

Statistical Learning for Biological and Information Systems

Problem Set #2

Chenye Yang cy2540@columbia.edu

October 11, 2019

P1

(a)

Ans:

Ridge regression optimization problem:

$$\begin{aligned}\min \quad & RSS + \lambda \sum_{j=1}^p \beta_j^2 \\ &= \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= (y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 + \lambda \beta_1^2 + \lambda \beta_2^2\end{aligned}\tag{1}$$

(b)

Ans:

From the statement, we know that $x_{11} = x_{12} = a$, $x_{21} = x_{22} = b$, $a + b = 0$, $\beta_0 = 0$.

The equation (1) write as:

$$\begin{aligned}\min \quad & [y_1 - (\beta_1 + \beta_2)a]^2 + [y_2 - (\beta_1 + \beta_2)b]^2 + \lambda(\beta_1^2 + \beta_2^2) \\ &= y_1^2 + y_2^2 - (\beta_1 + \beta_2)(2y_1a + 2y_2b) + (a^2 + b^2)(\beta_1 + \beta_2)^2 + \lambda(\beta_1^2 + \beta_2^2) \\ &= f(\beta_1, \beta_2)\end{aligned}\tag{2}$$

To solve this optimization problem, let:

$$\begin{aligned}\frac{\partial f(\beta_1, \beta_2)}{\partial \beta_1} &= 0 \\ \frac{\partial f(\beta_1, \beta_2)}{\partial \beta_2} &= 0\end{aligned}\tag{3}$$

We have:

$$\begin{aligned}-(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + 2\lambda\hat{\beta}_1 &= 0 \\ -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + 2\lambda\hat{\beta}_2 &= 0\end{aligned}\tag{4}$$

Therefore,

$$2\lambda\hat{\beta}_1 = 2\lambda\hat{\beta}_2 \quad (5)$$

Because in ridge regression, $\lambda > 0$. Thus:

$$\hat{\beta}_1 = \hat{\beta}_2$$

(c)

Ans:

Lasso optimization problem:

$$\begin{aligned} \min \quad & RSS + \lambda \sum_{j=1}^p |\beta_j| \\ &= \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| \\ &= (y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 + \lambda|\beta_1| + \lambda|\beta_2| \end{aligned} \quad (6)$$

(d)

Ans:

Same as **P1 (b)**, we write the lasso optimization equation (6) as:

$$\begin{aligned} \min \quad & y_1^2 + y_2^2 - (\beta_1 + \beta_2)(2y_1a + 2y_2b) + (a^2 + b^2)(\beta_1 + \beta_2)^2 + \lambda(|\beta_1| + |\beta_2|) \\ &= g(\beta_1, \beta_2) \end{aligned} \quad (7)$$

Then

$$\begin{aligned} \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_1} &= \begin{cases} -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + \lambda, & \hat{\beta}_1 > 0 \\ -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) - \lambda, & \hat{\beta}_1 < 0 \end{cases} \\ \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_2} &= \begin{cases} -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + \lambda, & \hat{\beta}_2 > 0 \\ -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) - \lambda, & \hat{\beta}_2 < 0 \end{cases} \end{aligned} \quad (8)$$

(1) When $\hat{\beta}_1, \hat{\beta}_2 > 0$

$$\begin{cases} \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_1} = 0 \\ \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_2} = 0 \end{cases}$$

\Rightarrow

$$-(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + \lambda = 0$$

\Rightarrow

$$\begin{aligned} \hat{\beta}_1 + \hat{\beta}_2 &= \frac{(2y_1a + 2y_2b) - \lambda}{2(a^2 + b^2)} \\ &= \frac{y_1a + y_2b - \lambda/2}{(a + b)^2 - 2ab} \\ &= -\frac{y_1}{2b} - \frac{y_2}{2a} + \frac{\lambda}{4ab} \end{aligned} \quad (9)$$

Therefore, any $\hat{\beta}_1, \hat{\beta}_2$ satisfies equation (9) and $\hat{\beta}_1, \hat{\beta}_2 > 0$ will be a solution to the lasso optimization problem in **P1 (c)**.

(2) Similarly, when $\hat{\beta}_1, \hat{\beta}_2 < 0$

$$\begin{cases} \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_1} = 0 \\ \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_2} = 0 \end{cases}$$

\Rightarrow

$$-(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) - \lambda = 0$$

\Rightarrow

$$\hat{\beta}_1 + \hat{\beta}_2 = -\frac{y_1}{2b} - \frac{y_2}{2a} - \frac{\lambda}{4ab} \quad (10)$$

Therefore, any $\hat{\beta}_1, \hat{\beta}_2$ satisfies equation (10) and $\hat{\beta}_1, \hat{\beta}_2 < 0$ will be a solution to the lasso optimization problem in **P1 (c)**.

(3) When $\hat{\beta}_1 < 0, \hat{\beta}_2 > 0$ or $\hat{\beta}_1 > 0, \hat{\beta}_2 < 0$

$$\begin{cases} \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_1} = 0 \\ \frac{\partial g(\beta_1, \beta_2)}{\partial \beta_2} = 0 \end{cases}$$

\Rightarrow

$$-(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) - \lambda = -(2y_1a + 2y_2b) + 2(a^2 + b^2)(\hat{\beta}_1 + \hat{\beta}_2) + \lambda = 0$$

$\Rightarrow \lambda = 0$, which is not Lasso.

