

Understanding Black-Box Models

Yifei Sun

Contents

Variable importance	2
Partial dependence plots (PDPs)	3
Individual conditional expectation (ICE) curves	4
lime	5

```
library(ISLR)
library(caret)
library(vip)
library(pdp)
library(lime)
library(gridExtra)
library(tidymodels)
```

Predict a baseball player's salary on the basis of various statistics associated with performance in the previous year.

```
data(Hitters)
Hitters <- na.omit(Hitters)

set.seed(2)
data_split <- initial_split(Hitters, prop = 0.8)

training_data <- training(data_split)
testing_data <- testing(data_split)

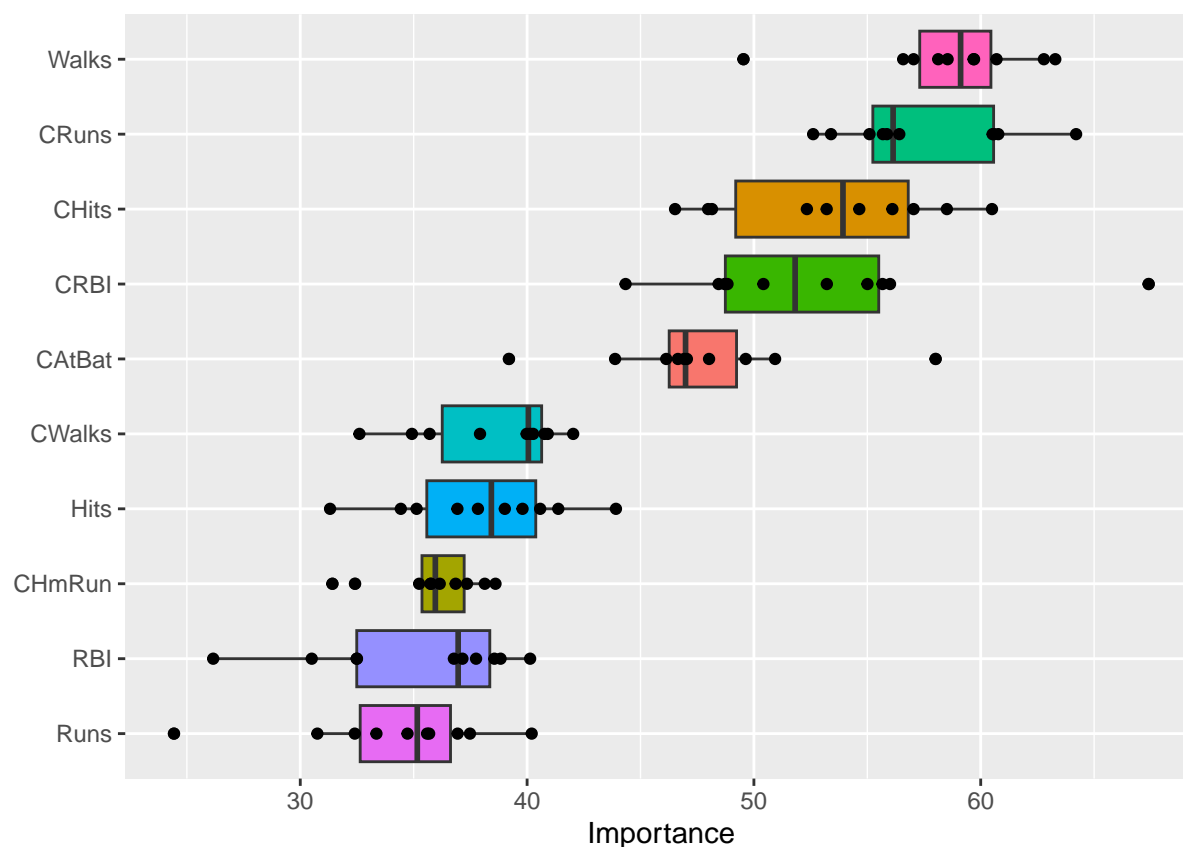
ctrl <- trainControl(method = "cv")

rf.grid <- expand_grid(mtry = 1:6,
                      splitrule = "variance",
                      min.node.size = 1:6)

set.seed(1)
rf.fit <- train(Salary~., training_data,
                method = "ranger",
                tuneGrid = rf.grid,
                trControl = ctrl)
```

