

Decision Trees and Random Forests

Laura Cline

23/10/2021

Contents

Fitting Classification Trees	1
Algorithm 8.1: Pruning Trees	4
Fitting Regression Trees	6
Bagging and Random Forests	9
Boosting	11

```
rm(list = ls(all=TRUE))
```

Fitting Classification Trees

The `tree` library is used to construct classification and regression trees.

```
#install.packages("tree")
#install.packages("gbm")
library(tree)
library(tidyverse)

## Registered S3 method overwritten by 'cli':
##   method      from
##   print.tree  tree

## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.3    v purrr   0.3.4
## v tibble  3.1.0    v dplyr   1.0.5
## v tidyr   1.1.3    v stringr 1.4.0
## v readr   1.4.0    v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(stringr)
library(MASS)

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##   select
```

```
library(ISLR)
library(randomForest)
```

```
## randomForest 4.6-14
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:dplyr':
##
##      combine
## The following object is masked from 'package:ggplot2':
##
##      margin
```

```
library(gbm)
```

```
## Loaded gbm 2.1.8
```

We first use classification trees to analyze the `Carseats` dataset. In this data, `Sales` is a continuous variable, and so we begin by recoding it as a binary variable. We use the `ifelse()` function to create a variable, called `High`, which takes on a value of `Yes` if the `Sales` variable exceeds 8, and takes on a value of `No` otherwise.

```
data <- Carseats %>%
  mutate(high = factor(if_else(Sales > 8, 1, 0)))
```

Finally, we use the `data.frame()` function to merge `High` with the rest of the `Carseats` data.

```
names(data) = str_to_lower(names(data))
```

We now use the `tree()` function to fit a classification tree in order to predict `High` using all variables but `Sales`. The syntax of the `tree()` function is quite similar to that of the `lm()` function.

```
# Set up initial tree
tree = tree(high ~ . - sales, data)
```

The `summary()` function lists the variables that are used as internal nodes in the tree, the number of terminal nodes, and the training error rate.

```
summary(tree)
```

```
##
## Classification tree:
## tree(formula = high ~ . - sales, data = data)
## Variables actually used in tree construction:
## [1] "shelve" "price" "income" "compprice" "population"
## [6] "advertising" "age" "us"
## Number of terminal nodes: 27
## Residual mean deviance: 0.4575 = 170.7 / 373
## Misclassification error rate: 0.09 = 36 / 400
```

We see that the training error rate is 9%. For classification trees, the deviance reported in the output of `summary()` is given by:

$$-2 \sum_m \sum_k n_{mk} \log \hat{p}_{mk},$$

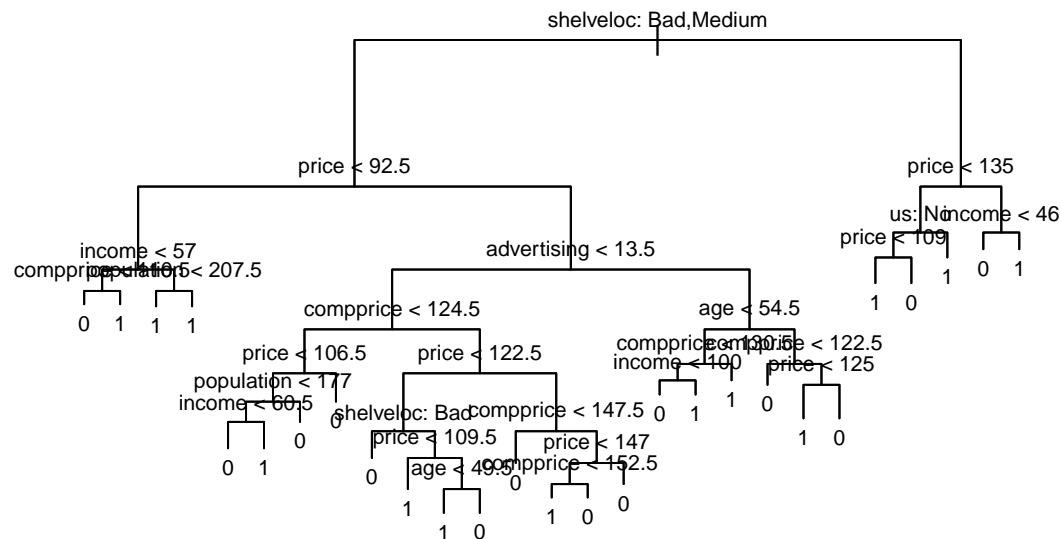
where n_{mk} is the number of observations in the m th terminal node that belong to the k th class. A small deviance indicates a tree that provides a good fit to the training data. The *residual mean deviance* reported

is simply the deviance divided by $n - |T_0|$, which in this case is $400 - 27 = 373$.

One of the most attractive properties of trees is that they can be graphically displayed. We use the `plot()` function to display the tree structure, and the `text()` function to display the node labels. The argument `pretty = 0` instructs R to include the category names for any qualitative predictors, rather than simply displaying a letter for each category.

We know that trees can easily be visually interpreted, so let's check out what our test case looks like. Be wary of doing any plotting because the tree can easily grow out of control and make the graph impossible to read. We see that in the following plot of our current results:

```
plot(tree)
text(tree, pretty = 0, cex = 0.65, digits = 1)
```



The most important indicator of **Sales** appears to be shelving location since the first branch differentiates 'Good' locations from 'Bad' and 'Medium' locations.

