

Brightbox: a short introduction

Benjamin Rollert and Ronny Li

2017-06-22

Brightbox contains functions that tackle the problem of inspecting internals for any blackbox supervised learner. The package is designed to work with the caret package as well as any model that is an ensemble of caret learners. As of writing, the approaches implemented in the package are **partial dependency plots** (`run_partial_dependency`) and **marginal variable importance** (`calculate_marginal_vimp`).

Installation

```
> devtools::install_github('breather/brightbox')
```

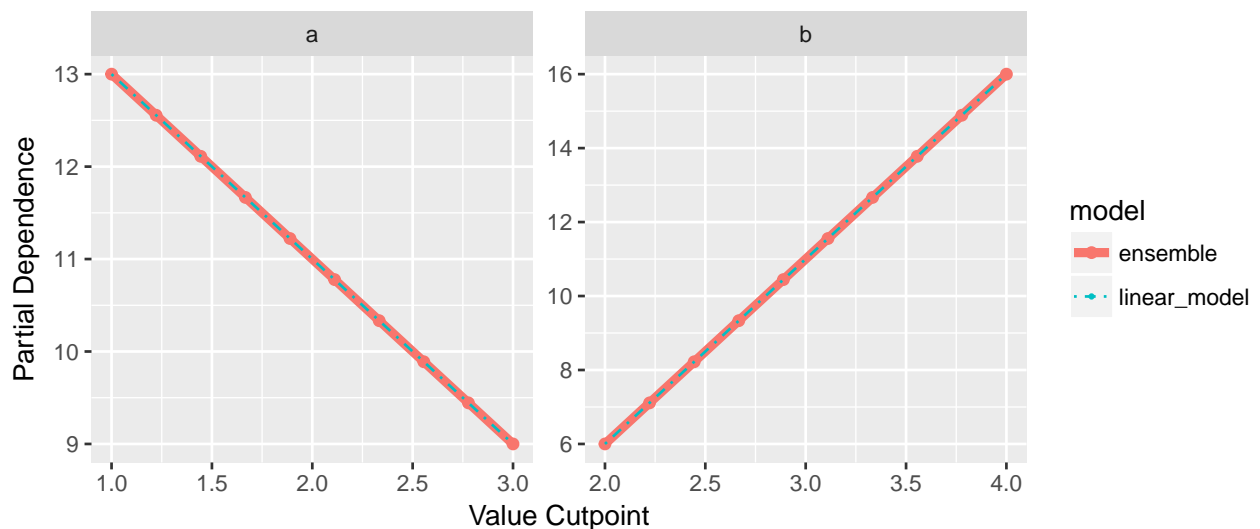
Partial Dependency Plots

Introduction

Partial dependency plots are a technique for visualizing the effect of a single feature on the response, marginalizing over the values of all other features. They are the visual equivalent of a coefficient in linear regression (and in fact, the partial dependency plot for a coefficient in a linear model will be a straight line).

```
# Example of partial dependency plots for a linear model
library(data.table)
dt <- data.table(a = 1:3, b = 2:4, c = c(8, 11, 14))
lm1 <- lm(c ~ a + b - 1, dt)
lm1$coefficients
#> a b
#> -2 5

# Note that the slopes of the plotted lines match the coefficients
library(brightbox)
pd <- run_partial_dependency(feature_dt = dt[, c("a", "b")],
                             model_list = list(linear_model = lm1))
```



The advantage of partial dependency plots shine when there is a non-linear relationship between the feature in question and the response. In such cases, a linear model will hide the true relationship due to its underlying assumptions. Ideally we would like to use a more flexible model with fewer assumptions but maintain interpretability.

The next section will work through an example dataset to demonstrate **brightbox**'s features with respect to partial dependency plots.

