

Unit 4 Lecture 3: Random forests

November 2, 2023

Today, we will learn how to train and tune random forests using the `randomForest` package.

First, let's load some libraries:

```
library(randomForest)
library(tidyverse)
```

Random forests for regression

Like last time, we will be using the `hitters` data, splitting into training and testing:

```
hitters_data <- read_csv("hitters-data.csv")

set.seed(1) # set seed for reproducibility
train_samples <- sample(1:nrow(hitters_data), round(0.8 * nrow(hitters_data)))
hitters_train <- hitters_data |> filter(row_number() %in% train_samples)
hitters_test <- hitters_data |> filter(!(row_number() %in% train_samples))
```

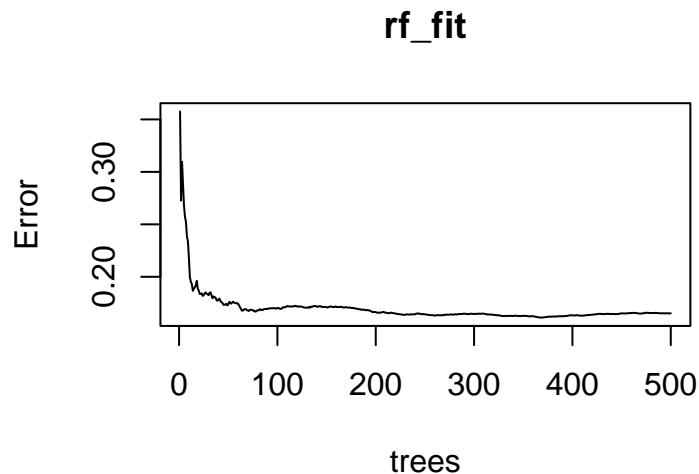
Training a random forest

To train a random forest with default settings, we use the following syntax:

```
rf_fit <- randomForest(Salary ~ ., data = hitters_train)
?randomForest
```

We can get a quick visualization by using `plot`, which shows us the OOB error as a function of the number of trees.

```
plot(rf_fit)
```



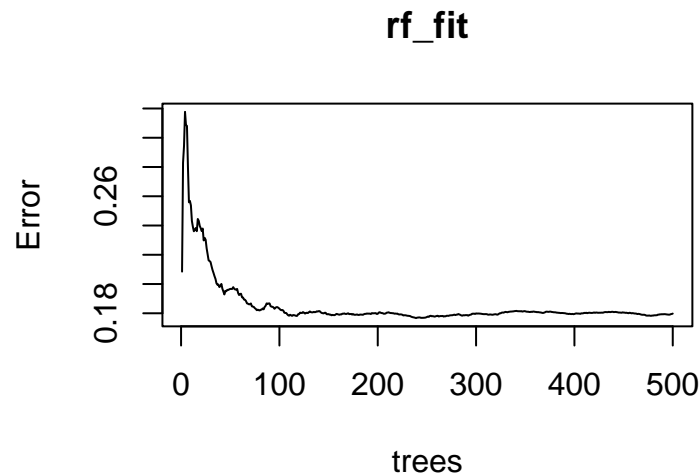
We see that this error stays flat as soon as `B` is large enough (in this case stabilizing around 100).

The key parameters controlling the random forest fit are the following:

- **mtry**: number of variables to sample for each split (called m in lecture), default $\text{floor}(p/3)$ for regression and $\text{sqrt}(p)$ for classification
- **nodesize**: minimum size of terminal nodes, default 1 for classification and 5 for regression
- **maxnodes**: maximum number of terminal nodes trees in the forest can have, default no maximum
- **ntree**: number of trees (called B in lecture), default 500

We might want to specify the **mtry** parameter manually. For example, to get the bagging predictions we can set **mtry** = 19, since 19 is the total number of features:

```
rf_fit <- randomForest(Salary ~ ., mtry = 19, data = hitters_train)
plot(rf_fit)
```