Variable importance

```
# You can use built-in VIMP in random forests
# The following code is to illustrate vip() in calculating VIMP in a general case
set.seed(1)
vip(rf.fit,
    method = "permute",
    train = training_data,
    target = "Salary",
    metric = "RMSE",
    nsim = 10,
    pred_wrapper = predict,
    geom = "boxplot",
    all_permutations = TRUE,
    mapping = aes_string(fill = "Variable"))
```



Partial dependence plots (PDPs)

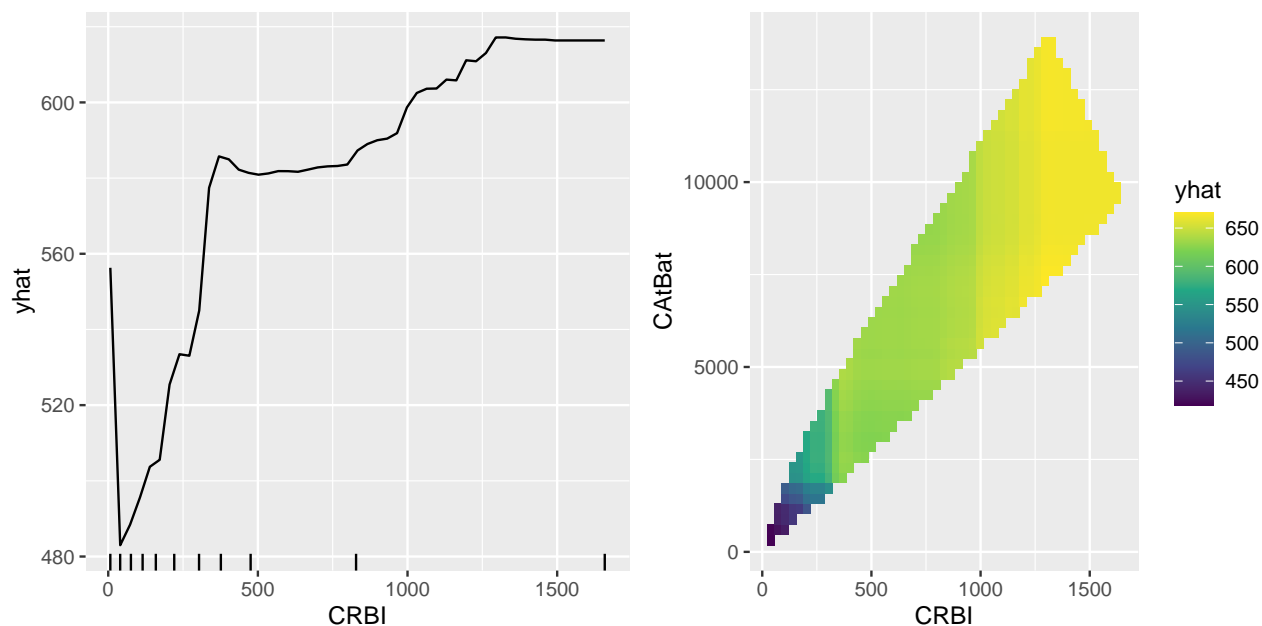
After the most relevant variables have been identified, the next step is to attempt to understand how the response variable changes based on these variables.

