

# FE590. Assignment #3 - Gang Ping Zhu

Enter Your Name Here, or “Anonymous” if you want to remain anonymous..

2017-11-10

## Instructions

In this assignment, you should use R markdown to answer the questions below. Simply type your R code into embedded chunks as shown above.

When you have completed the assignment, knit the document into a PDF file, and upload *both* the .pdf and .Rmd files to Canvas.

Note that you must have LaTeX installed in order to knit the equations below. If you do not have it installed, simply delete the questions below.

```
library(ISLR)

## Warning: package 'ISLR' was built under R version 3.4.2
library(boot)
library(leaps)

## Warning: package 'leaps' was built under R version 3.4.2
library(glmnet)

## Warning: package 'glmnet' was built under R version 3.4.2
## Loading required package: Matrix
## Loading required package: foreach
## Warning: package 'foreach' was built under R version 3.4.2
## Loaded glmnet 2.0-13
library(caTools)
library(pls)

## Warning: package 'pls' was built under R version 3.4.2
##
## Attaching package: 'pls'
## The following object is masked from 'package:stats':
##   loadings
library(tree)

## Warning: package 'tree' was built under R version 3.4.2
library(randomForest)

## Warning: package 'randomForest' was built under R version 3.4.2
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
```

```

library(gbm)

## Warning: package 'gbm' was built under R version 3.4.2
## Loading required package: survival
##
## Attaching package: 'survival'
## The following object is masked from 'package:boot':
## 
##      aml
## Loading required package: lattice
##
## Attaching package: 'lattice'
## The following object is masked from 'package:boot':
## 
##      melanoma
## Loading required package: splines
## Loading required package: parallel
## Loaded gbm 2.1.3

```

## Question 1 (based on JWHT Chapter 5, Problem 8)

In this problem, you will perform cross-validation on a simulated data set.

Generate a simulated data set as follows:

```

set.seed(1)
y <- rnorm(100)
x <- rnorm(100)
y <- x - 2*x^2 + rnorm(100)

```

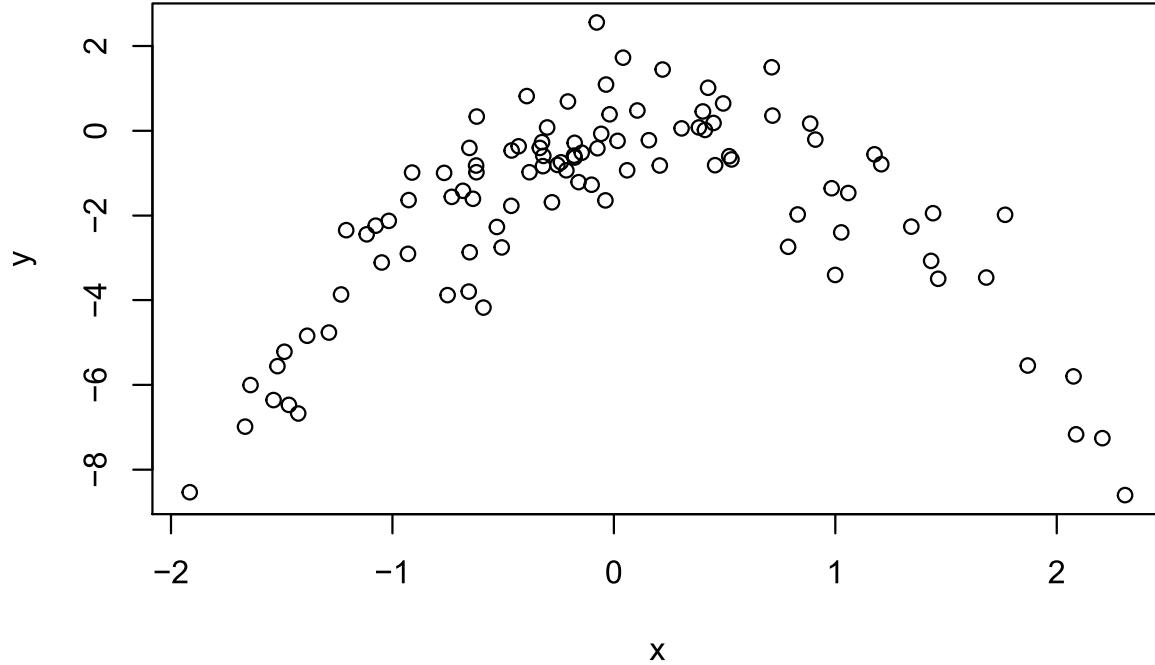
- (a) In this data set, what is  $n$  and what is  $p$ ?

```
df <- cbind.data.frame(x, y)
```

$n = 100$   $p = 1$

- (b) Create a scatterplot of  $x$  against  $y$ . Comment on what you find.

```
plot(x,y)
```



There is heavy curvature to the plot where we see an increase in Y when X gets closer to 0.

- (c) Set a random seed of 2, and then compute the LOOCV errors that result from fitting the following four models using least squares:

1.  $Y = \beta_0 + \beta_1 X + \epsilon$
2.  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$
3.  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon$
4.  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \epsilon$

```
set.seed(2)
glm.fit1 <- glm(y ~ x,data=df)
cv.err1 <- cv.glm(df, glm.fit1)
cv.err1$delta

## [1] 5.890979 5.888812

glm.fit2 <- glm(y ~ x + poly(x,2),data=df)
cv.err2 <- cv.glm(df, glm.fit2)

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading
```

```

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

cv.err4$delta

## [1] 1.114772 1.114334

```

- (d) Which of the models in (c) had the smallest LOOCV error? Is this what you expected? Explain your answer. The second model had the smallest LOOCV error. This is as expected as it's the model that closely resembles the  $y$  equation.
- (e) Comment on the statistical significance of the coefficient estimates that results from fitting each of the models in (c) using least squares. Do these results agree with the conclusions drawn based on the cross-validation results? Yes, these results agree with the cross-validation results as our p-value indicates. Our second model, which mostly resembled  $y$ , showed the lowest test error.

## Question 2 (based on JWHT Chapter 6, Problem 8)

In this exercise, we will generate simulated data, and will then use this data to perform best subset selection.

- (a) Set the random seed to be 10. Use the `rnorm()` function to generate a predictor  $X$  of length  $n = 100$ , as well as a noise vector  $\epsilon$  of length  $n = 100$ .

```

set.seed(10)
x <- rnorm(100)
epsil <- rnorm(100)

```

- (b) Generate a response vector  $Y$  of length  $n = 100$  according to the model

$$Y = 4 + 3X + 2X^2 + X^3 + \epsilon.$$

```

y <- 4 + 3*x + 2*x^2 + x^3 + epsil

```

- (c) Use the `regsubsets()` function to perform best subset selection in order to choose the best model containing the predictors  $X, X^2, \dots, X^{10}$ . What is the best model obtained according to  $C_p$ , BIC, and adjusted  $R^2$ ? Show some plots to provide evidence for your answer, and report the coefficients of the best model obtained. Note you will need to use the `data.frame()` function to create a single data set containing both  $X$  and  $Y$ .

```

df2 <- data.frame(x, y)
r.fit <- regsubsets(y ~ poly(x, 10), df2, nvmax = 10)
r.sum <- summary(r.fit)
names(r.sum)

## [1] "which"   "rsq"     "rss"      "adjr2"    "cp"       "bic"      "outmat"   "obj"

r.sum$rsq

## [1] 0.7286553 0.8860418 0.9648043 0.9658630 0.9661574 0.9663587 0.9664334
## [8] 0.9664580 0.9664635 0.9664680

```

According to the  $R^2$  from the summary, the last model (the one with 10 variables) is the best model with a value of 0.9664680.

```

r.sum$cp

## [1] 624.1972453 208.4653464 1.4155758 0.6057486 1.8243969
## [6] 3.2899388 5.0917595 7.0265590 9.0117735 11.0000000

```

According to the  $C_p$ , from the summary, the third model (containing three variables) is the best one with a value of 1.4155758.

```

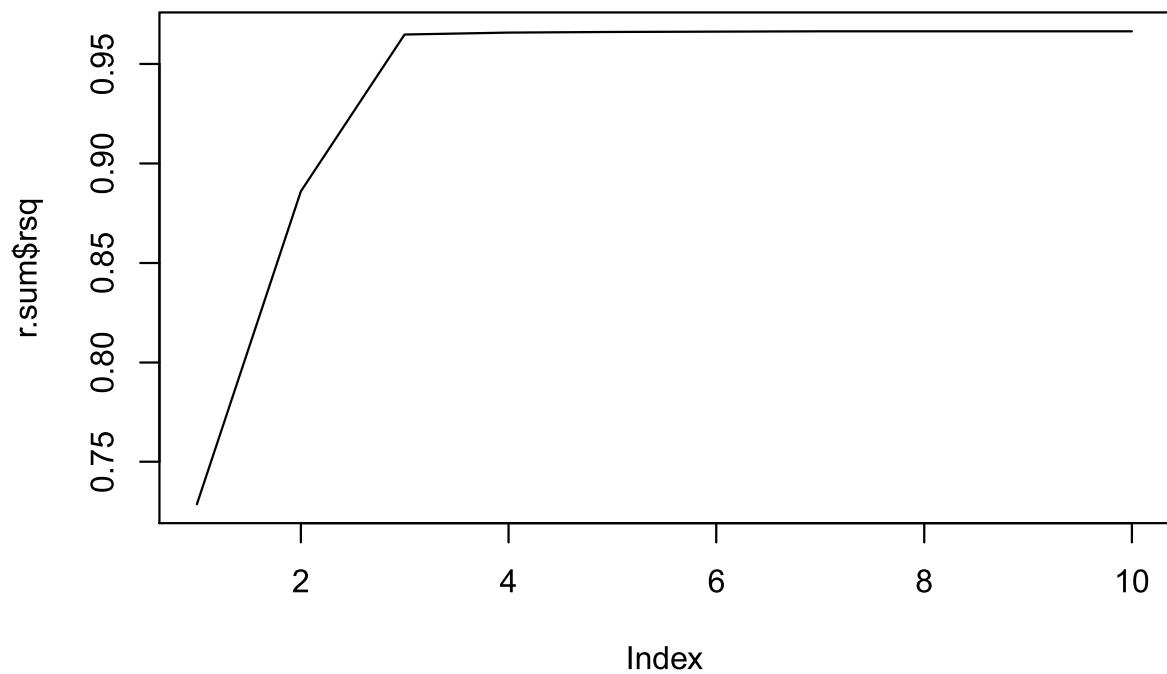
r.sum$bic

## [1] -121.2262 -203.3769 -316.2625 -314.7114 -310.9724 -306.9640 -302.5810
## [8] -298.0490 -293.4605 -288.8685

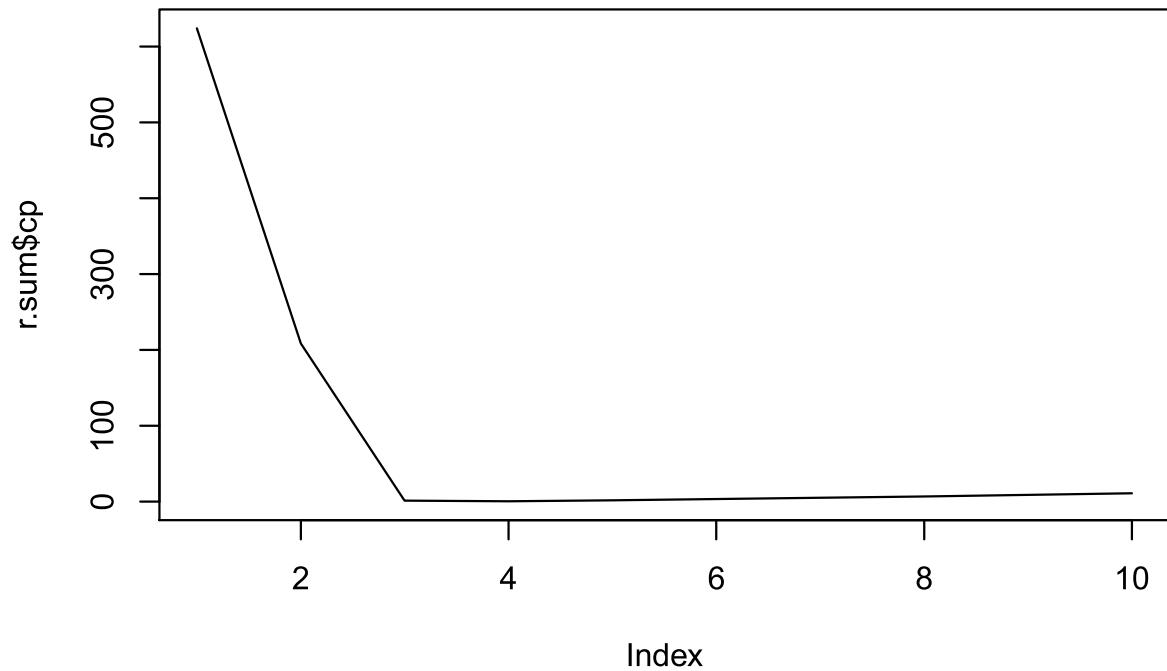
```

