

TMA4268 Statistical Learning V2018

Module 8: TREE-BASED METHODS

Mette Langaas, Thea Roksvåg, Julia Debik, Department of
Mathematical Sciences, NTNU

week 10, 2018 (version 06.03.2018)

Before we start

Learning material for this module:

- ▶ James et al (2013): An Introduction to Statistical Learning. Chapter 8.
- ▶ Classnotes 05.03.2018
- ▶ Classnotes 07.03.2018

Some of the figures in this presentation are taken (or are inspired) from “An Introduction to Statistical Learning, with applications in R” (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.

Topics

- ▶ Introduction
- ▶ Regression tree
- ▶ Classification tree
- ▶ Pruning
- ▶ Bagging
- ▶ Random forest
- ▶ Boosting
- ▶ Summing up
- ▶ Recommended exercises
- ▶ Packages to install before knitting this R Markdown file
- ▶ References

Introduction

What will you learn?

- ▶ Decision tree - idea and example
- ▶ Regression trees
 - ▶ what is a tree
 - ▶ how to grow a tree
- ▶ Classification trees - any changes to the above
- ▶ Pruning a tree
- ▶ Bagging
 - ▶ Variable importance plots
- ▶ Random forests
- ▶ Boosting

Running example: Detection of Minor Head Injury

Remark: this data is artificial.

Assume that you work as a statistician at a hospital, and the head of department asks you to develop a

simple method for detecting whether a patient is at risk for having a minor head injury.

The method should be easy to

- ▶ interpret for the medical personell that are not skilled in statistics, and
- ▶ the method should be fast, such that the medical personell quickly can identify a patient that needs treatment.

- ▶ The variable *clinically.important.brain.injury* will be the response of our model: It has value 1 if a person has an acute brain injury, and 0 otherwise.
- ▶ In this dataset, 250 (19%) of the patients have a clinically important brain injury.
- ▶ 10 variables are as explanatory variables. These variables describe the state of the patient:
 - ▶ Is he/she vomiting?
 - ▶ Is the Glasgow Coma Scale measure after 2 hours 15?
 - ▶ Has he/she an open skull fracture?
 - ▶ Has he/she had a loss of consciousness?
 - ▶ and so on.

These are questions the medical staff ask themselves when diagnosing the patient. Our job as a statistician is to systemize these questions in a good way so we can use them to estimate the probability of a brain injury.

This can be done by using tree-based methods.

The dataset includes data about 1321 patients and is a modified and smaller version of the (simulated) dataset *headInjury* from the *DAAG* library. We have made *age* instead of *age.65* giving us the “exact” age of the patient.

##	amnesia	bskullf	GCSdecr	GCS.13	GCS.15	risk	consc	oskullf	vomit
## 3	0	0	0	0	0	0	0	0	0
## 9	0	0	0	0	0	1	0	0	0
## 11	0	0	0	0	0	0	0	0	0
## 12	1	0	0	0	0	0	0	0	0
## 14	0	0	0	0	0	0	0	0	0
## 16	0	0	0	0	0	0	0	0	0

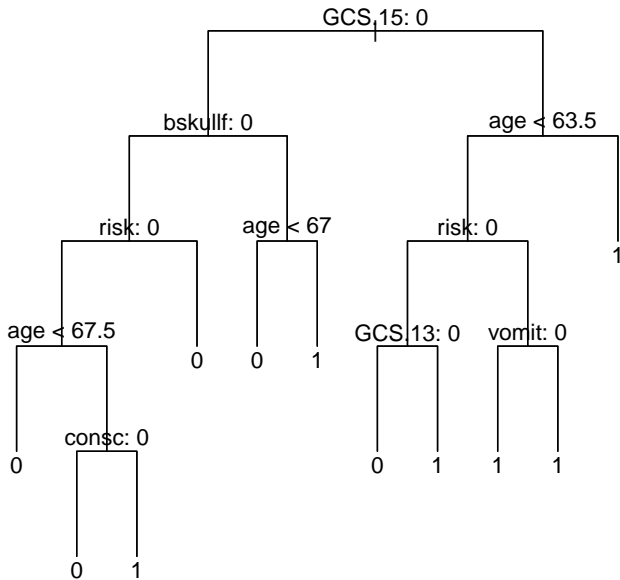
##	clinically.important.brain.injury	age
## 3	0	44
## 9	0	67
## 11	0	62
## 12	0	1
## 14	0	55
## 16	0	63

Main idea

The main idea behind tree-based methods is to

- ▶ derive a set of decision rules for segmenting the predictor space into a number of regions.
- ▶ In order to make a prediction for a new observation, we classify the observation into one of these regions by applying the decision rules.
- ▶ Then we typically use the mean or a majority vote of the training observations in this region as the prediction.

Below is a classification tree made from a training set of 850 randomly drawn observations.



```
## node), split, n, deviance, yval, (yprob)
##      * denotes terminal node
##
## 1) root 850 816.00 0 ( 0.8141 0.1859 )
##    2) GCS.15: 0 717 530.00 0 ( 0.8787 0.1213 )
##      4) bskullf: 0 664 407.00 0 ( 0.9081 0.0919 )
##        8) risk: 0 503 222.00 0 ( 0.9423 0.0577 )
##          16) age < 67.5 455 145.00 0 ( 0.9626 0.0374 ) *
##            17) age > 67.5 48  54.00 0 ( 0.7500 0.2500 )
##              34) consc: 0 42  37.80 0 ( 0.8333 0.1667 ) *
##                35) consc: 1 6   5.41 1 ( 0.1667 0.8333 ) *
##              9) risk: 1 161 161.00 0 ( 0.8012 0.1988 ) *
##            5) bskullf: 1 53  73.50 0 ( 0.5094 0.4906 )
##              10) age < 67 46  62.40 0 ( 0.5870 0.4130 ) *
##                11) age > 67 7   0.00 1 ( 0.0000 1.0000 ) *
##            3) GCS.15: 1 133 184.00 1 ( 0.4662 0.5338 )
##              6) age < 63.5 94 128.00 0 ( 0.5745 0.4255 )
##                12) risk: 0 62  72.80 0 ( 0.7258 0.2742 )
##                  24) GCS.13: 0 55  55.00 0 ( 0.8000 0.2000 ) *
##                    25) GCS.13: 1 7   5.74 1 ( 0.1429 0.8571 ) *
##                  13) risk: 1 32  38.00 1 ( 0.2812 0.7188 )
##                    26) vomit: 0 22  29.80 1 ( 0.4091 0.5909 ) *
##                      27) vomit: 1 10   0.00 1 ( 0.0000 1.0000 ) *
##                    7) age > 63.5 39  39.60 1 ( 0.2051 0.7949 ) *
```

By using simple decision rules related to the most important explanatory variables the medical staff can now assess the probability of a brain injury.

For example if the Glasgow Coma Scale of the patient is 15 ($GCS.15.2hours=1$) and the patient is older than 63 year then we should predict a minor injury according to the fitted tree, and the probability of that is estimated to be 0.7949 (node 7 in printout).

- ▶ The advantage of such decisions trees is that they are easier to interpret than many of the classification (and regression) methods that we have studied so far
- ▶ and they provide an easy way to visualize the data for non-statisticians.

Glossary

- ▶ Classification and regression trees are usually drawn upside down, where the top node is called the *root*.
- ▶ The *terminal nodes* or *leaf nodes* are the nodes at the bottom, with no splitting criteria. These represent the final predicted class (for classification trees) or predicted response value (for regression trees) and are written symbolically as R_j for $j = 1, 2, \dots, J$.
- ▶ *Internal nodes* are all nodes between the root and the terminal nodes. These nodes correspond to the partitions of the predictor space.
- ▶ *Branches*: segment of the tree connecting the nodes.

We will consider only binary splits on one variable, but multiway splits and linear combination of variables are possible - but not so common.

Constructing a decision tree

You can construct decision trees for both classification and regression problems. How do we construct a regression tree?

Regression tree

Assume that we have a dataset consisting of n pairs (\mathbf{x}_i, Y_i) , $i = 1, \dots, n$, and each predictor is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$.

Two steps:

1. Divide the predictor space into non-overlapping regions R_1, R_2, \dots, R_J .
2. For every observation that falls into region R_j we make the same prediction - which is the mean of the responses for the training observations that fall into R_j .

How to divide the predictor space R_1, R_2, \dots, R_J ?

Could try to minimize the RSS on the training set given by

$$\text{RSS} = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2,$$

where \hat{y}_{R_j} is the mean response for the training observations in region j . The mean \hat{y}_{R_j} is also the predicted value for a new observations that falls into region j .

To do this we need to consider every partition of the predictor space, and compute the RSS for each partition. This is computationally unfeasible (so no exhaustive search over possible trees).

Ripley (1996, p 216): Two types of optimality. a) Optimality of the partitioning of the predictor space : only feasible for small dimensions. b) Given partitioning of predictor space, how to represent this by a tree in the best possible way (=minimal expected number of tests) is a NP-complete problem.

A *greedy* approach is taken (aka top-down) - called *recursive binary splitting*.

