

# DSE 6011 - Module 6 - Ch Exercises

Nathan Monges

2024-08-12

Lab 8.3.1 Fitting Classification Trees Lab 8.3.2 Fitting Regression Trees Chapter 8 Exercises: Problems 9a-f

```
library(tree)
library(ISLR2)
```

## 8.3.1 Fitting Classification Trees

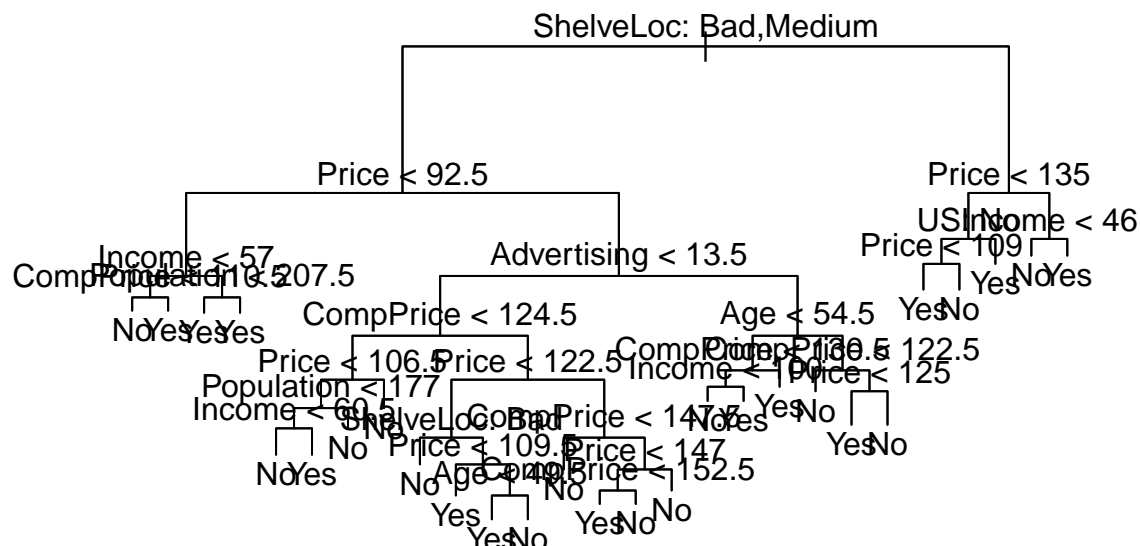
```
attach(Carseats)
```

```
High <- factor(ifelse(Sales <= 8, "No", "Yes"))
Carseats <- data.frame(Carseats, High)
```

```
tree.carseats <- tree(High ~ . - Sales, Carseats)
summary(tree.carseats)
```

```
##
## Classification tree:
## tree(formula = High ~ . - Sales, data = Carseats)
## Variables actually used in tree construction:
## [1] "ShelveLoc" "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

plot(tree.carseats)
text(tree.carseats, pretty = 0)
```