If we just type the name of the tree object, 'R' prints output corresponding to each branch of the tree. R displays the split criterion (e.g., **Price < 92.5**), the number of observations in that branch, the deviance, the overall prediction of that branch (**Yes** or **No**), and the fraction of observations in that branch that take on the values of **Yes** and **No**. Branches that lead to terminal nodes are indicated using asterisks.

```
#tree
```

In order to properly evaluate the performance of a classification tree on this data, we must estimate the test error rather than simply computing the training error. We split the observations into a training set and a test set, build the tree using the training set, and evaluate its performance on the test data. The `predict()` function can be used for this purpose. In the case of a classification tree, the argument `type="class"` instructs R to return the actual class prediction. This approach leads to correct predictions for around 77% of the locations in the test dataset.

In order to evaluate a classification tree, we need to use training and testing sets. Let's now repeat what we did above, this time including the calculation for the test error rate.

```
# Define our training/testing sets
set.seed(2)
train <- sample_n(data, 200)
test <- setdiff(data, train)
```

```
# Run the recursive partitioning algorithm
ttree <- tree(high ~. -sales, data=train)
```

```
# Make predictions and display the confusion matrix
test_predictions <- predict(ttree, test, type='class')
table(test_predictions, test$high)
```

```
##
## test_predictions    0    1
##                   0 104  33
##                   1  13  50
```

```
(104 + 50) / 200
```

```
## [1] 0.77
```

Next, we consider whether pruning the tree might lead to improved results. The function `cv.tree()` performs cross-validation in order to determine the optimal level of tree complexity; cost complexity pruning is used in order to select a sequence of trees for consideration. We use the argument `FUN=prune.misclass` in order to indicate that we want the classification error rate to guide the cross-validation and pruning process, rather than the default for `cv.tree()` function, which is deviance. The `cv.tree()` function reports the number of terminal nodes of each tree considered (`size`) as well as the corresponding error rate and the value of the cost-complexity parameter used (`k`, which corresponds to α).

We now add another layer of complexity by pruning our results. Recall that unpruned trees are prone to overfitting the data, so our method will be to watch variation in the test error rates as we increase the penalty in the number of terminal nodes. To refresh your memory, we summarize **Algorithm 8.1** below:

Algorithm 8.1: Pruning Trees

1. Grow your original tree T_0 using your training data.
2. As a function of α (the penalty parameter), define the sequence of best subtrees.
3. Use K-fold cross-validation to find the α that minimizes the average mean squared prediction error of the k th fold of the training data.
4. Find the best subtree from Step 2 using the α found in the previous step.

Luckily, `tree::cv.tree` will be doing most of the work for us. It will perform cross-validation required to determine the optimal tree size. It also allows us to choose the function by which the tree is pruned. In this case, pruning will be guided by the classification error rate.

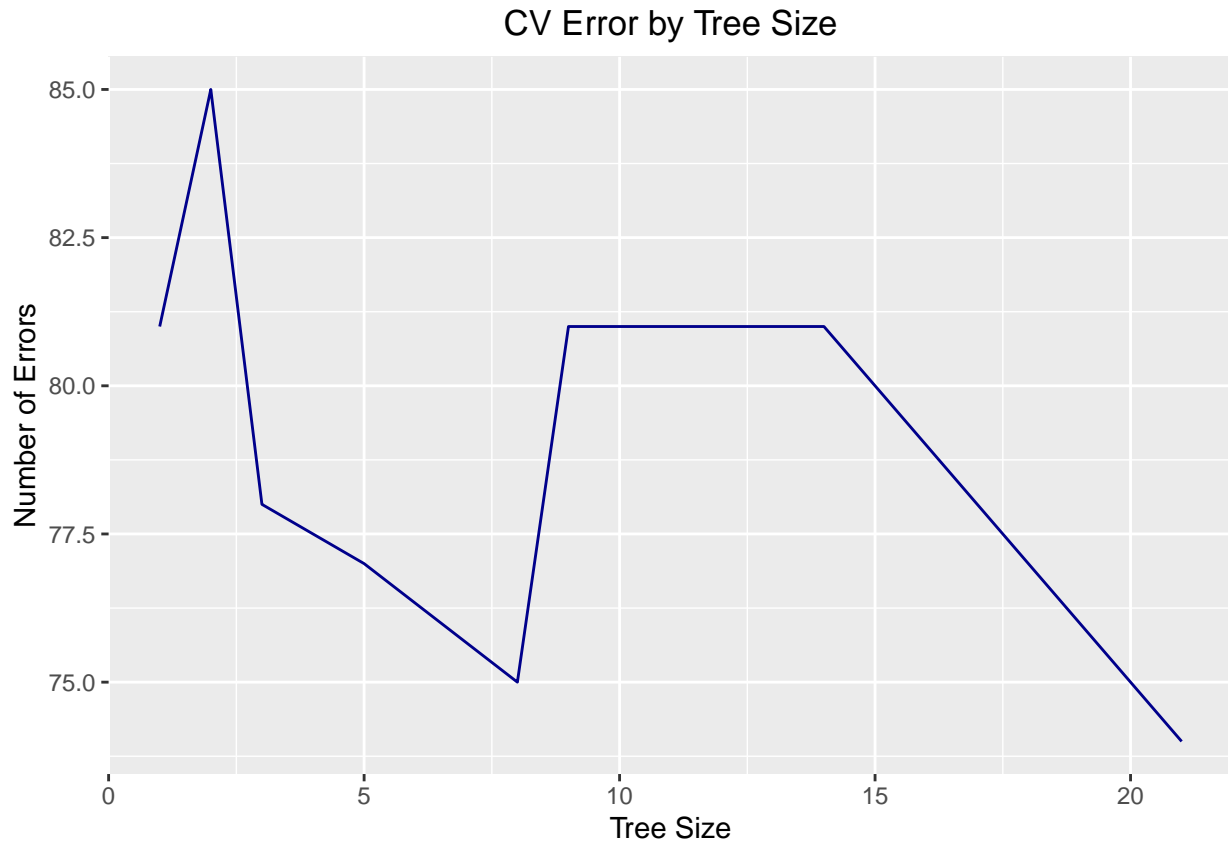
```
set.seed(3)
cv_tree <- cv.tree(ttree, FUN = prune.misclass)
cv_tree

## $size
## [1] 21 19 14  9  8  5  3  2  1
##
## $dev
## [1] 74 76 81 81 75 77 78 85 81
##
## $k
## [1] -Inf  0.0  1.0  1.4  2.0  3.0  4.0  9.0 18.0
##
## $method
## [1] "misclass"
##
## attr(,"class")
## [1] "prune"          "tree.sequence"
```