According to the BIC from the summary, the third model (containing three variables) is the best one with a value of -316.2625.

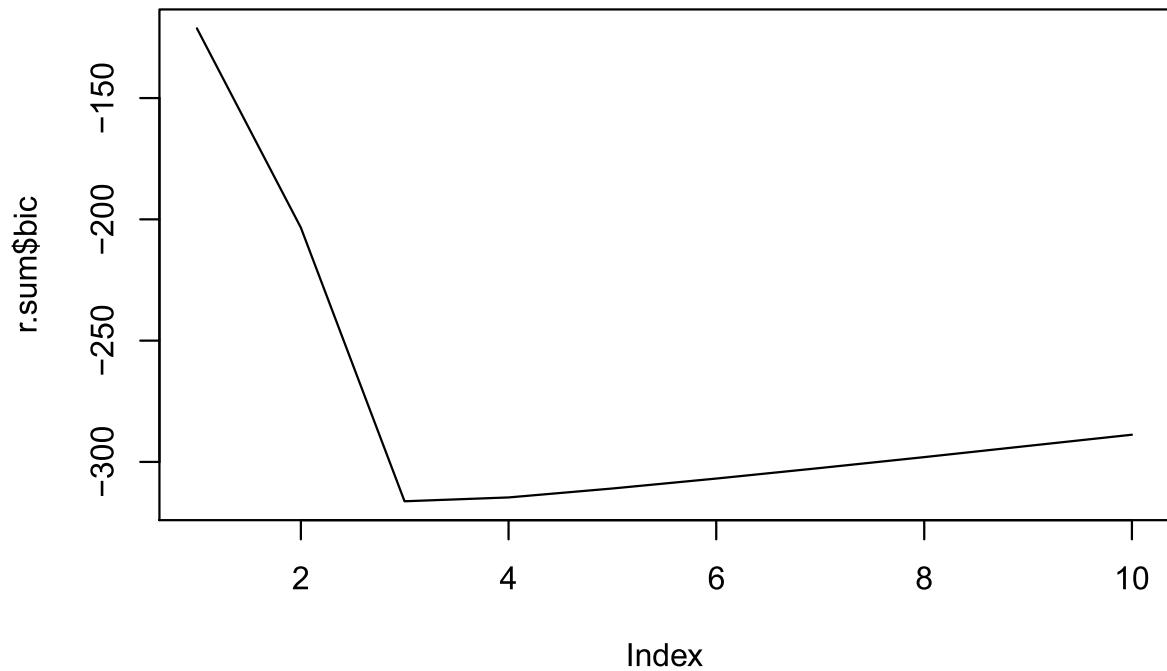
```
plot(r.sum$rsq, type = "l")
```



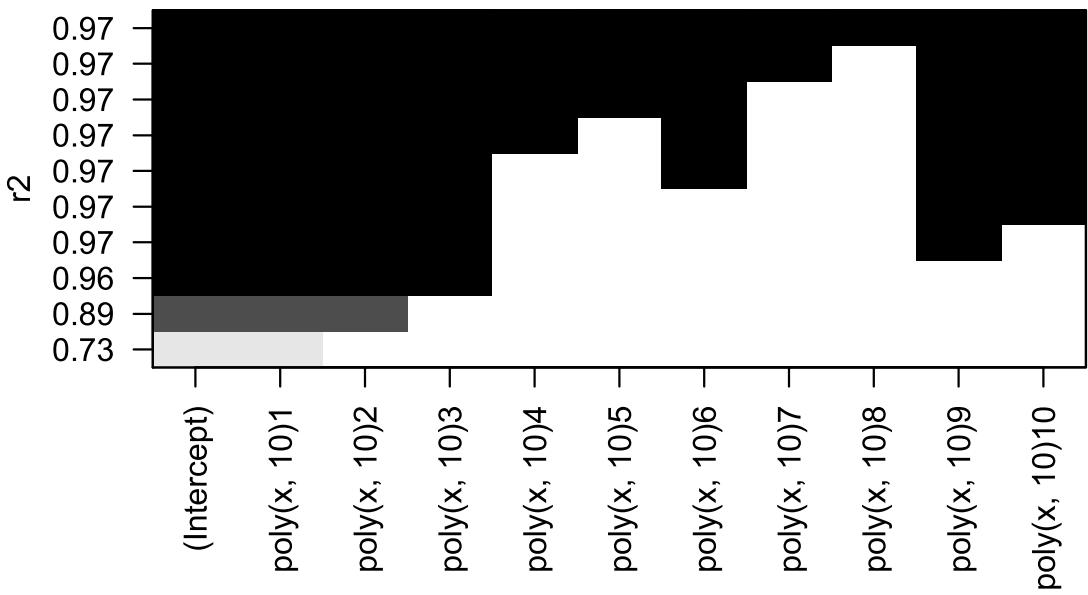
```
plot(r.sum$cp, type = "l")
```



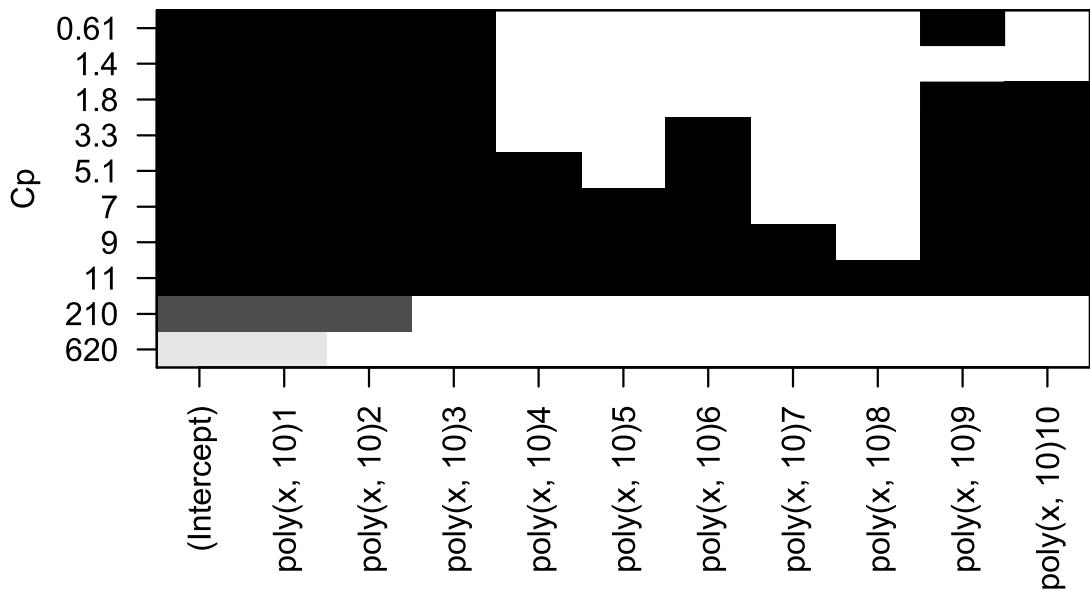
```
plot(r.sum$bic, type = "l")
```



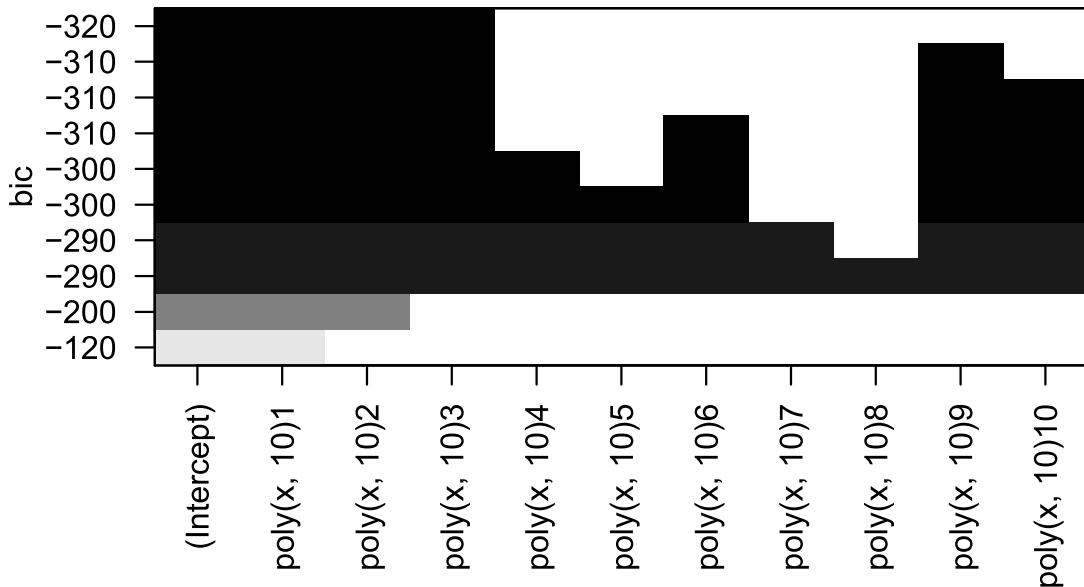
```
plot(r.fit, scale = "r2")
```



```
plot(r.fit, scale = "Cp")
```



```
plot(r.fit, scale = "bic")
```



