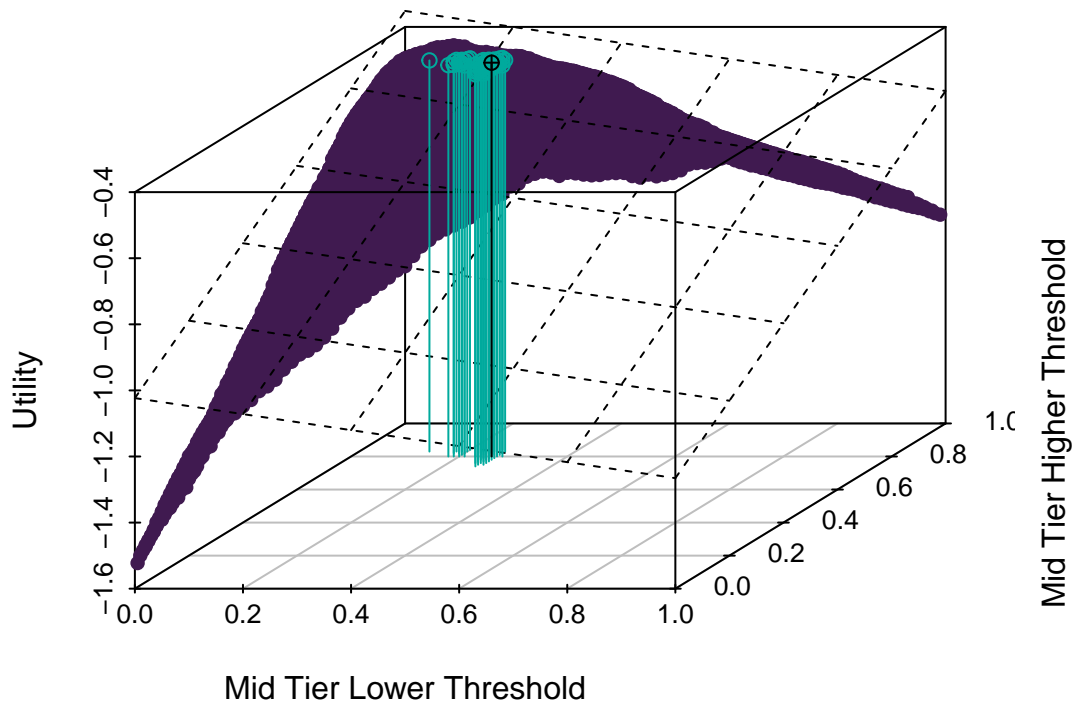
aa[, x1 := qnorm(target)]
aa[, x2 := runif(1000)]
aa[, predict := pnorm(Correl * x1 +
                      sqrt(1 - Correl ^2) *
                      qnorm(x2))]
aa[, target := as.numeric(ifelse(target < 0.5, 0, 1))]
data <- RemixAutoML::RedYellowGreen(
  aa,
  PredictColNumber = 4,
  ActualColNumber = 1,
  TruePositiveCost = 0,
  TrueNegativeCost = 0,
  FalsePositiveCost = -3,
  FalseNegativeCost = -2,
  MidTierCost = -0.5,
  Cores = 1
)
#> Loading required namespace: doParallel

knitr::kable(data[order(-Utility)][1:10])

```

TPP	TNP	FPP	FNP	MTDN	MTC	Threshold	MTLT	MHTT	Utility
0	0	-3	-2	TRUE	-0.5	0.80	0.26	0.80	-0.4082353
0	0	-3	-2	TRUE	-0.5	0.82	0.26	0.82	-0.4096833
0	0	-3	-2	TRUE	-0.5	0.83	0.26	0.83	-0.4105501
0	0	-3	-2	TRUE	-0.5	0.77	0.26	0.77	-0.4110934
0	0	-3	-2	TRUE	-0.5	0.79	0.26	0.79	-0.4112244
0	0	-3	-2	TRUE	-0.5	0.76	0.26	0.76	-0.4113786
0	0	-3	-2	TRUE	-0.5	0.78	0.26	0.78	-0.4113878
0	0	-3	-2	TRUE	-0.5	0.81	0.26	0.81	-0.4117595
0	0	-3	-2	TRUE	-0.5	0.80	0.25	0.80	-0.4120618
0	0	-3	-2	TRUE	-0.5	0.75	0.26	0.75	-0.4121054

Utility Maximizer – Main Threshold at 0.8



Lower Thresh = 0.26 and Upper Thresh = 0.8

Automated Feature Engineering Functions

This suite of functions are what will take your models to the next level. The core functions are the generalized distributed lag and rolling statistics functions. I have four of them.

Functions include:

- `GDL_Feature_Engineering()`
- `DT_GDL_Feature_Engineering()`
- `FAST_GDL_Feature_Engineering()`
- `Scoring_GDL_Feature_Engineering()`

The first three are used for building out lags and rolling statistics from target variables (numeric type; including classification models (0|1) and multinomial models with a little bit of work) and numeric features over your entire data set (no aggregation is done) with the option for creating the rolling statistics on the main variable or the lag1 version of the main variable. You can also compute time between records (by group) and add their lags and rolling statistics as well (really useful for transactional data). They can be generated using a single grouping variable (for multiple grouping variables you can concatenate them) and you can feed in a list of grouping variables to generate them by. The first function (**GDL__**) has the largest variety of rolling statics options but runs the slowest. The second function (**DT_GDL__**) runs the fastest but only generates moving averages. The third function (**FAST_GDL__**) is used for cases where you don't

need to generate the features across the entire data set. Suppose you have a limited number of target variable instance but a rich history of data. You can use the FAST_GDL_ version to create lags and rolling statistics for N number of records previous to each target instance (i.e. not the entire historical data set). The fourth function (**Scoring_GDL_**) is for use in a production setting where you need to generate single instances of the feature set quickly. You basically feed in the same arguments as you used for the other versions and out the other end is the same set of features, identically named.

DT_GDL_Feature_Engineering and Scoring_GDL_Feature_Engineering Demo (simulated data)

Find more demos at <https://www.remixcourses.com/course?courseid=intro-to-remixautoml-in-r>