Recursive binary splitting

We start at the top of the tree and divide the predictor space into two regions, R_1 and R_2 by making a decision rule for one of the predictors x_1, x_2, \dots, x_p . If we define the two regions by $R_1(j, s) = \{x | x_j < s\}$ and $R_2(j, s) = \{x | x_j \geq s\}$, it means that we need to find the (predictor) j and (splitting point) s that minimize

$$\sum_{i: x_i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2,$$

where \hat{y}_{R_1} and \hat{y}_{R_2} are the mean responses for the training observations in $R_1(j, s)$ and $R_2(j, s)$ respectively. This way we get the two first branches in our decision tree.

- ▶ We repeat the process to make branches further down in the tree.
- ▶ For every iteration we let each single split depend on only one of the predictors, giving us two new branches.
- ▶ This is done successively and in each step we choose the split that gives the best split at that particular step, i.e the split that gives the smallest RSS.
- ▶ We don't consider splits that further down the tree might give a tree with a lower overall RSS.

We continue splitting the predictor space until we reach some stopping criterion. For example we stop when no region contains more than 5 observations or when the reduction in the RSS is smaller than a specified limit.

Regression tree: ozone example

Consider the ozone data set from the `ElemStatLearn` library. The data set consists of 111 observations on the following variables:

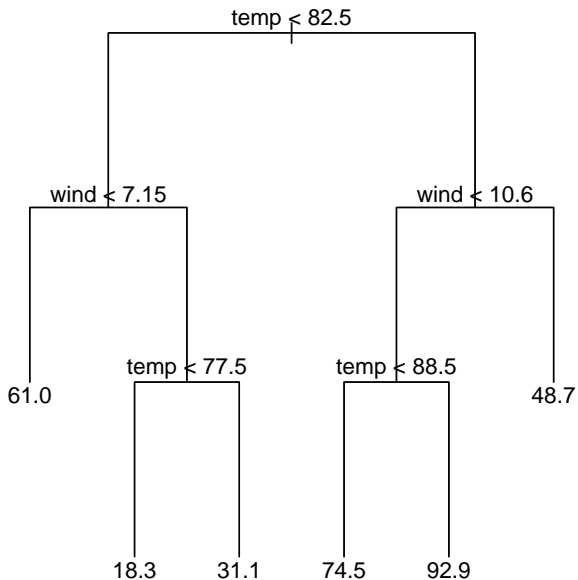
- ▶ `ozone` : the concentration of ozone in ppb
- ▶ `radiation`: the solar radiation (langleys)
- ▶ `temperature` : the daily maximum temperature in degrees F
- ▶ `wind` : wind speed in mph

Suppose we want to get an estimate of the ozone concentration based on the measurement of wind speed and the daily maximum temperature.

ozone	radiation	temperature	wind
41	190	67	7.4
36	118	72	8.0
12	149	74	12.6
18	313	62	11.5
23	299	65	8.6
19	99	59	13.8

We can fit a regression tree to the data, with ozone as our response variable and temperature and wind as predictors (not including radiation to make this easier to see). This gives us the following regression tree - what do you see?

```
ozone.tree = tree(ozone~temp+wind, data=myozone)
plot(ozone.tree,type="uniform")
text(ozone.tree)
```

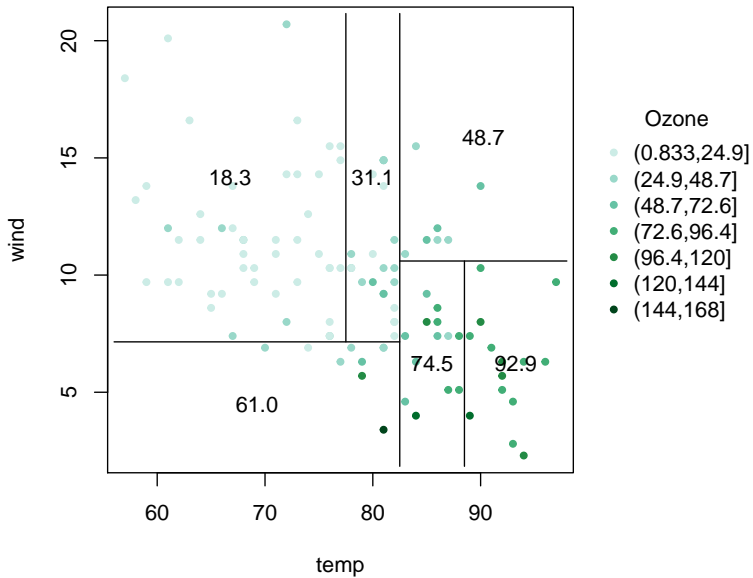


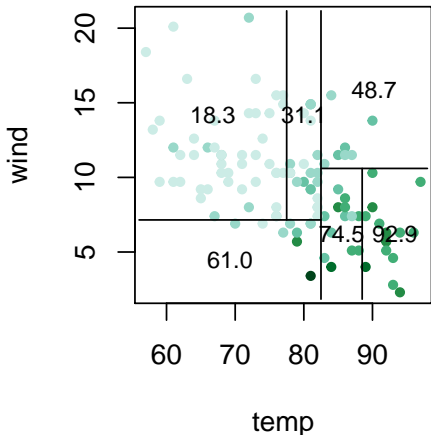
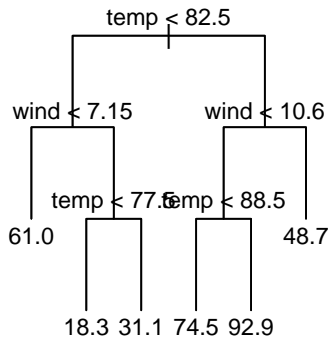
We see that temperature is the most important predictor for predicting the ozone concentration. Given a temperature below 82.5° F, and a low wind speed, the ozone level is high. For a low temperature and a high wind speed, the ozone level is lower. The highest ozone level is with a high temperature in combination with a wind speed below 10.6.

Now, we focus on the regions R_j , $j = 1, \dots, J$. What is J here?

Below we see the partition of the ozone data set from the `ozone.tree`, where the ozone levels have been color-coded, where a darker color corresponds to a higher ozone concentration.

We recognize the numbers from the leaf nodes on the figure. Each rectangle corresponds to one leaf node. Each number corresponding to a leaf node have been found by taking an average of all observations in the corresponding region (rectangle). Each line corresponds to an internal node, a binary partition of the predictor space.





- ▶ Explain the connection between the tree and the region plot.
- ▶ Why is recursive binary splitting classified as a greedy algorithm?
- ▶ Discuss the advantages and disadvantages of letting each single split depend on only one of the predictors.

Tree performance

To test the predictive performance of our regression tree, we can randomly divide our observations into a test and a training set (here 1/3 test).

```
set.seed(200)
ozone.trainID = sample(1:111, 75)
ozone.train = ozone[ozone.trainID, ]
ozone.test = ozone[-ozone.trainID,]

ozone.full = tree(ozone~temperature+wind,
                  data=ozone.train)

ozone.pred = predict(ozone.full, newdata=ozone.test)
ozone.MSE = mean((ozone.pred-ozone.test$ozone)^2)
ozone.MSE
```

```
## [1] 297.452
```

R: function tree in library tree

by Brian D. Ripley: Fit a Classification or Regression Tree

Description: A tree is grown by binary recursive partitioning using the response in the specified formula and choosing splits from the terms of the right-hand-side.

```
tree(formula, data, weights, subset, na.action = na.pass, control =  
tree.control(nobs, ...), method = "recursive.partition", split =  
c("deviance", "gini"), model = FALSE, x = FALSE, y = TRUE, wts  
= TRUE, ...)
```

- ▶ Details: A tree is grown by binary recursive partitioning using the response in the specified formula and choosing splits from the terms of the right-hand-side. Numeric variables are divided into $X < a$ and $X > a$; the levels of an unordered factor are divided into two non-empty groups. The split which maximizes the reduction in impurity is chosen, the data set split and the process repeated. Splitting continues until the terminal nodes are too small or too few to be split.
- ▶ Tree growth is limited to a depth of 31 by the use of integers to label nodes.
- ▶ Factor predictor variables can have up to 32 levels. This limit is imposed for ease of labelling, but since their use in a classification tree with three or more levels in a response involves a search over $2^{(k-1)} - 1$ groupings for k levels, the practical limit is much less.

Classification trees

remember the minor head injury example?

Let K be the number of classes for the response.

Building a decision tree in this setting is similar to building a regression tree for a quantitative response, but there are two main differences: *the prediction and the splitting criterion*

1) The prediction:

- ▶ In the regression case we use the mean value of the responses in R_j as a prediction for an observation that falls into region R_j .
- ▶ For the classification case however, we have two possibilities:
 - ▶ Majority vote: Predict that the observation belongs to the most commonly occurring class of the training observations in R_j .
 - ▶ Estimate the probability that an observation x_i belongs to a class k , $\hat{p}_{jk}(x_i)$, and then classify according to a threshold value. This estimated probability is the proportion of class k observations in region R_j with N_j observations:

$$\hat{p}_{jk} = \frac{1}{N_j} \sum_{i: x_i \in R_j} I(y_i = k).$$

2) The splitting criterion: We can not use RSS as a splitting criterion for a qualitative variable. Instead we can use some *measure of impurity* of the node.