Walkthrough

```
# First, load the data
library(data.table)
library(mlbench)
data(BostonHousing, package = "mlbench")
head(BostonHousing)

#>      crim zn indus chas   nox    rm  age    dis rad tax ptratio    b
#> 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900    1 296    15.3 396.90
#> 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671    2 242    17.8 396.90
#> 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671    2 242    17.8 392.83
#> 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622    3 222    18.7 394.63
#> 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622    3 222    18.7 396.90
#> 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622    3 222    18.7 394.12

#>    lstat medv
#> 1   4.98 24.0
#> 2   9.14 21.6
#> 3   4.03 34.7
#> 4   2.94 33.4
#> 5   5.33 36.2
#> 6   5.21 28.7

# Split into features and response
boston_dt <- data.table(BostonHousing)
x <- boston_dt[, -"medv", with = FALSE]
y <- boston_dt$medv
# Prep the data (numeric columns are friendlier)
x[, chas := as.numeric(chas)]
```

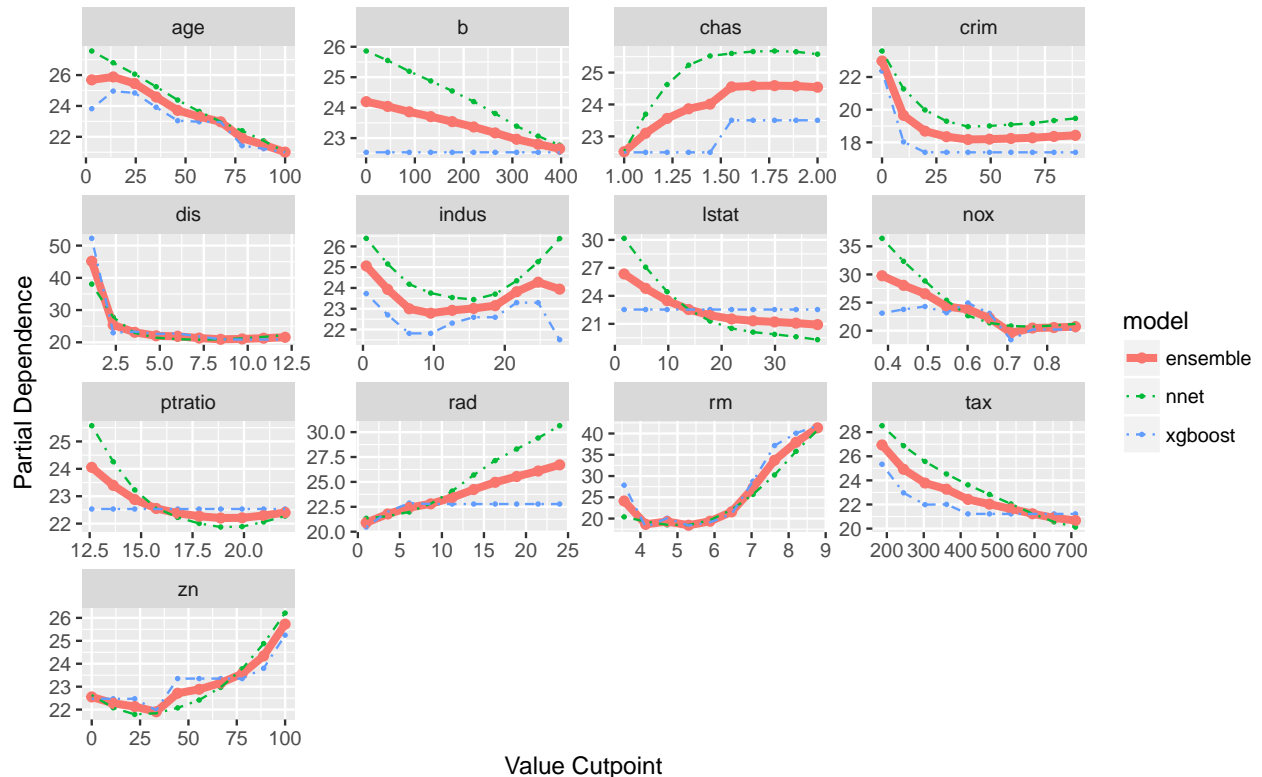
The dataset `BostonHousing` contains housing data for 506 census tracts of Boston from the 1970 census. `medv` is the response variable representing the median value of houses in each census tract (in USD 1000's). See `??BostonHousing` for additional details.