```
library(RemixAutoML)

# Build data for feature engineering for modeling
N <- 25116
ModelData <-
  data.table::data.table(GroupVariable = sample(
    x = c(letters,LETTERS,paste0(letters, letters),
          paste0(LETTERS, LETTERS),
          paste0(letters, LETTERS),
          paste0(LETTERS, letters))),
    DateTime = base::as.Date(Sys.time()),
    Target = stats::filter(rnorm(N,
                                mean = 50,
                                sd = 20),
                           filter = rep(1, 10),
                           circular = TRUE))
ModelData[, temp := seq(1:161), by = "GroupVariable"] [
  , DateTime := DateTime - temp] [
  , temp := NULL]
ModelData <- ModelData[order(DateTime)]
ModelData <- RemixAutoML::DT_GDL_Feature_Engineering(
  ModelData,
  lags      = c(seq(1, 5, 1)),
  periods   = c(3, 5, 10, 15, 20, 25),
  statsNames = c("MA"),
  targets    = c("Target"),
  groupingVars = "GroupVariable",
  sortDateName = "DateTime",
  timeDiffTarget = c("Time_Gap"),
  timeAgg      = c("days"),
  WindowingLag = 1,
  Type         = "Lag",
  Timer        = FALSE,
  SkipCols     = FALSE,
  SimpleImpute = TRUE)
#> [1] 22

# Build data for feature engineering for scoring
N <- 25116
ScoringData <-
  data.table::data.table(GroupVariable = sample(
```

```

x = c(letters,LETTERS,paste0(letters, letters),
      paste0(LETTERS, LETTERS),
      paste0(letters, LETTERS),
      paste0(LETTERS, letters)),
DateTime = base::as.Date(Sys.time()),
Target = stats::filter(rnorm(N,
                             mean = 50,
                             sd = 20),
                       filter = rep(1, 10),
                       circular = TRUE))
ScoringData[, temp := seq(1:161),
             by = "GroupVariable"][, DateTime := DateTime - temp]
ScoringData <- ScoringData[order(DateTime)]

# Use WindowingLag = 1 to build moving averages off of the lag1 Target Variable to eliminate forward leakage
ScoringData <- RemixerAutoML::Scoring_GDL_Feature_Engineering(
  ScoringData,
  lags          = c(seq(1, 5, 1)),
  periods       = c(3, 5, 10, 15, 20, 25),
  statsFUNs     = c(function(x) mean(x, na.rm = TRUE)),
  statsNames    = c("MA"),
  targets       = c("Target"),
  groupingVars  = c("GroupVariable"),
  sortDateName  = c("DateTime"),
  timeDiffTarget = c("Time_Gap"),
  timeAgg       = "days",
  WindowingLag  = 1,
  Type          = "Lag",
  Timer         = FALSE,
  SkipCols      = FALSE,
  SimpleImpute  = TRUE,
  AscRowByGroup = "temp",
  RecordsKeep   = 1
)

# View some of new features
knitr::kable(ModelData[order(GroupVariable,-DateTime)][1:10,c(3,4,14)])

```

Target	GroupVariable_LAG_1_Target	GroupVariableMA_3_GroupVariable_LAG_1_Target
453.0552	461.7106	495.1132
461.7106	419.2195	497.8995
419.2195	604.4096	534.8744
604.4096	470.0695	474.6354
470.0695	530.1442	479.4543
530.1442	423.6925	483.9627
423.6925	484.5260	476.0348
484.5260	543.6696	496.5063
543.6696	399.9087	498.3027
399.9087	545.9405	559.3506