Note that, despite the name, `dev` corresponds to the cross-validation error rate in this instance. The tree with 21 terminal nodes results in the lowest cross-validation error rate, with 74 cross validation errors. We plot the error rate as a function of both `size` and `k`.

```
library(ggplot2)

ggplot(data = data.frame(cv_tree$size, cv_tree$dev),
  aes(x = cv_tree$size, y = cv_tree$dev)) +
  geom_line(color = "darkblue") +
  labs(x = "Tree Size", y = "Number of Errors", title = "CV Error by Tree Size") +
  theme(plot.title = element_text(hjust = .5))
```



We now apply the `prune.misclass()` function in order to prune the tree to obtain the nine-node tree.

Now that we know exactly how many terminal nodes we want, we prune our tree with `prune.misclass()` to obtain the optimal tree. Then check to see if this tree performs any better on the testing set than the base tree T_0 did.

```
pruned <- prune.misclass(ttree, best=21)
```

How well does this pruned tree perform on the test dataset? Once again, we apply the `predict()` function.

```
test_predictions <- predict(pruned, data=test, type='class')
table(test_predictions, test$high)
```

```
##
## test_predictions  0  1
##                  0 79 49
##                  1 38 34
```

```
(77 + 38)/200
```

```
## [1] 0.575
```

Now, 57% of the test observations are correctly classified, so in this case, the pruning process has produced a more interpretable tree, but the classification accuracy has decreased.

Fitting Regression Trees

Here we fit a regression tree to the `Boston` dataset. First we create a training set and fit the tree to the training data.

Not much changes in terms of code when we switch to regression trees, so this section will pretty much be a recap of the previous one, just using different data. We pull the `Boston` dataset from the `MASS` library for this section.

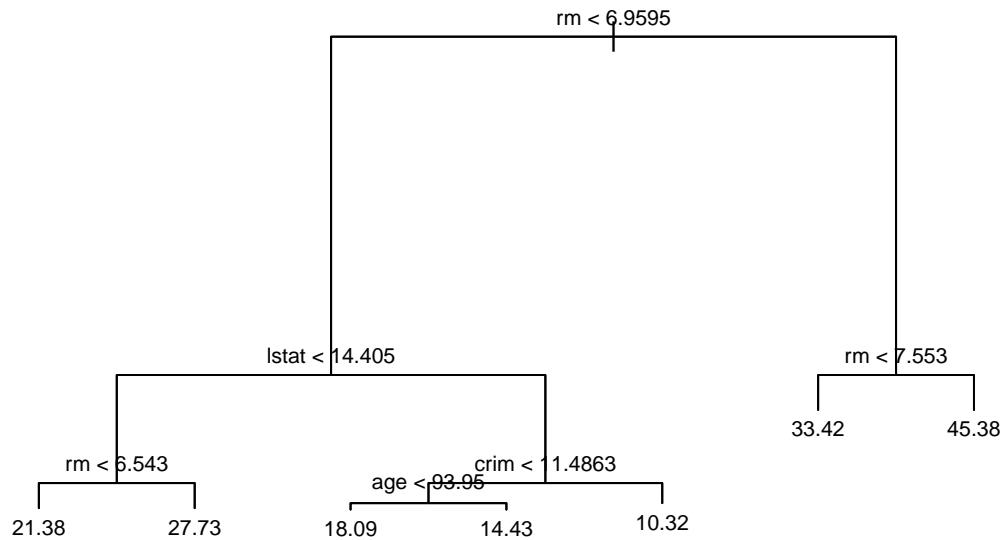
```
Boston <- MASS::Boston
set.seed(1)
train <- sample_frac(Boston, 0.5)
test <- setdiff(Boston, train)

tree_train <- tree(medv ~ ., data=train)
summary(tree_train)
```

```
##
## Regression tree:
## tree(formula = medv ~ ., data = train)
## Variables actually used in tree construction:
## [1] "rm"    "lstat" "crim"  "age"
## Number of terminal nodes: 7
## Residual mean deviance: 10.38 = 2555 / 246
## Distribution of residuals:
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -10.1800 -1.7770 -0.1775  0.0000  1.9230 16.5800
```

Notice that the output of `summary()` indicates that only four of the variables have been used in constructing the tree. In the context of a regression tree, the deviance is simply the sum of squared errors for the tree. We now plot the tree.

```
plot(tree_train)
text(tree_train, pretty = 0, cex = 0.65)
```



As you can see, $rm < 6.959$ is the first partition in the tree. The variable means the average number of rooms per dwelling so lower values of rm means less average rooms on the left side of the tree, this suggests that houses with more rooms end up with much larger median house prices.