Gini index:

$$G = \sum_{k=1}^K \hat{p}_{jk}(1 - \hat{p}_{jk}),$$

Cross entropy:

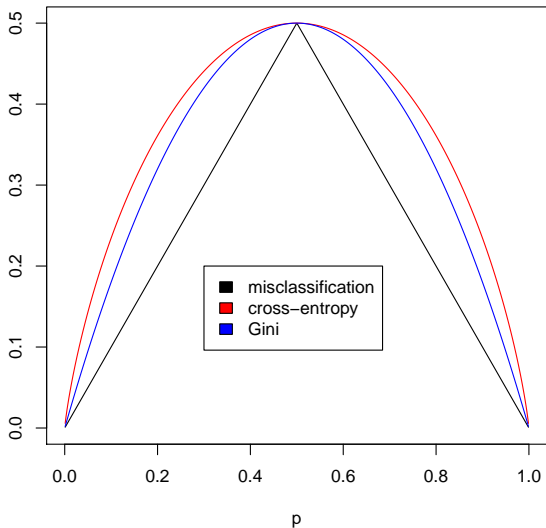
$$D = - \sum_{k=1}^K \hat{p}_{jk} \log \hat{p}_{jk}$$

Here \hat{p}_{jk} is the proportion of training observation in region j that are from class k .

When making a split in our classification tree, we want to minimize the Gini index or the cross-entropy.

⊕: Why these splitting criteria? Measure of impurity? How? See classnotes.

K=2



Minor head injury - continued

```
tree.HIClass=tree(clinically.important.brain.injury~.,  
                  data=headInjury2,  
                  subset=train,split="deviance")  
summary(tree.HIClass)
```

```
##
```

```
## Classification tree:
```

```
## tree(formula = clinically.important.brain.injury ~ ., data = headInj
```

```
##   subset = train, split = "deviance")
```

```
## Variables actually used in tree construction:
```

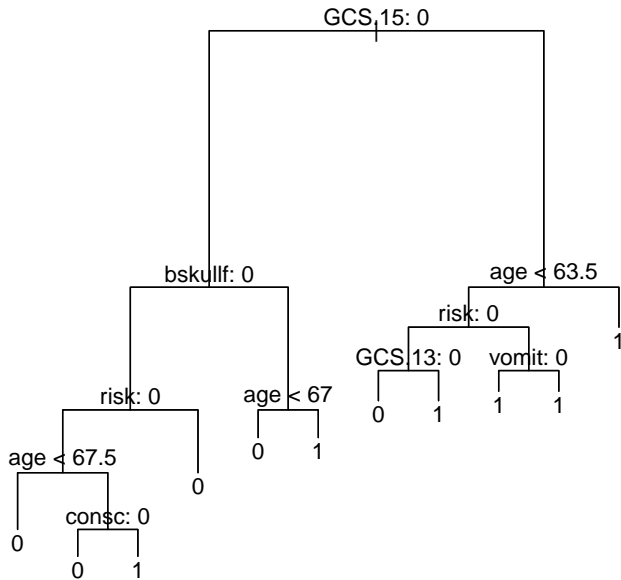
```
## [1] "GCS.15" "bskullf" "risk" "age" "consc" "GCS.13" "vom
```

```
## Number of terminal nodes: 11
```

```
## Residual mean deviance: 0.645 = 541 / 839
```

```
## Misclassification error rate: 0.124 = 105 / 850
```

$$\text{Deviance} = -2 \sum_j \sum_k n_{jk} \log(\hat{p}_{jk})$$



Length of branches are now proportional to the decrease in impurity.

Minor head injury - with Gini index

```
tree.HIClassG=tree(clinically.important.brain.injury~.,headInjury2,  
                  subset=train,split="gini")  
summary(tree.HIClassG)
```

```
##
```

```
## Classification tree:
```

```
## tree(formula = clinically.important.brain.injury ~ ., data = headInj
```

```
##   subset = train, split = "gini")
```

```
## Variables actually used in tree construction:
```

```
## [1] "GCS.15"  "bskullf" "age"      "risk"     "GCS.13"  "vomit"   "amn
```

```
## [8] "consc"   "oskullf"
```

```
## Number of terminal nodes: 78
```

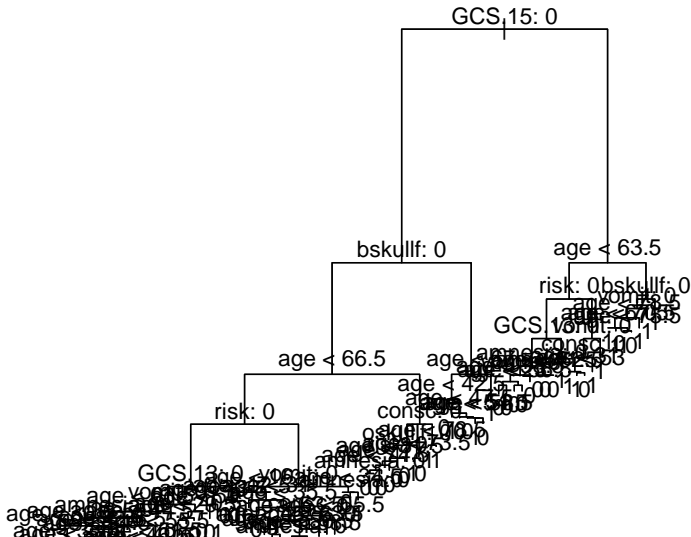
```
## Residual mean deviance: 0.488 = 377 / 772
```

```
## Misclassification error rate: 0.108 = 92 / 850
```

```

tree.HIClassG=tree(clinically.important.brain.injury~.,head
                    subset=train,split="gini")
plot(tree.HIClassG)
text(tree.HIClassG,pretty=1)

```



We also use the classification tree to predict the status of the patients in the test set (the one grown with deviance)

```
tree.pred=predict(tree.HIClass,headInjury2[test,],type="class")
misclass=table(tree.pred,headInjury2[test,]$clinically.important.brain)
print(misclass)
```

```
##
## tree.pred  0  1
##           0 361 50
##           1  18 42
```

The table shows how often the classification tree provides a correct prediction of the status of the patient. For example the number in row 1, column 2 is the number of times the tree predicts “0” when the real answer is “1”. The misclassification rate is given by

```
1-sum(diag(misclass))/(sum(misclass))
```

```
## [1] 0.144374
```

And for the Gini-grown tree

```
tree.predG=predict(tree.HIClassG,headInjury2[test,],type="class")
misclassG=table(tree.predG,headInjury2[test,]$clinically.important.brai)
print(misclassG)
1-sum(diag(misclassG))/(sum(misclassG))
```

```
##
## tree.predG    0    1
##           0 357  52
##           1  22  40
## [1] 0.157113
```

Questions:

- ▶ The classification tree has two terminal nodes with factor “1” originating from the same branch. Why do we get this “unnecessary” split?
- ▶ What if we have x_1 and x_2 and the true class boundary (two classes) is linear in x_1, x_2 space. How can we do that with our binary recursive splits?

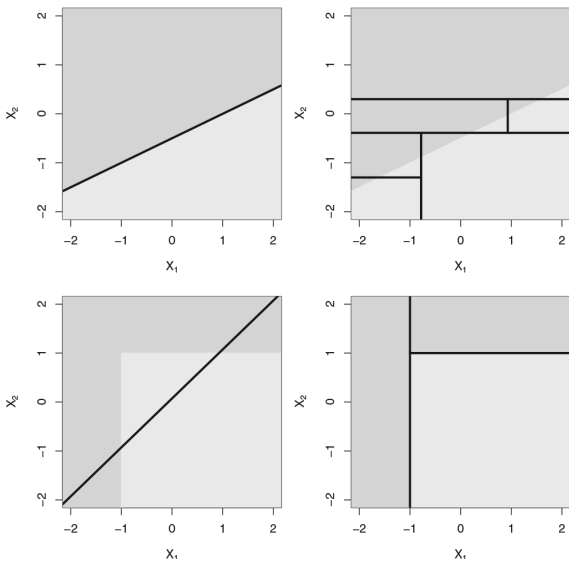


Figure 1: Linear boundary ISL Figure 8.7

- ▶ What about a rectangular boundary (figure above)?
- ▶ Study the above confusion matrices. One type of mistake is more severe than the other. Discuss if it is possible to change the algorithm in order to decrease the number of severe mistakes.

Pruning

Imagine that we have a data set with many predictors, and that we fit a large tree. Then, the number of observations from the training set that falls into some of the regions R_j may be small, and we may be concerned that we have overfitted the training data.

Pruning is a technique for solving this problem.

By *pruning* the tree we reduce the size or depth of the decision tree. When we reduce the number of terminal nodes and regions R_1, \dots, R_J , each region will probably contain more observations. This way we reduce the probability of overfitting, and we may get better predictions for test data.

If we have a large dataset with many explanatory variables and terminal nodes, we can also prune the tree if we want to create a simpler tree and increase the interpretability of the model.

In the classification tree (grown with deviance) we saw that we got several unnecessary splits. This is also something that can be avoided by pruning.

In the classification tree (grown with Gini index) there were 78 leaves - maybe we want a more easy interpretable tree?

But, there are many possible pruned versions of our full tree. Can we investigate all of these?