tree.carseats

```
## node), split, n, deviance, yval, (yprob)
##      * denotes terminal node
##
## 1) root 400 541.500 No ( 0.59000 0.41000 )
##    2) ShelfLoc: Bad,Medium 315 390.600 No ( 0.68889 0.31111 )
##      4) Price < 92.5 46 56.530 Yes ( 0.30435 0.69565 )
##        8) Income < 57 10 12.220 No ( 0.70000 0.30000 )
##          16) CompPrice < 110.5 5 0.000 No ( 1.00000 0.00000 ) *
##          17) CompPrice > 110.5 5 6.730 Yes ( 0.40000 0.60000 ) *
##          9) Income > 57 36 35.470 Yes ( 0.19444 0.80556 )
##            18) Population < 207.5 16 21.170 Yes ( 0.37500 0.62500 ) *
##            19) Population > 207.5 20 7.941 Yes ( 0.05000 0.95000 ) *
##          5) Price > 92.5 269 299.800 No ( 0.75465 0.24535 )
##            10) Advertising < 13.5 224 213.200 No ( 0.81696 0.18304 )
##              20) CompPrice < 124.5 96 44.890 No ( 0.93750 0.06250 )
##                40) Price < 106.5 38 33.150 No ( 0.84211 0.15789 )
##                  80) Population < 177 12 16.300 No ( 0.58333 0.41667 )
##                    160) Income < 60.5 6 0.000 No ( 1.00000 0.00000 ) *
##                    161) Income > 60.5 6 5.407 Yes ( 0.16667 0.83333 ) *
##                  81) Population > 177 26 8.477 No ( 0.96154 0.03846 ) *
##                41) Price > 106.5 58 0.000 No ( 1.00000 0.00000 ) *
##              21) CompPrice > 124.5 128 150.200 No ( 0.72656 0.27344 )
##                42) Price < 122.5 51 70.680 Yes ( 0.49020 0.50980 )
##                  84) ShelfLoc: Bad 11 6.702 No ( 0.90909 0.09091 ) *
##                  85) ShelfLoc: Medium 40 52.930 Yes ( 0.37500 0.62500 )
##                    170) Price < 109.5 16 7.481 Yes ( 0.06250 0.93750 ) *
##                    171) Price > 109.5 24 32.600 No ( 0.58333 0.41667 )
##                      342) Age < 49.5 13 16.050 Yes ( 0.30769 0.69231 ) *
##                      343) Age > 49.5 11 6.702 No ( 0.90909 0.09091 ) *
##                  43) Price > 122.5 77 55.540 No ( 0.88312 0.11688 )
##                    86) CompPrice < 147.5 58 17.400 No ( 0.96552 0.03448 ) *
##                    87) CompPrice > 147.5 19 25.010 No ( 0.63158 0.36842 )
##                      174) Price < 147 12 16.300 Yes ( 0.41667 0.58333 )
##                        348) CompPrice < 152.5 7 5.742 Yes ( 0.14286 0.85714 ) *
##                        349) CompPrice > 152.5 5 5.004 No ( 0.80000 0.20000 ) *
##                      175) Price > 147 7 0.000 No ( 1.00000 0.00000 ) *
##            11) Advertising > 13.5 45 61.830 Yes ( 0.44444 0.55556 )
##              22) Age < 54.5 25 25.020 Yes ( 0.20000 0.80000 )
##                44) CompPrice < 130.5 14 18.250 Yes ( 0.35714 0.64286 )
##                  88) Income < 100 9 12.370 No ( 0.55556 0.44444 ) *
##                  89) Income > 100 5 0.000 Yes ( 0.00000 1.00000 ) *
##                45) CompPrice > 130.5 11 0.000 Yes ( 0.00000 1.00000 ) *
##              23) Age > 54.5 20 22.490 No ( 0.75000 0.25000 )
##                46) CompPrice < 122.5 10 0.000 No ( 1.00000 0.00000 ) *
##                47) CompPrice > 122.5 10 13.860 No ( 0.50000 0.50000 )
##                  94) Price < 125 5 0.000 Yes ( 0.00000 1.00000 ) *
##                  95) Price > 125 5 0.000 No ( 1.00000 0.00000 ) *
##    3) ShelfLoc: Good 85 90.330 Yes ( 0.22353 0.77647 )
##      6) Price < 135 68 49.260 Yes ( 0.11765 0.88235 )
##        12) US: No 17 22.070 Yes ( 0.35294 0.64706 )
##          24) Price < 109 8 0.000 Yes ( 0.00000 1.00000 ) *
##          25) Price > 109 9 11.460 No ( 0.66667 0.33333 ) *
```

```
##      13) US: Yes 51 16.880 Yes ( 0.03922 0.96078 ) *
##      7) Price > 135 17 22.070 No ( 0.64706 0.35294 )
##      14) Income < 46 6 0.000 No ( 1.00000 0.00000 ) *
##      15) Income > 46 11 15.160 Yes ( 0.45455 0.54545 ) *
```

```
set.seed(2)

train <- sample(1:nrow(Carseats), 200)
carseats.test <- Carseats[-train,]
high.test <- High[-train]
tree.carseats <- tree(High ~ . - Sales, Carseats,
                      subset = train)

tree.pred <- predict(tree.carseats, carseats.test,
                    type = "class")

table(tree.pred, high.test)
```

```
##      high.test
## tree.pred No Yes
##      No 104 33
##      Yes 13 50
(104 + 50) / 200
```