We will train two blackbox learners on this data and see if we can interpret the patterns each model learned.

```
library(caret)
# Train an xgboost model
xgb <- train(x = x, y = y,
            method = "xgbTree", metric = "RMSE")
# Train a single-layer neural net
nn <- train(x = x, y = y,
            method = "nnet", metric = "RMSE",
            preProc = c("center", "scale"),
            tuneGrid = expand.grid(size = c(10, 15, 20), decay = c(0, 5e-4, 0.05)),
            linout = TRUE,
            trace = FALSE)
```

From the two trained models we can use Brightbox to inspect their internals and trivially, the internals of an ensemble of these models. In this example, we construct an ensemble that is the median of the xgboost and neural net models.

```
library(brightbox)
# Generate partial dependency plots for each feature
pd <- run_partial_dependency(x, model_list = list(xgboost = xgb, nnet = nn),
                             ensemble_colname = "ensemble", ensemble_fcn = median)
```



In the returned plot we can see how `medv` changes with respect to each feature for every model. As a result, there are quite a few things we can learn by inspecting the plot.

1. We can easily see which features have a positive or negative relationship with the response and that the patterns are more complex than in the earlier example with linear regression (`rm` and `nox` for example).
2. There are some features and some ranges within features where the two models are roughly in agreement. From this we can conclude that the plotted signal in those ranges is rather strong. (`age`, `rm`, `dis`)
3. Some features have a wider range in the y-axis than others. This means that the feature had a larger impact on the response. In general, if a partial dependency plot is completely flat (zero slope) then the feature does not affect the model's predictions at all (an exception being when interaction effects are present).

As was hinted in points 2 and 3 above, partial dependency plots can help us determine feature importance by inspecting the partial dependence range and variance. `run_partial_dependency` returns a data.table with such calculations pre-computed.

```
# pd was previously saved from run_partial_dependency
head(pd)
#>   feature feature_val  model prediction  vimp
#> 1:    dis    1.129600 xgboost    52.25150 24.22594
#> 2:    dis    2.351478 xgboost    22.96461 24.22594
#> 3:    dis    3.573356 xgboost    23.30251 24.22594
#> 4:    dis    4.795233 xgboost    22.58686 24.22594
```

```
#> 5:    dis    6.017111 xgboost    22.73083 24.22594
#> 6:    dis    7.238989 xgboost    21.81226 24.22594
```

pd contains the data necessary to reconstruct the partial dependency plots (columns `feature`, `feature_val`, `model`, and `prediction`) as well as a variable importance column (`vimp`). Variable importance is calculated as the y-axis range for a given feature and model, with the ensemble model chosen by default (TODO: incorporate variance into the variable importance calculation).

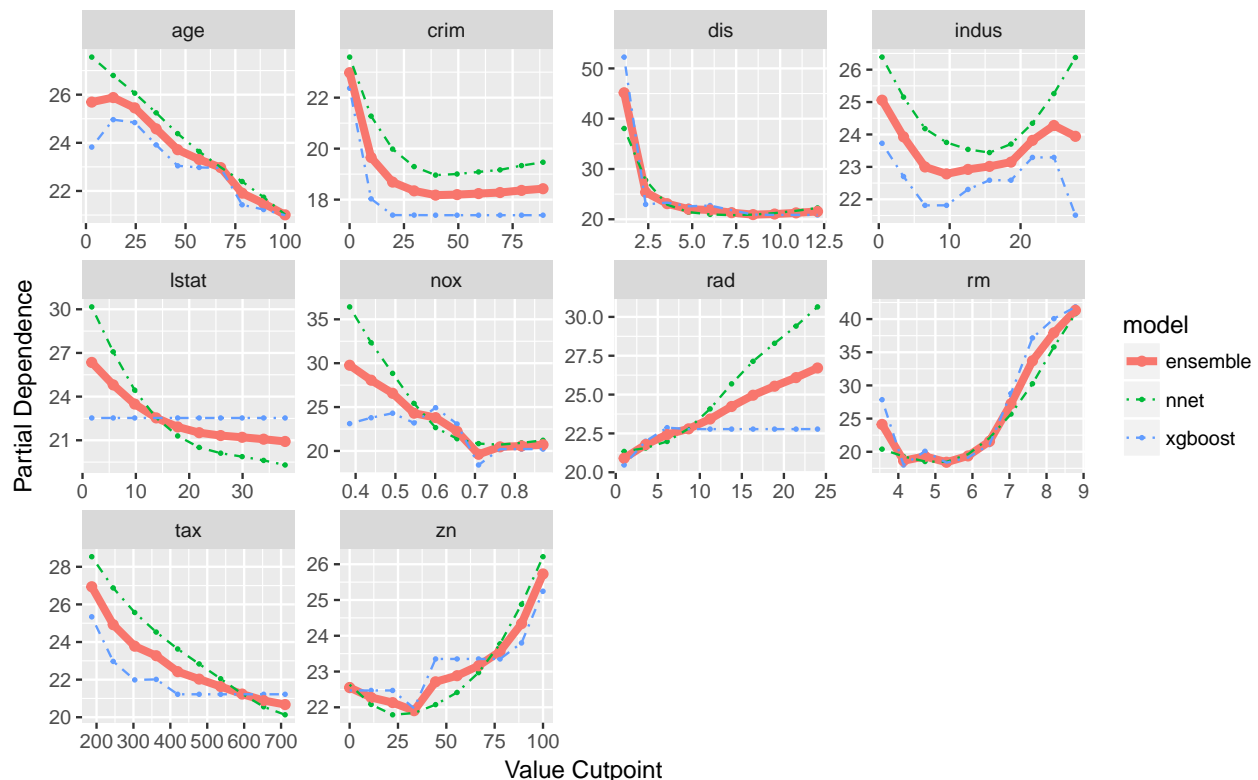
```
# Inspect the variable importance of each feature
```

```
unique(pd[, list(feature, vimp)])
```

```
#>    feature    vimp
#> 1:    dis 24.225936
#> 2:     rm 22.904621
#> 3:    nox 10.148432
#> 4:    tax  6.263277
#> 5:    rad  5.806844
#> 6:   lstat  5.430853
#> 7:    age  4.872998
#> 8:    crim  4.804854
#> 9:     zn  3.822874
#> 10:  indus  2.273781
#> 11:   chas  2.071185
#> 12: ptratio 1.849185
#> 13:     b  1.544938
```