- (d) Repeat (c), using forward stepwise selection and also using backwards stepwise selection. How does your answer compare to the results in (c)?

```
r.fitfwd <- regsubsets(y ~ poly(x, 10), df2, nvmax = 10, method = "forward")
r.sumfwd <- summary(r.fitfwd)
r.sumfwd$rsq

## [1] 0.7286553 0.8860418 0.9648043 0.9658630 0.9661574 0.9663587 0.9664334
## [8] 0.9664580 0.9664635 0.9664680

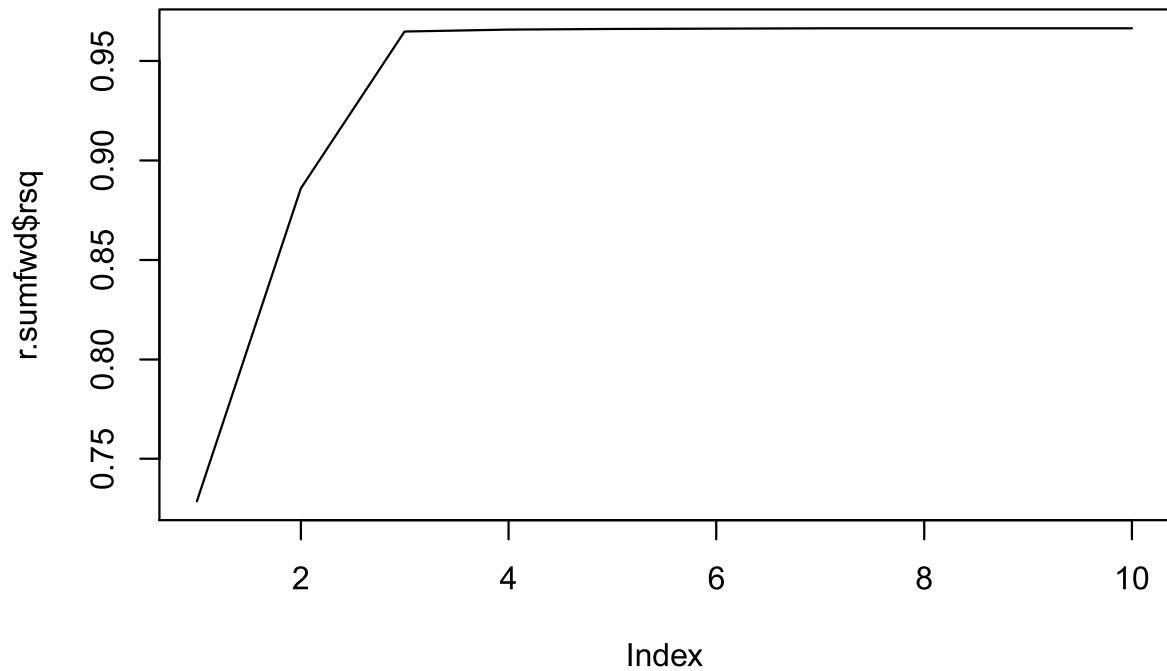
r.sumfwd$cp

## [1] 624.1972453 208.4653464 1.4155758 0.6057486 1.8243969
## [6] 3.2899388 5.0917595 7.0265590 9.0117735 11.0000000

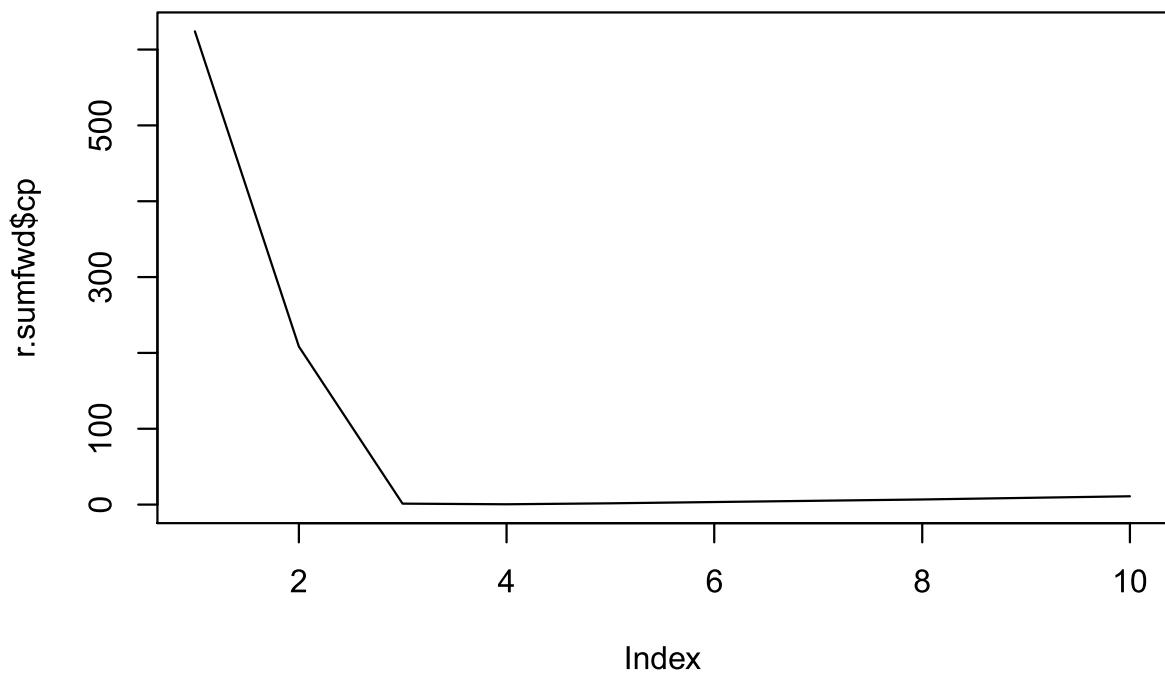
r.sumfwd$bic

## [1] -121.2262 -203.3769 -316.2625 -314.7114 -310.9724 -306.9640 -302.5810
## [8] -298.0490 -293.4605 -288.8685

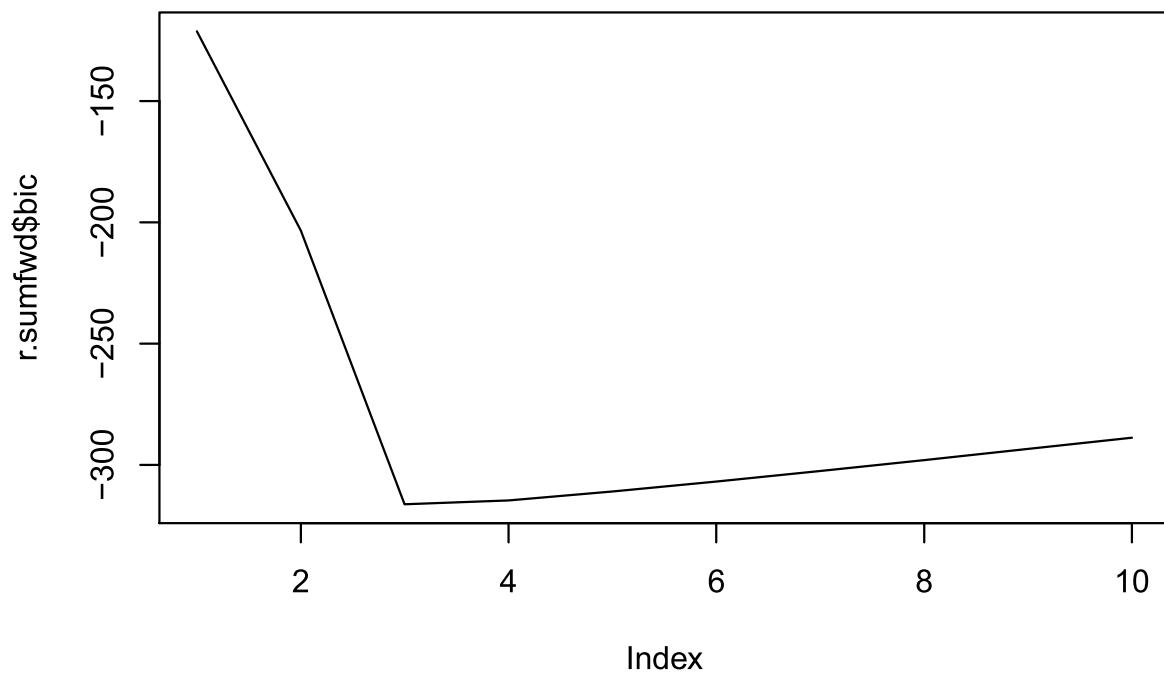
plot(r.sumfwd$rsq, type = "l")
```