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

```
set.seed(7)
cv.carseats <- cv.tree(tree.carseats, FUN = prune.misclass)

names(cv.carseats)
```

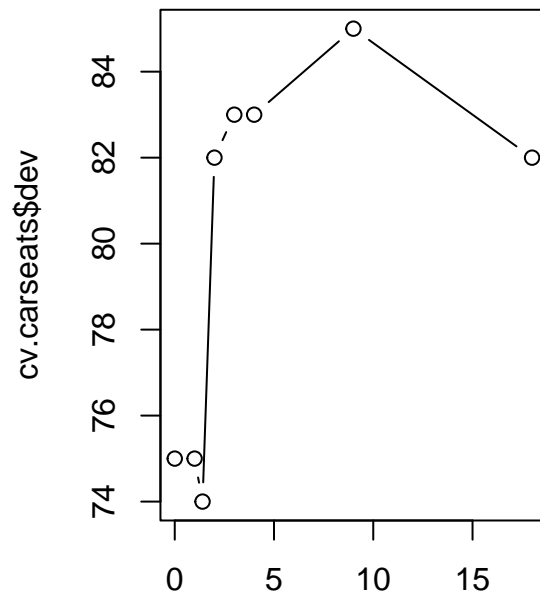
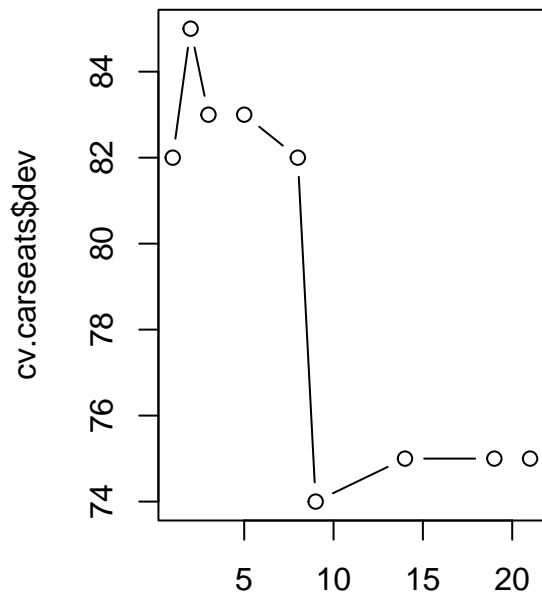
```
## [1] "size" "dev" "k" "method"
```

```
cv.carseats
```