From here we may be interested in plotting just the 10 most important features.

```
pd_top <- run_partial_dependency(x, model_list = list(xgboost = xgb, nnet = nn),
                                feature_cols = unique(pd[["feature"]])[1:10],
                                ensemble_colname = "ensemble", ensemble_fcn = median)
```



Partial dependency functions

While the function `run_partial_dependency` is the primary interface for partial dependency plots, it is composed of the following functions which can be useful if you want to avoid redundant calculations.

- `calculate_partial_dependency`
- `loop_calculate_partial_dependency`
- `facet_plot_fcn`
- `loop_plot_fcn`
- `plot_partial_dependency`
- `calculate_pd_vimp`

Marginal Variable Importance

(TODO: change this to a walkthrough instead)

```
calculate_marginal_vimp(x, y, method, loss_metric, resampling_indices, tuneGrid, trControl,
                        vars = names(x), allow_parallel = FALSE, ...)
```

Arguments `x`: data.table containing predictor variables

`y`: vector containing target variable

`method`: character string defining method to pass to caret

`loss_metric`: character. Loss metric to evaluate accuracy of model

`resampling_indices`: a list of integer vectors corresponding to the row indices used for each resampling iteration

`tuneGrid`: a data.frame containing hyperparameter values for caret. Should only contain one value for each hyperparameter. Set to NULL if caret method does not have any hyperparameter values.

`trControl`: trainControl object to be passed to caret `train`.

`vars`: character vector specifying variables for which to determine marginal importance Defaults to all predictor variables in `x`.

`allow_parallel`: boolean for parallel execution. If set to TRUE, user must specify parallel backend in their R session to take advantage of multiple cores. Defaults to FALSE.

`...`: additional arguments to pass to caret `train`

Description

A caret model is trained on training data according to resampling indices specified by the user, and error is calculated on out of sample data. Variable importance is determined by calculating the change in out of sample model performance when a variable is removed relative to baseline out of sample performance when all variables are included.

- **for** $b = 1, \dots, B$ **do**
 1. Draw a bootstrap sample of the data
 2. Fit the model and calculate its prediction error Err_b , using the OOB data
 3. Fit a second model, but without variable v , and calculate its prediction error, $Err_{v,b}^{marg}$
- **end for**
- Calculate the marginal VIMP by averaging: $\Delta_v^{marg} = \sum_{b=1}^B [Err_{v,b}^{marg} - Err_b] / B$

The workhorse function for Marginal Variable Importance is implemented in `calculate_marginal_vimp` (source).