

Statistical Learning (5454) - Assignment 3

Matthias Hochholzer, Lukas Pirnbacher, Anne Valder

Due: 2024-05-20

Exercise 1

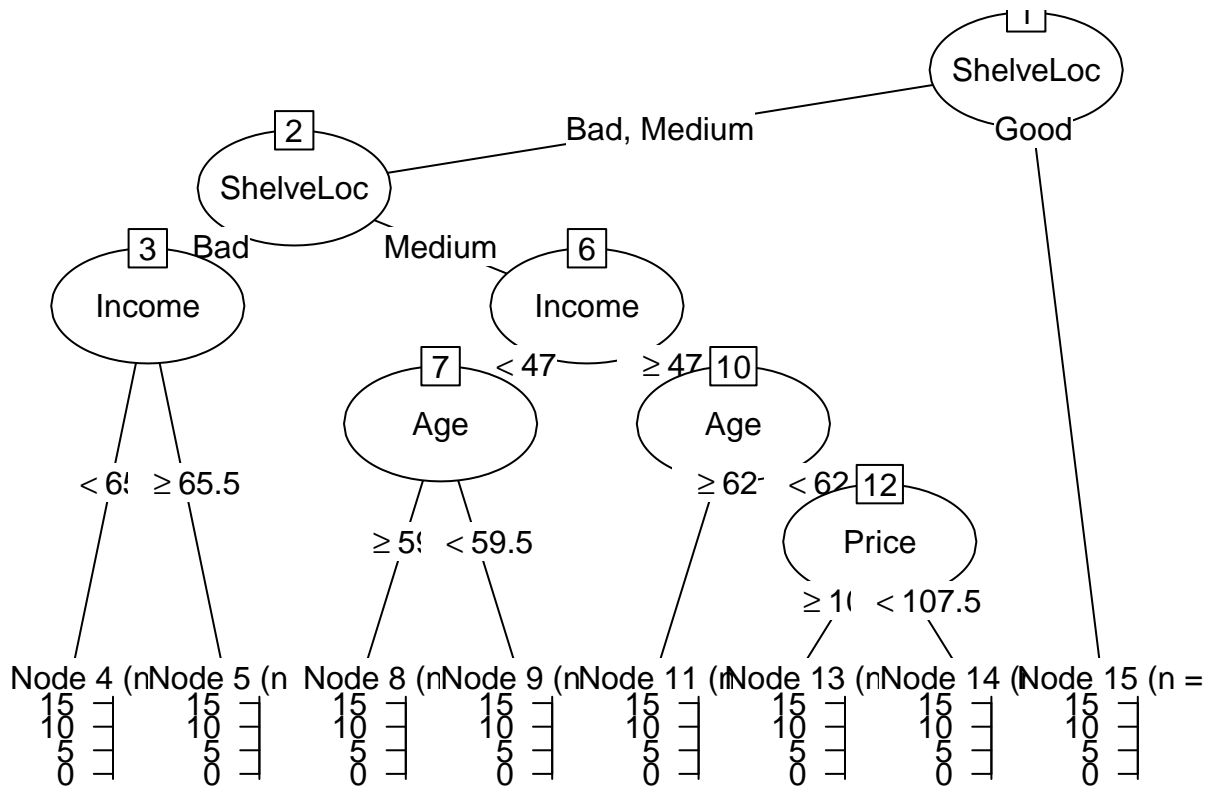
We load the data set *Carseats* from package *ISLR2*. We then split the data set into a training set and a test set.

We fit a regression tree to the training set.

```
##
## Regression tree:
## rpart(formula = Sales ~ ., data = train, method = "anova", parms = list(split = "gini"),
##       control = list(cp = 1e-04))
##
## Variables actually used in tree construction:
## [1] Age      Income   Price    ShelfLoc
##
## Root node error: 751/100 = 7.5
##
## n= 100
##
##      CP nsplit rel error xerror xstd
## 1 0.2574     0     1.00   1.02 0.15
## 2 0.1082     1     0.74   0.90 0.14
## 3 0.0657     2     0.63   0.84 0.12
## 4 0.0613     3     0.57   0.84 0.12
## 5 0.0325     4     0.51   0.77 0.11
## 6 0.0194     5     0.47   0.72 0.11
## 7 0.0163     6     0.46   0.73 0.12
## 8 0.0001     7     0.44   0.71 0.12
```