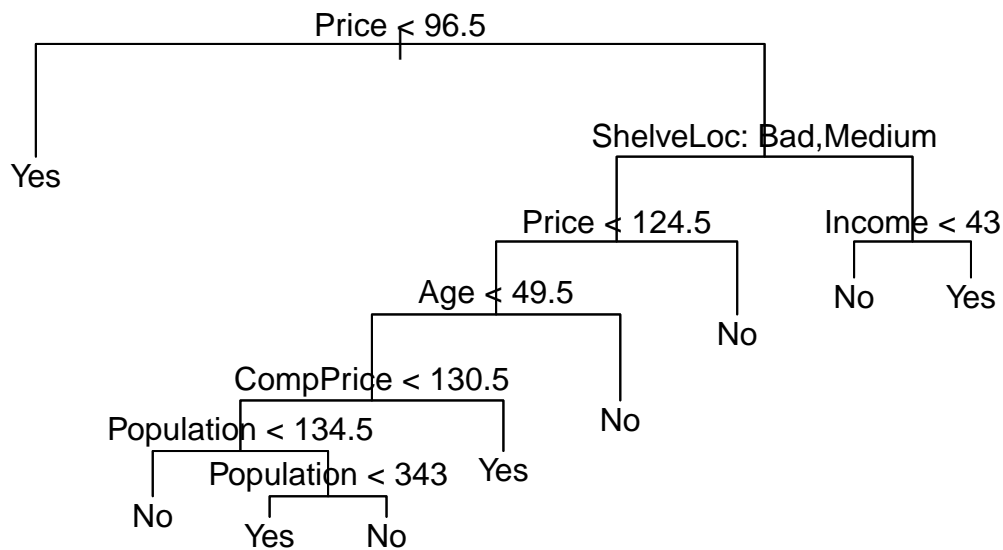
```
## $size
## [1] 21 19 14 9 8 5 3 2 1
##
## $dev
## [1] 75 75 75 74 82 83 83 85 82
##
## $k
## [1] -Inf 0.0 1.0 1.4 2.0 3.0 4.0 9.0 18.0
##
## $method
## [1] "misclass"
##
## attr("class")
## [1] "prune" "tree.sequence"
```

```
par(mfrow = c(1,2))
```

```
plot(cv.carseats$size, cv.carseats$dev, type = "b")
plot(cv.carseats$k, cv.carseats$dev, type = "b")
```



```
prune.carseats <- prune.misclass(tree.carseats, best = 9)
plot(prune.carseats)
text(prune.carseats, pretty = 0)
```



```
tree.pred <- predict(prune.carseats, carseats.test, type = "class")
table(tree.pred, high.test)
```

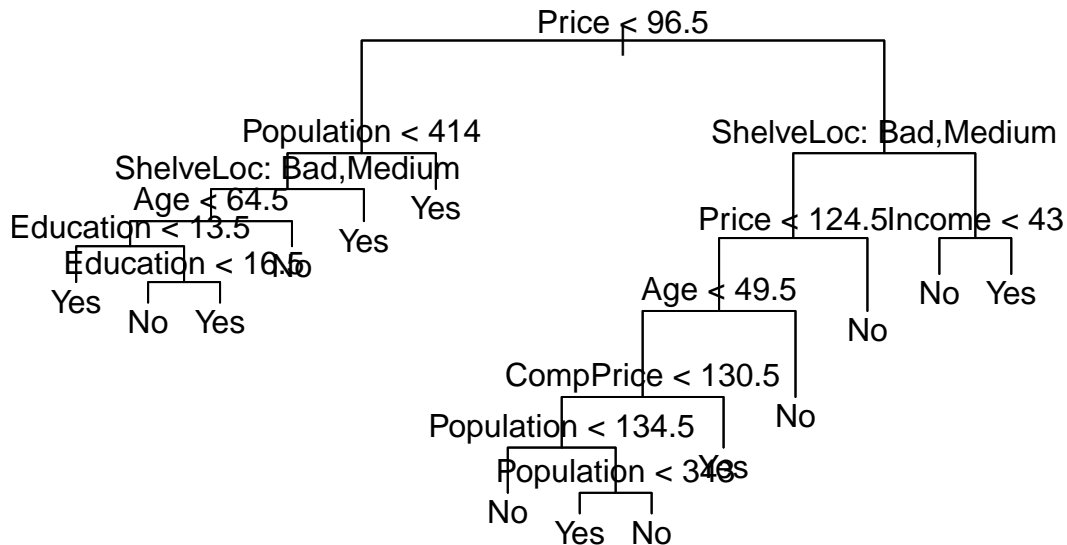
```
##           high.test
## tree.pred No  Yes
##           No  97  25
##           Yes 20  58
```

$$(97 + 58) / 200$$

```
## [1] 0.775
```

```
prune.carseats <- prune.misclass(tree.carseats, best = 14)
```

```
plot(prune.carseats)
text(prune.carseats, pretty = 0)
```



```
tree.pred <- predict(prune.carseats, carseats.test, type = "class")
table(tree.pred, high.test)
```

```
##           high.test
## tree.pred  No Yes
##           No 102 31
##           Yes  15 52
```