```
# Ensure names equal
knitr::kable(
  data.table::as.data.table(
    cbind(ModelData_Names = sort(names(ModelData)),
          ScoringData_Names = sort(names(ScoringData[, temp := NULL])))))
```

ModelData_Names	ScoringData_Names
DateTime	DateTime
GroupVariable	GroupVariable
GroupVariable_LAG_1_Target	GroupVariable_LAG_1_Target
GroupVariable_LAG_2_Target	GroupVariable_LAG_2_Target
GroupVariable_LAG_3_Target	GroupVariable_LAG_3_Target
GroupVariable_LAG_4_Target	GroupVariable_LAG_4_Target
GroupVariable_LAG_5_Target	GroupVariable_LAG_5_Target
GroupVariableMA_10_GroupVariable_LAG_1_Target	GroupVariableMA_10_GroupVariable_LAG_1_Target
GroupVariableMA_10_GroupVariableTime_Gap1	GroupVariableMA_10_GroupVariableTime_Gap1
GroupVariableMA_15_GroupVariable_LAG_1_Target	GroupVariableMA_15_GroupVariable_LAG_1_Target
GroupVariableMA_15_GroupVariableTime_Gap1	GroupVariableMA_15_GroupVariableTime_Gap1
GroupVariableMA_20_GroupVariable_LAG_1_Target	GroupVariableMA_20_GroupVariable_LAG_1_Target
GroupVariableMA_20_GroupVariableTime_Gap1	GroupVariableMA_20_GroupVariableTime_Gap1
GroupVariableMA_25_GroupVariable_LAG_1_Target	GroupVariableMA_25_GroupVariable_LAG_1_Target
GroupVariableMA_25_GroupVariableTime_Gap1	GroupVariableMA_25_GroupVariableTime_Gap1
GroupVariableMA_3_GroupVariable_LAG_1_Target	GroupVariableMA_3_GroupVariable_LAG_1_Target
GroupVariableMA_3_GroupVariableTime_Gap1	GroupVariableMA_3_GroupVariableTime_Gap1
GroupVariableMA_5_GroupVariable_LAG_1_Target	GroupVariableMA_5_GroupVariable_LAG_1_Target
GroupVariableMA_5_GroupVariableTime_Gap1	GroupVariableMA_5_GroupVariableTime_Gap1
GroupVariableTime_Gap1	GroupVariableTime_Gap1
GroupVariableTime_Gap2	GroupVariableTime_Gap2
GroupVariableTime_Gap3	GroupVariableTime_Gap3
GroupVariableTime_Gap4	GroupVariableTime_Gap4
GroupVariableTime_Gap5	GroupVariableTime_Gap5
Target	Target

Automated Feature Engineering Functions (Continued)

Functions include:

- `AutoWord2VecModeler()`
- `ModelDataPrep()`
- `DummifyDT()`

The **AutoWord2VecModeler** function converts your text features into numerical vector representations. You supply the function with your data set and all the text column names you want converted, and out the other end you have a data set with all the features merged on. The models can be saved to file and metadata saves their paths for scoring purposes in a production setting. The models built are based on H2O's word2vec algorithm and has done an excellent job at extracting high quality information out of those text columns. The **ModelDataPrep** function is used to prepare your data for modeling with the **AutoH2OModeler** function. It will convert character columns to factors, replace inf values to NA, and impute missing values (both numeric and factor based on supplied values). The **DummifyDT** function will turn your character (or factor) columns into dummy variable columns. You can specify one-hot encoding or not in which you will get N+1 columns for one-hot or N columns otherwise.

Miscellaneous Functions

Functions include:

- `AutoWordFreq()`
- `AutoH2OTextPrepScoring()`
- `RecomDataCreate()`
- `ChartTheme()`
- `RemixTheme()`
- `multiplot()`
- `PrintObjectsSize()`
- `percRank()`

The **AutoWordFreq** function will go through a process of cleaning your text column, doing some other text operations, and output a table with word frequencies and a word cloud plot. The **AutoH2OTextPrepScoring** will automatically prepare your text data for scoring. This function is run internally in the **AutoH2OScoring** function but you can utilize it outside for other purposes. The **ChartTheme** and **RemixTheme** functions will turn your ggplots into nicely formatted and colored charts, worthy of presentation. The **multiplot** function are for those who have had a terrible time plotting multiple graphs onto a single image. The **PrintObjectsSize** function is more of a debugging function for inspecting the size of variables in your environment (useful in looping functions). The **percRank** is simply a function to compute the percentile rank of every value in a column of data. **RecomDataCreate** will turn your transactional data set into a binary ratings matrix fast.