

RemixAutoML Library Introduction

Adrian Antico

2019-05-23

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](#) 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

There are several core packages RemixAutoML relies on which include. From a data management standpoint, I utilize data.table exclusively for data wrangling of all internal functions due to its ability to handle big data with a minimal memory footprint, along with the speed at which their functions process data. For the machine learning functions, I utilize H2O, Catboost, XGBoost, and forecast due to their quality results and ability to handle big data (some can run on GPU as well). I use these functions routinely for machine learning projects and they continue to outperform every other method 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. There are several unique functions to this package that help to optimize the machine learning process, such as feature engineering functions, model evaluation functions, model interpretation functions, and model output optimization functions.

Machine learning methodology

The package is designed to give your models 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, categorical, and text features. The functions in this package will help you squeeze as much information out

of your data set as possible. Then you simply supply those data sets into the automated machine learning functions for state of the art output with minimal effort on your part. With the extra time saved you can try out significantly more experiments to ensure you are putting the best models and frameworks possible for your machine learning use-cases.

Handling numerical features

Numerical features can hide relationships in a variety of ways and it's our goal as modeling professionals (or researchers) to capture as many of those relationships as we can (or enough to provide a sufficient return on investment). We have linear and nonlinear relationships, linear and nonlinear interactions, along with threshold effects. Those are what I call column-based model effects. Nowadays, you don't really have to bother with any of those when using tree-based ensembles (boosting and bagging) as those relationships can be captured without the need to manually create features to account for their peculiarities. However, there is still a significant amount of potential information to be extracted from your data if you don't account for time-based effects. In the business world, many applications require modeling data that is collected across time (think transactional data). With time series data (or panel data), where data is collected (or manipulated to be) across equally spaced time periods (hourly, daily, etc.), you tend to look at lags and moving averages of your target variable for predicting the present and future events. Well, why aren't we using the same type of features in our transaction data? That's where we offer our '*GDL*' suite of functions. They are designed to create lags and rolling statistics off of numeric target variables and numeric independent variables, by groups. They can also generate lags and rolling statistics off of the time between events, by groups. This is what I call row-based model effects. Any nonlinearities and interactions will be captured by the tree-based ensembles.

Taking this a step further, there is a concept out there called target encoding. What's done, essentially, is that the mean of the target is used as a replacement for factor levels, thus converting your factor variables into numeric variables. I generally have two problems with this. For one thing, there is inherent forward leak because you are using all values across time to predict values that occur historically. Second, it fails to account for all the other information in your data to be used in the transformation process. If you used the '*GDL*' suite of functions, you can do your target encoding without any forward leak (they are forced to prevent this), but you will also be able to utilize various window sizes (opposed to full history) to capture recent trends or cyclical effects, and do the same thing based on your independent variables that have effects distributed across time.

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 automated modeling functions. Turning your categorical features into dummy variables is problematic for high cardinality factor variables. There are other approaches and the automated modeling functions 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. Of course, if you're working with time-based data, I would recommend using the '*GDL*' suite of functions to convert your factor variables to numeric to prevent forward leak and to not miss out on all that sweet information to be gained!

Handling text data

We are living in a world now where text data is becoming more readily available and your machine learning models aren't suited to handle them without first managing them. With this package, you can simply run the **AutoWord2VecModeler()** function, which builds skip-gram models, for all your text columns, thus replacing your text data with numerical vectors suitable for modeling. The function will save the models and recreate them on the fly with the **AutoH2OScoring()** function in your production setting.

Machine learning algorithms

The machine learning algorithms available have been demonstrated to provide optimal performance on a wide-range of business use-cases over the years. They are all intended to remove the coding aspects behind tuning, evaluation, and interpretation, along with consistent output for comparison amongst them.

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](#)

How to install catboost

```
devtools::install_github('catboost/catboost', subdir = 'catboost/R-package')
```

[Review catboost](#)

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 Feature Engineering
- Automated Cost Sensitive Optimization
- A Few Miscellaneous Functions

Automated Supervised Learning Functions

Functions include:

- `AutoCatBoostClassifier()`
- `AutoXGBoostClassifier()`
- `AutoH2oGBMClassifier()`
- `AutoH2oDRFClassifier()`

- `AutoCatBoostMultiClass()`
- `AutoXGBoostMultiClass()`
- `AutoH2oGBMMultiClass()`
- `AutoH2oDRFMultiClass()`
- `AutoCatBoostRegression()`
- `AutoH2oGBMRegression()`
- `AutoH2oDRFRegression()`
- `AutoXGBoostRegression()`
- `AutoH2OModeler()`
- `AutoH2OScoring()`
- `AutoTS()`
- `AutoNLS()`
- `AutoRecomDataCreate()`
- `AutoRecommender()`
- `AutoRecommenderScoring()`

`AutoCatBoostRegression()`

`AutoXGBoostRegression()`

`AutoH2oGBMRegression()`

`AutoH2oDRFRegression()`

The `Auto__Regression()` set are automated regression modeling function that runs a variety of steps. First, the functions will run a random grid tune over N number of models and find which model is the best on holdout test data (a default model is always included in that set). Once the model is identified and built, several other outputs are generated: validation data with predictions, evaluation calibration plot, evaluation calibration boxplot, evaluation model metrics, variable importance, partial dependence calibration plots, partial dependence calibration box plots, grid metrics, grid arguments, and column names used in model fitting. You can fit standard expected value regression (all of them) along with quantile regression (catboost and h2o gbm).

`AutoCatBoostClassifier()`

`AutoXGBoostClassifier()`

`AutoH2oGBMClassifier()`

`AutoH2oDRFClassifier()`

The `Auto__Classifier()` set are automated binary classification modeling functions that runs a variety of steps. First, the function will run a random grid tune over N number of models and find which model is the best on holdout test data (a default model is always included in that set). Once the model is identified and built, several other outputs are generated: validation data with predictions, ROC plot, evaluation calibration plot, evaluation metrics, variable importance, partial dependence calibration plots, grid metrics, grid arguments, and column names used in model fitting.

`AutoCatBoostMultiClass()`

AutoXGBoostMultiClass()

AutoH2oGBMMultiClass()

AutoH2oDRFMultiClass()

The `Auto__MultiClass()` set are automated multiclass modeling functions that runs a variety of steps. First, the function will run a random grid tune over N number of models and find which model is the best on holdout test data (a default model is always included in that set). Once the best model is identified and built, several other outputs are generated: validation data with predictions, evaluation metrics, variable importance, grid metrics, grid arguments, and column names used in model fitting.

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, and AutoML (for Windows users XGBoost and LightGBM are not available).

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, and 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).

AutoRecomDataCreate()

This function will automatically turn your transactional data into a binary ratings matrix that you supply to the AutoRecommender() function for model building and the AutoRecommenderScoring() function for scoring.

AutoRecommender()

This function builds out several variations of collaborative filtering models on a binary ratings matrix. To automatically build the binary ratings matrix, see **AutoRecomDataCreate**. 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. This process is parallelized for fast scoring.

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, and Polynomial regression or monotonic regression.

Example of AutoXGBoostRegression()