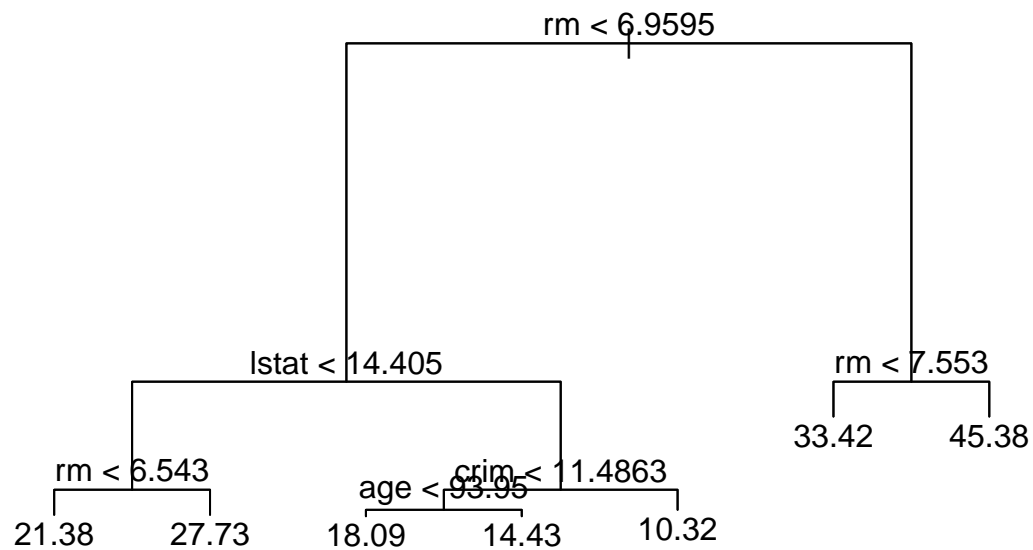
### 8.3.2 Fitting Regression Trees

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

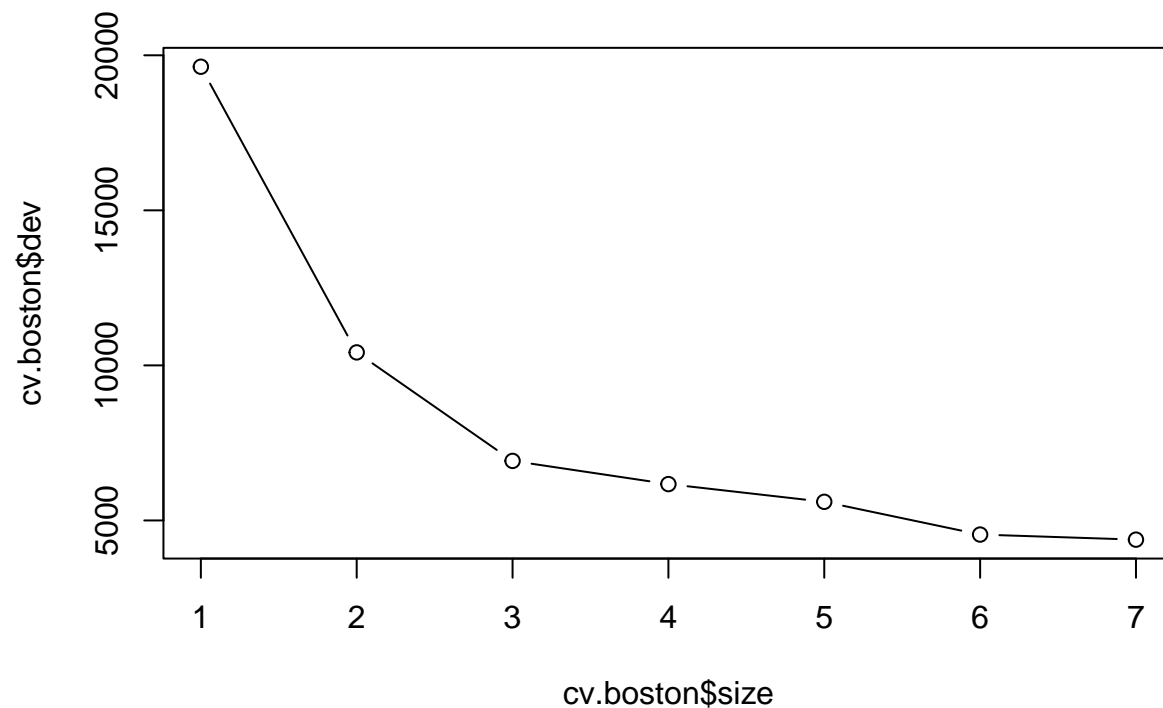
```
train <- sample(1:nrow(Boston), nrow(Boston) / 2)
tree.boston <- tree(medv ~ ., Boston, subset = train)
summary(tree.boston)
```

```
##
## Regression tree:
## tree(formula = medv ~ ., data = Boston, subset = 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
```