Cost complexity pruning

We can prune the tree by using an algorithm called *cost complexity pruning*. We first build a large tree T_0 by recursive binary splitting. Then we try to find a subtree $T \subset T_0$ that (for a given value of α) minimizes

$$C_\alpha(T) = Q(T) + \alpha|T|,$$

where $Q(T)$ is our cost function, $|T|$ is the number of terminal nodes in tree T . The parameter α is then a parameter penalizing the number of terminal nodes, ensuring that the tree does not get too many branches.

We proceed by repeating the process for the best subtree T , and this way we get a sequence of smaller and smaller subtrees where each tree is the best subtree of the previous tree.

For regression trees we choose $Q(T) = \sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2$, and for classification trees the entropy (deviance), Gini or misclassification rate.

Given a value of α we get a pruned tree (but the same pruned tree for ranges of α).

For $\alpha = 0$ we get T_0 and as α increases we get smaller and smaller trees.

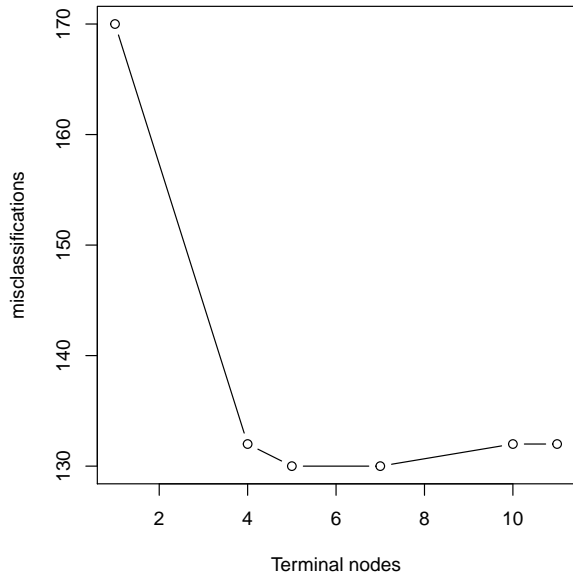
Please study this note from Bo Lindqvist in MA8702 - PhD version of Statistical Learning course for an example of how we perform cost complexity pruning in detail.

Building a regression (classification) tree: Algorithm 8.1

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of α .
3. Use K-fold cross-validation to choose α . That is, divide the training observations into K folds. For each $k = 1, \dots, K$:
 - ▶ Repeat Steps 1 and 2 on all but the kth fold of the training data.
 - ▶ Evaluate the mean squared prediction (misclassification, gini, cross-entropy) error on the data in the left-out kth fold, as a function of α .
 - ▶ Average the results for each value of α , and pick α to minimize the average error.
4. Return the subtree from Step 2 that corresponds to the chosen

Combining pruning and cross-validation to find optimal tree

We continue using the classification tree.



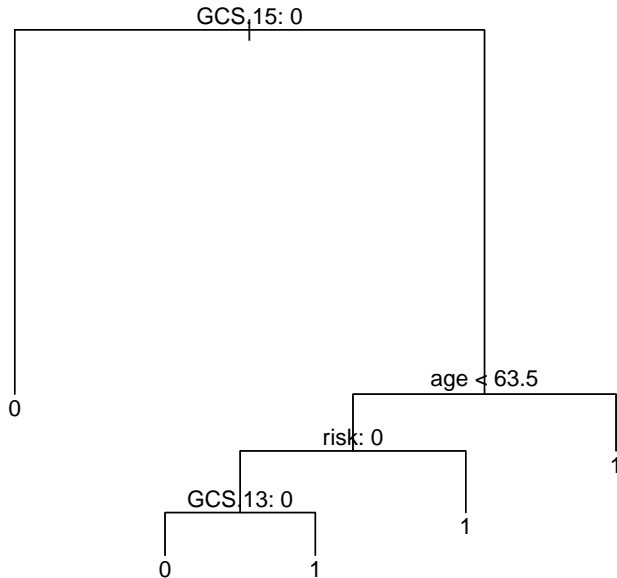
The function `cv.tree` automatically does 10-fold cross-validation. `dev` is here the number of misclassifications.

```
print(cv.head)
```

```
## $size
## [1] 11 10 7 5 4 1
##
## $dev
## [1] 132 132 130 130 132 170
##
## $k
## [1]      -Inf  0.00000  1.33333  3.50000  5.00000 12.33333
##
## $method
## [1] "misclass"
##
## attr(,"class")
## [1] "prune"          "tree.sequence"
```

We have done cross-validation on our training set of 850 observations. According to the plot, the number of misclassifications is lowest if we use 5 terminal nodes. Next, we prune the classification tree according to this value:

```
prune.HIClass=prune.misclass(tree.HIClass,best=5)  
#Five node tree.
```

We see that the new tree doesn't have any unnecessary splits, and we have a simple and interpretable decision tree. How is the predictive performance of the model affected?

```
tree.pred.prune=predict(prune.HIClass,headInjury2[test,],ty  
misclass.prune=table(tree.pred,headInjury2[test,]$clinical  
print(misclass.prune)
```

```
1-sum(diag(misclass.prune))/(sum(misclass.prune))
```

```
##
```

```
## tree.pred    0    1
```

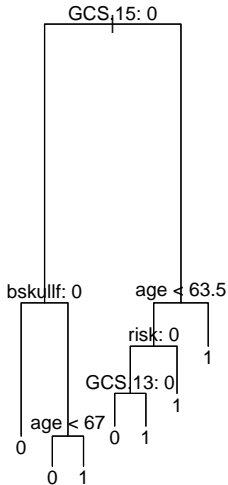
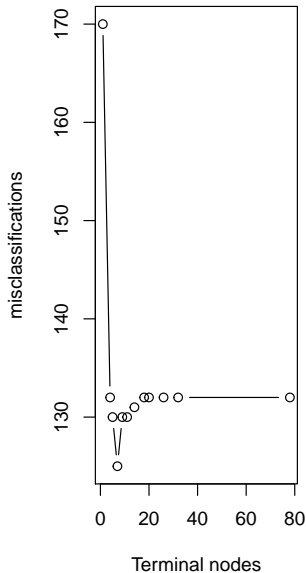
```
##           0 361  50
```

```
##           1  18  42
```

```
## [1] 0.144374
```

We see that the misclassification rate is as small as before indicating that the pruned tree is as good as the original tree for the test data.

The same repeated for the Gini-grown tree - comment on what is done.



```
print(cv.headG)
```

```
## $size
## [1] 78 32 26 20 18 14 11 9 7 5 4 1
##
## $dev
## [1] 132 132 132 132 132 131 130 130 125 130 132 170
##
## $k
## [1] -Inf 0.000000 0.166667 0.333333 0.500000
## [8] 1.500000 2.000000 3.500000 5.000000 12.333333
##
## $method
## [1] "misclass"
##
## attr(,"class")
## [1] "prune" "tree.sequence"
```

⊕ Questions

Discuss the bias-variance tradeoff of a regression tree when increasing/decreasing the number of terminal nodes, i.e:

- ▶ What happens to the bias?
- ▶ What happens to the variance of a prediction if we reduce the tree size?

From trees to forests

Advantages (+)

- ▶ Trees automatically select variables
- ▶ Tree-growing algorithms scale well to large n , growing a tree greedily
- ▶ Trees can handle mixed features (continuous, categorical) seamlessly, and can deal with missing data
- ▶ Small trees are easy to interpret and explain to people
- ▶ Some believe that decision trees mirror human decision making
- ▶ Trees can be displayed graphically

Disadvantages (-)

- ▶ Large trees are not easy to interpret
- ▶ Trees do not generally have good prediction performance (high variance)
- ▶ Trees are not very robust, a small change in the data may cause a large change in the final estimated tree

What is next?

- ▶ **Bagging**: grow many trees (from bootstrapped data) and average - to get rid of the non-robustness and high variance by averaging
- ▶ Variable importance plot - to see which variables make a difference (now that we have many trees).
- ▶ **Random forest**: inject more randomness (and even less variance) by just allowing a random selection of predictors to be used for the splits at each node.
- ▶ **Boosting**: make one tree, then another based on the residuals from the previous, repeat. The final predictor is a weighted sum of these trees.

But first,

Leo Breiman - the inventor of CART, bagging and random forests

Quotation from Wikipedia

Leo Breiman (January 27, 1928 – July 5, 2005) was a distinguished statistician at the University of California, Berkeley. He was the recipient of numerous honors and awards, and was a member of the United States National Academy of Science.

Breiman's work helped to bridge the gap between statistics and computer science, particularly in the field of machine learning. His most important contributions were his work on classification and regression trees and ensembles of trees fit to bootstrap samples. Bootstrap aggregation was given the name bagging by Breiman. Another of Breiman's ensemble approaches is the random forest.

From Breimans obituary

BERKELEY – Leo Breiman, professor emeritus of statistics at the University of California, Berkeley.

“It is trite to say so, but Leo Breiman was indeed a Renaissance man, and we shall miss him greatly,” said Peter Bickel, professor of statistics and chairman this summer of UC Berkeley’s statistics department.

Breiman retired in 1993, but as a Professor in the Graduate School, he continued to get substantial National Science Foundation grants and supervised three Ph.D. students. Bickel said that some of Breiman’s best work was done after retirement.