Tuning the random forest

A quick-and-dirty way to tune a random forest is to try out a few different values of **mtry**:

```
rf_3 <- randomForest(Salary ~ ., mtry = 3, data = hitters_train)
rf_6 <- randomForest(Salary ~ ., mtry = 6, data = hitters_train)
rf_19 <- randomForest(Salary ~ ., mtry = 19, data = hitters_train)
```

We can extract the OOB errors from each of these objects by using the **mse** field:

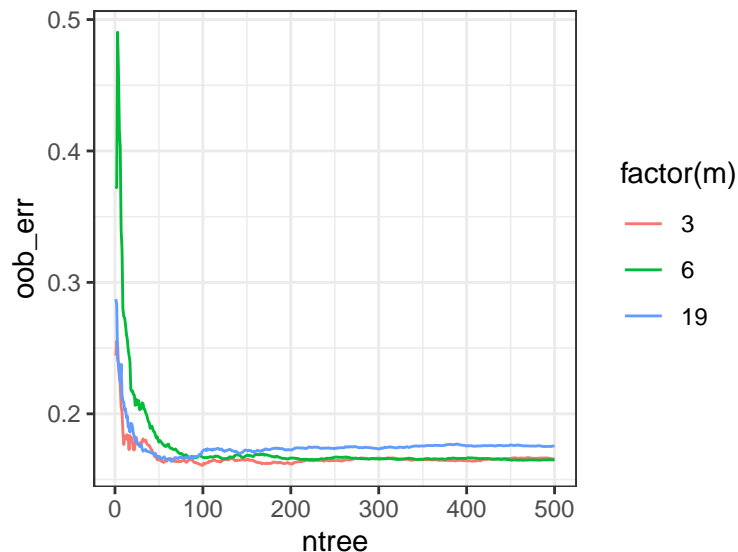
```
oob_errors <- bind_rows(
  tibble(ntree = 1:500, oob_err = rf_3$mse, m = 3),
  tibble(ntree = 1:500, oob_err = rf_6$mse, m = 6),
  tibble(ntree = 1:500, oob_err = rf_19$mse, m = 19)
)
oob_errors
```

```
## # A tibble: 1,500 x 3
##   ntree oob_err    m
##   <int> <dbl> <dbl>
## 1     1  0.244     3
## 2     2  0.255     3
## 3     3  0.253     3
## 4     4  0.236     3
## 5     5  0.228     3
## 6     6  0.222     3
## 7     7  0.205     3
## 8     8  0.201     3
```

```
## 9      9    0.187    3
## 10     10   0.177    3
## # i 1,490 more rows
```

We can then plot these as follows:

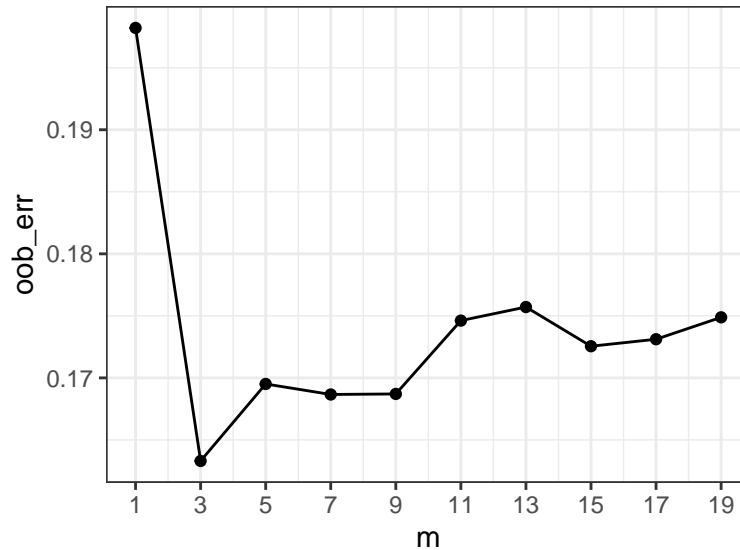
```
oob_errors |>
  ggplot(aes(x = ntree, y = oob_err, colour = factor(m))) +
  geom_line()
```



Which value of `mtry` seems to work the best here?