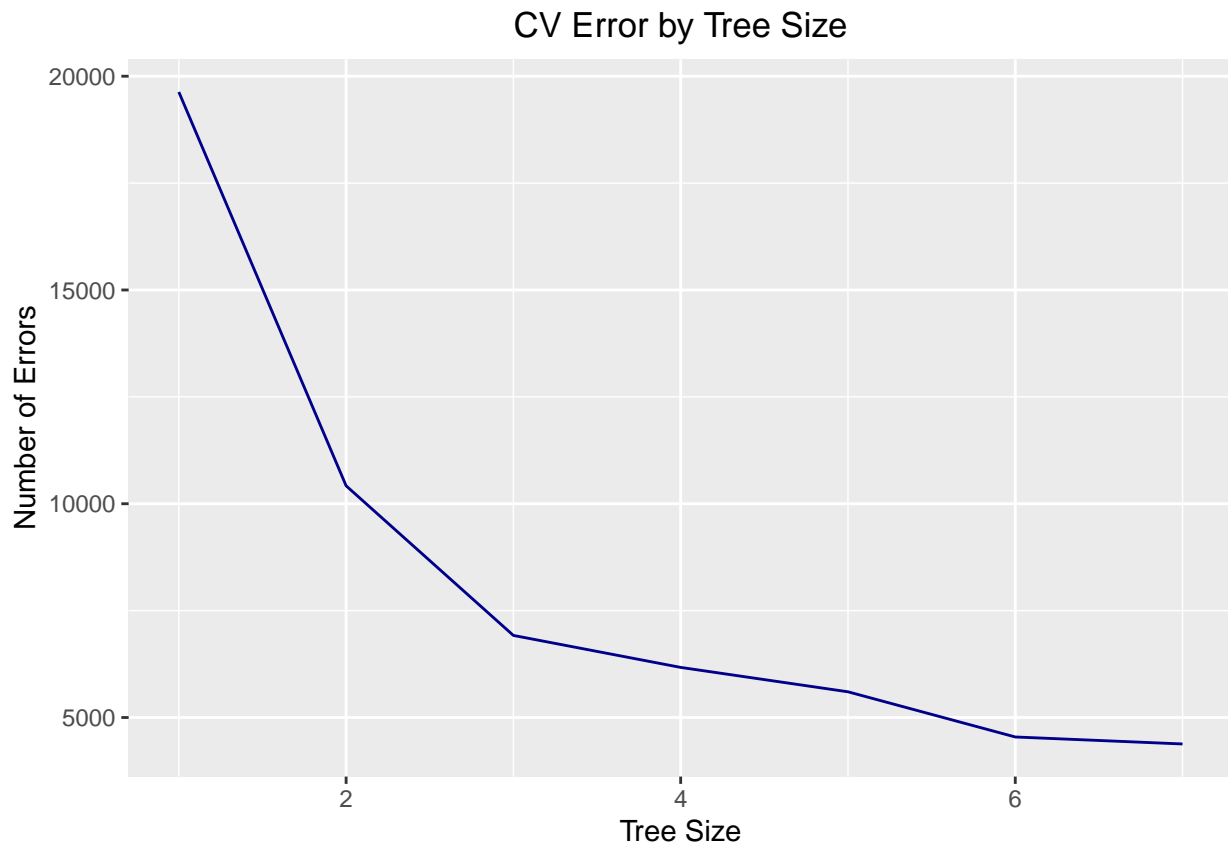
The tree predicts a median house price of \$25,380 for larger homes with more than seven rooms.

Now we use the `cv.tree()` function to see whether pruning the tree will improve performance.

```
cv_tree <- cv.tree(tree_train)
```

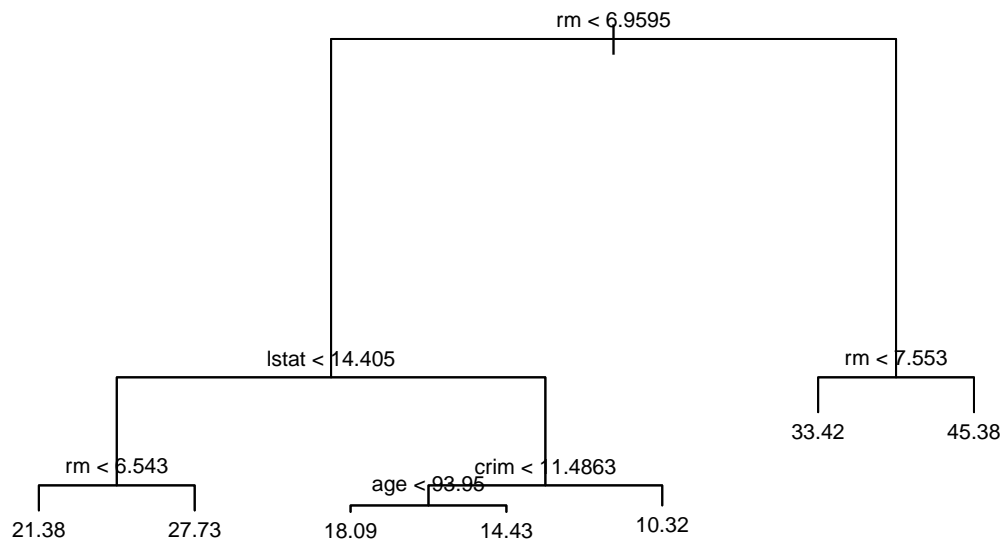
```
# Get an idea of change in error by changing tree size
```

```
ggplot(data = data.frame(cv_tree$size, cv_tree$dev), aes(x = cv_tree$size, y = cv_tree$dev)) +
  geom_line(color = "darkblue") +
  labs(x = "Tree Size", y = "Number of Errors", title = "CV Error by Tree Size") +
  theme(plot.title = element_text(hjust = 0.5))
```