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

```
library(RemixAutoML)

# Create Simulated Data to Demonstrate Modeling
Correl <- 0.85
N <- 10000
data <- data.table::data.table(Target = runif(N))
data[, x1 := qnorm(Target)]
data[, x2 := runif(N)]
data[, Independent_Variable1 := log(pnorm(Correl * x1 +
                                          sqrt(1-Correl^2) * qnorm(x2)))]
data[, Independent_Variable2 := (pnorm(Correl * x1 +
                                       sqrt(1-Correl^2) * qnorm(x2)))]
data[, Independent_Variable3 := exp(pnorm(Correl * x1 +
                                          sqrt(1-Correl^2) * qnorm(x2)))]
data[, Independent_Variable4 := exp(exp(pnorm(Correl * x1 +
                                              sqrt(1-Correl^2) * qnorm(x2))))]
data[, Independent_Variable5 := sqrt(pnorm(Correl * x1 +
                                           sqrt(1-Correl^2) * qnorm(x2)))]
data[, Independent_Variable6 := (pnorm(Correl * x1 +
                                       sqrt(1-Correl^2) * qnorm(x2)))^0.10]
data[, Independent_Variable7 := (pnorm(Correl * x1 +
                                       sqrt(1-Correl^2) * qnorm(x2)))^0.25]
```

```

data[, Independent_Variable8 := (pnorm(Correl * x1 +
                                     sqrt(1-Correl^2) * qnorm(x2)))^0.75]
data[, Independent_Variable9 := (pnorm(Correl * x1 +
                                     sqrt(1-Correl^2) * qnorm(x2)))^2]
data[, Independent_Variable10 := (pnorm(Correl * x1 +
                                     sqrt(1-Correl^2) * qnorm(x2)))^4]
data[, Independent_Variable11 := as.factor(
  ifelse(Independent_Variable2 < 0.20, "A",
    ifelse(Independent_Variable2 < 0.40, "B",
      ifelse(Independent_Variable2 < 0.6, "C",
        ifelse(Independent_Variable2 < 0.8, "D", "E")))))]
data[, ':= ' (x1 = NULL, x2 = NULL)]

# Create Train and Test Data
DataSets <- RemixAutoML::AutoDataPartition(data = data,
                                           NumDataSets = 3,
                                           Ratios = c(0.70,0.20,0.10),
                                           PartitionType = "random",
                                           StratifyColumnNames = NULL,
                                           TimeColumnName = NULL)

# Store data sets
train <- DataSets$TrainData
valid <- DataSets$ValidationData
test <- DataSets$TestData

# Build Models
TestModel <- RemixAutoML::AutoXGBoostRegression(data = train,
                                                  ValData = valid,
                                                  TestData = test,
                                                  TargetColumnName = 1,
                                                  FeatureColNames = 2:12,
                                                  CatFeatures = 12,
                                                  IDcols = NULL,
                                                  eval_metric = "RMSE",
                                                  Trees = 50,
                                                  GridTune = TRUE,
                                                  grid_eval_metric = "mae",
                                                  MaxModelsInGrid = 10,
                                                  NThreads = 8,
                                                  TreeMethod = "hist",
                                                  model_path = NULL,
                                                  ModelID = "FirstModel",
                                                  NumOfParDepPlots = 3,
                                                  Verbose = 0,
                                                  ReturnModelObjects = TRUE,
                                                  SaveModelObjects = FALSE)

#> [1] 1
#> [1] 2
#> [1] 3
#> [1] 4
#> [1] 5
#> [1] 6

```

```

#> [1] 7
#> [1] 8
#> [1] 9
#> [1] 10
#> [1] 11

# Evaluation metrics
TestModel$EvaluationMetrics

# Variable Importance
TestModel$VariableImportance
#>           Feature      Gain      Cover Frequency
#> 1: Independent_Variable1 0.7465 0.8053    0.8388
#> 2: Independent_Variable2 0.2410 0.0828    0.0812
#> 3: Independent_Variable11_E 0.0072 0.0114    0.0037
#> 4: Independent_Variable11_A 0.0037 0.0076    0.0081
#> 5: Independent_Variable3 0.0005 0.0320    0.0273
#> 6: Independent_Variable11_B 0.0004 0.0183    0.0124
#> 7: Independent_Variable11_C 0.0003 0.0239    0.0136
#> 8: Independent_Variable11_D 0.0003 0.0187    0.0149

# Grid List of Arguments Tested
TestModel$GridList
#>      eta max_depth min_child_weight subsample colsample_bytree
#> 1: 0.35          8                3          0.9              0.8
#> 2: 0.35          8                1          0.8              0.8
#> 3: 0.30          8                1          1.0              0.9
#> 4: 0.35         10                2          0.8              0.9
#> 5: 0.35          6                2          1.0              1.0
#> 6: 0.30          6                3          1.0              1.0
#> 7: 0.35          8                1          1.0              0.8
#> 8: 0.30          6                2          1.0              0.9
#> 9: 0.25         10                1          1.0              1.0
#> 10: 0.25         8                3          0.8              0.9
#> 11: 0.35         6                1          0.8              0.8

# Metrics from Grid Tuning
TestModel$GridMetrics
#>      ParamRow EvalStat
#> 1:          1  0.1219
#> 2:          2  0.1209
#> 3:          3  0.1216
#> 4:          4  0.1208
#> 5:          5  0.1216
#> 6:          6  0.1219
#> 7:          7  0.1216
#> 8:          8  0.1219
#> 9:          9  0.1216
#> 10:         10  0.1226
#> 11:         11  0.1213

# Evaluation Calibration Plot
TestModel$EvaluationPlot

