

Homework 4

Yanhao Li

Contents

Question 1	2
Question 2	8

```
library(tidyverse)
library(ISLR)
library(lasso2)
library(ISLR)
library(rpart)
library(rpart.plot)
library(ranger)
library(caret)
library(gbm)
```

Question 1

Load, clean, and tidy data

```
data("Prostate")

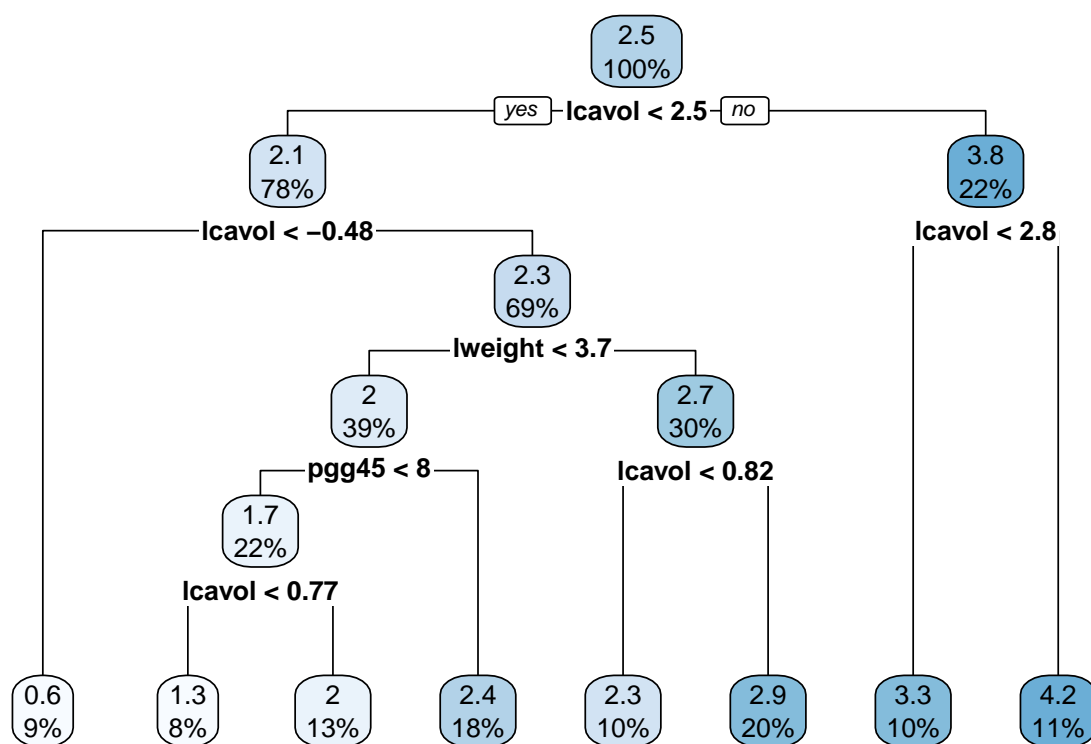
prostate = Prostate %>%
  janitor::clean_names()
```