In this case, the most complex tree is selected by cross-validation. However, if we wish to prune the tree, we could do so as follows using the `prune.tree()` function:

```
prune.boston = prune.tree(tree_train, best= 7)
plot(prune.boston)
text(prune.boston, pretty = 0, cex = 0.65)
```



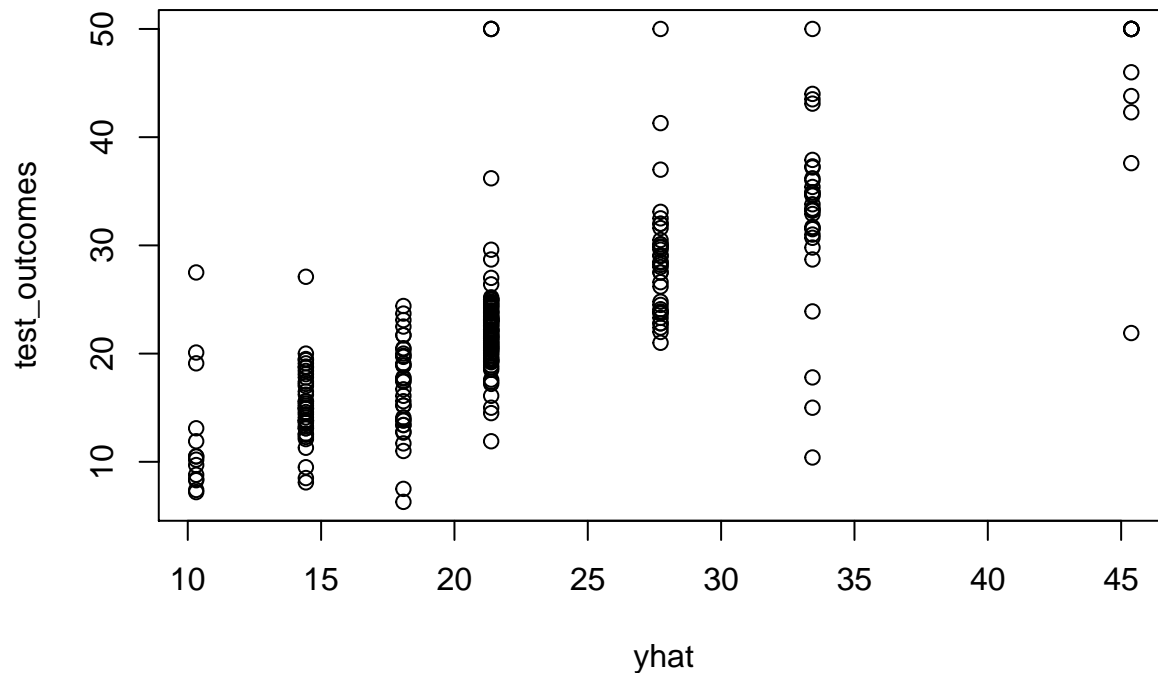
In keeping with cross-validation results, we use the unpruned tree to make predictions on the test set.

```
# Predict, plot and calculate MSE
yhat <- predict(tree_train, newdata=test)
```



```
test_outcomes <- test$medv
```

```
plot(yhat, test_outcomes)
```



```
mean((yhat - test_outcomes)^2)
```

```
## [1] 35.28688
```

In other words, the test set MSE associated with the regression tree is 35.29. The square root of the MSE is therefore around 5.94, indicating that this model leads to test predictions that are within \$5,940 of the true median home value for the suburb.

Bagging and Random Forests

Here we apply bagging and random forests to the **Boston** data, using the **randomForest** package in R. The exact results obtained in this section may depend on the version of R and the version of the **randomForest** package installed on your computer. Recall that bagging is simply a special case of a random forest with $m = p$. Therefore, the **randomForest()** function can be used to perform both random forests and bagging.

We'll be using the same data from the previous section and the **randomForest** package to help us accomplish some simple examples. We begin with a bagging example, where all predictors are used in each split. We perform bagging as follows:

```
set.seed(1)
```

```
train <- sample_frac(Boston, 0.5)
```

```
test <- setdiff(Boston, train)
```

```
# Set up the randomForest for the bagging case (all vars included)
```

```
bag <- randomForest(medv ~ ., data=train, mtry = 13, importance = TRUE)
```

```
bag
```

```
##
```

```
## Call:
```

```
## randomForest(formula = medv ~ ., data = train, mtry = 13, importance = TRUE)
```

```
##                Type of random forest: regression
##                Number of trees: 500
## No. of variables tried at each split: 13
##
##                Mean of squared residuals: 11.33119
##                % Var explained: 85.26
```

The argument `mtry = 14` indicates that all 13 predictors should be considered for each split of the tree - in other words, that bagging should be done. How does this bagged model perform on the test set?

```
# Calculate MSE of the testing set for the bagged regression tree
yhat <- predict(bag, test)
mean((yhat - test$medv)^2)
```

```
## [1] 23.4579
```

The test set MSE associated with the bagged regression tree is 23.46, almost half that obtained using an optimally-pruned single tree. We could change the number of trees grown by `randomForest()` using the `ntree` argument.

```
bag.boston <- randomForest(medv ~ ., data=train, mtry = 13, ntree = 25)
yhat <- predict(bag.boston, test)
mean((yhat - test$medv)^2)
```

```
## [1] 22.99145
```

Growing a random forest proceeds in exactly the same way, except that we use a smaller value of the `mtry` argument. By default, `randomForest()` uses $p/3$ variables when building a random forest of regression trees, and \sqrt{p} variables when building a random forest of classification trees. Here we use `mtry = 6`.