(4) When $\hat{\beta}_1 = 0$

Equation (7) changes to

$$\min \quad y_1^2 + y_2^2 - \beta_2(2y_1a + 2y_2b) + (a^2 + b^2)\beta_2^2 + \lambda(|\beta_2|) = g(\beta_2)$$

$$\frac{dg(\beta_2)}{d\beta_2} = -(2y_1a + 2y_2b) + 2(a^2 + b^2)\beta_2 \pm \lambda = 0$$

\Rightarrow

$$\hat{\beta}_2 = -\frac{y_1}{2b} - \frac{y_2}{2a} \pm \frac{\lambda}{4ab}, \quad \hat{\beta}_2 \geq 0 \quad (11)$$

Therefore, $\hat{\beta}_1 = 0$ and $\hat{\beta}_2$ satisfies equation (11) will be a solution to the lasso optimization problem in **P1 (c)**.

(5) When $\hat{\beta}_2 = 0$

Equation (7) changes to

$$\min \quad y_1^2 + y_2^2 - \beta_1(2y_1a + 2y_2b) + (a^2 + b^2)\beta_1^2 + \lambda(|\beta_1|) = g(\beta_1)$$

$$\frac{dg(\beta_1)}{d\beta_1} = -(2y_1a + 2y_2b) + 2(a^2 + b^2)\beta_1 \pm \lambda = 0$$

\Rightarrow

$$\hat{\beta}_1 = -\frac{y_1}{2b} - \frac{y_2}{2a} \pm \frac{\lambda}{4ab}, \quad \hat{\beta}_1 \geq 0 \quad (12)$$

Therefore, $\hat{\beta}_2 = 0$ and $\hat{\beta}_1$ satisfies equation (12) will be a solution to the lasso optimization problem in **P1 (c)**.

(6) When $\hat{\beta}_1 = \hat{\beta}_2 = 0$

Equation (7) changes to $\min y_1^2 + y_2^2$, which means the regression function is $y = 0$.

In summary, in this setting, the lasso coefficients $\hat{\beta}_1, \hat{\beta}_2$ are not unique.

P2

(a)

Ans:

When $p = 1$, equation (1) in the question changes to:

$$\begin{aligned} \min \sum_{j=1}^1 (y_j - \beta_j)^2 + \lambda \sum_{j=1}^1 \beta_j^2 \\ = (y_1 - \beta_1)^2 + \lambda \beta_1^2 = f(\beta_1) \end{aligned}$$

Let $y_1 = 1$, $\lambda = 0.5$, plot the ridge regression optimization function as black line in Figure 1. Also plot the ridge regression estimates equation (3) in question, $\hat{\beta}_1^R = y_1/(1 + \lambda)$, as a vertical red line in same figure. Two lines intersect at the small red circle, which is exactly on the lowest point of curve. Considering the lowest point of curve means the solve of ridge regression optimization problem, the estimated coefficient point. Therefore, the plot confirms that (1) is solved by (3)

Result:

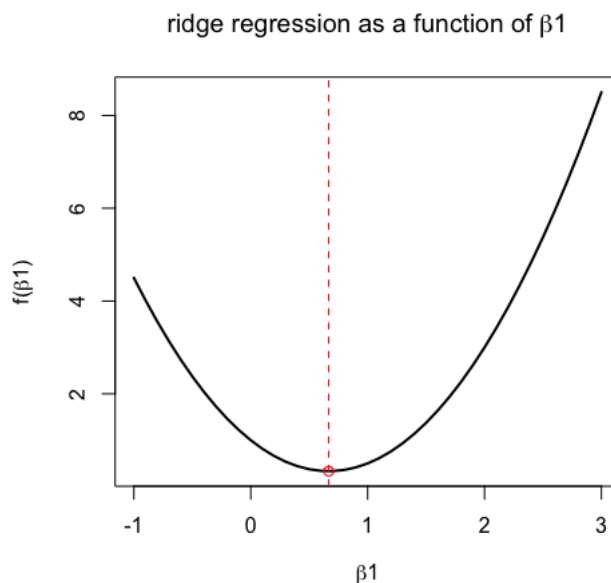


Figure 1: Ridge regression optimization function over one coefficient

Code:

```
1 # (a)
2 y1 = 1
3 lambda = 0.5
4 beta1 = seq(from=-1, to=3, by=0.1)
5 f = (y1 - beta1)^2 + lambda*(beta1^2)
```

```

6 plot(beta1, f, type="l", lwd=2, xlab=expression(paste(beta, '1')),
      ylab=expression(paste('f(', beta, '1)')), main=expression(paste("
      ridge_regression_as_a_function_of", beta, '1')))
7 abline(v=y1/(1+lambda), lty=2, col="red")
8 beta1 = y1/(1+lambda)
9 f = (y1 - beta1)^2 + lambda*(beta1^2)
10 points(beta1, f, col="red")

```

(b)**Ans:**

When $p = 1$, equation (2) in the question changes to:

$$\min \sum_{j=1}^1 (y_j - \beta_j)^2 + \lambda \sum_{j=1}^1 |\beta_j|$$

$$= (y_1 - \beta_1)^2 + \lambda |\beta_1| = g(\beta_1)$$

For $y_1 > \lambda/2$, let $y_1 = 1, \lambda = 0.5$, plot the lasso optimization function as black line in Fig.2. Also plot the lasso estimates equation (4) in question, $\hat{\beta}_1^L = y_j - \lambda/2$, as a vertical red line in same figure.

For $y_1 < -\lambda/2$, let $y_1 = -1, \lambda = 0.5$, plot the lasso optimization function as black line in Fig.3. Also plot the lasso estimates equation (4) in question, $\hat{\beta}_1^L = y_j + \lambda/2$, as a vertical red line in same figure.

For $|y_1| \leq \lambda/2$, let $y_1 = 0.1, \lambda = 0.5$, plot the lasso optimization function as black line in Fig.4. Also plot the lasso estimates equation (4) in question, $\hat{\beta}_1^L = 0$, as a vertical red line in same figure.