Question a

```
set.seed(1)

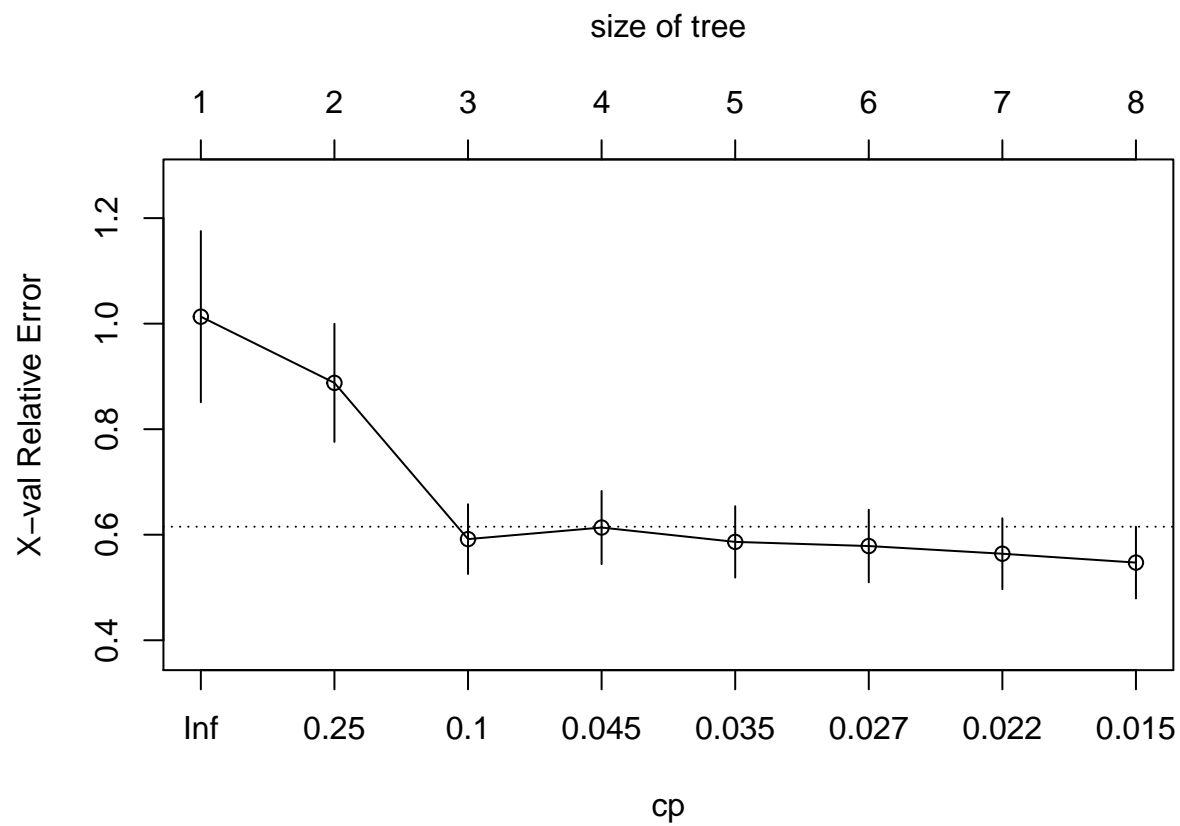
tree1 = rpart(formula = lpsa ~ .,
               data = prostate)

rpart.plot(tree1)
```



```
cpTable <- tree1$cptable
```