We can be a little more systematic in tuning the random forest by choosing a grid of values of `mtry` and plotting the OOB error for 500 trees versus `mtry`:

```
# might want to cache this chunk!
mvalues <- seq(1, 19, by = 2)
oob_errors <- numeric(length(mvalues))
ntree <- 500
for (idx in 1:length(mvalues)) {
  m <- mvalues[idx]
  rf_fit <- randomForest(Salary ~ ., mtry = m, data = hitters_train)
  oob_errors[idx] <- rf_fit$mse[ntree]
}
tibble(m = mvalues, oob_err = oob_errors) |>
  ggplot(aes(x = m, y = oob_err)) +
  geom_line() +
  geom_point() +
  scale_x_continuous(breaks = mvalues)
```



Variable importance

Let's go back to the default random forest fit:

```
rf_fit <- randomForest(Salary ~ ., data = hitters_train)
```

This object contains the purity-based feature importance in the `importance` field:

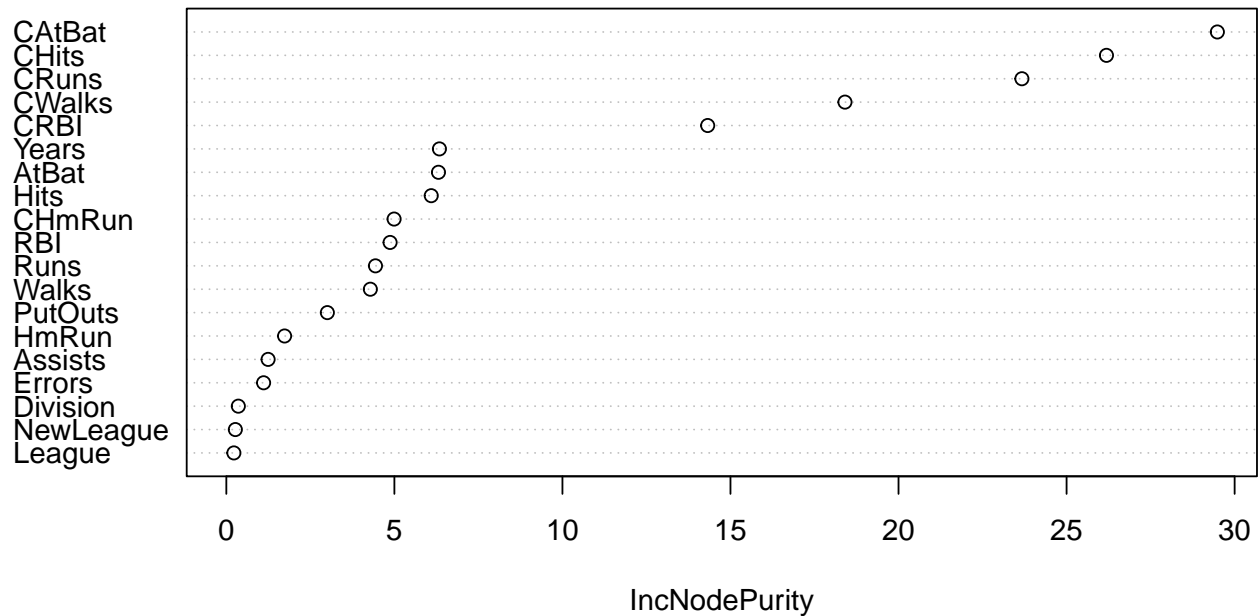
```
rf_fit$importance
```

```
##          IncNodePurity
## AtBat          6.3118423
## Hits           6.0965680
## HmRun          1.7326849
## Runs           4.4369543
## RBI            4.8722013
## Walks          4.2846233
## Years          6.3403421
## CAtBat        29.4857472
## CHits         26.1866816
## CHmRun         4.9909252
## CRuns         23.6713576
## CRBI          14.3213960
## CWalks        18.4065924
## League         0.2242469
## Division       0.3567897
## PutOuts        3.0054984
## Assists        1.2426939
## Errors         1.1086783
## NewLeague      0.2678145
```