In these figures, two lines intersect at the small red circle, which is exactly on the lowest point of curve. Considering the lowest point of curve means the solve of lasso optimization problem, the estimated coefficient point. In summary, the plot confirms that (2) is solved by (4)

Code:

```

1 # (b)
2 # y1 > lambda/2
3 y1 = 1
4 lambda = 0.5
5 beta1 = seq(from=-1, to=3, by=0.1)
6 g = (y1 - beta1)^2 + lambda*(abs(beta1))
7 plot(beta1, g, type="l", lwd=2, xlab=expression(paste(beta, '1')),
      ylab=expression(paste('g(', beta, '1)')), main=expression(paste("
      lasso_as_a_function_of", beta, '1')))
8 abline(v=y1-lambda/2, lty=2, col="red")
9 beta1 = y1-lambda/2
10 g = (y1 - beta1)^2 + lambda*(abs(beta1))
11 points(beta1, g, col="red")
12 # y1 < -lambda/2

```

```

13 y1 = -1
14 lambda = 0.5
15 beta1 = seq(from=-3, to=1, by=0.1)
16 g = (y1 - beta1)^2 + lambda*(abs(beta1))
17 plot(beta1, g, type="l", lwd=2, xlab=expression(paste(beta, '1')),
      ylab=expression(paste('g(', beta, ')')), main=expression(paste("
      lasso as a function of ", beta, '1')))
18 abline(v=y1+lambda/2, lty=2, col="red")
19 beta1 = y1+lambda/2
20 g = (y1 - beta1)^2 + lambda*(abs(beta1))
21 points(beta1, g, col="red")
22 # |y1| <= lambda/2
23 y1 = 0.1
24 lambda = 0.5
25 beta1 = seq(from=-2, to=2, by=0.1)
26 g = (y1 - beta1)^2 + lambda*(abs(beta1))
27 plot(beta1, g, type="l", lwd=2, xlab=expression(paste(beta, '1')),
      ylab=expression(paste('g(', beta, ')')), main=expression(paste("
      lasso as a function of ", beta, '1')))
28 abline(v=0, lty=2, col="red")
29 beta1 = 0
30 g = (y1 - beta1)^2 + lambda*(abs(beta1))
31 points(beta1, g, col="red")

```

Result:

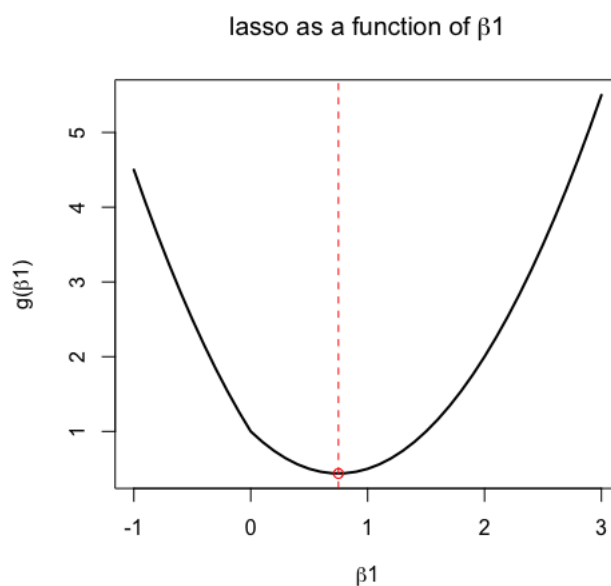


Figure 2: Lasso optimization function over one coefficient, $y_1 > \lambda/2$

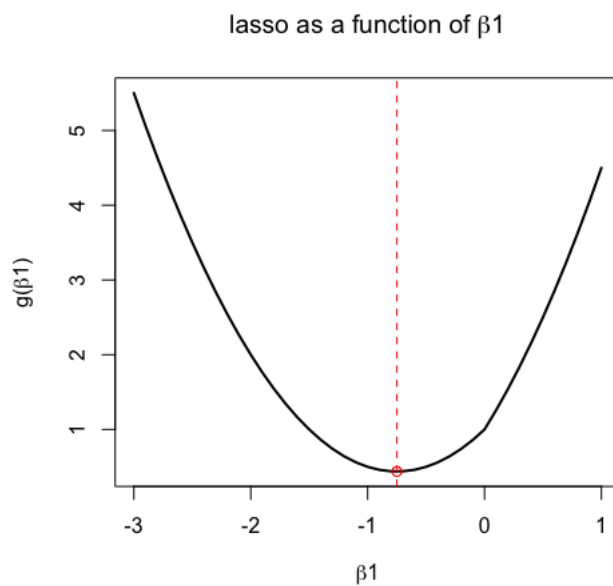


Figure 3: Lasso optimization function over one coefficient, $y_1 < -\lambda/2$

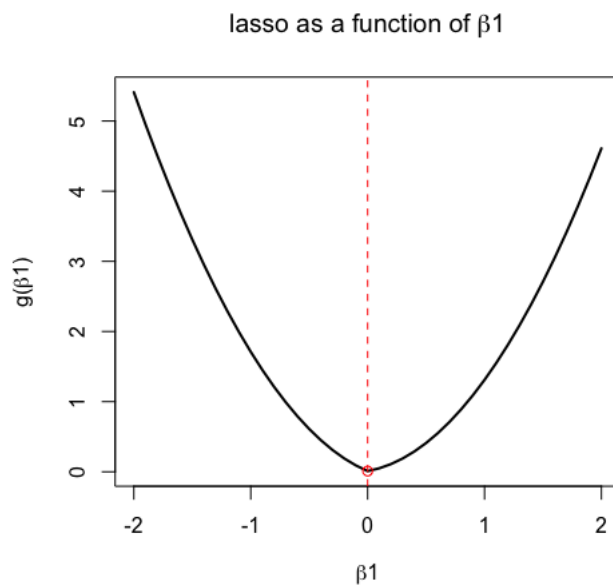


Figure 4: Lasso optimization function over one coefficient, $|y_1| \leq \lambda/2$

