# Random Forest

### Tree

We only illustrate it via a classification tree. Much of the followings are also true for regression tree.

```
library(kernlab) # for the data spam

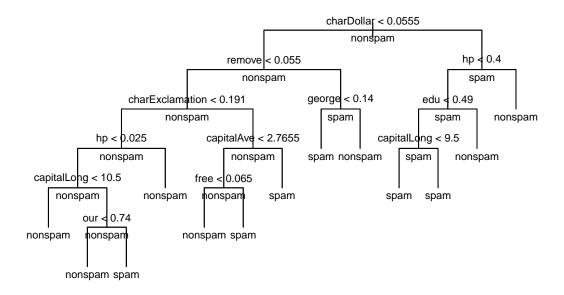
##
## Attaching package: 'kernlab'

## The following object is masked from 'package:purrr':
##
## cross

## The following object is masked from 'package:ggplot2':
##
## alpha

data(spam)

tree_spam <- tree(type ~ ., spam)
plot(tree_spam, type = "uniform")
text(tree_spam, pretty = 1, all = TRUE, cex = 0.7)</pre>
```



## Random forest

In R, there are packages like randomForest and ranger to perform random forest. However, we want to implement it from scratch.

The main idea bebind random forest is to resample the dataset by sampling the row with replacement and selecting the columns randomly.

For the spam data, suppose we want to predict the class of a random observation. (We randomly draw an observation from the original data and pretend that we do not know its class)

```
set.seed(141)
new_data <- spam %>% sample_n(1)
```

With the whole tree, we could do

```
predict(tree_spam, new_data)
```

```
## nonspam spam
## 1 0.812709 0.187291
```

that gives the probability of nonspam about 0.81.

However it is known that a single tree is not predictive. We want to use bootstrap to increase the predictivity.

```
r <- 1000 # in practise, we need a larger value, say 10000
m <- 8
n <- nrow(spam)
all_col_names <- names(spam)[1:57] # skip "type"

probs <- map_dbl(seq_len(r), function(i) {
   col_names <- c("type", sample(all_col_names, m))
   spam_boot <- spam[sample(n, n, replace = TRUE), col_names]
   tree_spam_boot <- tree(type ~ ., spam_boot)
   # we only need the probabliy of spam, because the sum of the two values is always 1
   predict(tree_spam_boot, new_data)[2]
})</pre>
```

There are two ways to yield the final predicted class, either by consensus or by averaging probablies. Either way, we need a baseline to compare with - using the prior proportion as the baseline is a simplest way (though may not be the best way). One may also use CV to select the baseline.

```
(baseline <- mean(spam$type == "spam"))
## [1] 0.3940448
```

#### Consensus

```
mean(probs > baseline)
```

## [1] 0.513

Since more than 50% of the trees predicted spam, by consensus, the predicted class for the new data is spam.

### By averaging probablies

```
mean(probs)
```

```
## [1] 0.4475824
```

The average probably acress all trees is 0.45 > baseline so the predicted class is "spam". For this new data, we have the same prediction using average probability.

In general, it is more stable to use average probability rather than consensus.

### Confidence interval

To construst a CI, you may be thinking of

```
quantile(probs, c(0.025, 0.975))
```

## 2.5% 97.5% ## 0.07925649 0.93187178

This confidence interval is essentially the bootstrap percentile interval for a tree model which randomly selects m predictors. It is not very reliable because the sampling of the columns introduces extra variability.

A correct way is to make use of the Jackknife, see https://arxiv.org/pdf/1311.4555.pdf

Well, it is too hard!? Use the package ranger!