```



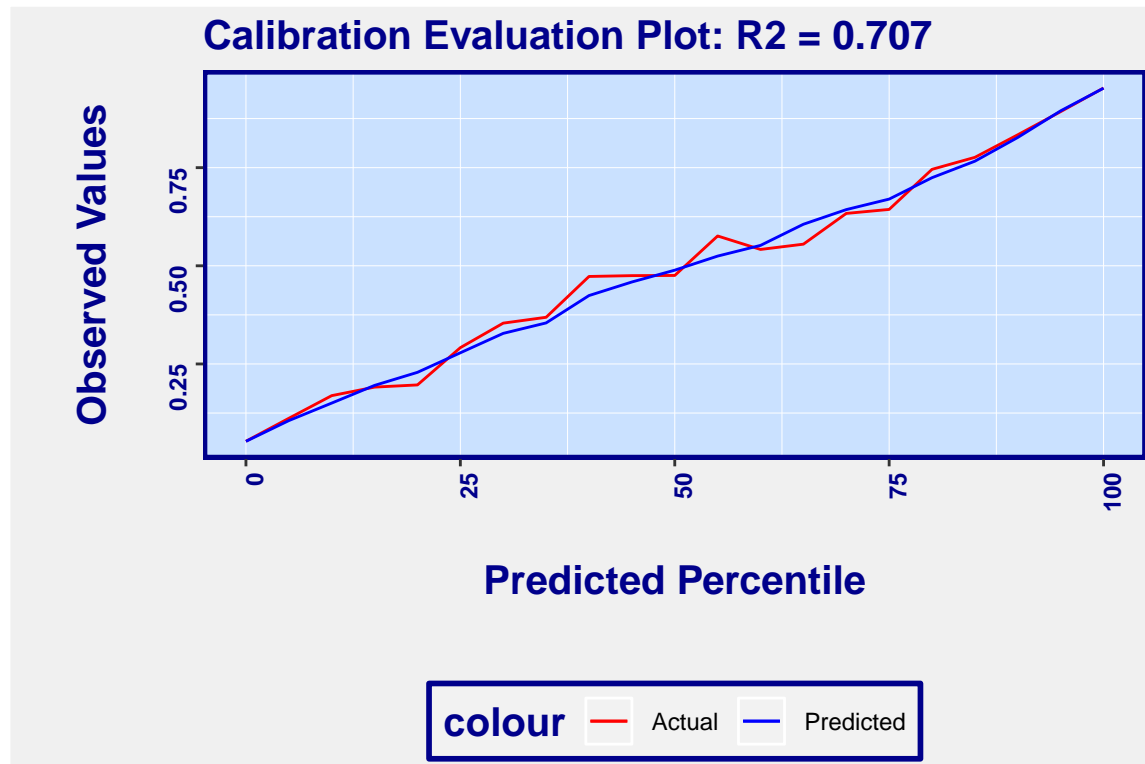
```

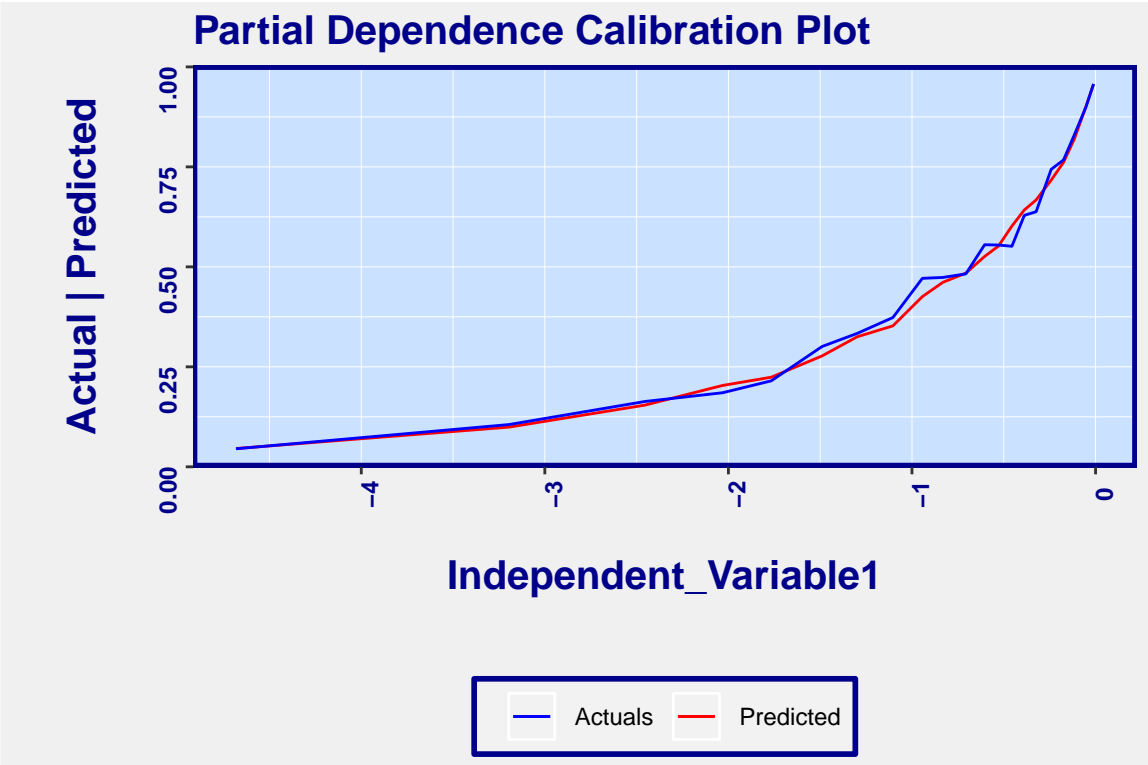
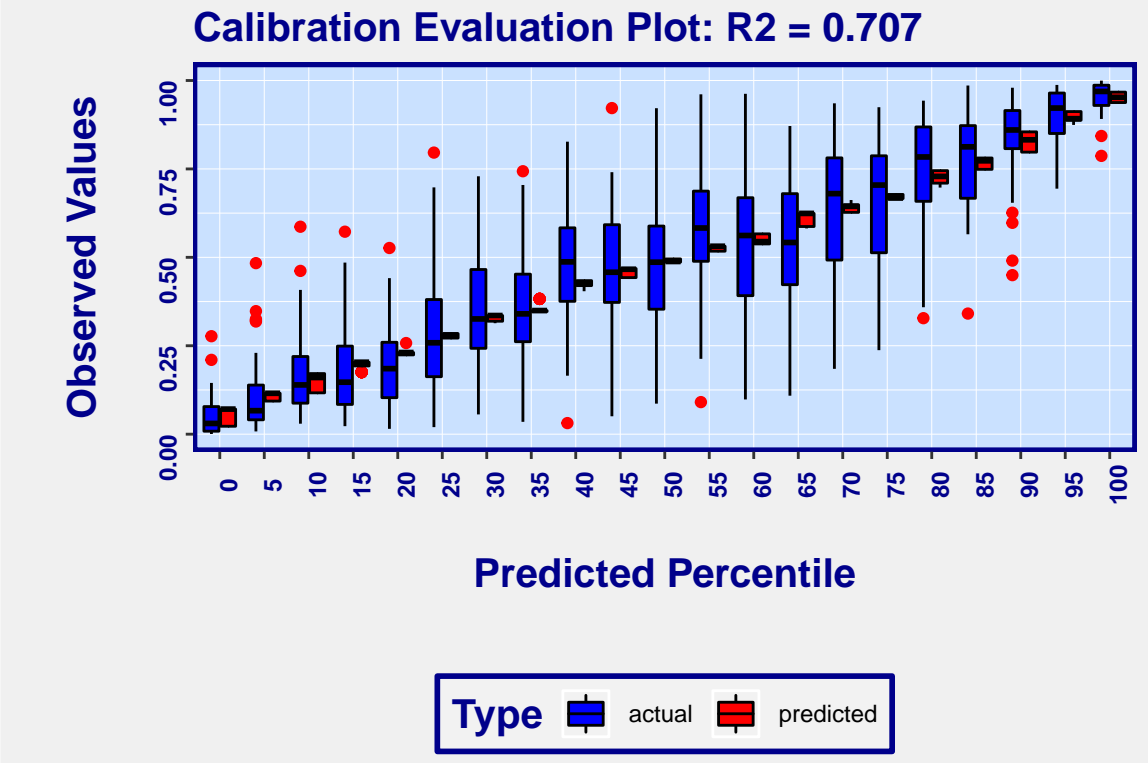
# Evaluation Calibration BoxPlot
TestModel$EvaluationBoxPlot

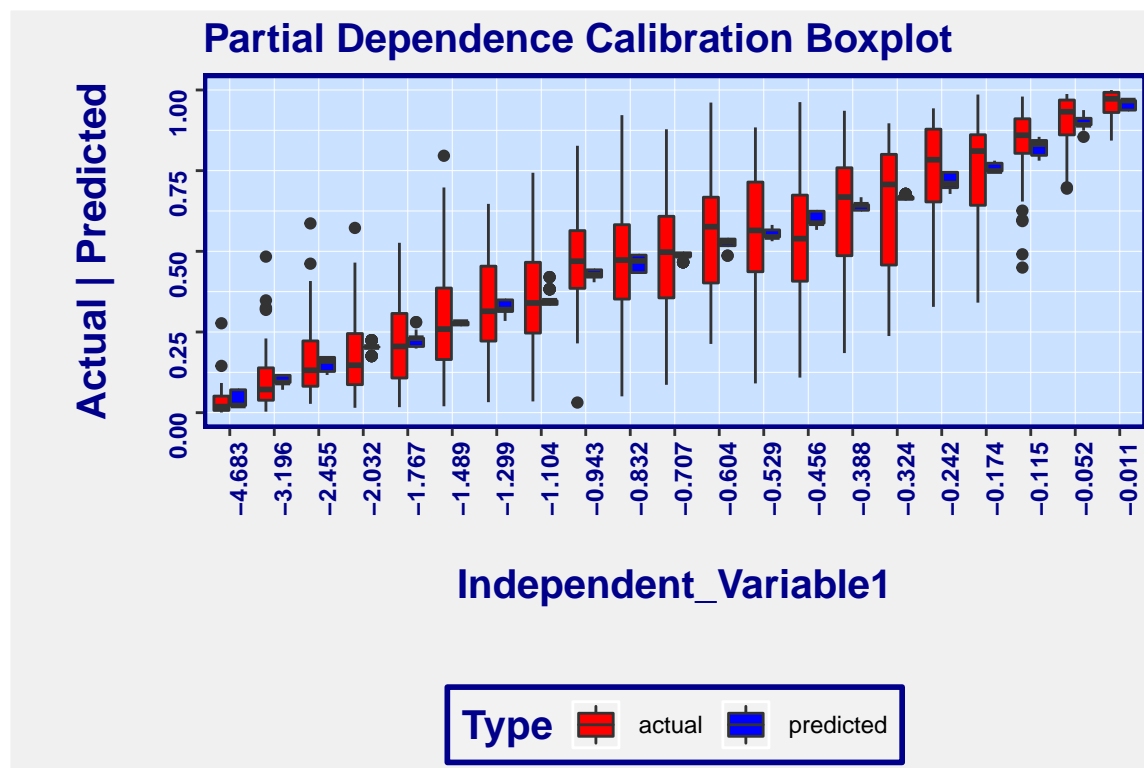
# Partial Dependence Calibration Plot
TestModel$PartialDependencePlots$Independent_Variable1

# Partial Dependence Calibration BoxPlot
TestModel$PartialDependenceBoxPlots$Independent_Variable1

```







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
)