P3

(a)

Ans:

Because $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are independent and identically distributed from a $N(0, \sigma^2)$ distribution, i.e., the data points are normally distributed about the regression line, the likelihood of the data can be write as follow:

$$\mathcal{L}(\boldsymbol{\beta}|\mathbf{y}) := P(\mathbf{y}|\boldsymbol{\beta}) = \prod_{i=1}^n P_Y(y_i|\boldsymbol{\beta}, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} \quad (13)$$

(b)

Ans:

Because $\beta_1, \beta_2, \dots, \beta_p$ are independent and identically distributed according to a double-exponential distribution, and $p(\boldsymbol{\beta}) = \frac{1}{2^p b} \exp(-|\boldsymbol{\beta}|/b)$, $|\boldsymbol{\beta}| = \sum_{j=1}^p |\beta_j|$, the posterior for $\boldsymbol{\beta}$ in this setting is:

$$\begin{aligned} P(\boldsymbol{\beta}|\mathbf{y}) &= \frac{P(\mathbf{y}|\boldsymbol{\beta})P(\boldsymbol{\beta})}{P(\mathbf{y})} \\ &= \frac{1}{\prod_{i=1}^n P(y_i)} \times \frac{1}{2^p b} e^{-\sum_{j=1}^p |\beta_j|/b} \times \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} \end{aligned} \quad (14)$$

(c)

Ans:

From the view of maximum a posterior probability estimate (usually abbreviated by MAP), we should find the maximum of the posterior in equation (14) in order to find the best estimate of predictors $\boldsymbol{\beta}$:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{MAP}} &= \arg \max_{\boldsymbol{\beta}} P(\boldsymbol{\beta}|\mathbf{y}) \\ &= \arg \max_{\boldsymbol{\beta}} \frac{1}{\prod_{i=1}^n P(y_i)} \times \frac{1}{2^p b} e^{-\sum_{j=1}^p |\beta_j|/b} \times \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} \\ &= \arg \max_{\boldsymbol{\beta}} \left[\log \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} + \log \frac{1}{2^p b} e^{-\sum_{j=1}^p |\beta_j|/b} \right] \\ &= \arg \max_{\boldsymbol{\beta}} \left[\sum_{i=1}^n -\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2} - \sum_{j=1}^p |\beta_j|/b \right] \\ &= \arg \min_{\boldsymbol{\beta}} \frac{1}{2\sigma^2} \left[\sum_{i=1}^n [y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2 + \frac{2\sigma^2}{b} \sum_{j=1}^p |\beta_j| \right] \\ &= \arg \min_{\boldsymbol{\beta}} \left[\sum_{i=1}^n [y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2 + \lambda \sum_{j=1}^p |\beta_j|, \quad \lambda = \frac{2\sigma^2}{b} \right] \end{aligned} \quad (15)$$

The above equation (15) is exactly the form of Lasso estimation. When b is small, meaning the density of double-exponential distribution has a sharp increase at its mean, β_j is more likely to have the mean value 0. So some of the coefficient of predictors can be 0, which is the variable selection. When b is large, β_j is more possible to have any value, just like the least squared estimation. Therefore, the lasso estimate is the mode for β under this posterior distribution.

(d)

Ans:

Because $\beta_1, \beta_2, \dots, \beta_p$ are independent and identically distributed according to a normal distribution with mean 0 and variance c , the posterior for β in this setting is:

$$\begin{aligned} P(\beta|\mathbf{y}) &= \frac{P(\mathbf{y}|\beta)P(\beta)}{P(\mathbf{y})} \\ &= \frac{1}{\prod_{i=1}^n P(y_i)} \times \prod_{j=1}^p \frac{1}{\sqrt{c}\sqrt{2\pi}} e^{-\frac{\beta_j^2}{2c}} \times \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} \end{aligned} \quad (16)$$

(e)

Ans: From the view of maximum a posterior probability estimate (usually abbreviated by MAP), we should find the maximum of the posterior in equation (16) in order to find the best estimate of predictors β :

$$\begin{aligned} \hat{\beta}_{\text{MAP}} &= \arg \max_{\beta} P(\beta|\mathbf{y}) \\ &= \arg \max_{\beta} \frac{1}{\prod_{i=1}^n P(y_i)} \times \prod_{j=1}^p \frac{1}{\sqrt{c}\sqrt{2\pi}} e^{-\frac{\beta_j^2}{2c}} \times \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} \\ &= \arg \max_{\beta} \left[\log \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2}} + \log \prod_{j=1}^p \frac{1}{\sqrt{c}\sqrt{2\pi}} e^{-\frac{\beta_j^2}{2c}} \right] \\ &= \arg \max_{\beta} \left[\sum_{i=1}^n -\frac{[y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2}{2\sigma^2} - \sum_{j=1}^p \frac{\beta_j^2}{2c} \right] \\ &= \arg \min_{\beta} \frac{1}{2\sigma^2} \left[\sum_{i=1}^n [y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2 + \frac{\sigma^2}{c} \sum_{j=1}^p \beta_j^2 \right] \\ &= \arg \min_{\beta} \left[\sum_{i=1}^n [y_i - (\beta_0 + \sum_{j=1}^p x_{ij}\beta_j)]^2 + \lambda \sum_{j=1}^p \beta_j^2 \right], \quad \lambda = \frac{\sigma^2}{c} \end{aligned} \quad (17)$$

The above equation (17) is exactly the form of ridge regression estimation. When c is small, meaning the density is more shrunken to its mean value 0, β_j is more likely to have a value near 0. So the coefficient of predictors are restricted small. When c is large, β_j is more possible to have any value, just like the least squared estimation. Therefore, the ridge regression estimate is both the mode and the mean for β under this posterior distribution.