Compare the MSE of the bagged random forest to the optimally pruned single tree found earlier - it's much lower. We manually change the amount of variables at each split in the above bagging example, but we might achieve better results using a general random forest. By default, `randomForest` uses $p/3$ variables when building a forest of regression trees and \sqrt{p} for classification trees. In the following example, we will use `mtry = 6` ($m \approx p/2$).

```
forest = randomForest(medv ~ ., data=train, mtry = 6, importance = TRUE)
```

```
yhat <- predict(forest, test)
mean((yhat - test$medv)^2)
```

```
## [1] 20.16422
```

We find that this approach worked - our MSE is now reduced to 19.88, lower than the previous two methods we tried.

Using the `importance()` function, we can view the importance of each variable.

```
importance(forest)
```

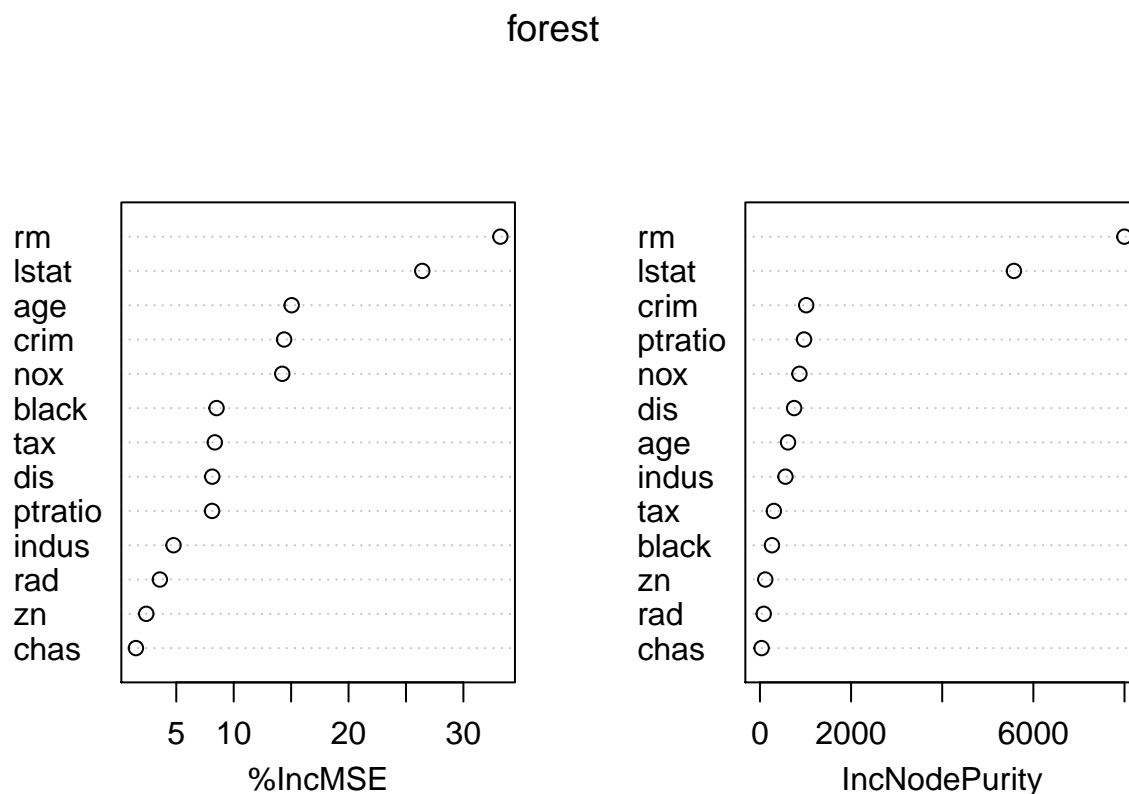
```
##          %IncMSE IncNodePurity
## crim      14.389520    1015.38819
## zn         2.383675     116.75818
## indus      4.759755     560.00649
## chas       1.485182      36.49459
## nox       14.227996     866.91472
## rm        33.219106    7997.53209
## age       15.046634     616.54222
## dis        8.131485     751.23862
## rad        3.571826      83.08033
```

```
## tax      8.356359    305.72575
## ptratio  8.115858    967.88590
## black    8.507384    264.38182
## lstat    26.418993   5573.29309
```

The first column represents the mean decrease in accuracy of the prediction when the variable is removed from the model, and the second column is a measure of the total decrease in node purity resulting from splits over that variable (averaged over all of the trees).

Two measures of variable importance are reported. The former is based upon the mean decrease of accuracy in predictions on the out of bag samples when a given variable is excluded from the model. The latter is a measure of the total decrease in node impurity that results from splits over that variable, averaged over all trees. In the case of regression trees, the node impurity is measured by the training RSS, and for classification trees by the deviance. Plots of these importance measures can be produced using the `varImpPlot()` function.

```
varImpPlot(forest)
```



The results indicate that across all of the trees considered in the random forest, the average number of rooms (`rm`) and the wealth level of the community (`lstat`) are by far the two most important variables.