We can visualize these importances using the built-in function called `varImpPlot`:

```
varImpPlot(rf_fit)
```

rf_fit



In lecture, we discussed that there were two variable importance measures. If we want to compute the second one (OOB-based importance), we need to explicitly specify this in the call to `randomForest`:

```
rf_fit <- randomForest(Salary ~ ., importance = TRUE, data = hitters_train)
```

Now let's see what the `importance` field looks like:

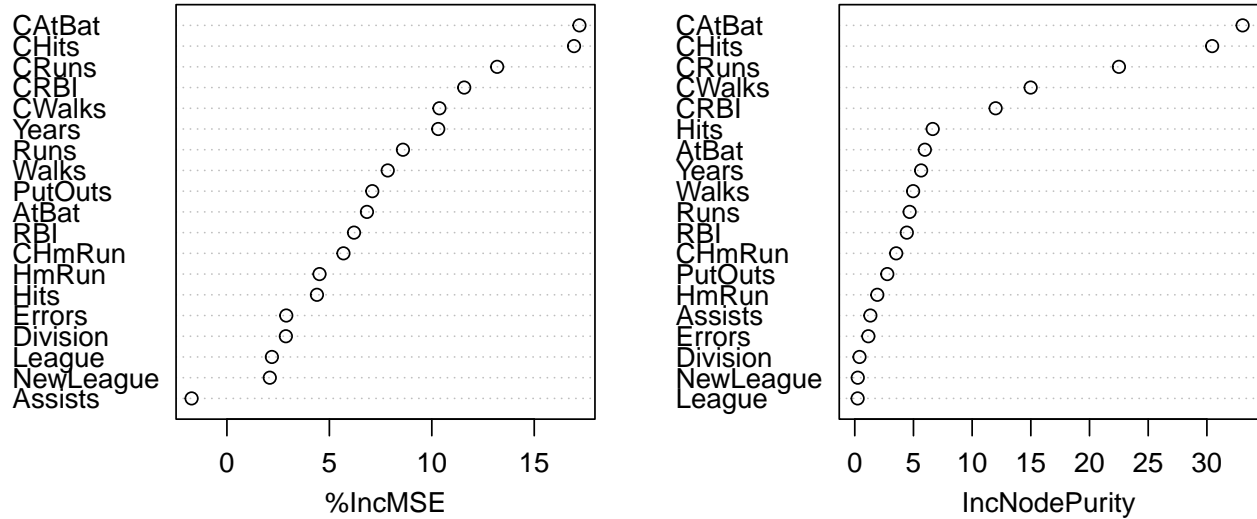
```
rf_fit$importance
```

```
##           %IncMSE  IncNodePurity
## AtBat      0.0160370055      5.9696004
## Hits       0.0105353096      6.6381789
## HmRun      0.0072632171      1.9189612
## Runs       0.0200710538      4.6750218
## RBI        0.0122292305      4.4341936
## Walks      0.0149404658      4.9715901
## Years      0.0252956559      5.6550609
## CAtBat     0.2082537275     33.0437649
## CHits      0.1851466785     30.4478334
## CHmRun     0.0133121878      3.5293561
## CRuns      0.1314162990     22.5059305
## CRBI       0.1027781783     11.9997599
## CWalks     0.0611970496     14.9848253
## League     0.0010101773      0.2441357
## Division   0.0014266008      0.4051892
## PutOuts    0.0081385982      2.7696280
## Assists    -0.0022045174      1.3285290
## Errors     0.0022181862      1.1570355
## NewLeague  0.0007964947      0.2555317
```