P4

(a)

Code:

```

1 # (a)
2 set.seed(1)
3 X = rnorm(100, mean = 0, sd = 1)
4 e = rnorm(100, mean = 0, sd = 1)

```

Result:

```

1 > summary(X)
2   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
3 -2.2147 -0.4942  0.1139  0.1089  0.6915  2.4016
4 > summary(e)
5   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
6 -1.91436 -0.65105 -0.17722 -0.03781  0.50090  2.30798

```

(b)

Code:

```

1 # (b)
2 beta0 = 1
3 beta1 = 2
4 beta2 = 3
5 beta3 = 5
6 Y = beta0 + beta1*X + beta2*X^2 + beta3*X^3 + e

```

Result:

```

1 > summary(Y)
2   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
3 -43.6798 -0.3189  1.9192  4.6648  5.7107  92.7915

```

(c)

Code:

```

1 # (c)
2 # create a single data set containing both X and Y
3 simulatedData = data.frame(
4   response = Y,
5   X_1 = X,
6   X_2 = X^2,
7   X_3 = X^3,
8   X_4 = X^4,

```

```

9   X_5 = X^5,
10  X_6 = X^6,
11  X_7 = X^7,
12  X_8 = X^8,
13  X_9 = X^9,
14  X_10 = X^10,
15  stringsAsFactors = FALSE
16 )
17 str(simulatedData) # Get the structure of the data frame.
18 library(leaps)
19 # perform best subset selection, using 10 predictors
20 regfit.full = regsubsets(response~., simulatedData, nvmax = 10)
21 reg.summary = summary(regfit.full)
22 # RSS
23 plot(reg.summary$rss, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="RSS", type = "l")
24 # Adjusted RSq
25 plot(reg.summary$adjr2, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="Adjusted_RSq", type = "l")
26 abline(v=which.max(reg.summary$adjr2), lty=2, col="red")
27 points(which.max(reg.summary$adjr2), reg.summary$adjr2[which.max(reg.
    summary$adjr2)], col="red")
28 coef(regfit.full, which.max(reg.summary$adjr2))
29 # Cp
30 plot(reg.summary$cp, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="Cp", type = "l")
31 abline(v=which.min(reg.summary$cp), lty=2, col="red")
32 points(which.min(reg.summary$cp), reg.summary$cp[which.min(reg.summary$
    cp)], col="red")
33 coef(regfit.full, which.min(reg.summary$cp))
34 # BIC
35 plot(reg.summary$bic, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="BIC", type = "l")
36 abline(v=which.min(reg.summary$bic), lty=2, col="red")
37 points(which.min(reg.summary$bic), reg.summary$bic[which.min(reg.
    summary$bic)], col="red")
38 coef(regfit.full, which.min(reg.summary$bic))

```

Result:

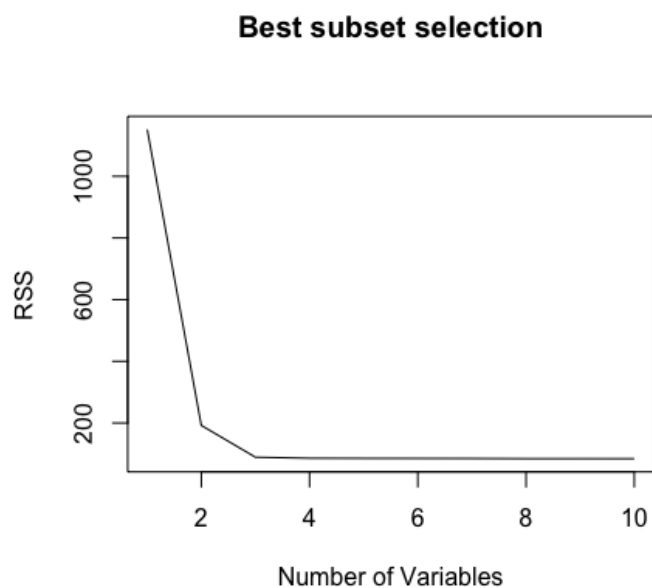
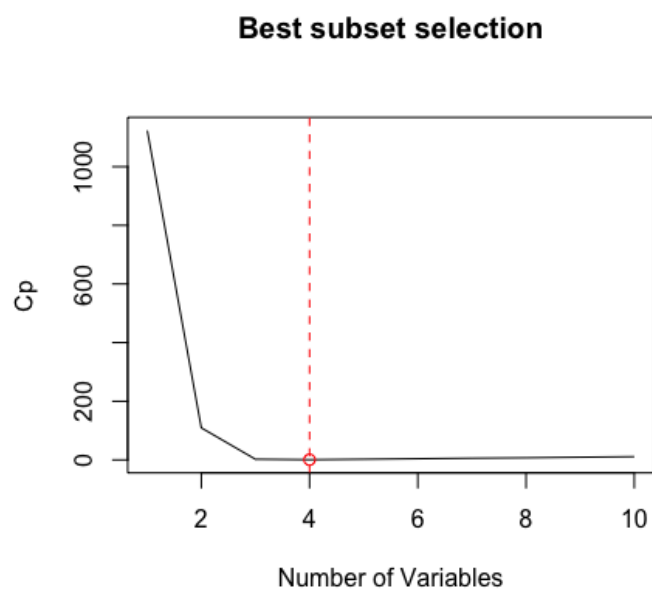