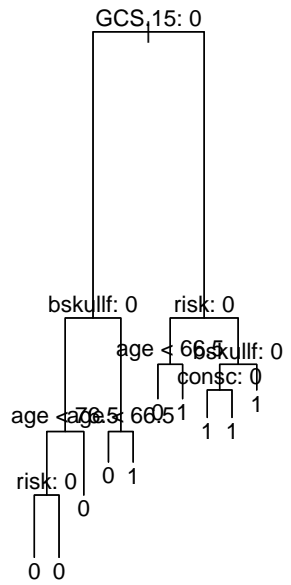
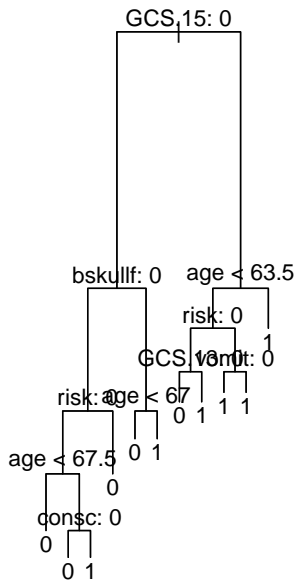
“In particular,” said Bickel, “he developed one of the most successful state-of-the-art classification programs, ‘Random Forest.’ This method was based on a series of new ideas that he developed in papers during the last seven years, and it is extensively used in government and industry.”

Breiman’s best known work is considered to be “Classification and Regression Trees,” a work in collaboration with three other scholars that facilitates practical applications, such as the diagnosis of diseases, from a multitude of symptoms.

Bagging

Decision trees often suffer from high variance. By this we mean that the trees are sensitive to small changes in the predictors: If we change the observation set, we may get a very different tree.

Let's draw a new training set for our data and see what happens if we fit our full classification tree (deviance grown).



This classification tree is constructed by using 850 observations, just like the tree in the classification trees section, but we get two different trees that will give different predictions for a test set.

To reduce the variance of decision trees we can apply *bootstrap aggregating (bagging)*, invented by Leo Breiman in 1996 (see references).

Independent data sets

Assume we have B i.i.d. observations of a random variable X each with the same mean and with variance σ^2 . We calculate the mean $\bar{X} = \frac{1}{B} \sum_{b=1}^B X_b$. The variance of the mean is

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{B} \sum_{b=1}^B X_b\right) = \frac{1}{B^2} \sum_{b=1}^B \text{Var}(X_b) = \frac{\sigma^2}{B}.$$

By averaging we get reduced variance. This is the basic idea!

But, we will not draw random variables - we want to fit decision trees: $\hat{f}_1(\mathbf{x}), \hat{f}_2(\mathbf{x}), \dots, \hat{f}_B(\mathbf{x})$ and average those.

$$\hat{f}_{avg}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(\mathbf{x})$$

However, we do not have many independent data set - so we use *bootstrapping* to construct B data sets.

Bootstrapping (from Module 5)

Problem: we want to draw samples from a population with distribution F .

But: we do not know F and do not have a population to draw from, we only have our one sample.

Solution: we may use our sample as an empirical estimate for the distribution F - by assuming that each sample point has probability $1/n$ for being drawn.

Therefore: we draw with replacement n observations from our sample - and that is our first *bootstrap sample*.

We repeat this B times and get B bootstrap samples - that we use as our B data sets.

Bootstrap samples and trees

For each bootstrap sample we construct a decision tree, $\hat{f}^{*b}(x)$ with $b = 1, \dots, B$, and we then use information from all of the trees to draw inference.

For a regression tree, we take the average of all of the predictions and use this as the final result:

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x).$$

For a classification tree we record the predicted class (for a given observation x) for each of the B trees and use the most occurring classification (majority vote) as the final prediction - or alternatively average posterior probabilities for each class.

Originally, Breiman (1996) suggested to prune each tree, but later research has found that it is better to leave the trees at maximal size (a bushy tree), to make the trees as different from each other as possible.

The number B is chosen to be as large as “necessary”. An increase in B will not lead to overfitting, and B is not regarded as a tuning parameter. If a goodness of fit measure is plotted as a function of B (soon) we see that (given that B is large enough) increasing B will not change the goodness of fit measure.

But first, a smart way to avoid doing cross-validation.

Out-of-bag error estimation

- ▶ We use a subset of the observations in each bootstrap sample. From Module 5 we know that the probability that an observation is in the bootstrap sample is approximately $1 - e^{-1} = 0.632121$ (approximately $2/3$).
- ▶ when an observation is left out of the bootstrap sample it is not used to build the tree, and we can use this observation as a part of a “test set” to measure the predictive performance and error of the fitted model, $f^{*b}(x)$.

An other words: Since each observation i has a probability of approximately $2/3$ to be in a bootstrap sample, and we make B bootstrap samples, then observation i will be outside the bootstrap sample in approximately $B/3$ of the fitted trees.

The observations left out are referred to as the *out-of-bag* observations, and the measured error of the $B/3$ predictions is called the *out-of-bag error*.

Example

We can do bagging by using the function `randomForest()` in the `randomForest` library.

```
library(randomForest)
set.seed(1)
bag=randomForest(clinically.important.brain.injury~.,
                 data=headInjury2,subset=train,
                 mtry=10,ntree=500,importance=TRUE)
bag$confusion
1-sum(diag(bag$confusion))/sum(bag$confusion[1:2,1:2])
```

```
##      0  1 class.error
## 0 642 50  0.0722543
## 1  82 76  0.5189873
## [1] 0.155294
```

The variable `mtry=10` because we want to consider all 10 predictors in each split of the tree. The variable `ntree = 500` because we want to average over 500 trees.

Predictive performance of the bagged tree on unseen test data:

```
yhat.bag=predict(bag,newdata=headInjury2[test,])
misclass.bag=table(yhat.bag,headInjury2[test,]$clinically.)
print(misclass.bag)
1-sum(diag(misclass.bag))/(sum(misclass.bag))
```

```
##
## yhat.bag    0    1
##           0 351  43
##           1  28  49
## [1] 0.150743
```

We note that the misclassification rate has increased slightly for the bagged tree (as compared to our previous full and pruned tree). **In other examples an improvement is very often seen.**

Prediction by majority vote vs. by averaging the probabilities

Consider the case when you have grown B classification tree with a binary response with classes 0 and 1. You might wonder which approach to choose to make a final prediction: majority vote or an average of the probabilities? Or would the prediction be the same in each case?

The difference between these two procedures can be compared to the difference between the mean value and median of a set of numbers. If we average the probabilities and make a classification thereafter, we have the mean value. If we sort all of our classifications, so that the classifications corresponding to one class would be lined up after each other, followed by the classifications corresponding to the other class we obtain the median value.

We examine this by an example:

Suppose we have $B = 5$ (no, B should be higher - this is only for illustration) classification tree and have obtained the following 5 estimated probabilities: $\{0.4, 0.4, 0.4, 0.4, 0.9\}$. If we average the probabilities, we get 0.5, and if we use a cut-off value of 0.5, our predicted class is 1. However, if we take a majority vote, using the same cut of value, the predicted classes will be $\{0, 0, 0, 0, 1\}$. The predicted class, based on a majority vote, would accordingly be 0.

The two procedures thus have their pros and cons: By averaging the predictions no information is lost. We do not only get the final classification, but the probability for belonging to the class 0 or 1. However, this method is not robust to outliers. By taking a majority vote, outliers have a smaller influence on the result.

When should we use bagging?

Bagging can be used for predictors (regression and classification) that are not trees, and according to Breiman (1996)

- ▶ the vital element is the instability of the prediction method
- ▶ if perturbing the learning set can cause significant changes in the predictor constructed, then bagging can improve accuracy.

Breiman (1996) suggests that these methods should be suitable for bagging:

- ▶ neural nets, classification and regression trees, subset selection in linear regression

however not nearest neighbours - since

- ▶ the stability of nearest neighbour classification methods with respect to perturbations of the data distinguishes them from competitors such as trees and neural nets.

Variable importance plots

Bagging is an example of an *ensemble method*, so is boosting and random forests (to come next). For all of these methods many trees are grown and combined, and the predictive power can be highly improved. However, this comes at a cost of interpretability. Instead of having one tree, the resulting model consists of B trees, where B often is 300 or 500 (or maybe even 5000 when boosting).

Variable importance plots show *the relative importance of the predictors*: the predictors are sorted according to their importance, such that the top variables have a higher importance than the bottom variables. There are in general two types of variable importance plots:

- ▶ variable importance based on decrease in node impurity and
- ▶ variable importance based on randomization.

Variable importance based on node impurity

The term *important* relates to *total decrease in the node impurity, over splits for a predictor*, and is defined differently for regression trees and classification trees.

Regression trees:

- ▶ The importance of each predictor is calculated using the RSS.
- ▶ The algorithm records the total amount that the RSS is decreased due to splits for each predictor (there may be many splits for one predictor for each tree).
- ▶ This decrease in RSS is then averaged over the B trees. The higher the decrease, the more important the predictor.

Classification trees:

- ▶ The importance of each predictor is calculated using the Gini index.
- ▶ The importance is the mean decrease (over all B trees) in the Gini index by splits of a predictor.