PDPs plot the change in the average predicted value as specified feature(s) vary over their marginal distribution.

```
pdp1.rf <- rf.fit %>%
  pdp::partial(pred.var = c("CRBI")) %>%
  autoplot(train = training_data, rug = TRUE)

pdp2.rf <- rf.fit %>%
  pdp::partial(pred.var = c("CRBI", "CAtBat"), chull = TRUE) %>%
  autoplot(train = training_data, rug = TRUE)

grid.arrange(pdp1.rf, pdp2.rf, nrow = 1)
```



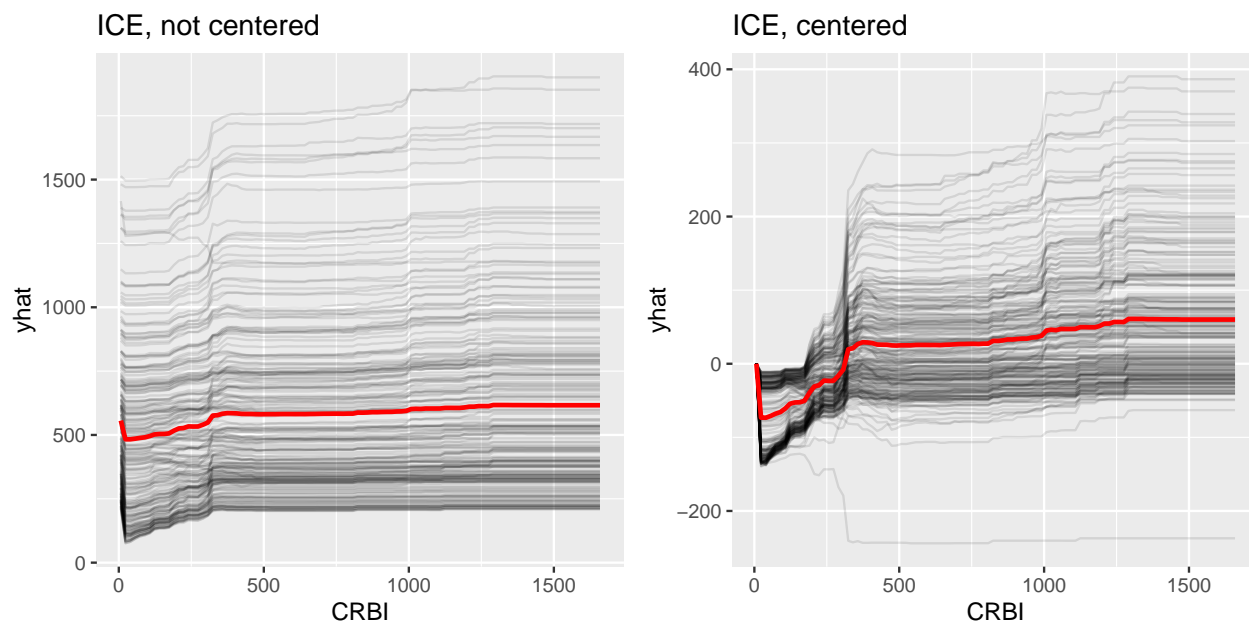
Individual conditional expectation (ICE) curves

ICE curves are an extension of PDP plots but, rather than plot the average marginal effect on the response variable, we plot the change in the predicted response variable for each observation as we vary each predictor variable.

```
ice1.rf <- rf.fit %>%
  pdp::partial(pred.var = "CRBI",
               grid.resolution = 100,
               ice = TRUE) %>%
  autoplot(train = training_data, alpha = .1) +
  ggtitle("ICE, not centered")

ice2.rf <- rf.fit %>%
  pdp::partial(pred.var = "CRBI",
               grid.resolution = 100,
               ice = TRUE) %>%
  autoplot(train = training_data, alpha = .1,
           center = TRUE) +
  ggtitle("ICE, centered")

grid.arrange(ice1.rf, ice2.rf, nrow = 1)
```