Figure 5: RSS in best subset selection

Figure 6: Best model obtained according to C_p

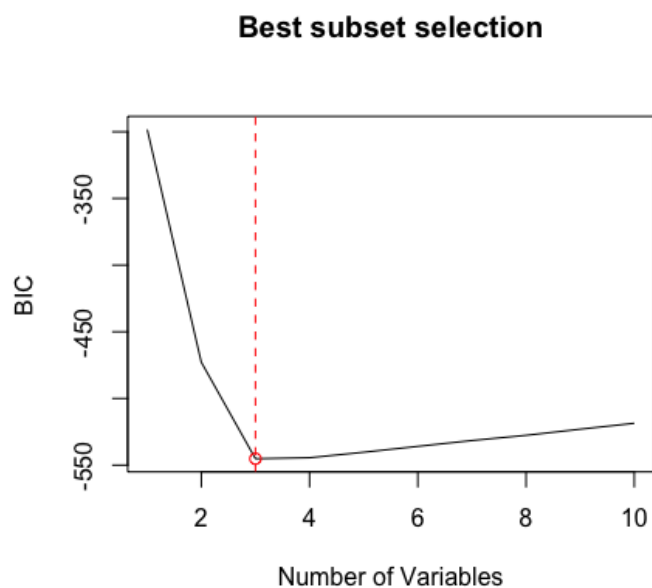


Figure 7: Best model obtained according to BIC

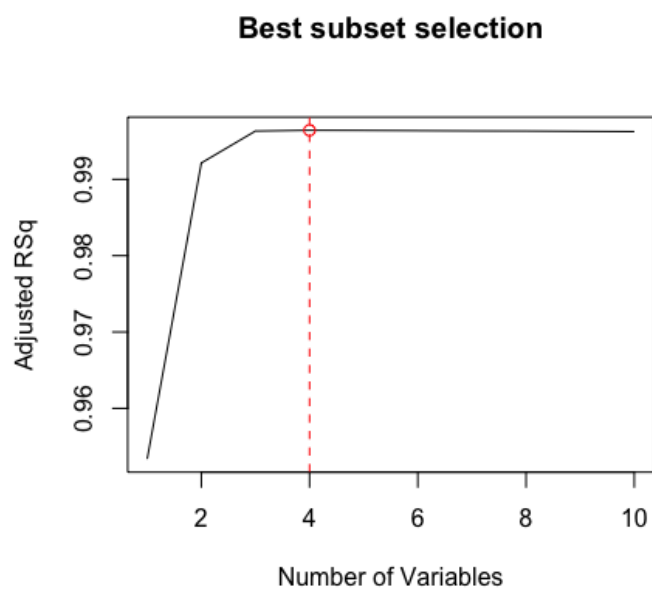


Figure 8: Best model obtained according to $Adjusted R^2$

The training data:

```

1 > str(simulatedData) # Get the structure of the data frame.
2 'data.frame': 100 obs. of 11 variables:
3 $ response: num -0.925 1.542 -2.405 32.283 1.509 ...
4 $ X_1 : num -0.626 0.184 -0.836 1.595 0.33 ...
5 $ X_2 : num 0.3924 0.0337 0.6983 2.5449 0.1086 ...
6 $ X_3 : num -0.24585 0.00619 -0.5835 4.05986 0.03578 ...
7 $ X_4 : num 0.15401 0.00114 0.48759 6.47662 0.01179 ...
8 $ X_5 : num -0.096482 0.000209 -0.407443 10.332031 0.003884 ...
9 $ X_6 : num 6.04e-02 3.84e-05 3.40e-01 1.65e+01 1.28e-03 ...
10 $ X_7 : num -3.79e-02 7.04e-06 -2.85e-01 2.63e+01 4.22e-04 ...
11 $ X_8 : num 2.37e-02 1.29e-06 2.38e-01 4.19e+01 1.39e-04 ...
12 $ X_9 : num -1.49e-02 2.38e-07 -1.99e-01 6.69e+01 4.58e-05 ...
13 $ X_10 : num 9.31e-03 4.36e-08 1.66e-01 1.07e+02 1.51e-05 ...

```

The coefficients of the best model obtained according to C_p :

```

1 > coef(regfit.full, which.min(reg.summary$cp))
2 (Intercept)      X_1      X_2      X_3      X_5
3 1.07200775 2.38745596 2.84575641 4.55797426 0.08072292

```

The coefficients of the best model obtained according to BIC :

```

1 > coef(regfit.full, which.min(reg.summary$bic))
2 (Intercept)      X_1      X_2      X_3
3 1.061507 1.975280 2.876209 5.017639

```

The coefficients of the best model obtained according to $Adjusted R^2$:

```

1 > coef(regfit.full, which.max(reg.summary$adjr2))
2 (Intercept)      X_1      X_2      X_3      X_5
3 1.07200775 2.38745596 2.84575641 4.55797426 0.08072292

```

Ans:

From Figure 6, the best model using best subset selection method according to C_p contains 4 predictors. The best model obtained is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_5 X^5 \\ &= 1.072 + 2.387X + 2.846X^2 + 4.558X^3 + 0.081X^5\end{aligned}$$

From Figure 7, the best model using best subset selection method according to BIC contains 3 predictors. The best model obtained is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 \\ &= 1.062 + 1.975X + 2.876X^2 + 5.018X^3\end{aligned}$$

From Figure 8, the best model using best subset selection method according to $Adjusted R^2$ contains 4 predictors. The best model obtained is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_5 X^5 \\ &= 1.072 + 2.387X + 2.846X^2 + 4.558X^3 + 0.081X^5\end{aligned}$$

The model from C_p and *Adjusted R^2* still have a noise variable, although whose coefficient is small. The estimated coefficients of signal variables are close to the true value.