Boosting

Here we use the `gbm` package, and within it the `gbm()` function, to fit boosted regression trees to the `Boston` dataset. we run `gbm()` with the option of `distribution="gaussian"` since this is a regression problem; if it were a binary classification problem, we would use `distribution=bernoulli`". The argument `n.trees=5000` indicates that we want 5,000 trees, and the option `interaction.depth=4` limits the depth of each tree.

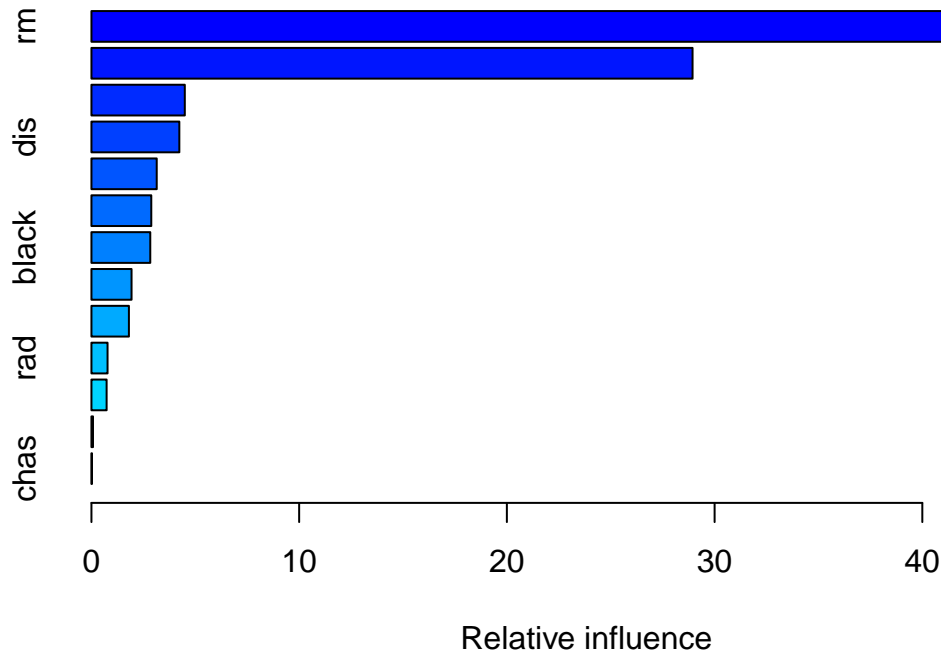
We'll be using the `gbm` package to help us fit boosted regression trees to the `Boston` dataset, which you should be familiar with by now.

```
set.seed(1)
train <- sample_frac(Boston, 0.5)
test <- setdiff(Boston, train)

# Regression => distr = gaussian
boosted <- gbm(medv ~ ., train, distribution="gaussian", n.trees=5000, interaction.depth = 4)
```

The `summary()` function produces a relative influence plot and also outputs the relative influence statistics.

```
# Summarize and produce a quick plot to highlight importance of variables
summary(boosted)
```

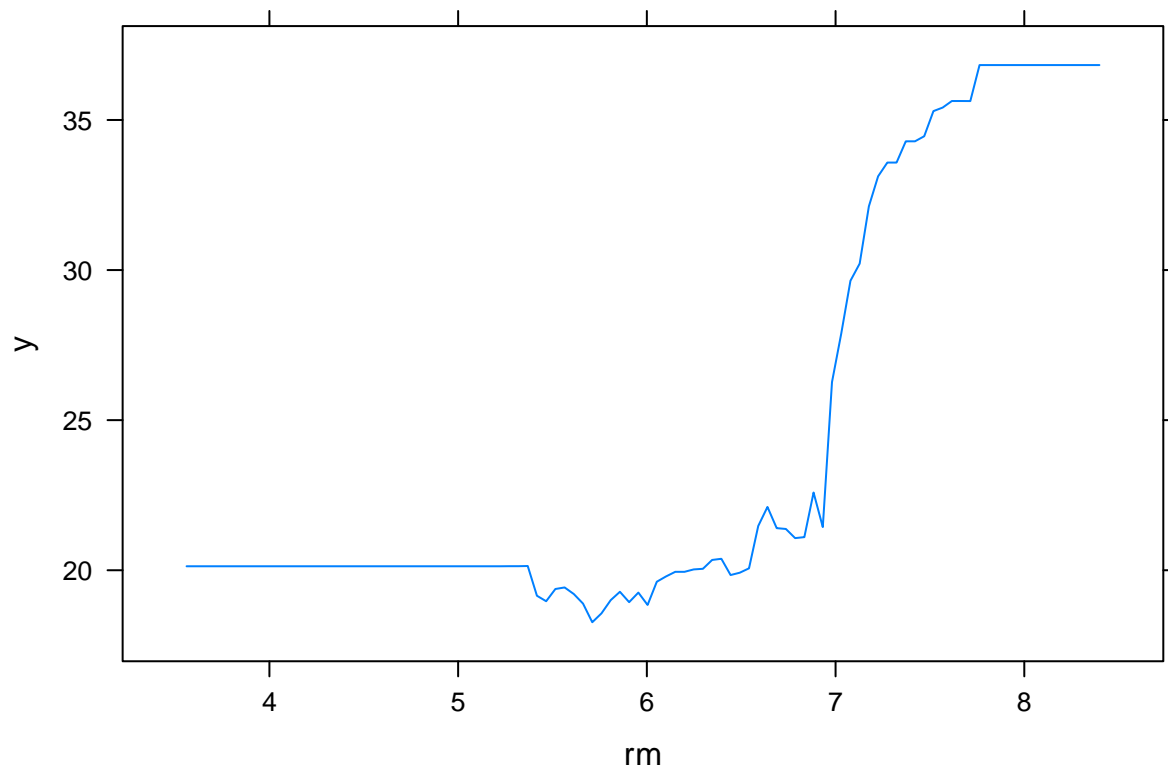


```
##      var      rel.inf
## rm      rm 48.13967682
## lstat    lstat 28.93851185
## crim     crim  4.49413146
## dis      dis  4.23182696
## age      age  3.14221169
## nox      nox  2.88094283
## black    black 2.83238772
## ptratio  ptratio 1.93050932
## tax      tax  1.80427054
## rad      rad  0.77569461
## indus    indus 0.73110525
## zn       zn   0.07442923
## chas     chas  0.02430170
```