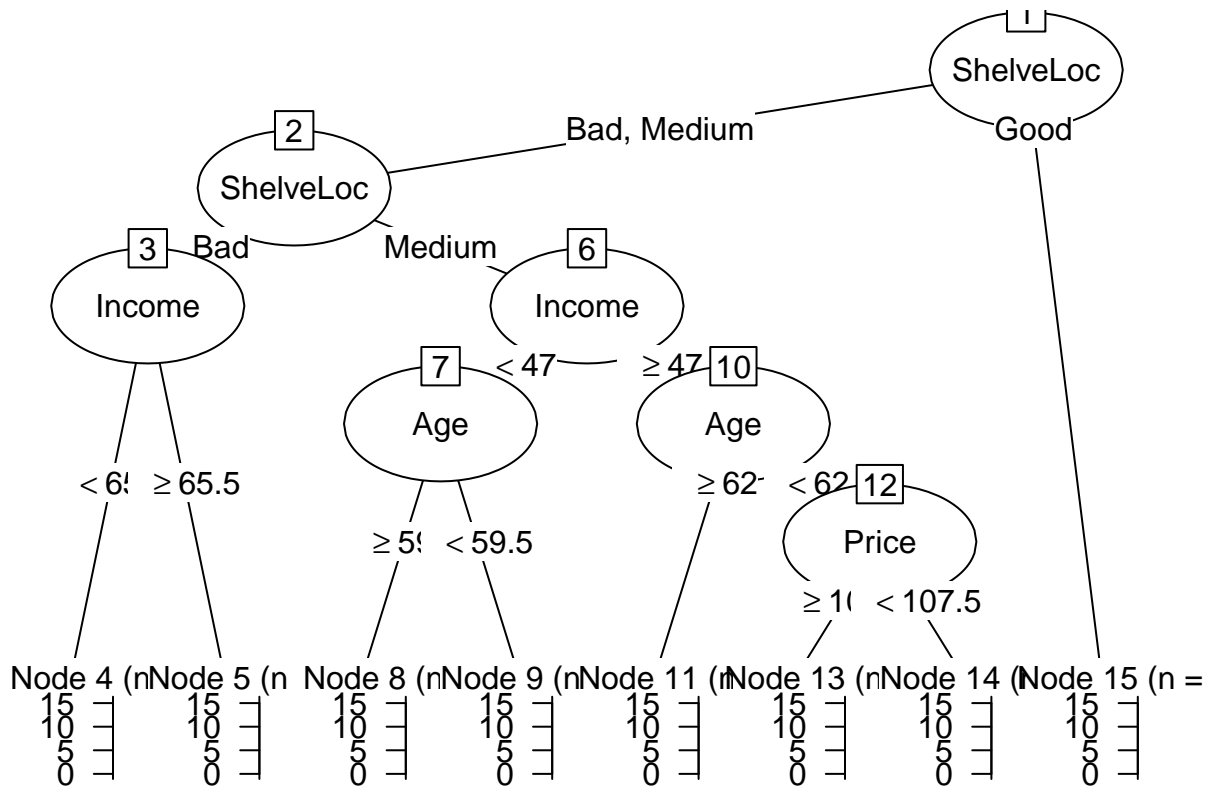
The plot of the tree is seen below.

```
## Loading required package: grid
## Loading required package: libcoin
## Loading required package: mvtnorm
```



The test MSE is

[1] 5.76



The test MSE with pruning is

[1] 6.198

Exercise 2

We draw 100 observations from four independent variables X_1, \dots, X_4 where

- X_1 follows a uniform distribution,
- X_2 follows a standard normal distribution,
- X_3 follows a Bernoulli distribution with success probability $\pi = 0.5$,
- X_4 follows a Bernoulli distribution with success probability $\pi = 0.1$.

We repeat 1000 times the following:

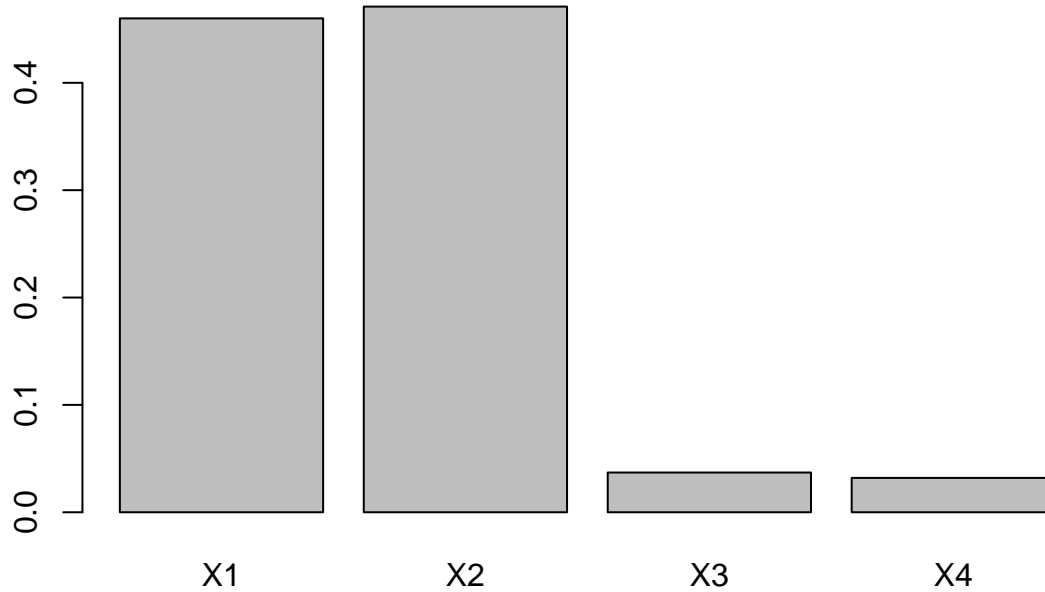
- Draw a dependent variable y from a standard normal distribution which is independent of the four independent variables.
- Fit a tree stump, i.e., a tree which contains only one split.
- Determine which variable was used for splitting.

Create the table of relative frequencies how often each of the variables was selected for splitting. Given that all independent variables are not associated with the dependent variable, is the probability of including them as a split variable the same? If not, why would they differ?

##

X1 X2 X3 X4

0.460 0.471 0.037 0.032



Exercise 3

Let's assume the following data generating process

$$Y = X + \epsilon,$$

with $X \sim N(0, 1)$ and $\epsilon \sim N(0, 1)$ independent. In addition 20 covariates Z_1, \dots, Z_{20} are given with

$$Z_i \sim \sqrt{0.9}X + \epsilon Z_i,$$

where $\epsilon Z_i \sim N(0, 0.1)$.

We draw a training data with 30 observations and a test data with 10,000 observations from the data generating process including the additional covariates \mathbf{Z} .

We sample 100 bootstrap samples of size 30 from the training data of the previous example by drawing with replacement.

Fit to each bootstrap sample:

1. a regression tree,
2. the null model with predicted value equal to the observed empirical mean of Y ,
3. a linear model including linear effects for X and all Z variables and
4. a linear model potentially including linear effects for X and all Z variables, but using model selection with the AIC to select a suitable model starting from the null model.

Determine the predicted values on the test data for the bagged model estimator by calculating the average predictions over the 100 trees fitted to the bootstrap samples, the 100 null models, the 100 linear models including all linear effects and the 100 linear models based on model selection.

```
##      Model      Mean
## 1   Tree 0.1912
## 2   Null 0.1847
## 3    LM 4.2230
## 4   AIC 0.3806
```

Last, we determine the mean squared error of the four bagged model estimators on the test sample of size 10,000.

```
##      Model      MSE
## 1   Tree 2.073
## 2   Null 2.070
## 3    LM 20.002
## 4   AIC 2.187
```

Exercise 4

We load the dataset `icu` in package **aplore3** which contains information on patients who were admitted to an adult intensive care unit (ICU). We develop a predictive model for the probability of survival to hospital discharge of these patients. To fit a predictive model to the data we use random forests.

We select a suitable number of bootstrap iterations.

Assess the influence of varying the hyperparameter `m` on the out-of-bag error obtained and select a suitable value.

Inspect the variable importance measures. Compare the mean decrease Gini and the mean decrease accuracy measures and assess if the observed differences in relative importance assigned might be related to the predictor variable being numeric or not.

Exercise 5

Assume that there are four predictor variables which have the following distributions:

$$X_1 \sim N(0, 1), \quad X_2 \sim U(0, 1), \quad (1)$$

$$X_3 \sim M(1, (0.5, 0.5)), \quad X_4 \sim M(1, (0.2, 0.2, 0.2, 0.2, 0.2)). \quad (2)$$

This means we have two continuous variables which follow either a standard normal or a standard uniform distribution ($U(0, 1)$) and two categorical variables with balanced categories with either 2 or 5 categories, i.e., $M(N, \pi)$ is the multinomial distribution for N trials and success probability vector π . The dependent variable y is assumed to be a binary categorical variable with equal-sized classes. The sample size is set to $N = 200$.

We generate 100 datasets for each setting and fit a random forest to each dataset and determine the mean decrease Gini and mean decrease accuracy values for each of the predictor variables.

```
## randomForest 4.7-1.1
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

Let's suitably visualize the results and interpret them.