(d)

Code:

```

1 # (d)
2 # perform forward stepwise selection , using 10 predictors
3 regfit.fwd = regsubsets(response~., simulatedData , nvmax = 10 , method = "
  forward")
4 reg.summary = summary(regfit.fwd)
5 # RSS
6 plot(reg.summary$rss , main="Forward_stepwise_selection" , xlab="Number_of
  Variables" , ylab="RSS" , type = "l")
7 # Adjusted RSq
8 plot(reg.summary$adjr2 , main="Forward_stepwise_selection" , xlab="Number_of
  Variables" , ylab="Adjusted_RSq" , type = "l")
9 abline(v=which.max(reg.summary$adjr2) , lty=2, col="red")
10 points(which.max(reg.summary$adjr2) , reg.summary$adjr2[which.max(reg.
  summary$adjr2)] , col="red")
11 coef(regfit.fwd , which.max(reg.summary$adjr2))
12 # Cp
13 plot(reg.summary$cp , main="Forward_stepwise_selection" , xlab="Number_of
  Variables" , ylab="Cp" , type = "l")
14 abline(v=which.min(reg.summary$cp) , lty=2, col="red")
15 points(which.min(reg.summary$cp) , reg.summary$cp[which.min(reg.summary$
  cp)] , col="red")
16 coef(regfit.fwd , which.min(reg.summary$cp))
17 # BIC
18 plot(reg.summary$bic , main="Forward_stepwise_selection" , xlab="Number_of
  Variables" , ylab="BIC" , type = "l")
19 abline(v=which.min(reg.summary$bic) , lty=2, col="red")
20 points(which.min(reg.summary$bic) , reg.summary$bic[which.min(reg.
  summary$bic)] , col="red")
21 coef(regfit.fwd , which.min(reg.summary$bic))
22 # perform backward stepwise selection , using 10 predictors
23 regfit.bwd = regsubsets(response~., simulatedData , nvmax = 10 , method = "
  backward")
24 reg.summary = summary(regfit.bwd)
25 # RSS
26 plot(reg.summary$rss , main="Backward_stepwise_selection" , xlab="Number_of
  Variables" , ylab="RSS" , type = "l")
27 # Adjusted RSq
28 plot(reg.summary$adjr2 , main="Backward_stepwise_selection" , xlab="Number
  of Variables" , ylab="Adjusted_RSq" , type = "l")
29 abline(v=which.max(reg.summary$adjr2) , lty=2, col="red")

```



```

30 points(which.max(reg.summary$adjr2), reg.summary$adjr2[which.max(reg.
    summary$adjr2)], col="red")
31 coef(regfit.bwd, which.max(reg.summary$adjr2))
32 # Cp
33 plot(reg.summary$cp, main="Backward_stepwise_selection", xlab="Number_of
    _Variables", ylab="Cp", type = "l")
34 abline(v=which.min(reg.summary$cp), lty=2, col="red")
35 points(which.min(reg.summary$cp), reg.summary$cp[which.min(reg.summary$
    cp)], col="red")
36 coef(regfit.bwd, which.min(reg.summary$cp))
37 # BIC
38 plot(reg.summary$bic, main="Backward_stepwise_selection", xlab="Number_
    of_Variables", ylab="BIC", type = "l")
39 abline(v=which.min(reg.summary$bic), lty=2, col="red")
40 points(which.min(reg.summary$bic), reg.summary$bic[which.min(reg.
    summary$bic)], col="red")
41 coef(regfit.bwd, which.min(reg.summary$bic))

```

Result:

The coefficients of the forward stepwise selection obtained according to C_p :

```

1 > coef(regfit.fwd, which.min(reg.summary$cp))
2 (Intercept)      X_1      X_2      X_3      X_5
3  1.07200775  2.38745596  2.84575641  4.55797426  0.08072292

```

The coefficients of the forward stepwise selection obtained according to BIC :

```

1 > coef(regfit.fwd, which.min(reg.summary$bic))
2 (Intercept)      X_1      X_2      X_3
3  1.061507  1.975280  2.876209  5.017639

```

The coefficients of the forward stepwise selection obtained according to $Adjusted R^2$:

```

1 > coef(regfit.fwd, which.max(reg.summary$adjr2))
2 (Intercept)      X_1      X_2      X_3      X_5
3  1.07200775  2.38745596  2.84575641  4.55797426  0.08072292

```

The coefficients of the backward stepwise selection obtained according to C_p :

```

1 > coef(regfit.bwd, which.min(reg.summary$cp))
2 (Intercept)      X_1      X_2      X_3      X_9
3  1.079236362  2.231905828  2.833494180  4.819555807  0.001290827