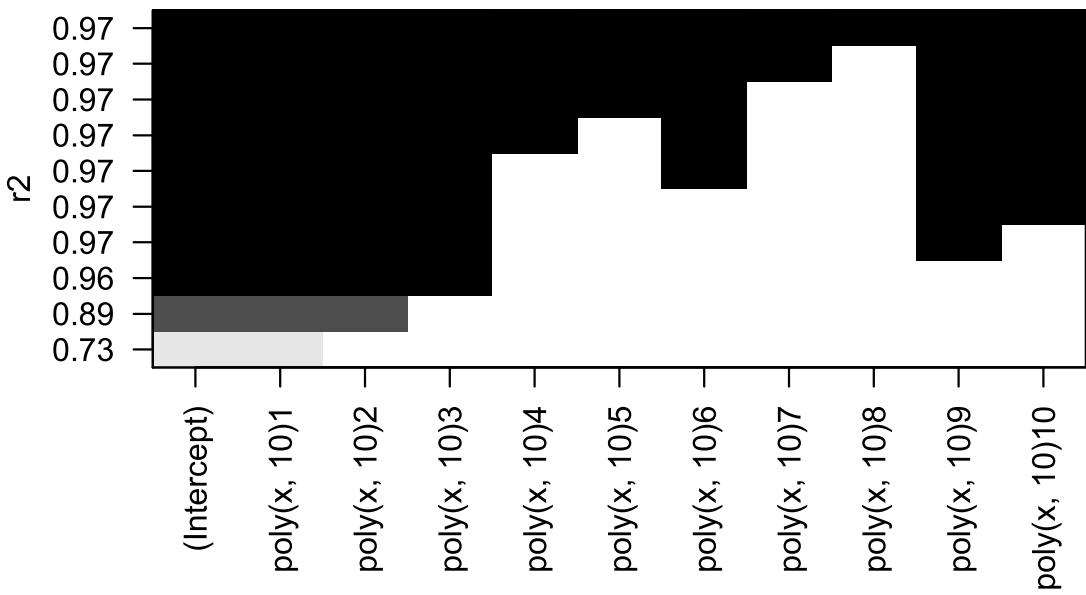
```
plot(r.sumfwd$cp, type = "l")
```



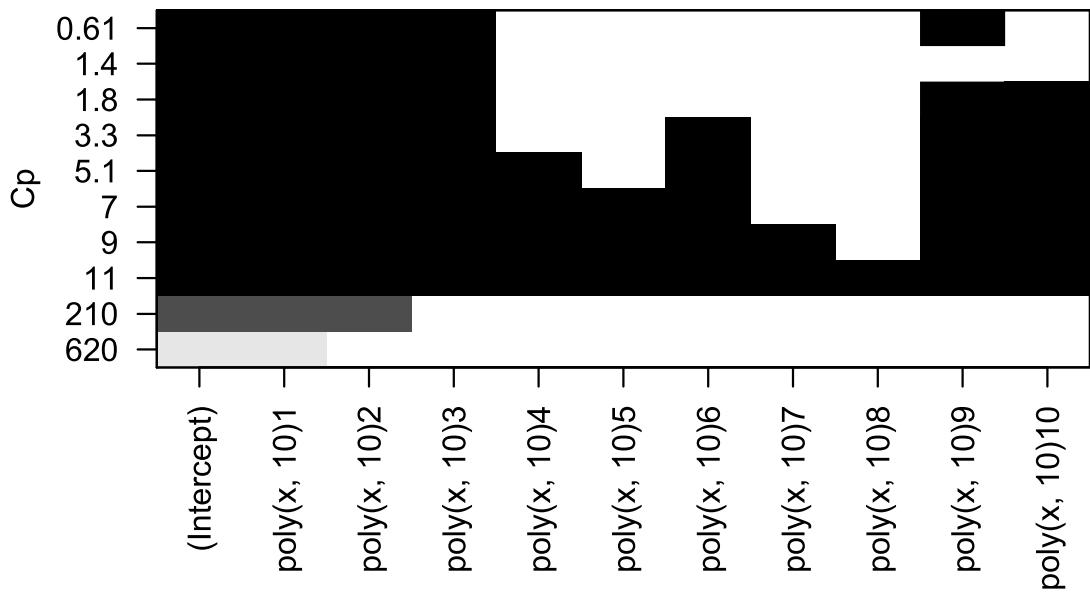
```
plot(r.sumfwd$bic, type = "l")
```



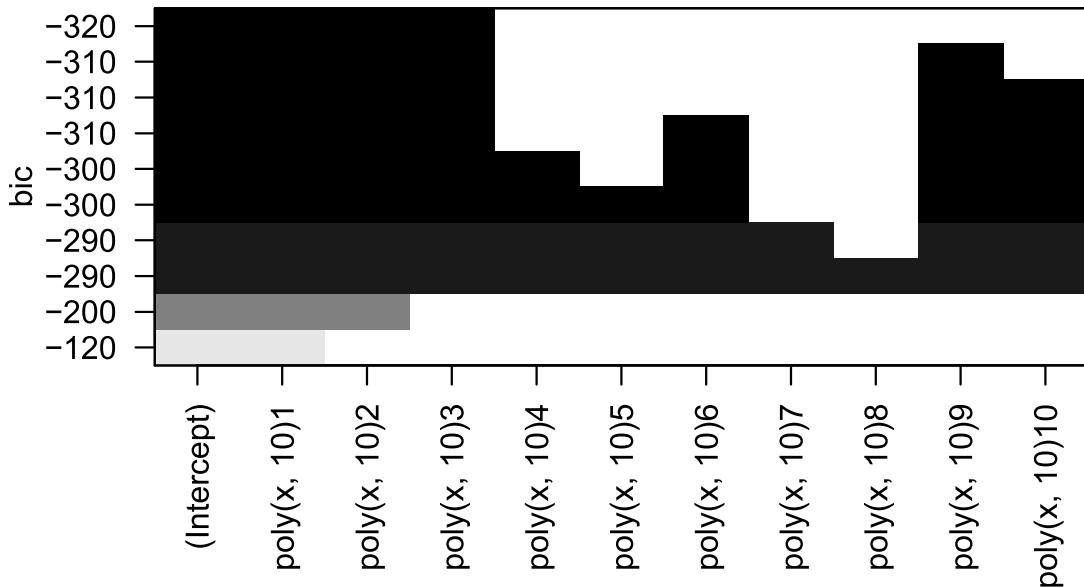
```
plot(r.fitfwd, scale = "r2")
```



```
plot(r.fitfwd, scale = "Cp")
```



```
plot(r.fitfwd, scale = "bic")
```



The forward step selection returned similar results where the model with 3 variables is the best for CP and BIC. For the  $R^2$ , the model with 10 variables is the best.

```
r.fitbwd <- regsubsets(y ~ poly(x, 10), df2, nvmax = 10, method = "backward")
r.sumbwd <- summary(r.fitbwd)
r.sumbwd$rsq

## [1] 0.7286553 0.8860418 0.9648043 0.9658630 0.9661574 0.9663587 0.9664334
## [8] 0.9664580 0.9664635 0.9664680

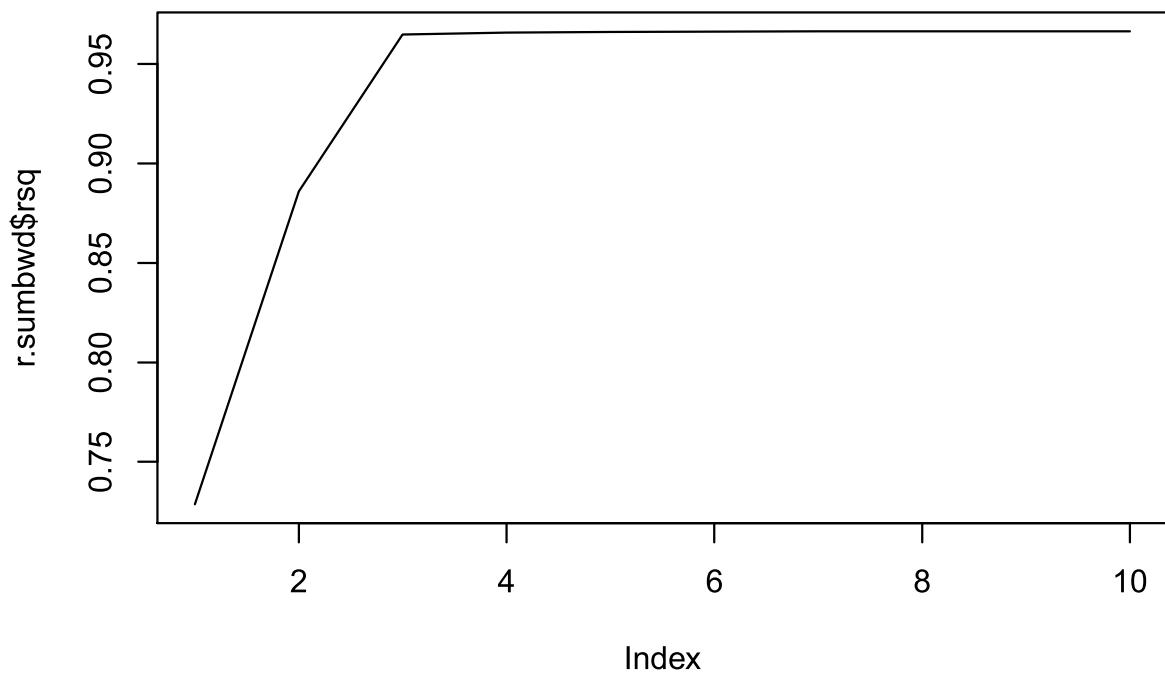
r.sumbwd$cp

## [1] 624.1972453 208.4653464 1.4155758 0.6057486 1.8243969
## [6] 3.2899388 5.0917595 7.0265590 9.0117735 11.0000000

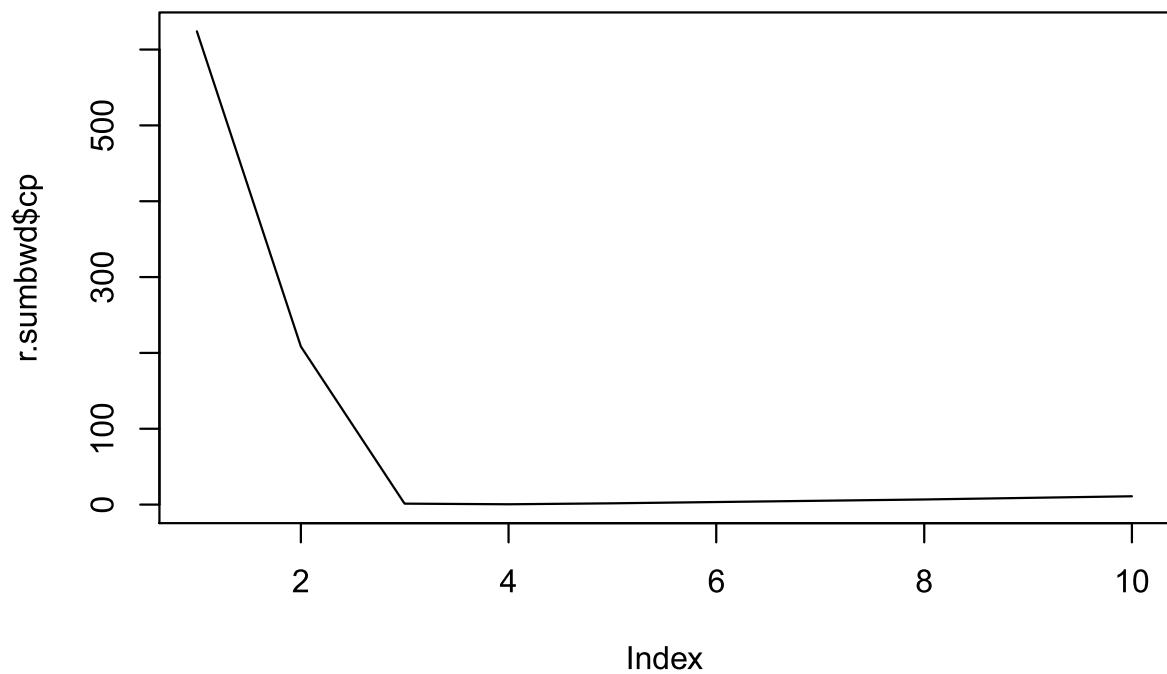
r.sumbwd$bic

## [1] -121.2262 -203.3769 -316.2625 -314.7114 -310.9724 -306.9640 -302.5810
## [8] -298.0490 -293.4605 -288.8685

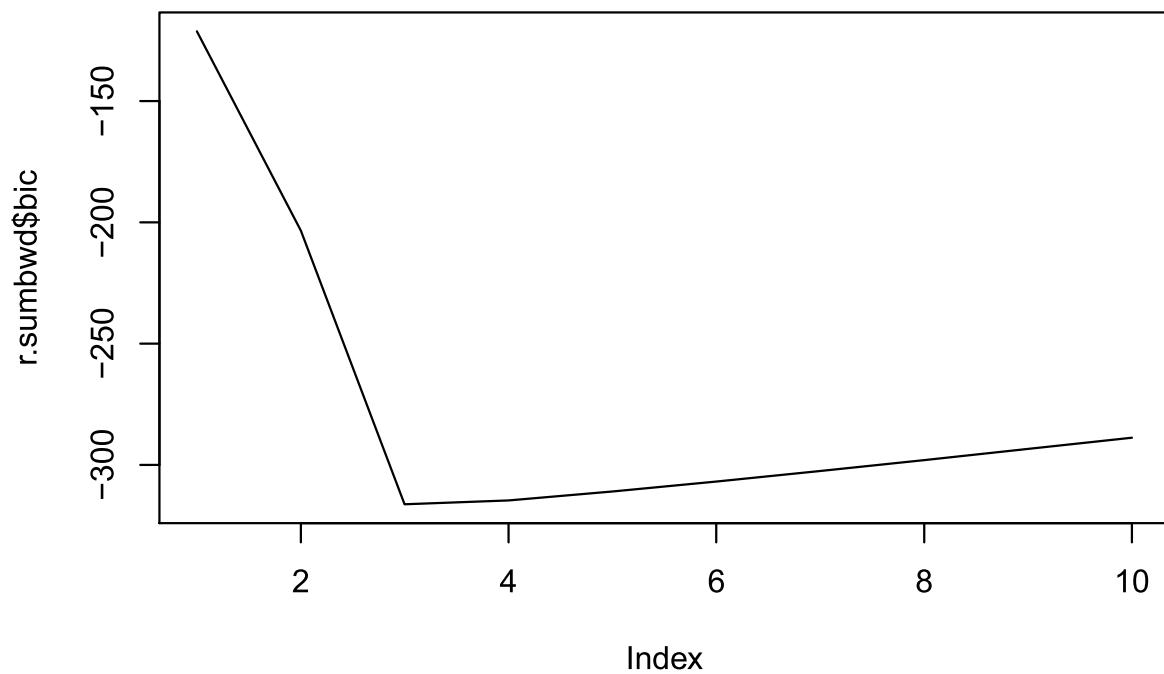
plot(r.sumbwd$rsq, type = "l")
```