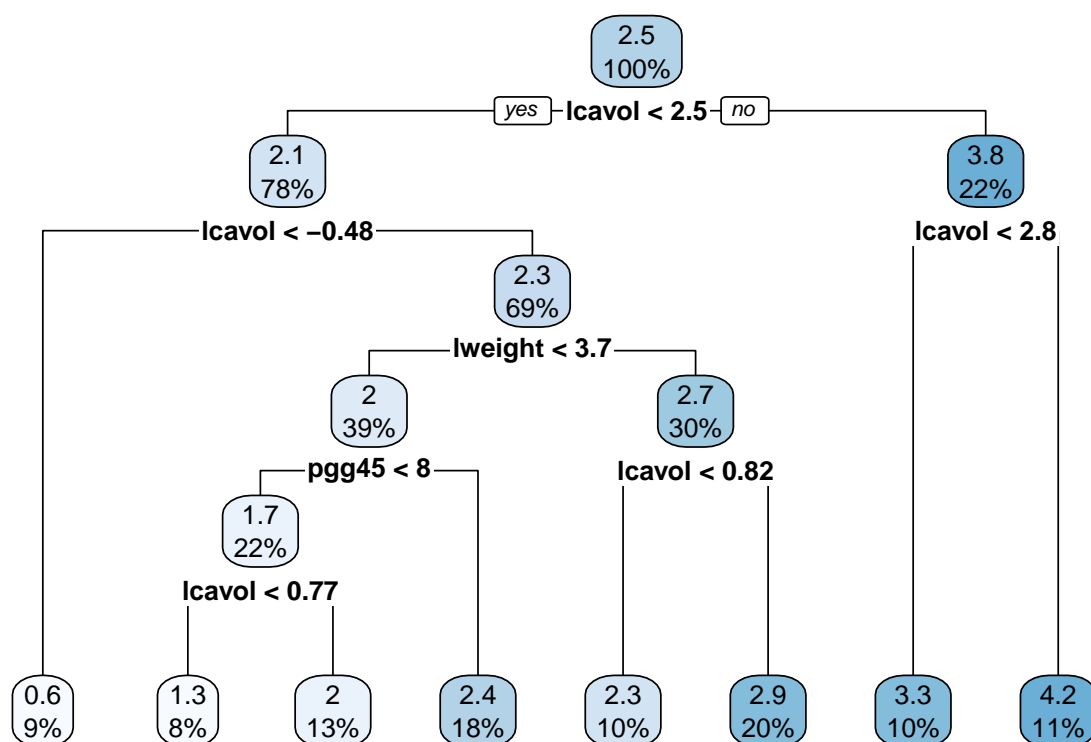
```
plotcp(tree1)
```



```
# minimum cross-validation error
minErr <- which.min(cpTable[,4])

tree2 <- prune(tree1, cp = cpTable[minErr,1])

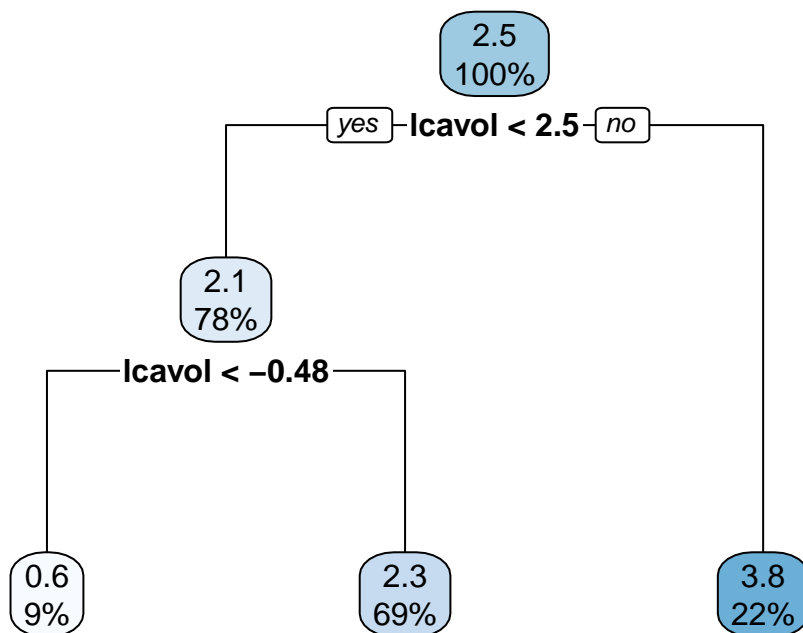
rpart.plot(tree2)
```



```

# 1SE rule
tree3 <- prune(tree1, cp = cpTable[cpTable[,4] < cpTable[minErr,4] + cpTable[minErr,5],1][1])
rpart.plot(tree3)

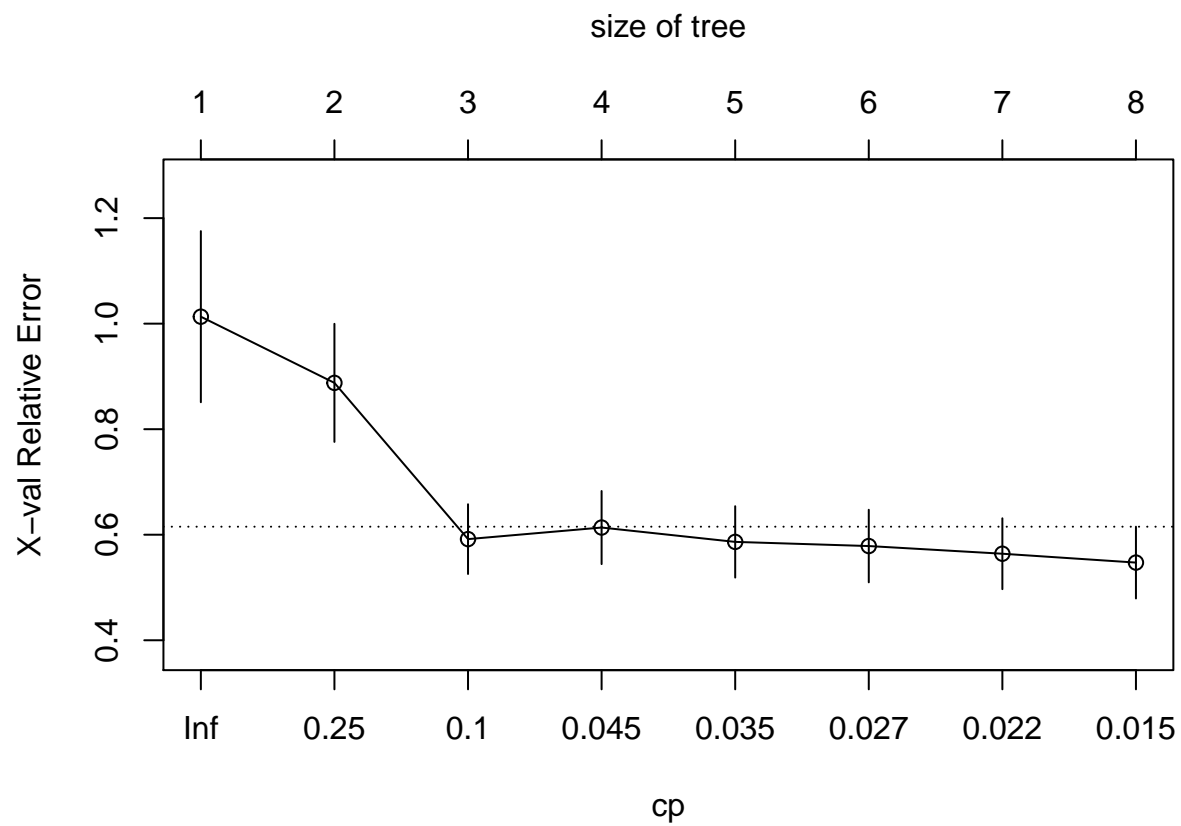
```