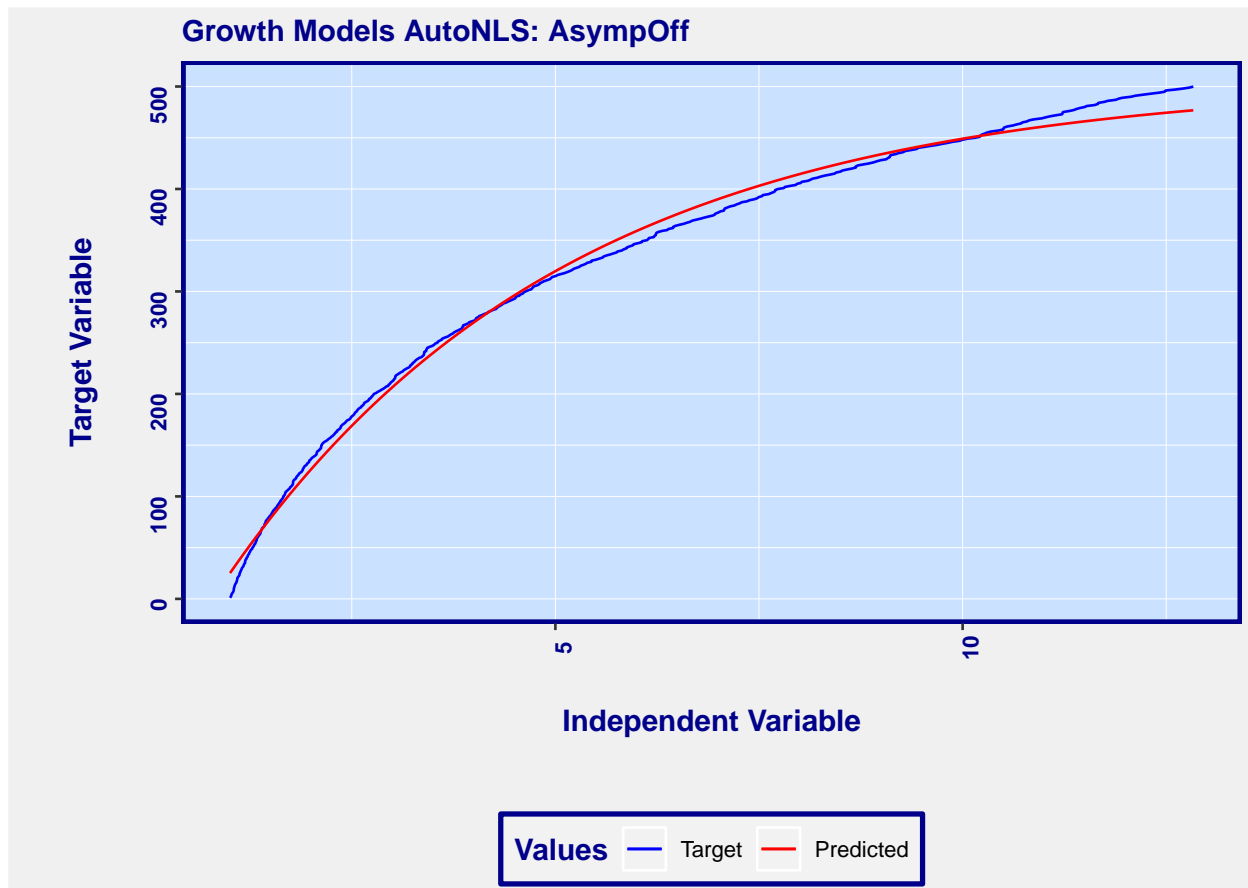
# Join predictions to source data
data2 <- merge(
  data1,
  data11$PredictionData,
  by = "Variable",
  all = FALSE
)

# 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")) +
  RemixerAutoML::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 ~ SSasymptOff(Variable, Asym, lrc, c0)
#>
#> Parameters:
#>      Estimate Std. Error t value Pr(>|t|)
#> Asym 505.63358    2.17609  232.36  <2e-16 ***
#> lrc  -1.43642     0.01115 -128.80  <2e-16 ***
#> c0    0.79079     0.01028   76.95  <2e-16 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 9.586 on 497 degrees of freedom
#>

```

```
#> Number of iterations to convergence: 0
#> Achieved convergence tolerance: 6.581e-07
data11$EvaluationMetrics
#>      ModelName MeanAbsError
#> 1:    AsympOff      8.200265
#> 2:      Asymp      8.200265
#> 3:    Gompertz     17.025430
#> 4:   AsympOrig     17.257745
#> 5: Michal_Menton     19.041491
#> 6:    Logistic     22.638652
#> 7:       Poly      82.344628
```



Example of AutoTS()

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

```
library(RemixAutoML)

# From WalMart data in vignette folder: Store 1 and Department 1
dates <- c("2010-02-05", "2010-02-12", "2010-02-19", "2010-02-26", "2010-03-05", "2010-03-12",
           "2010-03-19", "2010-03-26", "2010-04-02", "2010-04-09", "2010-04-16", "2010-04-23",
           "2010-04-30", "2010-05-07", "2010-05-14", "2010-05-21", "2010-05-28", "2010-06-04",
           "2010-06-11", "2010-06-18", "2010-06-25", "2010-07-02", "2010-07-09", "2010-07-16",
```

```

"2010-07-23", "2010-07-30", "2010-08-06", "2010-08-13", "2010-08-20", "2010-08-27",
"2010-09-03", "2010-09-10", "2010-09-17", "2010-09-24", "2010-10-01", "2010-10-08",
"2010-10-15", "2010-10-22", "2010-10-29", "2010-11-05", "2010-11-12", "2010-11-19",
"2010-11-26", "2010-12-03", "2010-12-10", "2010-12-17", "2010-12-24", "2010-12-31",
"2011-01-07", "2011-01-14", "2011-01-21", "2011-01-28", "2011-02-04", "2011-02-11",
"2011-02-18", "2011-02-25", "2011-03-04", "2011-03-11", "2011-03-18", "2011-03-25",
"2011-04-01", "2011-04-08", "2011-04-15", "2011-04-22", "2011-04-29", "2011-05-06",
"2011-05-13", "2011-05-20", "2011-05-27", "2011-06-03", "2011-06-10", "2011-06-17",
"2011-06-24", "2011-07-01", "2011-07-08", "2011-07-15", "2011-07-22", "2011-07-29",
"2011-08-05", "2011-08-12", "2011-08-19", "2011-08-26", "2011-09-02", "2011-09-09",
"2011-09-16", "2011-09-23", "2011-09-30", "2011-10-07", "2011-10-14", "2011-10-21",
"2011-10-28", "2011-11-04", "2011-11-11", "2011-11-18", "2011-11-25", "2011-12-02",
"2011-12-09", "2011-12-16", "2011-12-23", "2011-12-30", "2012-01-06", "2012-01-13",
"2012-01-20", "2012-01-27", "2012-02-03", "2012-02-10", "2012-02-17", "2012-02-24",
"2012-03-02", "2012-03-09", "2012-03-16", "2012-03-23", "2012-03-30", "2012-04-06",
"2012-04-13", "2012-04-20", "2012-04-27", "2012-05-04", "2012-05-11", "2012-05-18",
"2012-05-25", "2012-06-01", "2012-06-08", "2012-06-15", "2012-06-22", "2012-06-29",
"2012-07-06", "2012-07-13", "2012-07-20", "2012-07-27", "2012-08-03", "2012-08-10",
"2012-08-17", "2012-08-24", "2012-08-31", "2012-09-07", "2012-09-14", "2012-09-21",
"2012-09-28", "2012-10-05", "2012-10-12", "2012-10-19", "2012-10-26")

weekly_sales <- c(24924.50, 46039.49, 41595.55, 19403.54, 21827.90, 21043.39, 22136.64, 26229.21,
57258.43, 42960.91, 17596.96, 16145.35, 16555.11, 17413.94, 18926.74, 14773.04,
15580.43, 17558.09, 16637.62, 16216.27, 16328.72, 16333.14, 17688.76, 17150.84,
15360.45, 15381.82, 17508.41, 15536.40, 15740.13, 15793.87, 16241.78, 18194.74,
19354.23, 18122.52, 20094.19, 23388.03, 26978.34, 25543.04, 38640.93, 34238.88,
19549.39, 19552.84, 18820.29, 22517.56, 31497.65, 44912.86, 55931.23, 19124.58,
15984.24, 17359.70, 17341.47, 18461.18, 21665.76, 37887.17, 46845.87, 19363.83,
20327.61, 21280.40, 20334.23, 20881.10, 20398.09, 23873.79, 28762.37, 50510.31,
41512.39, 20138.19, 17235.15, 15136.78, 15741.60, 16434.15, 15883.52, 14978.09,
15682.81, 15363.50, 16148.87, 15654.85, 15766.60, 15922.41, 15295.55, 14539.79,
14689.24, 14537.37, 15277.27, 17746.68, 18535.48, 17859.30, 18337.68, 20797.58,
23077.55, 23351.80, 31579.90, 39886.06, 18689.54, 19050.66, 20911.25, 25293.49,
33305.92, 45773.03, 46788.75, 23350.88, 16567.69, 16894.40, 18365.10, 18378.16,
23510.49, 36988.49, 54060.10, 20124.22, 20113.03, 21140.07, 22366.88, 22107.70,
28952.86, 57592.12, 34684.21, 16976.19, 16347.60, 17147.44, 18164.20, 18517.79,
16963.55, 16065.49, 17666.00, 17558.82, 16633.41, 15722.82, 17823.37, 16566.18,
16348.06, 15731.18, 16628.31, 16119.92, 17330.70, 16286.40, 16680.24, 18322.37,
19616.22, 19251.50, 18947.81, 21904.47, 22764.01, 24185.27, 27390.81)