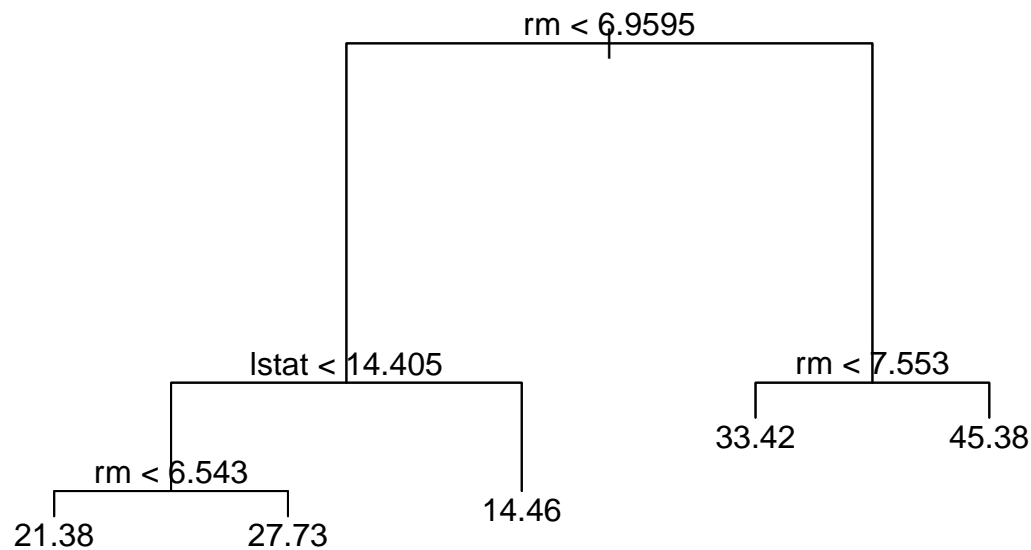
```
plot(tree.boston)
text(tree.boston, pretty = 0)
```



```
cv.boston <- cv.tree(tree.boston)
plot(cv.boston$size, cv.boston$dev, type = "b")
```