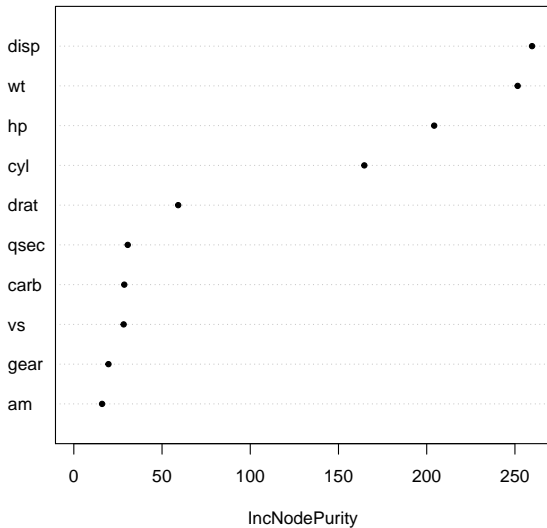
R: `varImpPlot` (or `importance`) in `randomForest` with `type=2`.

Auto data example

```
set.seed(4268)
data(mtcars)
mtcars.rf <- randomForest(mpg ~ ., data=mtcars, ntree=1000,
                          keep.forest=FALSE,
                          importance=TRUE)
```

```
varImpPlot(mtcars.rf, type=2, pch=20)
```

mtcars.rf



Variable importance based on randomization

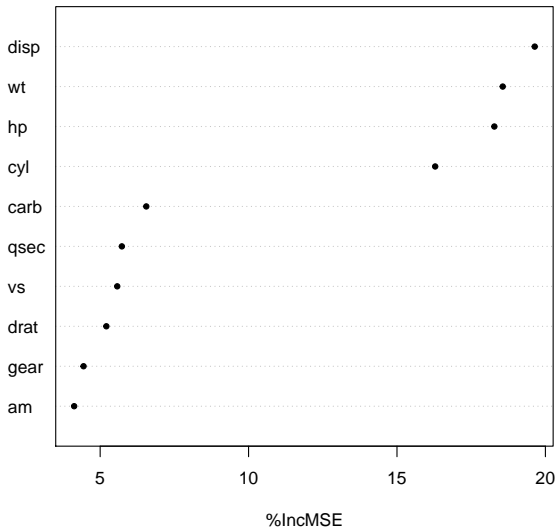
Variable importance based on randomization is calculated using the OOB sample.

- ▶ Computations are carried out for one bootstrap sample at a time.
- ▶ Each time a tree is grown the OOB sample is used to test the predictive power of the tree.
- ▶ Then for one predictor at a time, repeat the following:
 - ▶ permute the OOB observations for the j th variable x_j and calculate the new OOB error.
 - ▶ If x_j is important, permuting its observations will decrease the predictive performance.
- ▶ The difference between the two is averaged over all trees (and normalized by the standard deviation of the differences).

R: `varImpPlot` (or `importance`) in `randomForest` with `type=1`.

```
varImpPlot(mtcars.rf, type=1, pch=20)
```

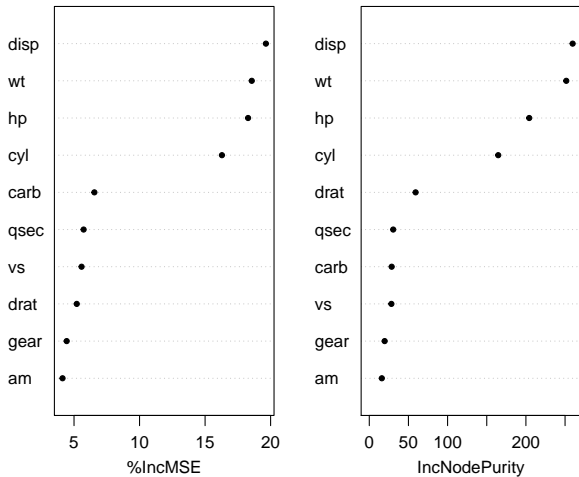
mtcars.rf



The two types together: do they agree?

```
varImpPlot(mtcars.rf, pch=20)
```

mtcars.rf



Random Forest

If there is a strong predictor in the dataset, the decision trees produced by each of the bootstrap samples in the bagging algorithm becomes very similar: Most of the trees will use the same strong predictor in the top split.

We have seen this for our example trees for the minor head injury example, the predictor *GCS.15.2hours* was chosen in the top split every time. This is probably the case for a large amount of the bagged trees as well.

This is not optimal, because we get B trees that are highly correlated. We don't get a large reduction in variance by averaging $\hat{f}^{*b}(x)$ when the correlation between the trees is high. In the previous section we actually saw a (marginal) decrease in the predictive performance for the bagged tree compared to the pruned tree and the full tree.

Random forests is a solution to this problem and a method for decorrelating the trees.

The effect of correlation on the variance of the mean

The variance of the average of B observations of i.i.d random variables X , each with variance σ^2 is $\frac{\sigma^2}{B}$. Now, suppose we have B observations of a random variable X which are identically distributed, each with mean μ and variance σ^2 , but not independent.

That is, suppose the variables have a positive correlation ρ

$$\text{Cov}(X_i, X_j) = \rho\sigma^2, \quad i \neq j.$$

The variance of the average is

$$\begin{aligned}\text{Var}(\bar{X}) &= \text{Var}\left(\frac{1}{B} \sum_{i=1}^B X_i\right) \\ &= \sum_{i=1}^B \frac{1}{B^2} \text{Var}(X_i) + 2 \sum_{i=2}^B \sum_{j=1}^{i-1} \frac{1}{B} \frac{1}{B} \text{Cov}(X_i, X_j) \\ &= \frac{1}{B} \sigma^2 + 2 \frac{B(B-1)}{2} \frac{1}{B^2} \rho \sigma^2 \\ &= \frac{1}{B} \sigma^2 + \rho \sigma^2 - \frac{1}{B} \rho \sigma^2 \\ &= \rho \sigma^2 + \frac{1-\rho}{B} \sigma^2 \\ &= \frac{1 - (1-B)\rho}{B} \sigma^2\end{aligned}$$

Check: $\rho = 0$ and $\rho = 1$? (Most negative values of ρ will not give a positive definite covariance matrix. The covariance matrix is positive definite if $\rho > -1/(B-1)$.)

The idea behind random forests is to *improve the variance reduction of bagging* by reducing the correlation between the trees.

The procedure is thus as in bagging, but with the important difference, that

- ▶ at each split we are only allowed to consider $m < p$ of the predictors.

A new sample of m predictors is taken at each split and

- ▶ typically $m \approx \sqrt{p}$ (classification) and $m = p/3$ (regression)

The general idea is that for very correlated predictors m is chosen to be small.

The number of trees, B , is not a tuning parameter, and the best is to choose it large enough.

If B is sufficiently large (three times the number needed for the random forest to stabilize), the OOB error estimate is equivalent to LOOCV (Efron and Hastie, 2016, p 330).

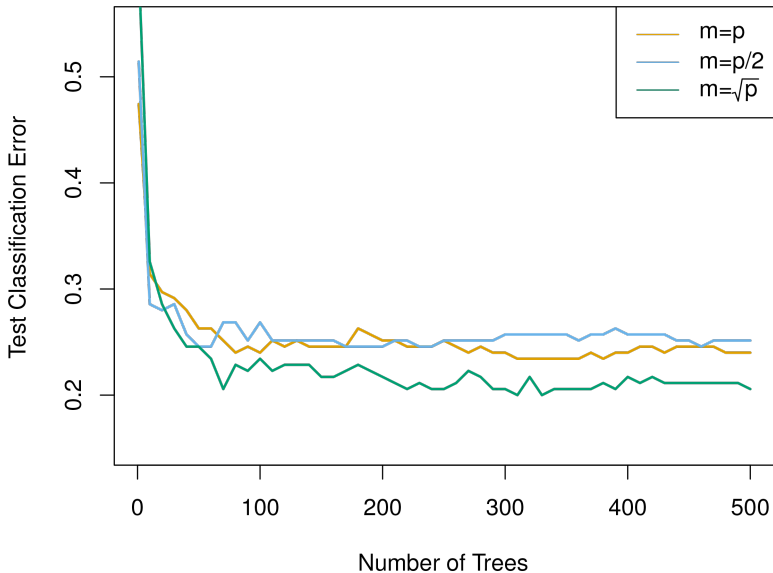


Figure 2: ISLR Figure 8.10, gene expression data set with 15 classes and 500 predictors

Example

We decorrelate the trees by using the *randomForest()* function again, but this time we set *mtry=3*. This means that the algorithm only considers three of the predictors in each split. We choose 3 because we have 10 predictors in total and $\sqrt{10} \approx 3$.