# Convert to data.table
data <- data.table::data.table(Date = dates, Weekly_Sales = weekly_sales)

# Names of data columns
names(data)
#> [1] "Date"          "Weekly_Sales"

# Build models and generate forecasts
output <- RemixAutoML::AutoTS(data,
                               TargetName = "Weekly_Sales",
                               DateName   = "Date",
                               FCPeriods  = 120,
                               HoldOutPeriods = 12,

```

```

TimeUnit      = "week",
Lags           = 5,
SLags          = 1,
NumCores       = 4,
SkipModels     = NULL,
StepWise       = TRUE,
TSClean        = TRUE,
PrintUpdates   = TRUE)

#> DSHW FITTING
#> ARFIMA FITTING
#> ARIMA FITTING
#> ETS FITTING
#> TBATS FITTING
#> TSLM FITTING
#> 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"
#> FIND WINNER
#> GENERATE FORECASTS
#> FULL DATA ARIMA FITTING

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

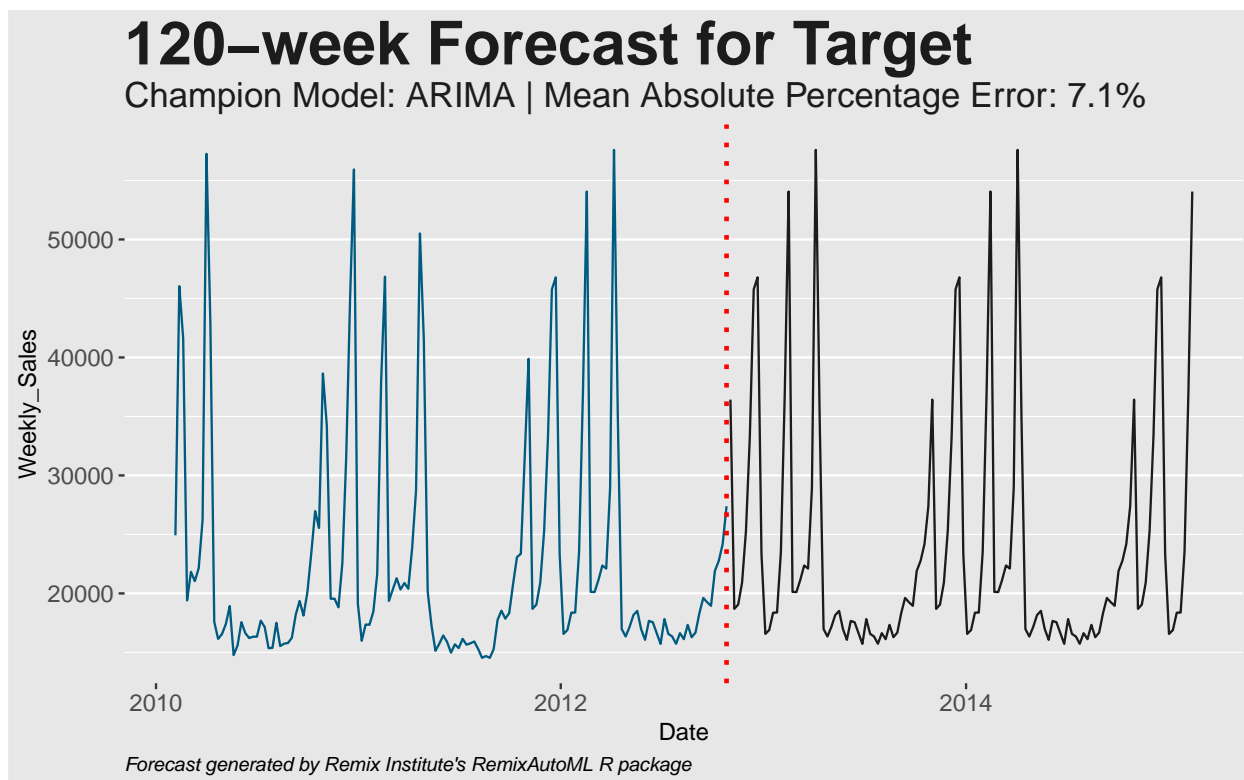
```

ModelName	MeanResid	MeanPercError	MAPE	ID
ARIMA	609.55	0.04692	0.07130	1
TBATS_TSC	1574.01	0.09270	0.11628	2
NN	2135.10	0.11308	0.12424	3
TBATS	688.19	0.03370	0.12458	4
TBATS_ModelFreq	688.19	0.03370	0.12458	5
TBATS_ModelFreqTSC	688.19	0.03370	0.12458	6
ARFIMA_TSC	431.74	0.02090	0.12866	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

ModelName	MeanResid	MeanPercError	MAPE	ID
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
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	3831.90	0.23878	0.24034	21
NN_ModelFreqTSC	3853.51	0.24064	0.24237	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

# Time Series Plot
output$TimeSeriesPlot
#> Warning: Removed 120 rows containing missing values (geom_path).
#> Warning: Removed 143 rows containing missing values (geom_path).
```

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")

```

```

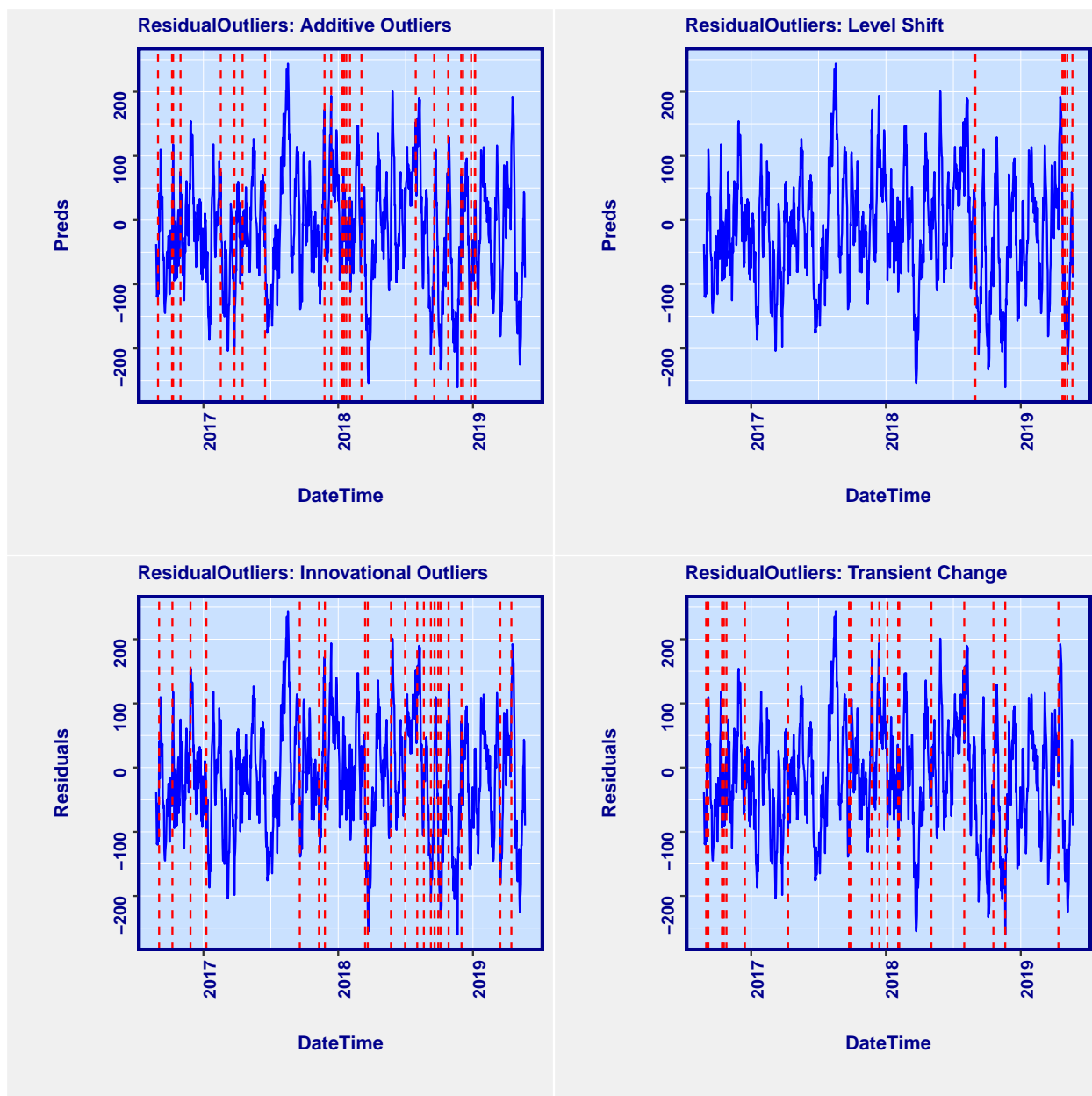
p22 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Residuals),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "IO", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "IO"][["DateTime"]]),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