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



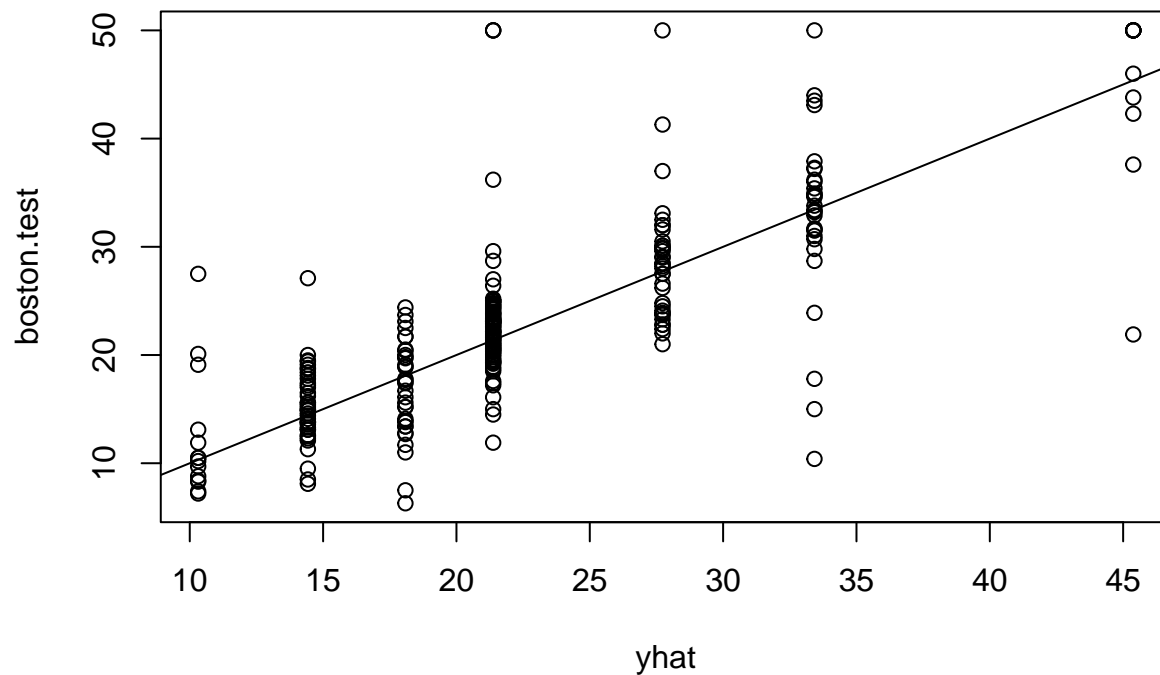
```

yhat <- predict(tree.boston, newdata = Boston[-train, ])

boston.test <- Boston[-train, "medv"]

plot(yhat, boston.test)
abline(0, 1)

```



```
mean((yhat - boston.test)^2)
```

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

## Exercise 9

- Create a training set containing a random sample of 800 observations, and a test set containing the remaining observations.

```
attach(OJ)

set.seed(1)

train <- sample(1:nrow(OJ), 800)
oj_train <- OJ[train, ]

oj_test <- OJ[-train, ]
```

- b) Fit a tree to the training data, with Purchase as the response and the other variables as predictors. Use the summary() function to produce summary statistics about the tree, and describe the results obtained. What is the training error rate? How many terminal nodes does the tree have?

```
tree.oj <- tree(Purchase ~ ., OJ, subset = train)

summary(tree.oj)

##
## Classification tree:
## tree(formula = Purchase ~ ., data = OJ, subset = train)
## Variables actually used in tree construction:
## [1] "LoyalCH"      "PriceDiff"    "SpecialCH"    "ListPriceDiff"
## [5] "PctDiscMM"
## Number of terminal nodes: 9
## Residual mean deviance: 0.7432 = 587.8 / 791
## Misclassification error rate: 0.1588 = 127 / 800
```

The training error rate is 0.1588 and there are 9 terminal nodes in the tree model.