We see there are now two columns instead of one! We can plot both of these feature importance measures using the same syntax as above:

```
varImpPlot(rf_fit)
```

rf_fit



Making predictions based on a random forest

We can make predictions using `predict`, as usual:

```
rf_predictions <- predict(rf_fit, newdata = hitters_test)
rf_predictions
```

```
##      1      2      3      4      5      6      7      8
## 6.730500 4.723214 4.524503 4.787757 5.949685 4.768927 7.093869 6.586140
##      9     10     11     12     13     14     15     16
## 5.910027 6.639294 7.190061 5.745520 6.401716 6.951979 5.820046 6.495988
##     17     18     19     20     21     22     23     24
## 4.461315 6.812642 6.217612 6.059059 6.605661 5.600980 6.732921 6.332523
##     25     26     27     28     29     30     31     32
## 6.206834 7.054491 4.653610 6.205098 6.679432 5.711332 5.143725 6.657193
##     33     34     35     36     37     38     39     40
## 5.233679 4.513696 6.194446 6.046111 6.239535 6.586366 5.988931 6.538774
##     41     42     43     44     45     46     47     48
## 6.369567 6.967355 6.543839 4.959653 6.406267 6.971199 4.866071 5.943582
##     49     50     51     52     53
## 5.916375 5.543489 5.043576 6.794856 6.081674
```

We can compute the mean-squared prediction error as usual too:

```
mean((rf_predictions - hitters_test$Salary)^2)
```

```
## [1] 0.2465996
```

Random forests for classification

Random forests work very similarly for classification. Let's continue with the heart disease data from last time:

```
heart_data <- read_csv("heart-data.csv")
```

As it turns out, this dataset contains missing values. While trees deal nicely with missing values, random forests do not. So let's drop the missing values

```
heart_data <- heart_data |> na.omit()
```

```
set.seed(1) # set seed for reproducibility
train_samples <- sample(1:nrow(heart_data), round(0.8 * nrow(heart_data)))
heart_train <- heart_data |> filter(row_number() %in% train_samples)
heart_test <- heart_data |> filter(!(row_number() %in% train_samples))
```

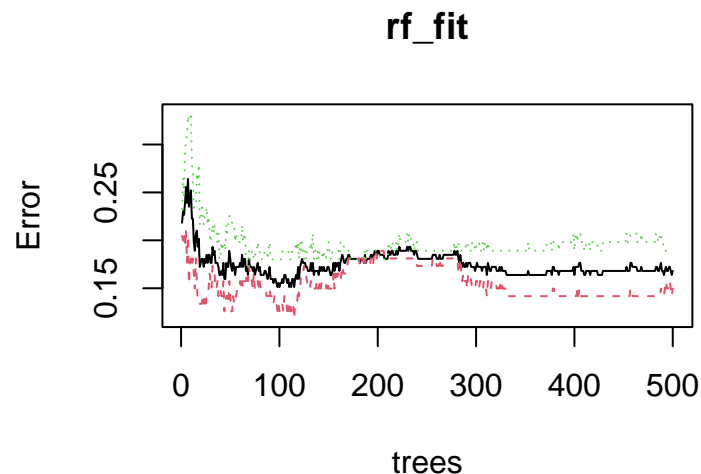
Fitting a random forest uses the same basic syntax:

```
# IMPORTANT: RESPONSE MUST BE CODED AS A FACTOR!
rf_fit <- randomForest(factor(AHD) ~ ., data = heart_train)
```

Note that for random forests the default value of `mtry` is the square root of the number of features, in this case `floor(sqrt(13)) = 3`.

When we go to make the random forest plot it looks slightly different though:

```
plot(rf_fit)
```



That is strange! Why does this happen? What's being plotted are three versions of the OOB error, which are stored in `rf_fit$err.rate`:

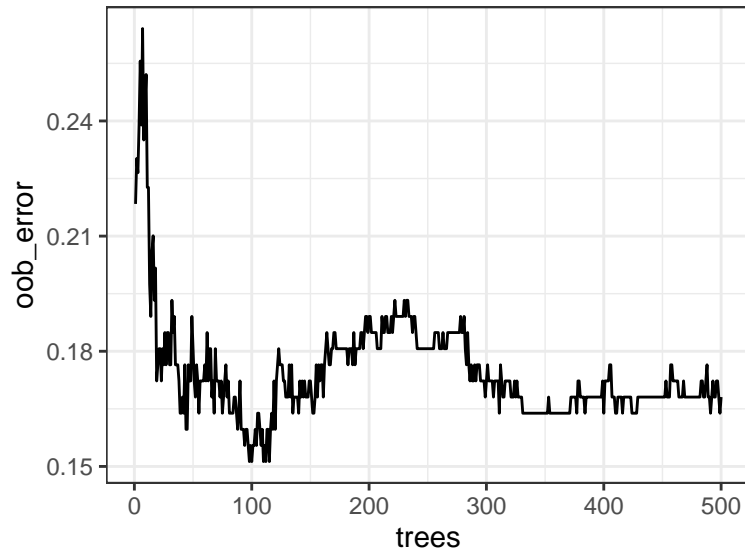
```
rf_fit$err.rate |> head()
```

```
##           OOB           No           Yes
## [1,] 0.2183908 0.2045455 0.2325581
## [2,] 0.2302158 0.2000000 0.2656250
## [3,] 0.2265193 0.2083333 0.2470588
## [4,] 0.2390244 0.1926606 0.2916667
## [5,] 0.2556054 0.2118644 0.3047619
## [6,] 0.2389381 0.1764706 0.3084112
```

We have the OOB error column as well as two other columns, which correspond to error rates specific to each value of the response. In this class we'll ignore the latter two and focus on the OOB error, which we can plot

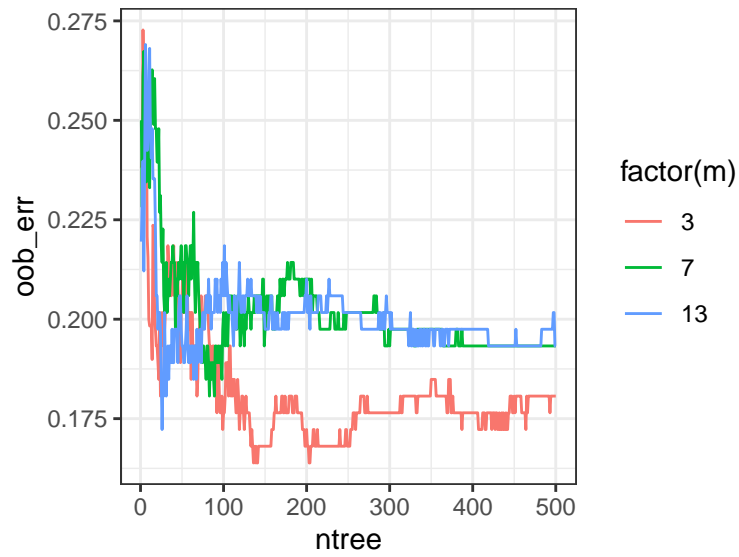
as follows:

```
tibble(  
  oob_error = rf_fit$err.rate[, "OOB"],  
  trees = 1:500  
) |>  
  ggplot(aes(x = trees, y = oob_error)) +  
  geom_line()
```



We can use the same parameters `ntree`, `mtry`, `nodesize`, and `maxnodes` as for regression random forests. For example, let's take a look at what happens when we vary `mtry`:

```
rf_3 <- randomForest(factor(AHD) ~ ., mtry = 3, data = heart_train)  
rf_7 <- randomForest(factor(AHD) ~ ., mtry = 7, data = heart_train)  
rf_13 <- randomForest(factor(AHD) ~ ., mtry = 13, data = heart_train)  
  
oob_errors <- bind_rows(  
  tibble(ntree = 1:500, oob_err = rf_3$err.rate[, "OOB"], m = 3),  
  tibble(ntree = 1:500, oob_err = rf_7$err.rate[, "OOB"], m = 7),  
  tibble(ntree = 1:500, oob_err = rf_13$err.rate[, "OOB"], m = 13)  
)  
  
oob_errors |>  
  ggplot(aes(x = ntree, y = oob_err, colour = factor(m))) +  
  geom_line() +  
  theme_bw()
```

We can make variable importance plots in the same way too:

```
rf_fit <- randomForest(factor(AHD) ~ ., importance = TRUE, data = heart_train)
varImpPlot(rf_fit)
```

rf_fit