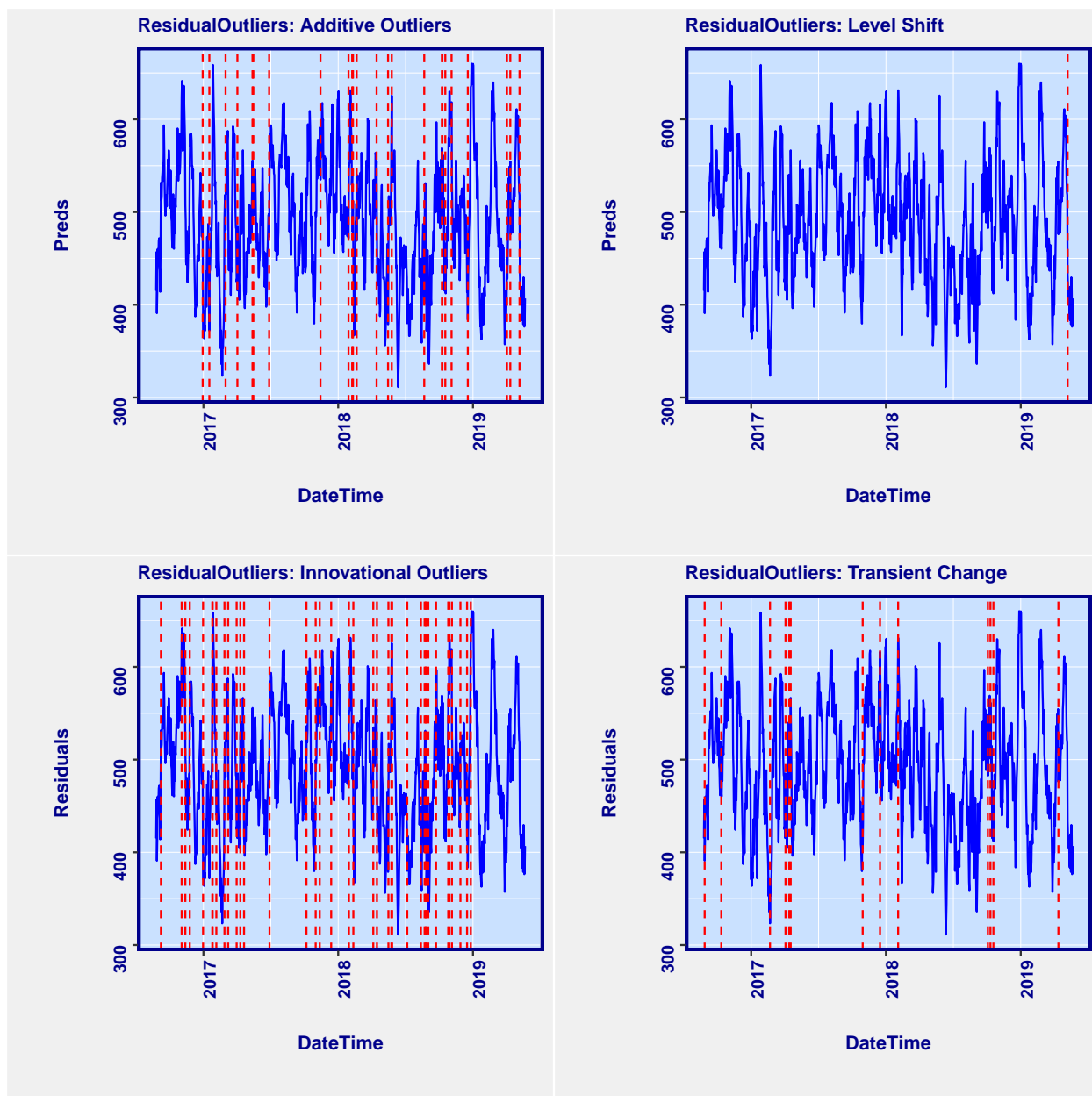
p33 <- ggplot2::ggplot(data, ggplot2::aes(x = DateTime)) +
  ggplot2::geom_line(ggplot2::aes(y = Preds),
    color = "blue") +
  RemixAutoML::ChartTheme(Size = 10) +
  ggplot2::geom_vline(data = outliers[type == "LS", "DateTime"],
    ggplot2::aes(xintercept = outliers[
      type == "LS"][["DateTime"]]),
    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()
- ParDepCalPlots()
- threshOptim()
- RedYellowGreen()

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)

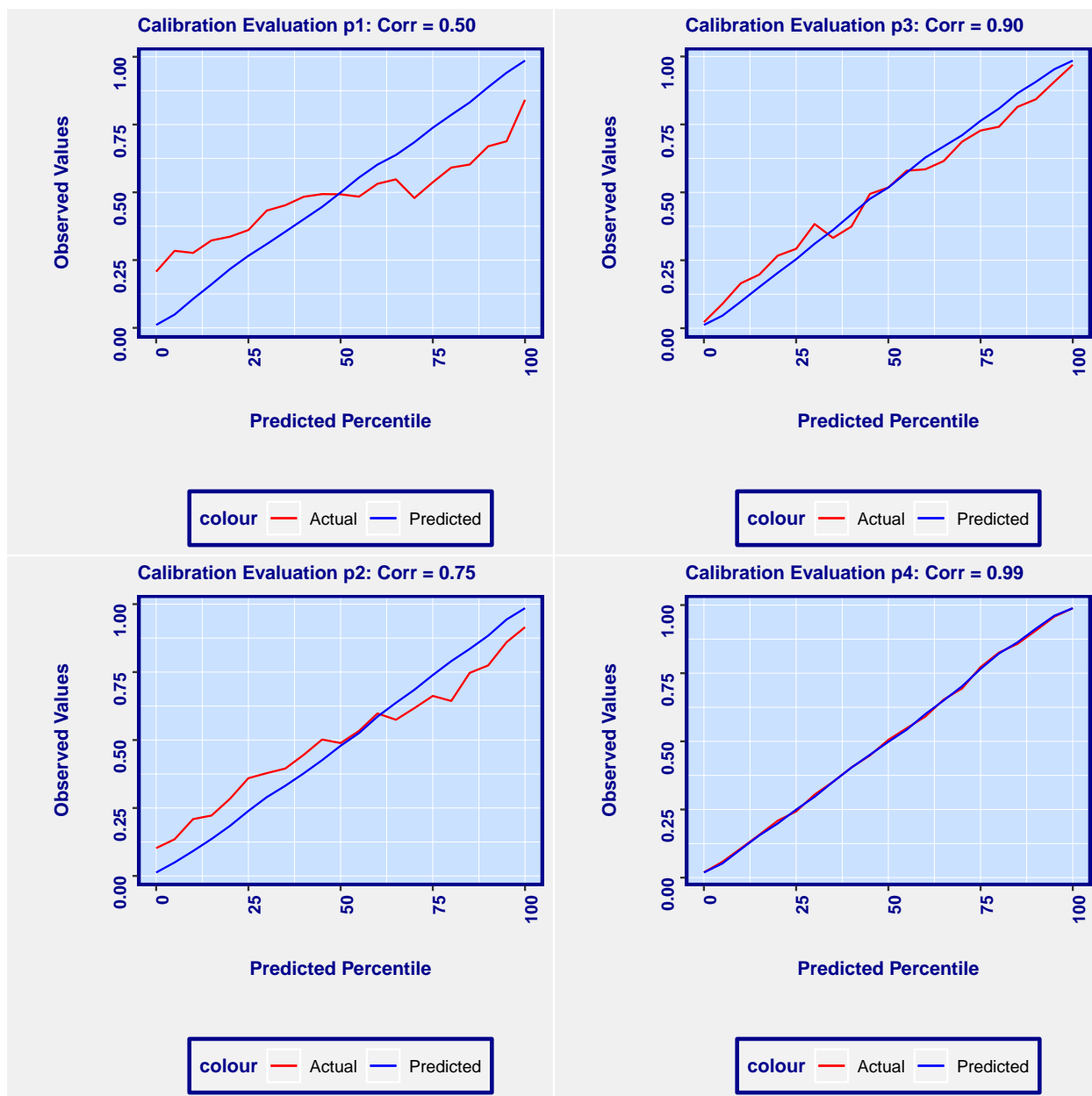
p2 <- RemixAutoML::EvalPlot(data = data2,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p2 <- p2 + ggplot2::ggtitle("Calibration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