- c) Type in the name of the tree object in order to get a detailed text output. Pick one of the terminal nodes, and interpret the information displayed.

```
tree.oj

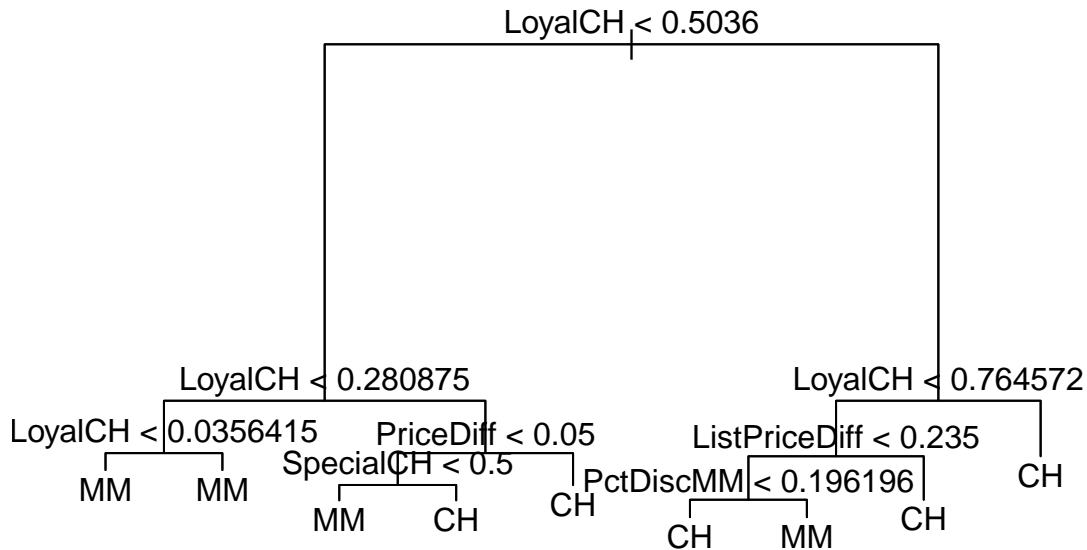
## node), split, n, deviance, yval, (yprob)
##      * denotes terminal node
##
## 1) root 800 1073.00 CH ( 0.60625 0.39375 )
##    2) LoyalCH < 0.5036 365 441.60 MM ( 0.29315 0.70685 )
##      4) LoyalCH < 0.280875 177 140.50 MM ( 0.13559 0.86441 )
##        8) LoyalCH < 0.0356415 59 10.14 MM ( 0.01695 0.98305 ) *
##        9) LoyalCH > 0.0356415 118 116.40 MM ( 0.19492 0.80508 ) *
##      5) LoyalCH > 0.280875 188 258.00 MM ( 0.44149 0.55851 )
##        10) PriceDiff < 0.05 79 84.79 MM ( 0.22785 0.77215 )
##          20) SpecialCH < 0.5 64 51.98 MM ( 0.14062 0.85938 ) *
##          21) SpecialCH > 0.5 15 20.19 CH ( 0.60000 0.40000 ) *
##        11) PriceDiff > 0.05 109 147.00 CH ( 0.59633 0.40367 ) *
##    3) LoyalCH > 0.5036 435 337.90 CH ( 0.86897 0.13103 )
##      6) LoyalCH < 0.764572 174 201.00 CH ( 0.73563 0.26437 )
##        12) ListPriceDiff < 0.235 72 99.81 MM ( 0.50000 0.50000 )
##          24) PctDiscMM < 0.196196 55 73.14 CH ( 0.61818 0.38182 ) *
##          25) PctDiscMM > 0.196196 17 12.32 MM ( 0.11765 0.88235 ) *
##      13) ListPriceDiff > 0.235 102 65.43 CH ( 0.90196 0.09804 ) *
##      7) LoyalCH > 0.764572 261 91.20 CH ( 0.95785 0.04215 ) *
```

From the information displayed in the tree mode, the second terminal node split is  $\text{LoyalCH} < 0.5036$  with 365 observation in the branch. The majority class is “MM” with 70.69% of observation being “MM”,



d) Create a plot of the tree, and interpret the results.

```
plot(tree.oj)
text(tree.oj, pretty = 0)
```



e) Predict the response on the test data, and produce a confusion matrix comparing the test labels to the predicted test labels. What is the test error rate?

```
oj.pred <- predict(tree.oj, oj_test, type = "class")
table(oj.pred, oj_test$Purchase)
```

```
##
## oj.pred  CH  MM
##      CH 160  38
##      MM   8  64
```

```
(160 + 64) / 270
```

```
## [1] 0.8296296
```

f) Apply the cv.tree() function to the training set in order to determine the optimal tree size.

```
cv_oj <- cv.tree(tree.oj, FUN = prune.misclass)
cv_oj
```

```
## $size
## [1] 9 8 7 4 2 1
##
## $dev
## [1] 150 150 149 158 172 315
##
## $k
## [1] -Inf 0.000000 3.000000 4.333333 10.500000 151.000000
##
## $method
## [1] "misclass"
##
## attr(,"class")
## [1] "prune" "tree.sequence"
```

From the cross validation on the `oj.tree` model it is evident that the optima tree size is one with 7 terminal nodes since it has the lwoest deviance of 149.