Tree size corresponds to the lowest cross-validation error is 8. It is different from the tree size obtained using the 1 SE rule, which is 3.

Question b

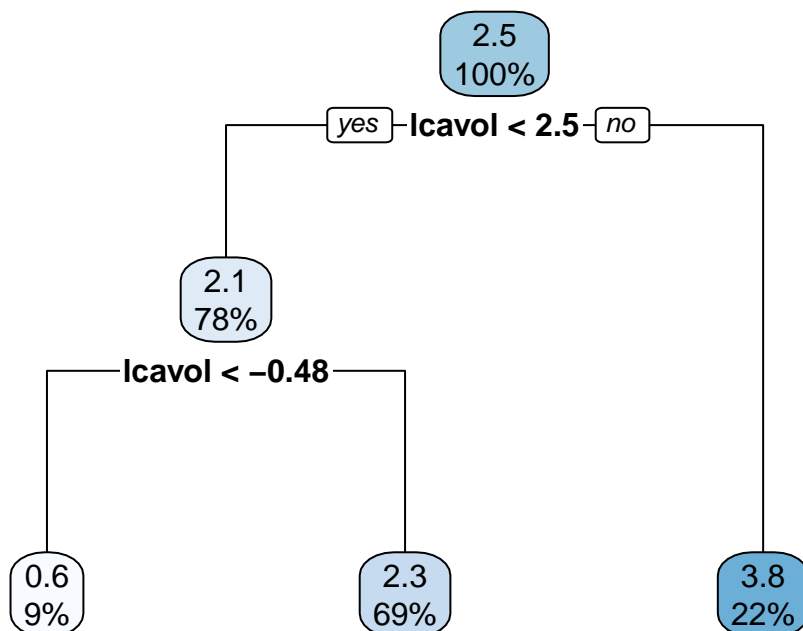
```
plotcp(tree1)
```



```
set.seed(1)

tree4 <- rpart(formula = lpsa ~ .,
               data = prostate,
               control = rpart.control(cp = 0.1))

rpart.plot(tree4)
```



A good choice of cp for pruning is often the leftmost value for which the mean lies below the horizontal line. According to the plot, I choose cp equals to 0.1 and size of tree equals to 3.

In terminal node where $lcavol$ is less than -0.48, the mean $lpsa$ is 0.6. This node contains 9% of the sample.

Question c

Question d

Question e

Question f

Question 2

Load, clean, and tidy data

Question a

Question b

Question c