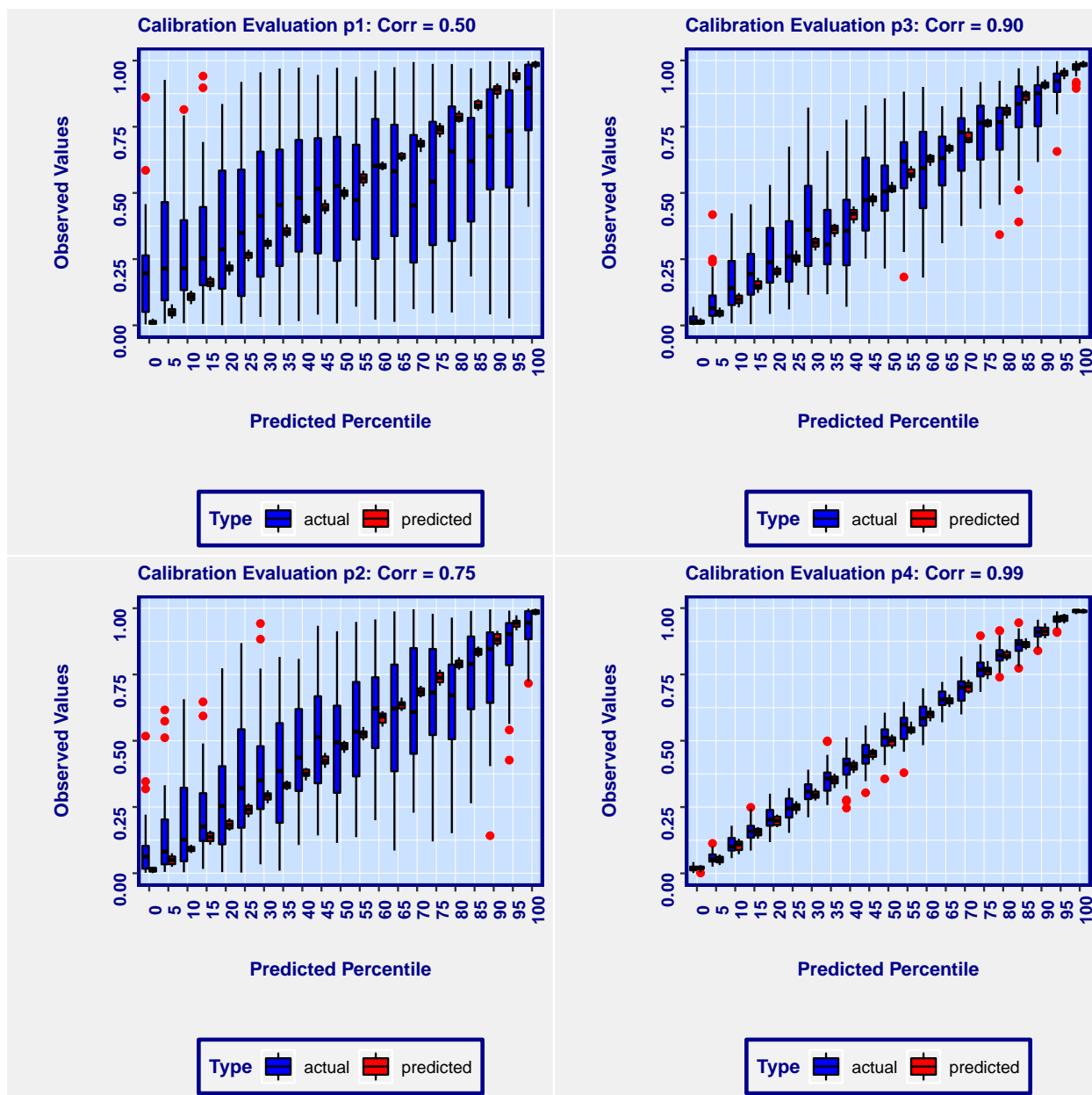
p3 <- RemixAutoML::EvalPlot(data = data3,
                             PredictionColName = "predict",
                             TargetColName = "target",
                             GraphType = "boxplot",
                             PercentileBucket = 0.05)
p3 <- p3 + ggplot2::ggtitle("Calibration 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.

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()

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))]
  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)

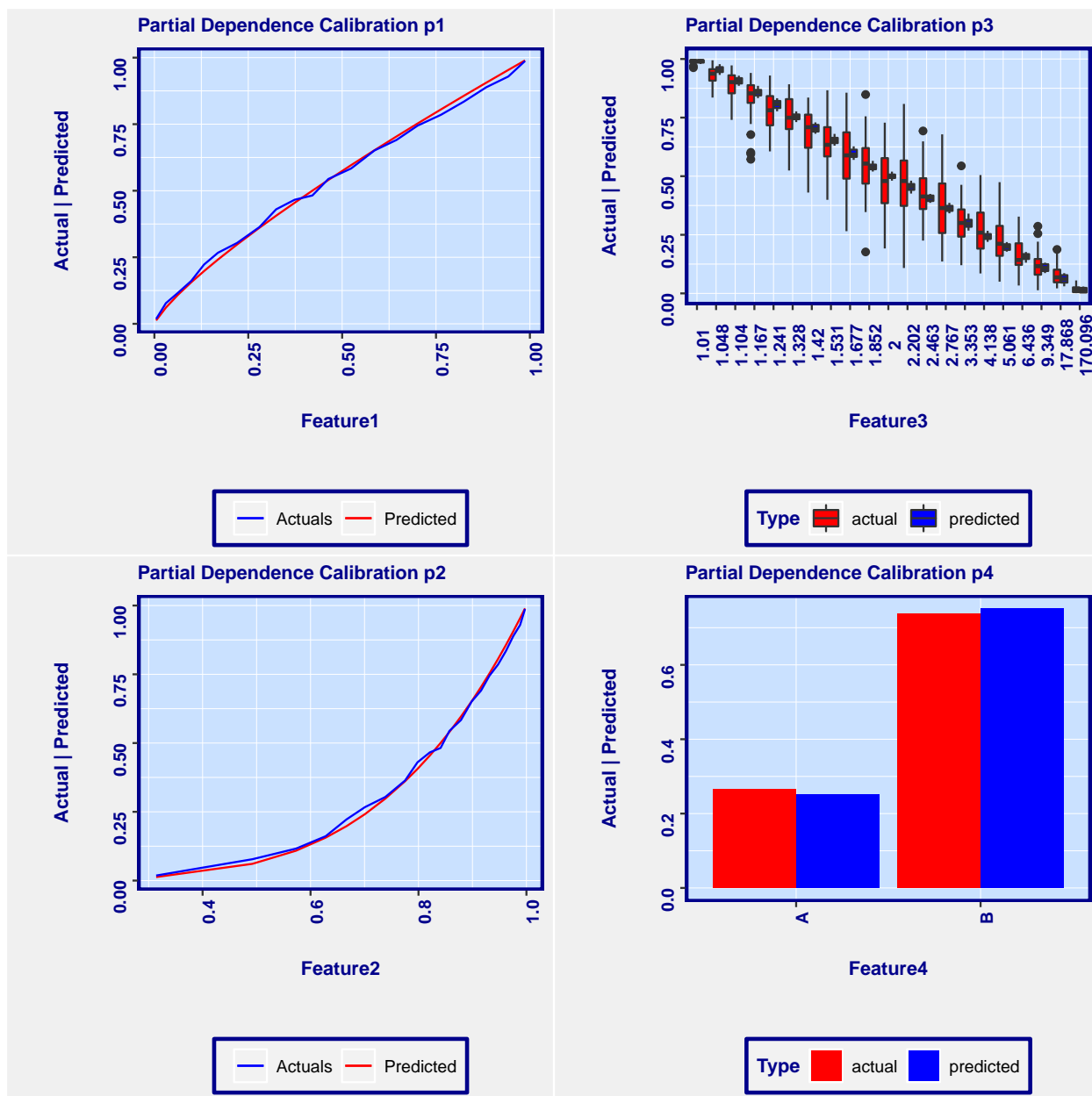
p2 <- RemixAutoML::ParDepCalPlots(data = data1,
  PredictionColName = "predict",
  TargetColName = "target",
  IndepVar = "Feature2",
  GraphType = "calibration",
  PercentileBucket = 0.05,
  Function = function(x) mean(x,
                                na.rm = TRUE),
  FactLevels = 10)
p2 <- p2 + ggplot2::ggtitle("Partial 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)

```



Demo of RedYellowGreen()

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