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

```
rf=randomForest(clinically.important.brain.injury~.,  
                data=headInjury2,subset=train,  
                mtry=3,ntree=500,importance=TRUE)
```

We check the predictive performance as before:

```
rf$confusion
1-sum(diag(rf$confusion[1:2,1:2]))/(sum(rf$confusion[1:2,1:2]))

yhat.rf=predict(rf,newdata=headInjury2[test,])

misclass.rf=table(yhat.rf,headInjury2[test,]$clinically.important.brain)
print(misclass.rf)

1-sum(diag(misclass.rf))/(sum(misclass.rf))
```

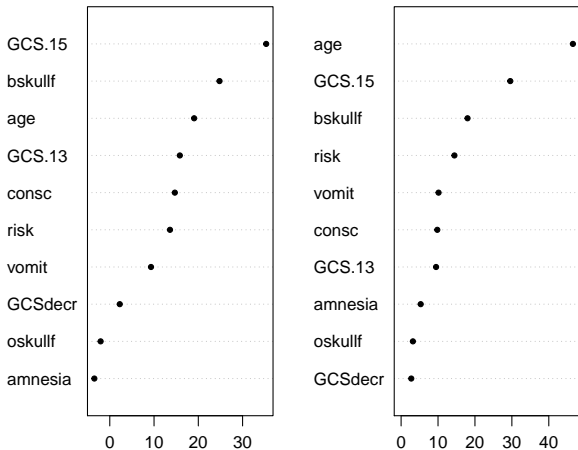
```
##      0  1 class.error
## 0 671 21  0.0303468
## 1  97 61  0.6139241
## [1] 0.138824
##
## yhat.rf   0   1
##          0 367 51
##          1  12 41
## [1] 0.133758
```

The misclassification rate is slightly decreased compared to the bagged tree (and to the pruned tree).

By using the `varImpPlot()` function we can study the importance of each predictor.

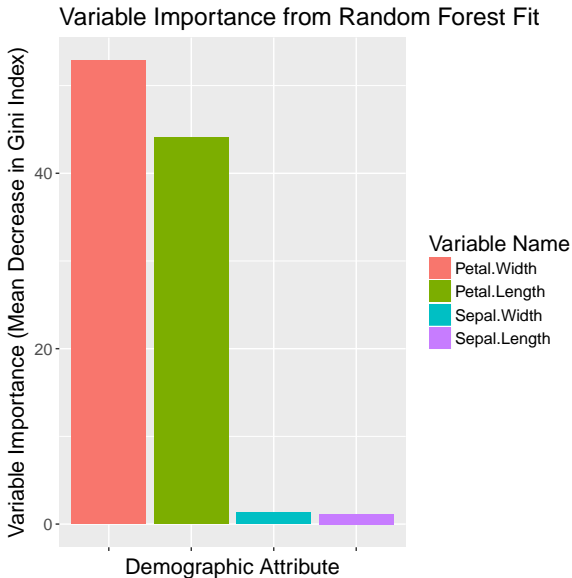
```
varImpPlot(rf, pch=20)
```

rf



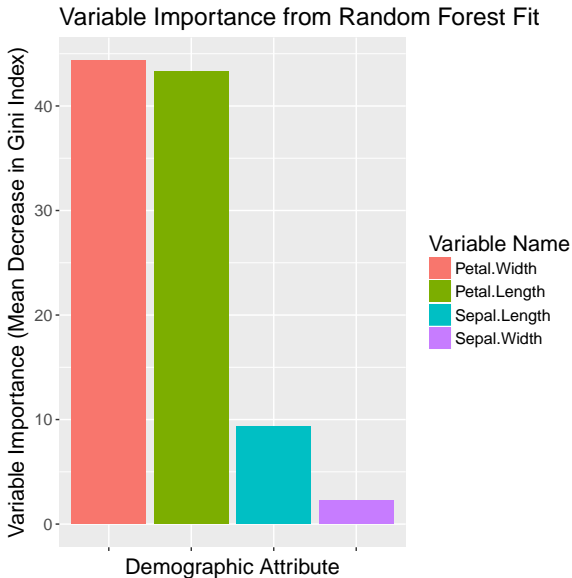
Iris example

Variable importance plot for bagging:



Iris example

Variable importance plot for random forest:



Boosting

Boosting is an alternative approach for improving the predictions resulting from a decision tree. We will only consider the description of boosting regression trees (and not classification trees) in this course.

In boosting the trees are grown *sequentially* so that each tree is grown using information from the previous tree.

- ▶ First build a decision tree with d splits (and $d + 1$ terminal nodes).
- ▶ Next, improve the model in areas where the model didn't perform well. This is done by fitting a decision tree to the *residuals of the model*. This procedure is called *learning slowly*.
- ▶ The first decision tree is then updated based on the residual tree, but with a weight.
- ▶ The procedure is repeated until some stopping criterion is reached. Each of the trees can be very small, with just a few terminal nodes (or just one split).

Algorithm 8.2: Boosting for regression trees

1. Set $\hat{f}(x) = 0$ and $r_i = y_i$ for all i in the training set.
2. For $b = 1, 2, \dots, B$, repeat:
 - 2.1 Fit a tree \hat{f}^b with d splits ($d + 1$ terminal nodes) to the training data.
 - 2.2 Update \hat{f} by adding in a shrunken version of the new tree:

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x).$$

- 2.3 Update the residuals,

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i).$$

3. The boosted model is $\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x)$.

Boosting has three tuning parameters which need to set, and can be found using cross-validation.

Tuning parameters

- ▶ The number of trees to be grown, B . The value of B could be chosen using cross-validation. A too small value of B would imply that much information is unused (remember that boosting is a slow learner), whereas a too large value of B may lead to overfitting.
- ▶ λ : This is a shrinkage parameter and controls the rate at which boosting learns. The role of λ is to scale the new information, when added to the existing tree. We add information from the b -th tree to our existing tree \hat{f} , but scaled by the λ . Choosing a small value for λ ensures that the algorithm learns slowly, but will require a larger tree ensemble. Typical values of λ is 0.1 or 0.01.
- ▶ Interaction depth d : The number of splits in each tree. This parameter controls the complexity of the boosted tree ensemble (the level of interaction between variables that we may estimate). By choosing $d = 1$ a tree stump will be fitted at each step and this gives an additive model.

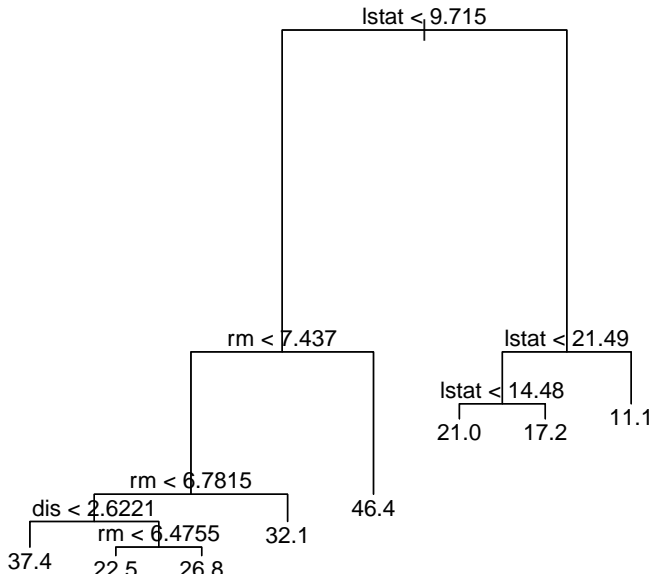
Example: Boston data set

First - to get acquainted with the data set we run through trees, bagging and random forests - before arriving at boosting. See also the ISLR book, Section 8.3.4.

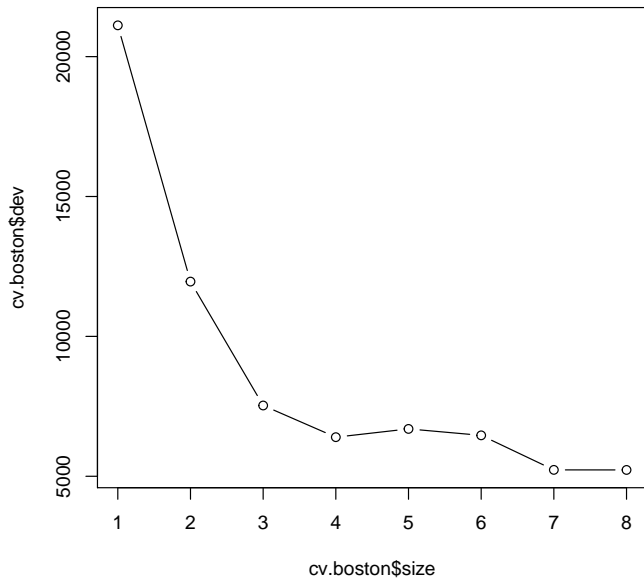
```
library(MASS)
library(tree)
set.seed(1)
train = sample(1:nrow(Boston), nrow(Boston)/2)
colnames(Boston)
head(Boston)
```

```
## [1] "crim"      "zn"        "indus"     "chas"      "nox"
## [8] "dis"       "rad"       "tax"       "ptratio"   "black"
##      crim zn indus chas nox  rm age  dis rad tax
## 1 0.00632 18  2.31   0 0.538 6.575 65.2 4.0900  1 296
## 2 0.02731  0  7.07   0 0.469 6.421 78.9 4.9671  2 242
## 3 0.02729  0  7.07   0 0.469 7.185 61.1 4.9671  2 242
## 4 0.03237  0  2.18   0 0.458 6.998 45.8 6.0622  3 222
```

```
tree.boston=tree(medv~.,Boston,subset=train)
summary(tree.boston); plot(tree.boston)
text(tree.boston,pretty=0)
```