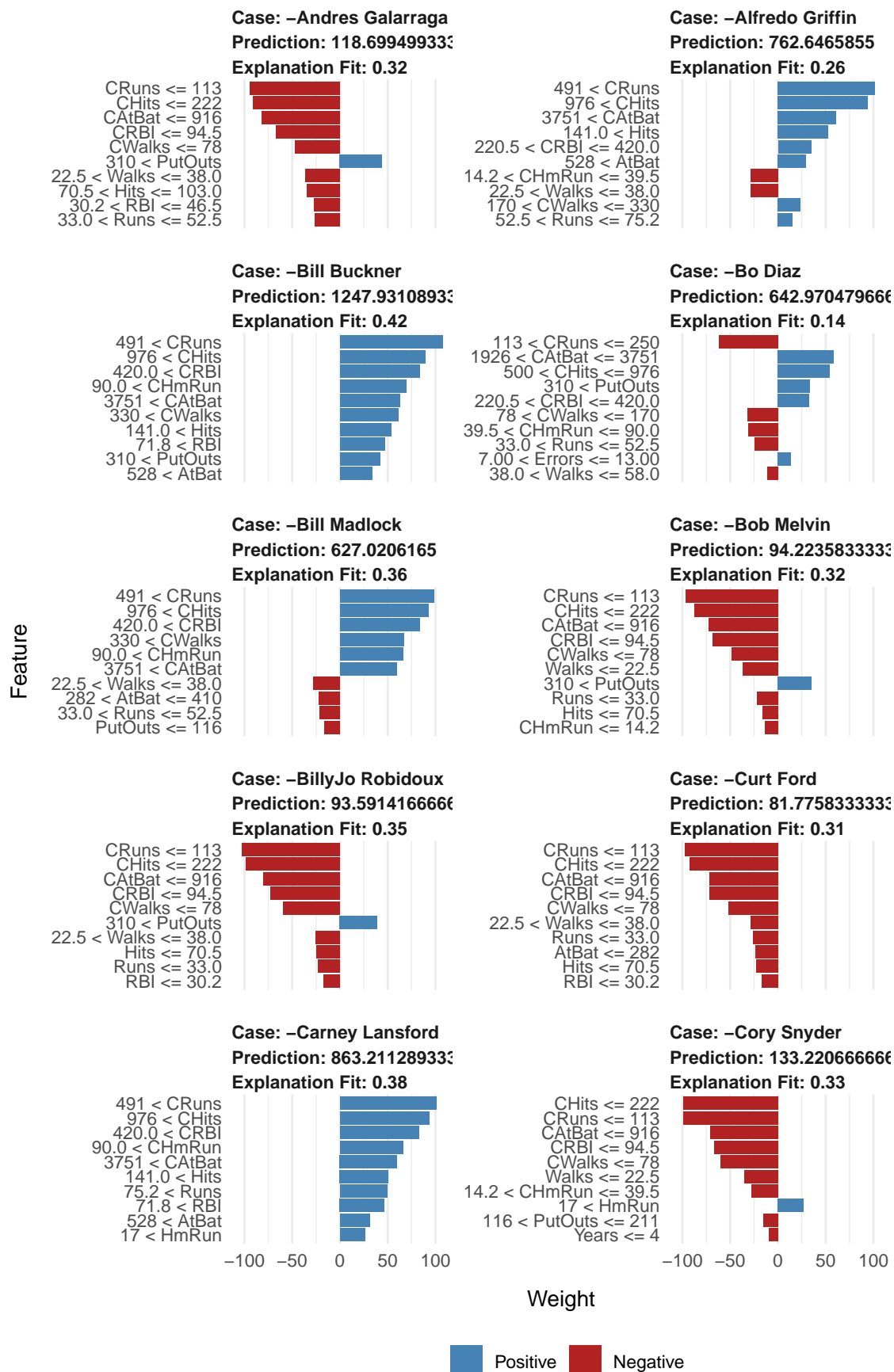
lime

Once an explainer has been created using the `lime()` function, it can be used to explain the result of the model on new observations. The `explain()` function takes new observation along with the explainer and returns a data.frame with prediction explanations, one observation per row. The function `plot_features()` creates a compact visual representation of the explanations for each case and label combination in an explanation. Each extracted feature is shown with its weight, thus giving the importance of the feature in the label prediction.

```
explainer.rf <- lime(training_data[, -19], rf.fit)

new_obs <- testing_data[, -19][1:10,]
explanation.obs <- lime::explain(new_obs,
                               explainer.rf,
                               n_features = 10)

plot_features(explanation.obs)
```



```
plot_explanations(explanation.obs)
```