```
library(RemixAutoML)
Correl <- 0.70
data <- data.table::data.table(target = runif(1000))
data[, x1 := qnorm(target)]
data[, x2 := runif(1000)]
data[, predict := pnorm(Correl * x1 +
  sqrt(1 - Correl ^2) *
  qnorm(x2))]
data[, target := as.numeric(ifelse(target < 0.5, 0, 1))]
```

```

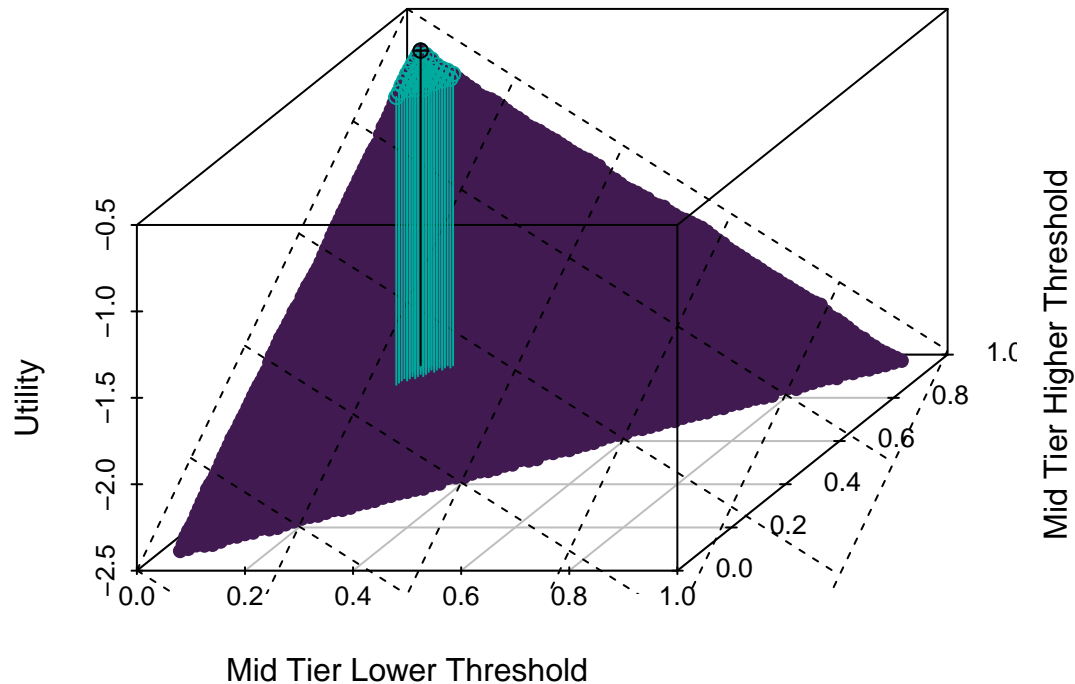
data <- RemixAutoML::RedYellowGreen(
  data,
  PredictColNumber = 4,
  ActualColNumber = 1,
  TruePositiveCost = 0,
  TrueNegativeCost = 0,
  FalsePositiveCost = -3,
  FalseNegativeCost = -2,
  MidTierCost = -0.5,
  Cores = 1,
  Boundaries = c(0.05,0.95)
)
#> Loading required namespace: doParallel

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

```

TPP	TNP	FPP	FNP	MTDN	MTC	Threshold	MTLT	MTHT	Utility
0	0	-3	-2	TRUE	-0.5	0.95	0.05	0.95	-0.678830
0	0	-3	-2	TRUE	-0.5	0.94	0.05	0.94	-0.692739
0	0	-3	-2	TRUE	-0.5	0.95	0.06	0.95	-0.704661
0	0	-3	-2	TRUE	-0.5	0.93	0.05	0.93	-0.708635
0	0	-3	-2	TRUE	-0.5	0.94	0.06	0.94	-0.718570
0	0	-3	-2	TRUE	-0.5	0.92	0.05	0.92	-0.722544
0	0	-3	-2	TRUE	-0.5	0.95	0.07	0.95	-0.730492
0	0	-3	-2	TRUE	-0.5	0.93	0.06	0.93	-0.734466
0	0	-3	-2	TRUE	-0.5	0.91	0.05	0.91	-0.738440
0	0	-3	-2	TRUE	-0.5	0.94	0.07	0.94	-0.744401

Utility Maximizer – Main Threshold at 0.95



Lower Thresh = 0.05 and Upper Thresh = 0.95

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

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.remixxcourses.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 <- RemixAutoML::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
505.5307	390.4048	436.3974
390.4048	449.7426	485.8290
449.7426	469.0448	528.2404
469.0448	538.6996	586.0295
538.6996	576.9768	514.7926
576.9768	642.4121	486.2020
642.4121	324.9888	440.0333
324.9888	491.2050	509.1019
491.2050	503.9060	506.9274

Target	GroupVariable_LAG_1_Target	GroupVariableMA_3_GroupVariable_LAG_1_Target
503.9060	532.1947	508.7755

```
# 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

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()`
- `ProblematicFeatures()`
- `ProblematicRecords()`
- `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 **ProblematicFeatures** identified problematic columns for machine learning. **ProblematicRecords** finds problematic rows in your data set that you should investigate further. 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. **AutoRecomDataCreate** will turn your transactional data set into a binary ratings matrix fast.

```
# Create Some Data
data <- data.table::data.table(
DESCR = c("Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru, Gru,
Urkle, Urkle, Urkle, Urkle, Urkle, Urkle, Urkle, Urkle, Gru, Gru, Gru,
bears, bears, bears, bears, bears, bears, bears, smug, smug, smug, smug,
smug, smug, smug, smug, smug, smug, smug, smug, smug, smug, smug,
eats, eats, eats, eats, eats, eats, beats, beats, beats, beats,
beats, beats, beats, beats, beats, beats, beats, science, science,
Dwigt, Dwigt, Dwigt, Dwigt, Dwigt, Dwigt, Dwigt, Dwigt, Dwigt, Dwigt,
Schrute, Schrute, Schrute, Schrute, Schrute, Schrute, Schrute,
James, James, James, James, James, James, James, James, James, James,
Halpert, Halpert, Halpert, Halpert, Halpert, Halpert, Halpert, Halpert"))

# Run function
data <- AutoWordFreq(data,
  TextColName = "DESCR",
  GroupColName = NULL,
  GroupLevel = NULL,
  RemoveEnglishStopwords = FALSE,
  Stemming = FALSE,
  StopWords = c("Bla"))

#>      word freq
#> 1:    gru   16
#> 2:    smug   15
#> 3:    beats  11
#> 4:    dwigt  10
```

```

#> 5: james 10
#> 6: halpert 8
#> 7: schrute 7
#> 8: urkle 7
#> 9: bears 6
#> 10: eats 6
#> NULL

# View word frequency table
print(data)
#>      word freq
#> 1: gru 16
#> 2: smug 15
#> 3: beats 11
#> 4: dwigt 10
#> 5: james 10
#> 6: halpert 8
#> 7: schrute 7
#> 8: urkle 7
#> 9: bears 6
#> 10: eats 6
#> 11: science 2

```