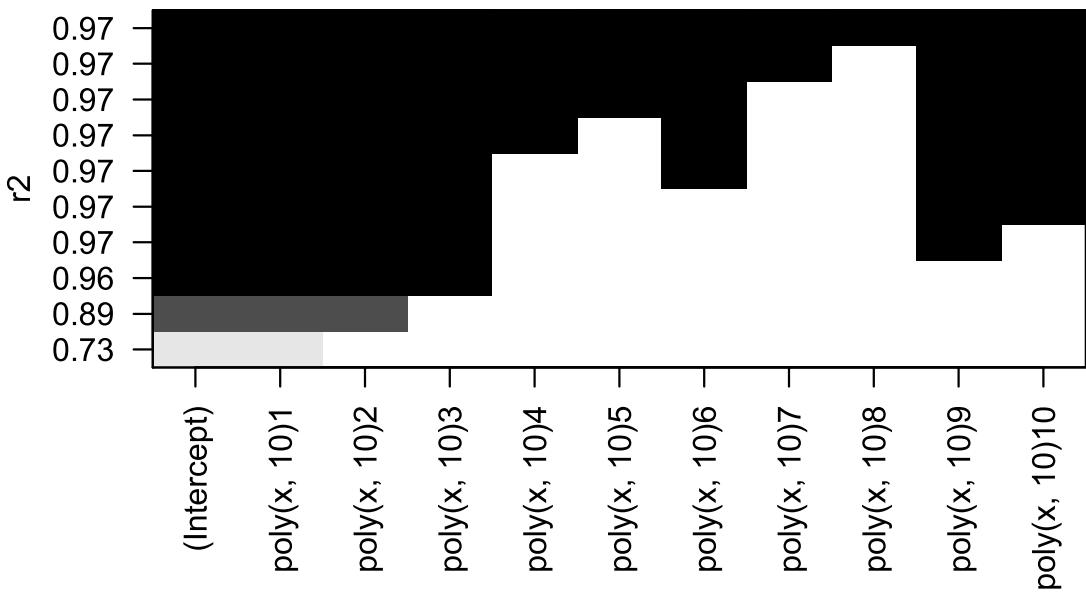
```
plot(r.sumbwd$cp, type = "l")
```



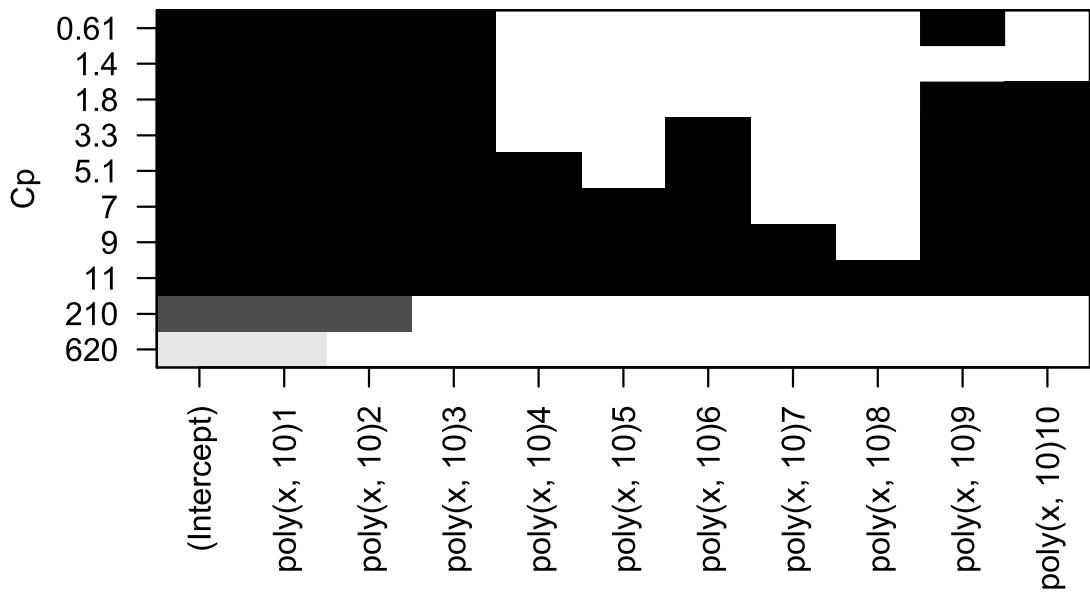
```
plot(r.sumbwd$bic, type = "l")
```



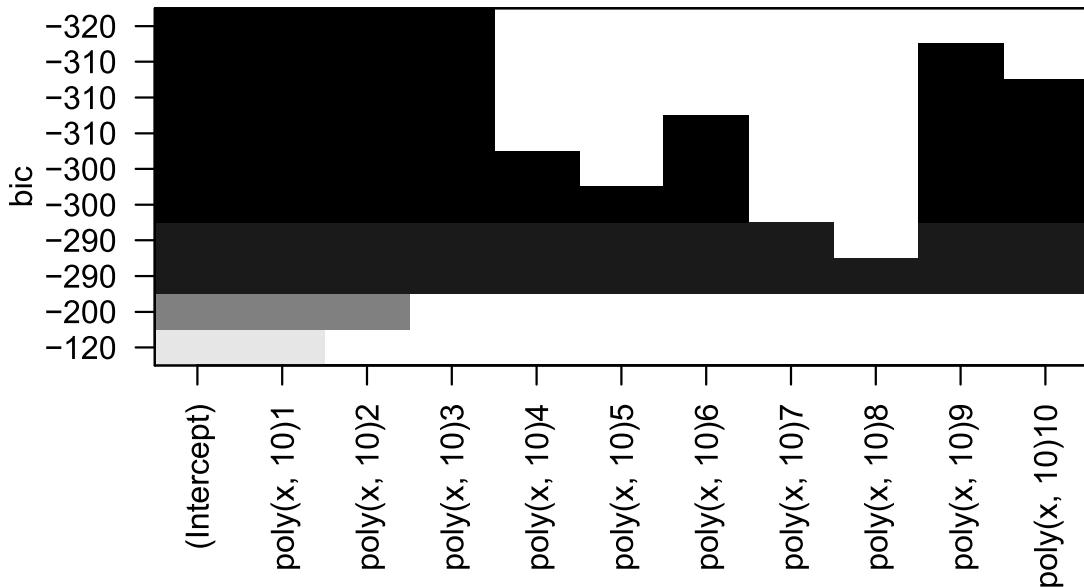
```
plot(r.fitbwd, scale = "r2")
```



```
plot(r.fitbwd, scale = "Cp")
```



```
plot(r.fitbwd, scale = "bic")
```



The backward step selection returned similar results to the forward step selection where the model with 3 variables is the best for CP and BIC. For the  $R^2$ , the model with 10 variables is the best.

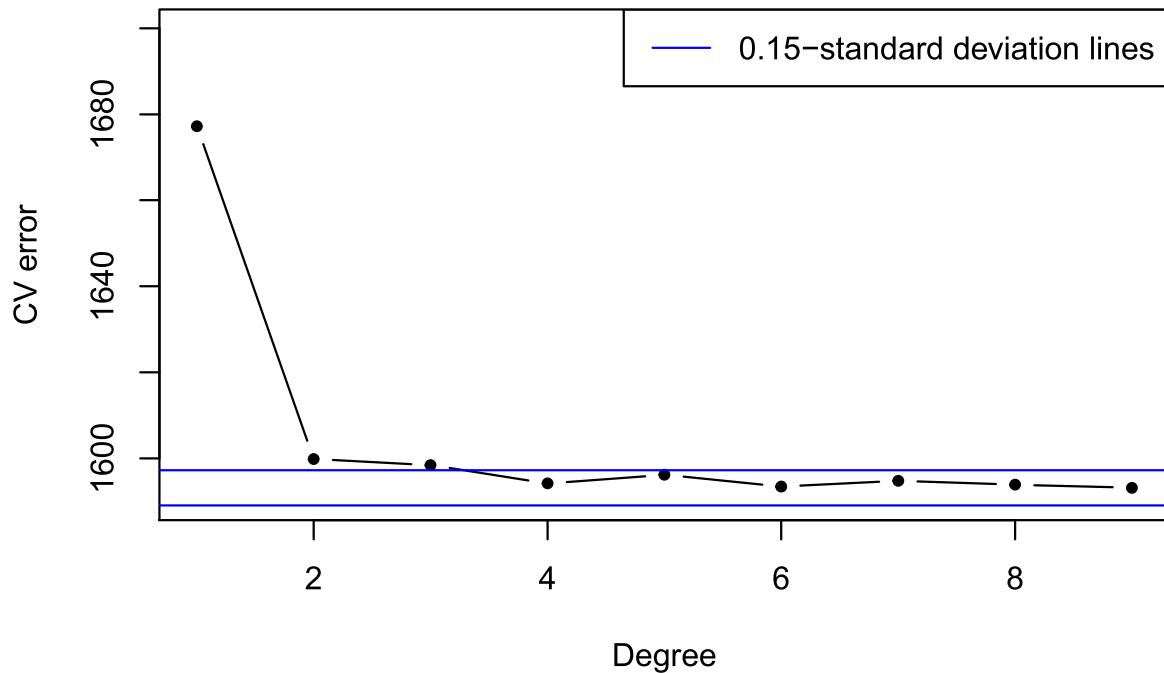
### Question 3 (based on JWHT Chapter 7, Problem 6)

In this exercise, you will further analyze the `Wage` data set.