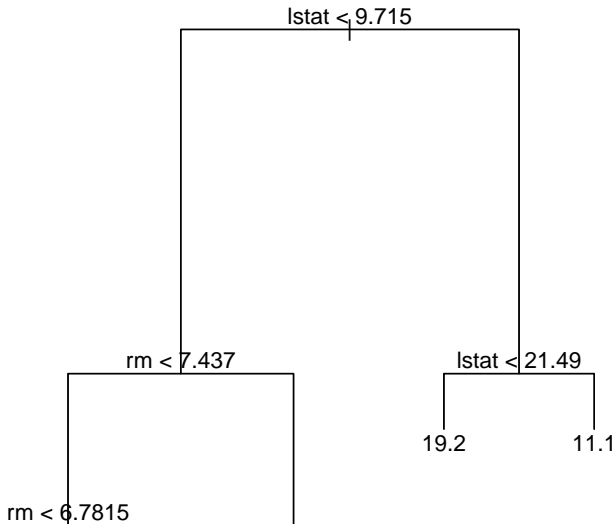


```
cv.boston=cv.tree(tree.boston)
plot(cv.boston$size,cv.boston$dev,type='b')
```

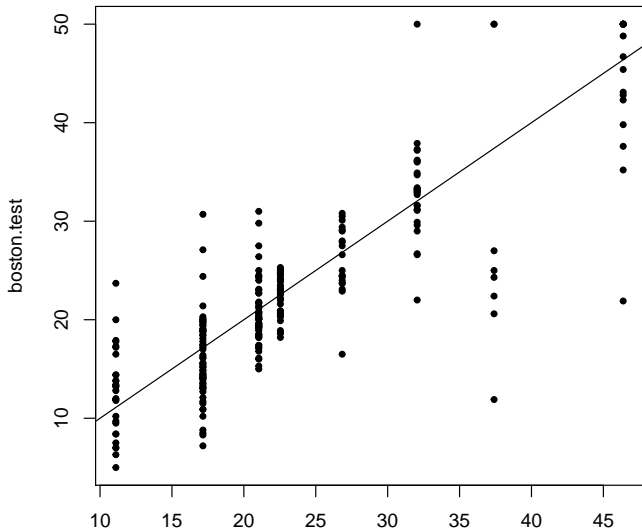


Just to show pruning (even if most complex tree was selected).

```
prune.boston=prune.tree(tree.boston,best=5)  
plot(prune.boston)  
text(prune.boston,pretty=0)
```



```
yhat=predict(tree.boston,newdata=Boston[-train,])  
boston.test=Boston[-train,"medv"]  
plot(yhat,boston.test, pch=20)  
abline(0,1)
```



```
library(randomForest)
set.seed(1)
bag.boston=randomForest(medv~.,data=Boston,subset=train,mtry=
bag.boston
```

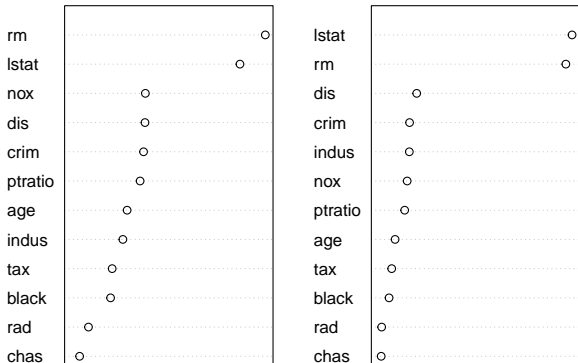
```
##
## Call:
## randomForest(formula = medv ~ ., data = Boston, mtry =
##           Type of random forest: regression
##           Number of trees: 500
## No. of variables tried at each split: 13
##
##           Mean of squared residuals: 11.0251
##           % Var explained: 86.65
```

```
yhat.bag = predict(bag.boston,newdata=Boston[-train,])
plot(yhat.bag, boston.test,pch=20)
abline(0,1)
```



```
set.seed(1)
rf.boston=randomForest(medv~.,data=Boston,subset=train,mtry
yhat.rf = predict(rf.boston,newdata=Boston[-train,])
mean((yhat.rf-boston.test)^2)
importance(rf.boston)
varImpPlot(rf.boston)
```

rf.boston



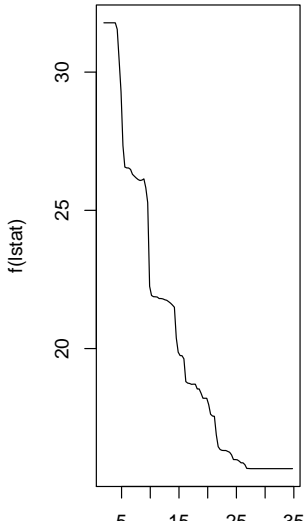
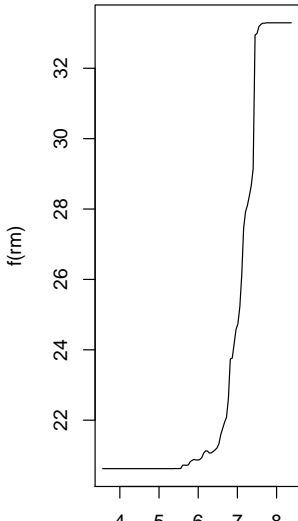
Finally, boosting Boston

```
library(gbm)
set.seed(1)
boost.boston=gbm(medv~.,data=Boston[train,],
                 distribution="gaussian",
                 n.trees=5000,interaction.depth=4)
summary(boost.boston,plotit=FALSE)
```

```
##           var      rel.inf
## lstat      lstat 45.9627334
## rm         rm    31.2238187
## dis        dis   6.8087398
## crim       crim   4.0743784
## nox        nox   2.5605001
## ptratio    ptratio 2.2748652
## black      black  1.7971159
## age        age   1.6488532
## tax        tax   1.3595005
```

Partial dependency plots - integrating out other variables

```
par(mfrow=c(1,2))  
plot(boost.boston,i="rm")  
plot(boost.boston,i="lstat")
```



Prediction on test set

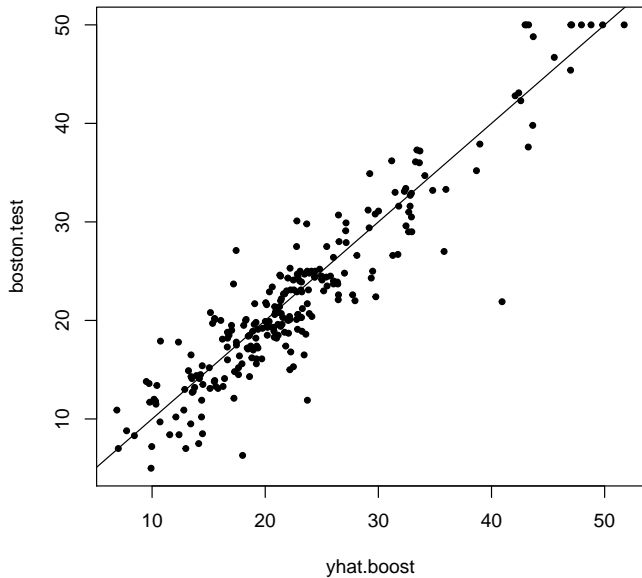
First for model with $\lambda = 0.001$ (default), then with $\lambda = 0.2$: MSE on test set. We could have done cross-validation to find the best λ over a grid.

```
yhat.boost=predict(boost.boston,newdata=Boston[-train,],n.trees=5000)
mean((yhat.boost-boston.test)^2)
boost.boston=gbm(medv~.,data=Boston[train,],distribution="gaussian",
                 n.trees=5000,interaction.depth=4,shrinkage=0.1)
yhat.boost=predict(boost.boston,newdata=Boston[-train,],n.trees=5000)
mean((yhat.boost-boston.test)^2)
```

```
## [1] 11.8443
```

```
## [1] 11.5111
```

```
plot(yhat.boost,boston.test,pch=20)  
abline(0,1)
```



Summing up

with a quiz on Module 8: Tree-based methods - also hosted in Kahoot!

Recommended exercises

1. Theoretical questions:

1. Show that each bootstrap sample will contain on average approximately $2/3$ of the observations.

2. Understanding the concepts and algorithms:

1. Do Exercise 1 in our book (page 331?)

Draw an example (of your own invention) of a partition of two-dimensional feature space that could result from recursive binary splitting. Your example should contain at least six regions. Draw a decision tree corresponding to this partition. Be sure to label all aspects of your figures, including the regions R_1, R_2, \dots , the cutpoints t_1, t_2, \dots , and so forth.

If the class border of the two dimensional space is linear, how can that be done with recursive binary splitting?

R packages

These packages needs to be install before knitting this R Markdown file.

```
install.packages("gamlss.data")
install.packages("tidyverse")
install.packages("GGally")
install.packages("Matrix")
install.packages("tree")
install.packages("randomForest")
install.packages("gbm")
```

References and further reading

- ▶ Videos on YouTube by the authors of ISL, Chapter 8, and corresponding slides
- ▶ Solutions to exercises in the book, chapter 8

Breiman, Leo. 1996. "Bagging Predictors." *Machine Learning* 24: 123–40.

———. 2001. "Random Forest." *Machine Learning* 45: 5–32.

Efron, Bradley, and Trevor Hastie. 2016. *Computer Age Statistical Inference - Algorithms, Evidence, and Data Science*. Cambridge University Press.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2001. *The Elements of Statistical Learning*. Vol. 1. Springer series in statistics New York.

James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2013. *An Introduction to Statistical Learning*. Vol. 112. Springer.

Rinley, Brian D. 1996. *Pattern Recognition and Neural Networks*