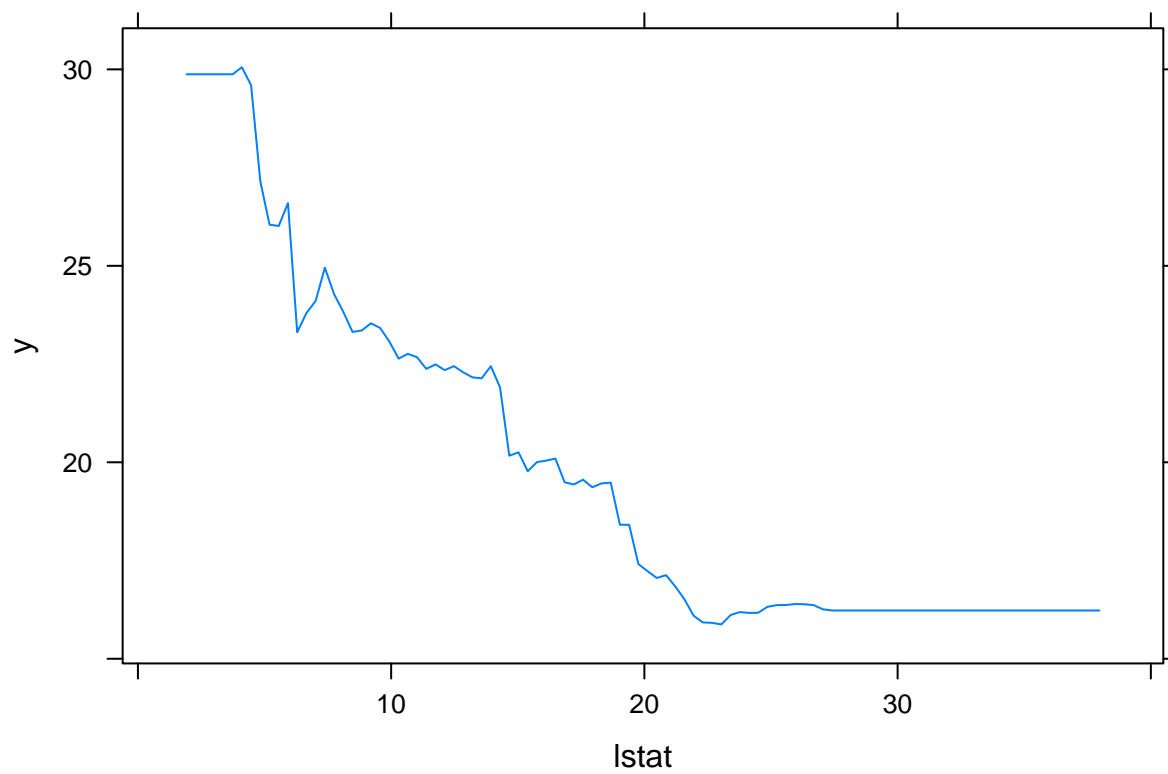
We see that `rm` and `lstat` are by far the most important variables. We can also produce *partial dependence plots* for these two variables. These plots illustrate the marginal effect of the selected variables on the response after *integrating* out the other variables. In this case, as we might expect, median house prices are increasing with `rm` and decreasing with `lstat`.

Let's plot the marginal effect of these two variables, `lstat` and `rm`.

```
par(mfrow = c(1,2))  
plot(boosted, i = 'rm')
```



```
plot(boosted, i = 'lstat')
```



Alright, this just confirmed what we should have already been expecting: median house values are decreasing with `lstat` and increasing in `rm`.

We now use the boosted model to predict `medv` on the test set.

Let's now test how well this boosted regression tree performs on the testing data.

```
yhat <- predict(boosted, newdata = test, n.trees=5000)
mean((yhat - test$medv)^2)
```

```
## [1] 19.37033
```

Not amazing, but not bad. The boosted model performed just about the same as random forests and superior to that of the bagging model, but we might be able to squeeze out some extra performance by changing the shrinkage parameter λ .

The test MSE obtained is 19.37; similar to the test MSE for random forests and superior to that for bagging. If we want to, we can perform boosting with a different value of the shrinkage parameter λ . The default value is 0.01, but this is easily modified. Here we take $\lambda = 0.2$.

```
boosted <- gbm(medv ~ ., train, distribution="gaussian", n.trees=5000, interaction.depth = 4, shrinkage = 0.2)
yhat <- predict(boosted, newdata=test, n.trees=5000)
mean((yhat - test$medv)^2)
```

```
## [1] 18.68911
```

Changing the shrinkage parameter actually made a difference - we are now just slightly under what we got from our previous model where it was equal to 0.001.