- (a) Perform polynomial regression to predict `wage` using `age`. Use cross-validation to select the optimal degree  $d$  for the polynomial. What degree was chosen? Make a plot of the resulting polynomial fit to the data.

```
data(Wage)
set.seed(999)
all errs <- rep(NA, 9)
for(i in 1:9){
  glm.fit <- glm(wage~poly(age, i), data=Wage)
  all errs[i] <- cv.glm(Wage, glm.fit, K=9)$delta[2]
}

plot(1:9, all errs, xlab="Degree", ylab="CV error", type="b", pch=20, ylim=c(1590, 1700))
sd.p <- sd(all errs)
min.p <- min(all errs)
abline(h=min.p + 0.15 * sd.p, col="blue", lty=1)
abline(h=min.p - 0.15 * sd.p, col="blue", lty=1)
legend("topright", "0.15-standard deviation lines", col="blue", lty=1)
```



From the plot above, the degree of 3 is the smallest degree with a small CV error.

```
f1 <- lm(wage~age, data=Wage)
f2 <- lm(wage~poly(age, 2), data=Wage)
f3 <- lm(wage~poly(age, 3), data=Wage)
f4 <- lm(wage~poly(age, 4), data=Wage)
f5 <- lm(wage~poly(age, 5), data=Wage)
f6 <- lm(wage~poly(age, 6), data=Wage)
f7 <- lm(wage~poly(age, 7), data=Wage)
f8 <- lm(wage~poly(age, 8), data=Wage)
f9 <- lm(wage~poly(age, 9), data=Wage)
f10 <- lm(wage~poly(age, 10), data=Wage)
anova(f1, f2, f3, f4, f5, f6, f7, f8, f9, f10)
```

```
## Analysis of Variance Table
##
## Model 1: wage ~ age
## Model 2: wage ~ poly(age, 2)
## Model 3: wage ~ poly(age, 3)
## Model 4: wage ~ poly(age, 4)
## Model 5: wage ~ poly(age, 5)
## Model 6: wage ~ poly(age, 6)
## Model 7: wage ~ poly(age, 7)
## Model 8: wage ~ poly(age, 8)
## Model 9: wage ~ poly(age, 9)
## Model 10: wage ~ poly(age, 10)
##   Res.Df      RSS Df Sum of Sq      F    Pr(>F)
```

```

## 1 2998 5022216
## 2 2997 4793430 1 228786 143.7638 < 2.2e-16 ***
## 3 2996 4777674 1 15756 9.9005 0.001669 **
## 4 2995 4771604 1 6070 3.8143 0.050909 .
## 5 2994 4770322 1 1283 0.8059 0.369398
## 6 2993 4766389 1 3932 2.4709 0.116074
## 7 2992 4763834 1 2555 1.6057 0.205199
## 8 2991 4763707 1 127 0.0796 0.777865
## 9 2990 4756703 1 7004 4.4014 0.035994 *
## 10 2989 4756701 1 3 0.0017 0.967529
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

The anova shows that all the polynomials above 3 are not significant

```
plot(wage~age, data=Wage)
```

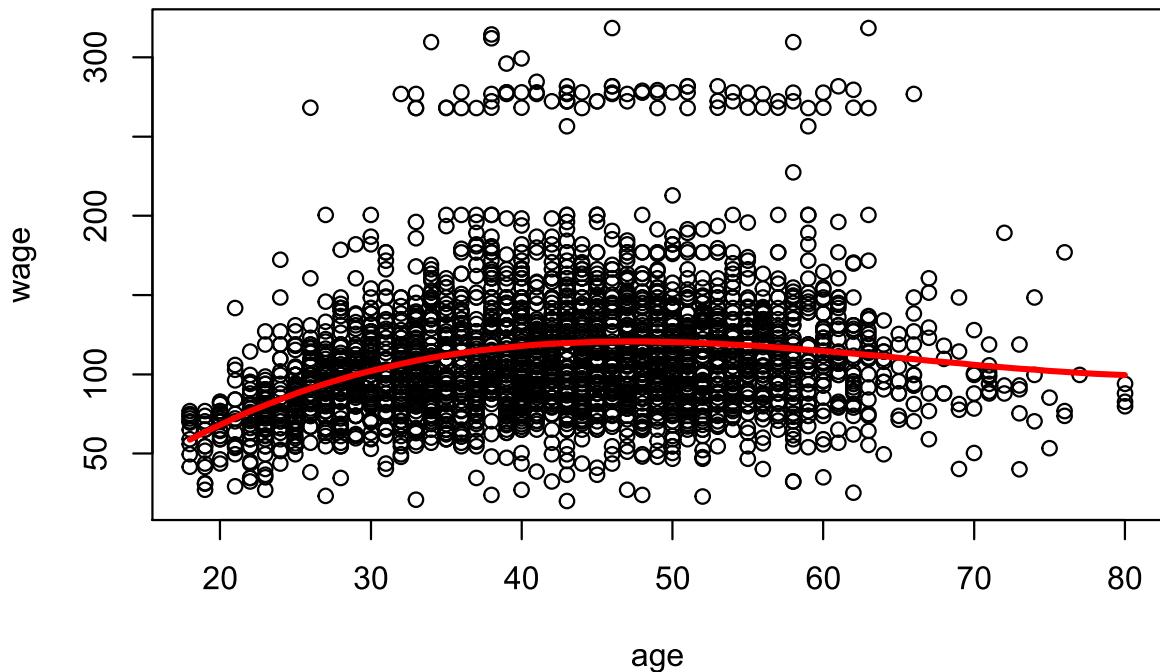
```
lmfit <- lm(wage~poly(age, 3), data=Wage)
```

```
agelms <- range(Wage$age)
```

```
agegrid <- seq(from=agelms[1], to=agelms[2])
```

```
lmpred <- predict(lmfit, data.frame(age=agegrid))
```

```
lines(agegrid, lmpred, col="red", lwd=3)
```



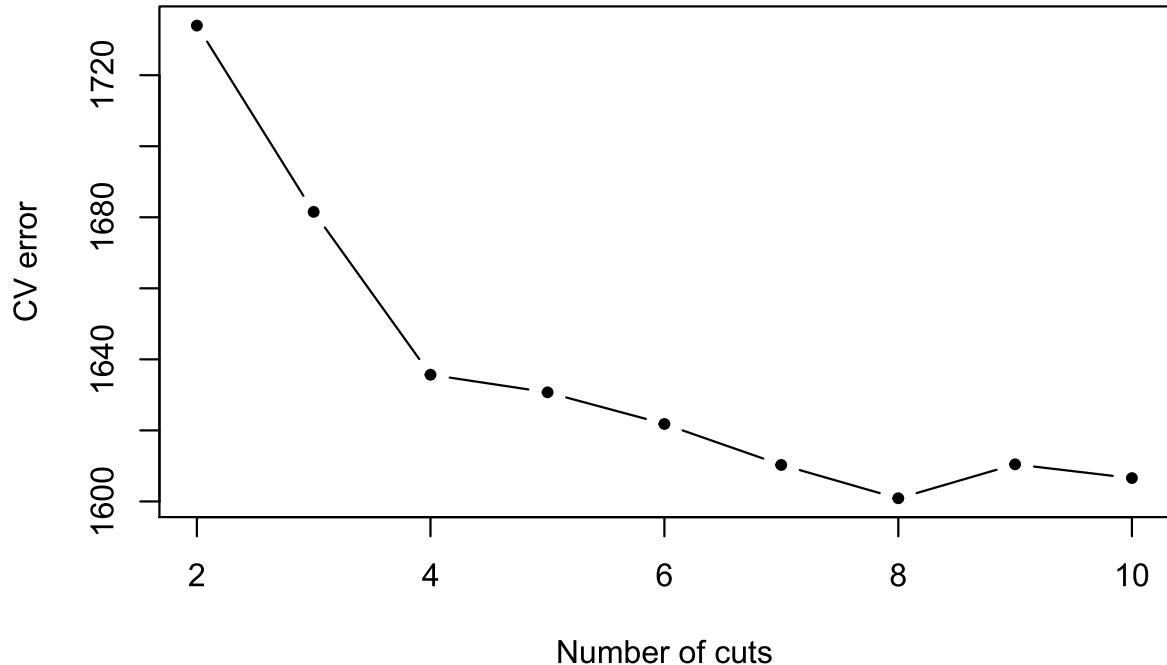
- (b) Fit a step function to predict `wage` using `age`, and perform cross-validation to choose the optimal

number of cuts. Make a plot of the fit obtained.

```
cvseerrs <- rep(NA, 10)
dfage <- Wage$age

for (i in 2:10) {
  Wage$age.cut <- cut(dfage, i)
  lm.fit <- glm(wage~age.cut, data=Wage)
  cvseerrs[i] <- cv.glm(Wage, lm.fit, K=10)$delta[2]
}

plot(2:10, cvseerrs[-1], xlab="Number of cuts", ylab="CV error", type="b", pch=20)
```



Based on the chart above, we'll use 8 cuts as our optimal number to plot.

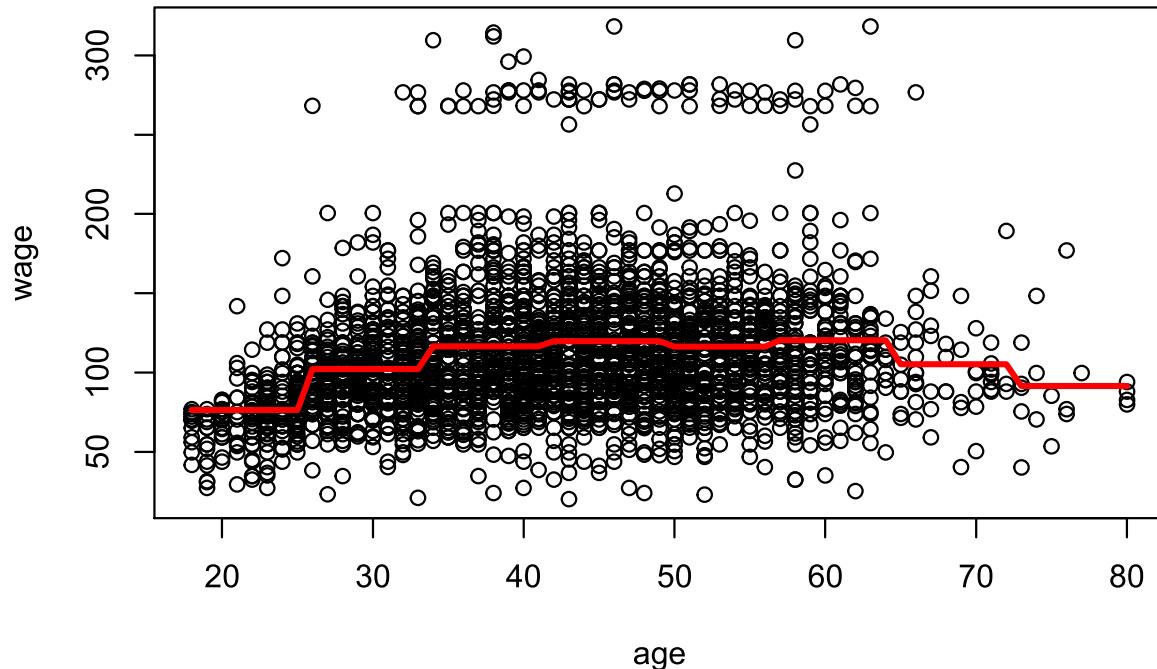
```
plot(wage~age, data=Wage)

lm.fit <- glm(wage~cut(age, 8), data=Wage)

agelms <- range(dfage)
agegrid <- seq(from=agelms[1], to=agelms[2])

lmpred <- predict(lm.fit, data.frame(age=agegrid))

lines(agegrid, lmpred, col="red", lwd=3)
```



## Question 4 (based on JWHT Chapter 8, Problem 8)

In the lab, a classification tree was applied to the `Carseats` data set after converting Sales into a qualitative response variable. Now we will seek to predict Sales using regression trees and related approaches, treating the response as a quantitative variable.

- (a) Split the data set into a training set and a test set.

```
dfcs <- Carseats
dim(dfcs)

## [1] 400 11
set.seed(55)
n <- 1:floor(nrow(dfcs)/2)
TrainingSet <- dfcs[n, ]
TestSet <- dfcs[-n, ]
dim(TrainingSet)

## [1] 200 11
dim(TestSet)

## [1] 200 11
```

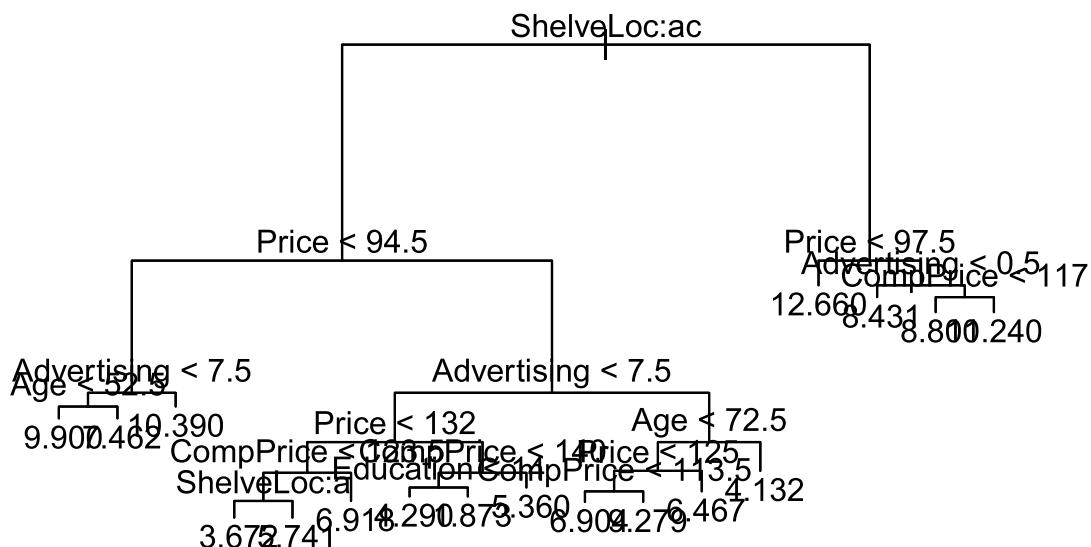
- (b) Fit a regression tree to the training set. Plot the tree, and interpret the results. What test MSE do you obtain?

```

treecs <- tree(Sales ~ ., data = TrainingSet)
summary(treecs)

##
## Regression tree:
## tree(formula = Sales ~ ., data = TrainingSet)
## Variables actually used in tree construction:
## [1] "ShelveLoc"      "Price"          "Advertising"    "Age"           "CompPrice"
## [6] "Education"
## Number of terminal nodes:  17
## Residual mean deviance:  2.283 = 417.9 / 183
## Distribution of residuals:
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## -4.9900 -0.8437  0.1344  0.0000  0.9115  3.3500
plot(treecs)
text(treecs)

```



```

mean((TestSet$Sales - predict(treecs, TestSet))^2)

```

```

## [1] 4.547366

```

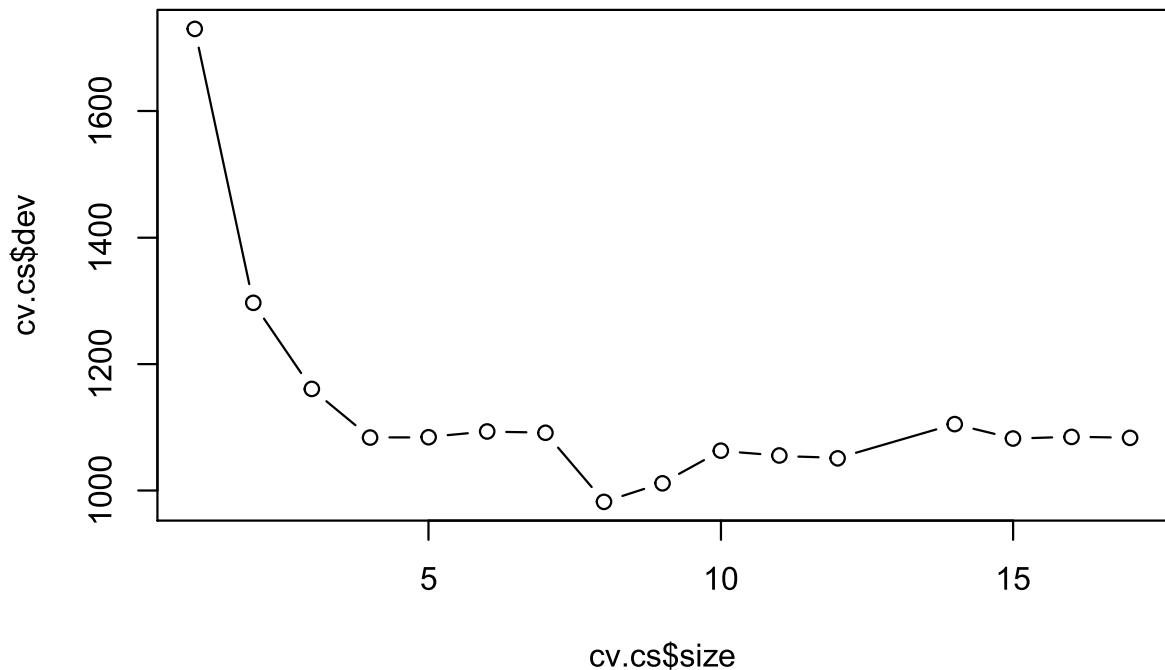
- (c) Use cross-validation in order to determine the optimal level of tree complexity. Does pruning the tree improve the test MSE?

```

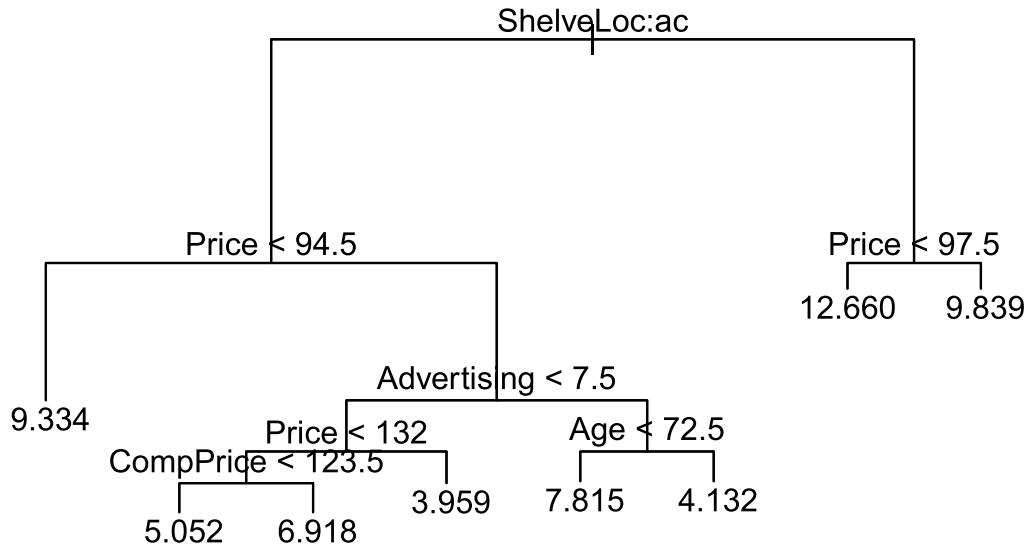
set.seed(77)
cv.cs <- cv.tree(treecs, FUN = prune.tree)

```

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



```
prunedcs <- prune.tree(treecs, best = 8)
plot(prunedcs)
text(prunedcs)
```



```
mean((TestSet$Sales - predict(prunedcs, TestSet))^2)
```

```
## [1] 4.919946
```

Pruning improved it to 4.92.

- (d) Use the bagging approach in order to analyze this data. What test MSE do you obtain? Use the `importance()` function to determine which variables are most important.

```
set.seed(111)
```

```
bcs= randomForest(Sales ~ ., data = TrainingSet, ntree = 1000, mtry = 5, importance = TRUE)
mean((TestSet$Sales - predict(bcs, TestSet))^2)
```

```
## [1] 2.579638
```

```
importance(bcs)
```

	%IncMSE	IncNodePurity
## CompPrice	24.6383726	163.627681
## Income	2.4041882	107.579418
## Advertising	24.8895406	173.588911
## Population	-1.8399948	88.974481
## Price	62.1443853	474.262763
## ShelveLoc	63.2162936	419.148421
## Age	13.0479944	149.897709
## Education	-0.9657177	49.671722
## Urban	-3.0906280	8.831572
## US	2.3295792	17.817760

The MSE is 2.57 and the four important predictors for Sale are CompPrice, Advertising, Price, and Shelve Location.

## Question 5 (based on JWTH Chapter 8, Problem 10)

Use boosting (and bagging) to predict Salary in the Hitters data set

- (a) Remove the observations for which salary is unknown, and then log-transform the salaries

```
dfh <- Hitters
dfh <- na.omit(dfh)
dfh$LogSalary <- log(dfh$Salary)
```

- (b) Split the data into training and testing sets for cross validation purposes.

```
set.seed(50)
n <- 1:floor(nrow(dfh)/2)
HitTrainingSet <- dfh[n, ]
HitTestSet <- dfh[-n, ]
dim(HitTrainingSet)

## [1] 131 21
dim(HitTestSet)

## [1] 132 21
```

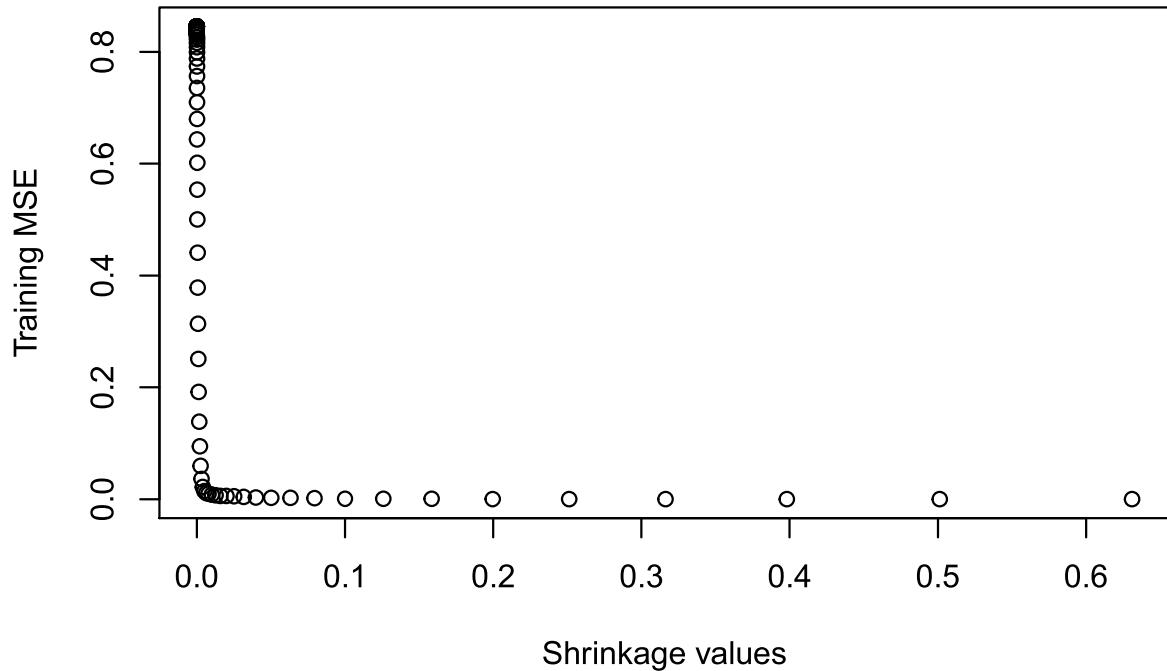
- (c) Perform boosting on the training set with 1000 trees for a range of values of the shrinkage parameter  $\lambda$ . Produce a plot with different shrinkage parameters on the x-axis and the corresponding training set MSE on the y-axis

```
set.seed(66)
lambda <- 10^(seq(-10, -0.2, by = 0.1))
length(lambda)

## [1] 99

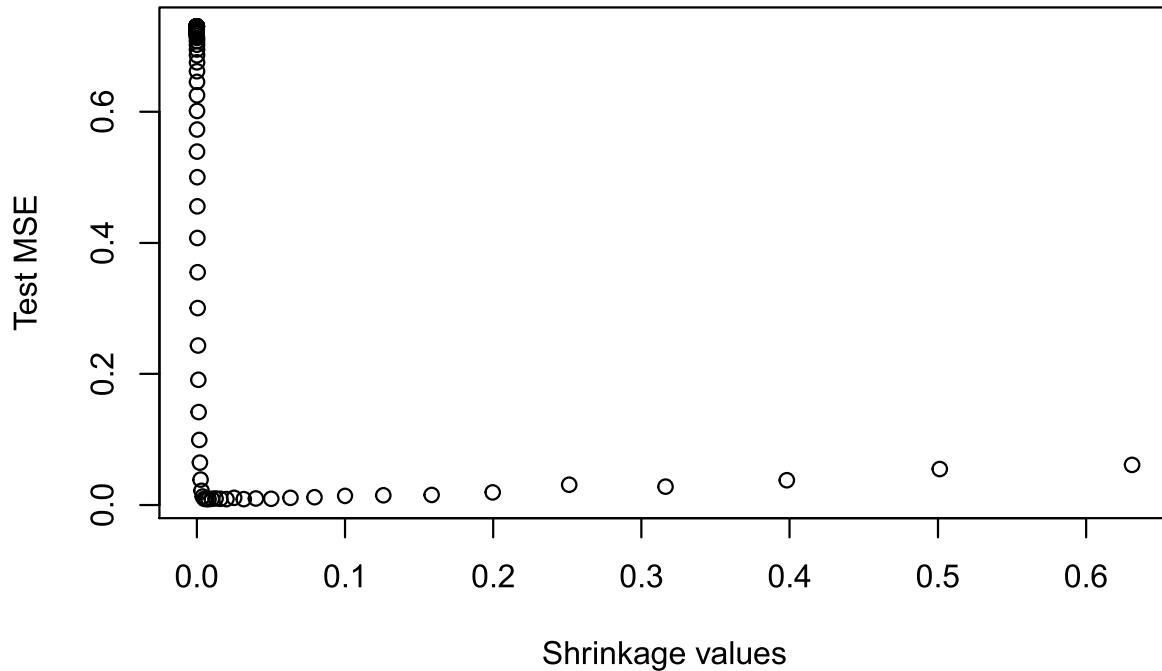
trerr <- rep(NA, 99)
set.seed(66)
for (i in 1:99) {
  boost.hitters <- gbm(LogSalary ~ ., data = HitTrainingSet, distribution = "gaussian", n.trees = 1000)
  pred.train <- predict(boost.hitters, HitTrainingSet, n.trees = 1000)
  trerr[i] <- mean((pred.train - HitTrainingSet$LogSalary)^2)
}

plot(lambda, trerr, type = "p", xlab = "Shrinkage values", ylab = "Training MSE")
```



(d) Produce a plot similar to the last one, but this time using the test set MSE

```
testerr <- rep(NA, 99)
set.seed(22)
for (i in 1:99) {
  boost.hitters <- gbm(LogSalary ~ ., data = HitTrainingSet, distribution = "gaussian", n.trees = 1000)
  pred.test <- predict(boost.hitters, HitTestSet, n.trees = 1000)
  testerr[i] <- mean((pred.test - HitTestSet$LogSalary)^2)
}
plot(lambda, testerr, type ="p", xlab = "Shrinkage values", ylab = "Test MSE")
```



- (e) Fit the model using two other regression techniques (from previous classes) and compare the MSE of those techniques to the results of these boosted trees.

```
fitlm <- lm(LogSalary ~ ., data = HitTrainingSet)
mean((HitTestSet$LogSalary - predict(fitlm, HitTestSet))^2)

## [1] 0.1144346

x.train <- model.matrix(LogSalary ~ ., data = HitTrainingSet)

x.test <- model.matrix(LogSalary ~ ., data = HitTestSet)

y <- HitTrainingSet$LogSalary

lassofit <- glmnet(x.train, y, alpha = 1)

mean((HitTestSet$LogSalary - predict(lassofit, s = 0.01, newx = x.test))^2)

## [1] 0.09814853
```

- (f) Reproduce (c) and (d), but this time use bagging instead of boosting and compare to the boosted MSE's and the MSE's from (e)

```
set.seed(333)
lambda <- 10^(seq(-10, -0.2, by = 0.1))
length(lambda)

## [1] 99
```

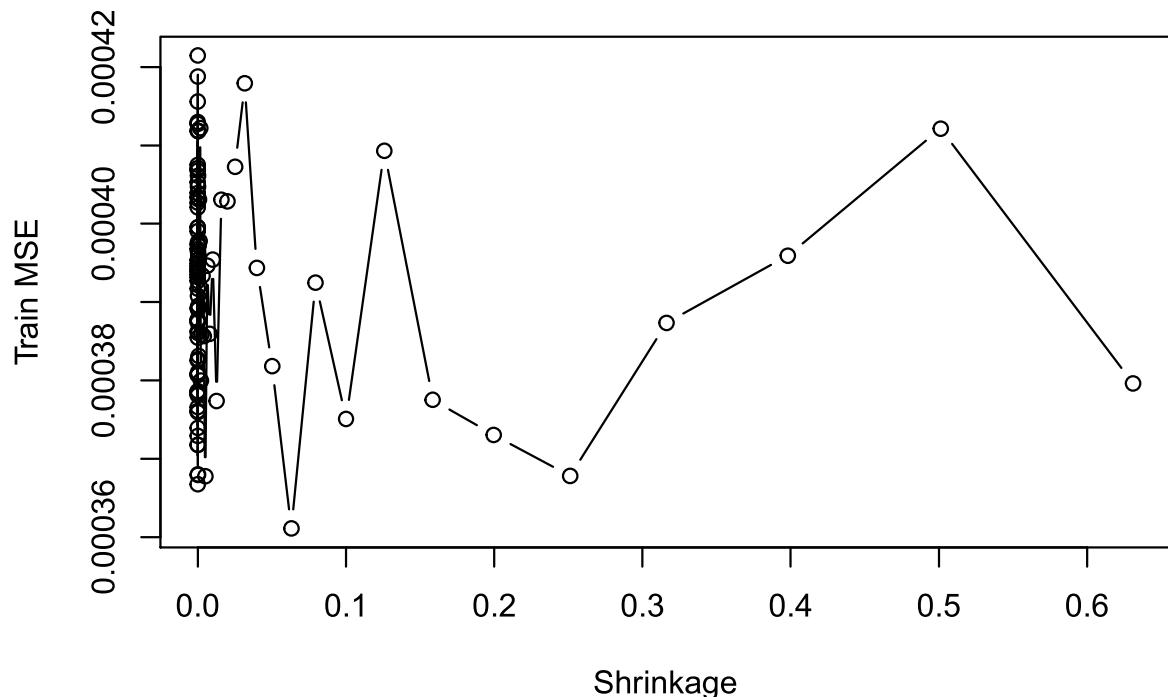
```

trainerror <- rep(NA, 99)
testerror <- rep(NA, 99)

for (i in 1:99) {
  baghit <- randomForest(LogSalary ~ ., data = HitTrainingSet, ntree = 1000, mtry = 20)
  trainerror[i] <- mean((HitTrainingSet$LogSalary - predict(baghit, HitTrainingSet, n.trees = 1000))^2)
  testerror[i] <- mean((HitTestSet$LogSalary - predict(baghit, HitTestSet, n.trees = 1000))^2)
}

plot(lambda, trainerror, type = "b", xlab = "Shrinkage", ylab = "Train MSE")

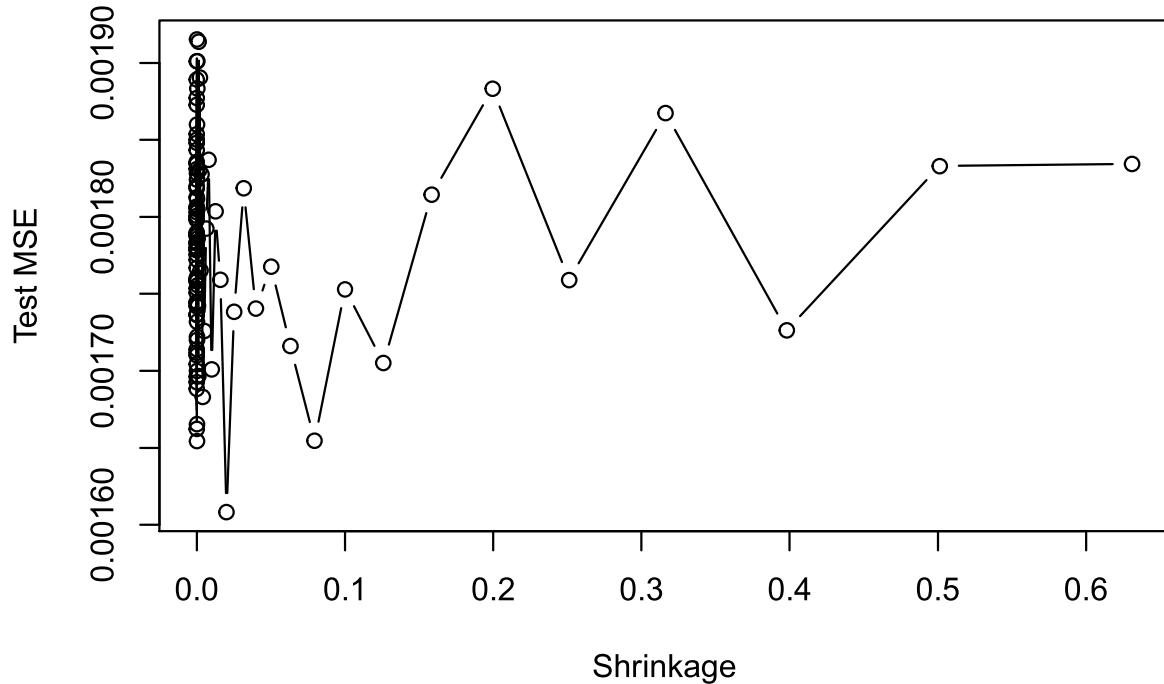
```



```

plot(lambda, testerror, type = "b", xlab = "Shrinkage", ylab = "Test MSE")

```



```
mean((HitTestSet$LogSalary - predict(baghit, newdata = HitTestSet))^2)
## [1] 0.001834307
```

```
mean((HitTrainingSet$LogSalary - predict(baghit, newdata = HitTrainingSet))^2)
```

```
## [1] 0.0003796239
```

The test MSE for bagging is lower than the MSE for boosting.