```

The coefficients of the backward stepwise selection obtained according to BIC :

```

1 > coef(regfit.bwd, which.min(reg.summary$bic))
2 (Intercept)      X_1      X_2      X_3
3  1.061507  1.975280  2.876209  5.017639

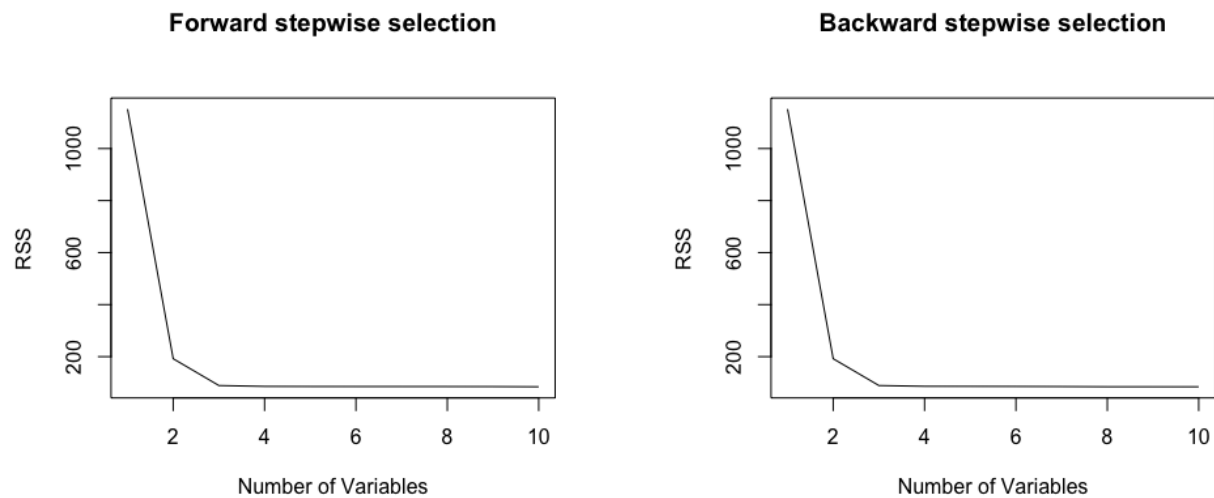
```

The coefficients of the backward stepwise selection obtained according to $Adjusted R^2$:

```

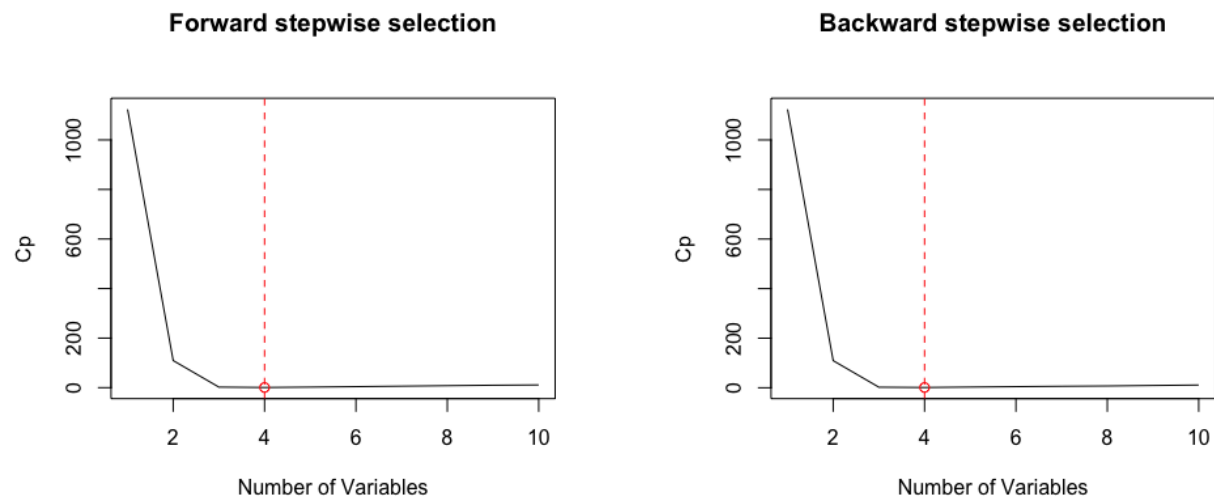
1 > coef(regfit.bwd, which.max(reg.summary$adjr2))
2 (Intercept)      X_1      X_2      X_3      X_9
3 1.079236362  2.231905828  2.833494180  4.819555807  0.001290827

```



(a) Forward stepwise selection

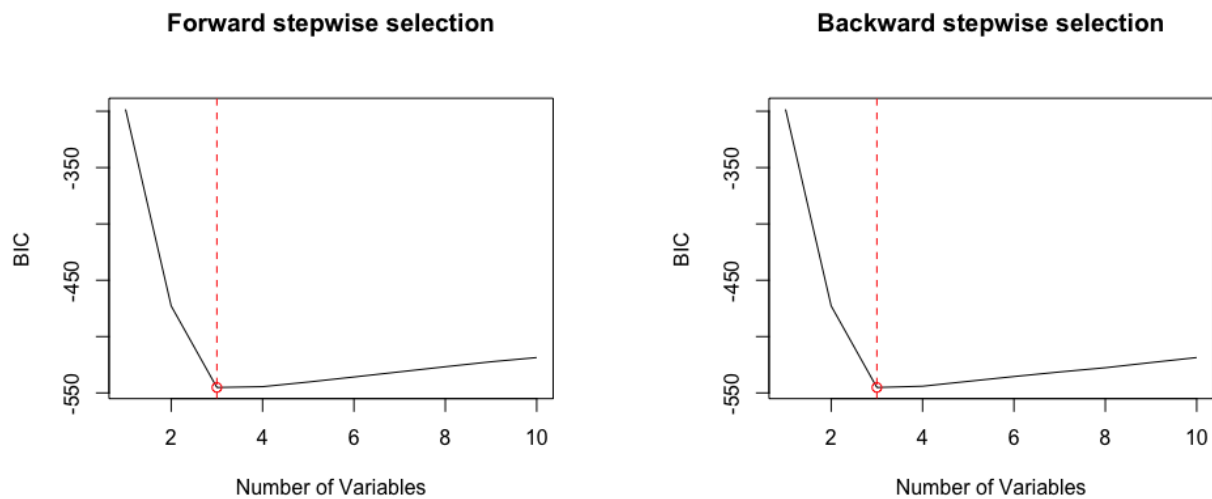
(b) Backward stepwise selection

Figure 9: RSS 

(a) Forward stepwise selection

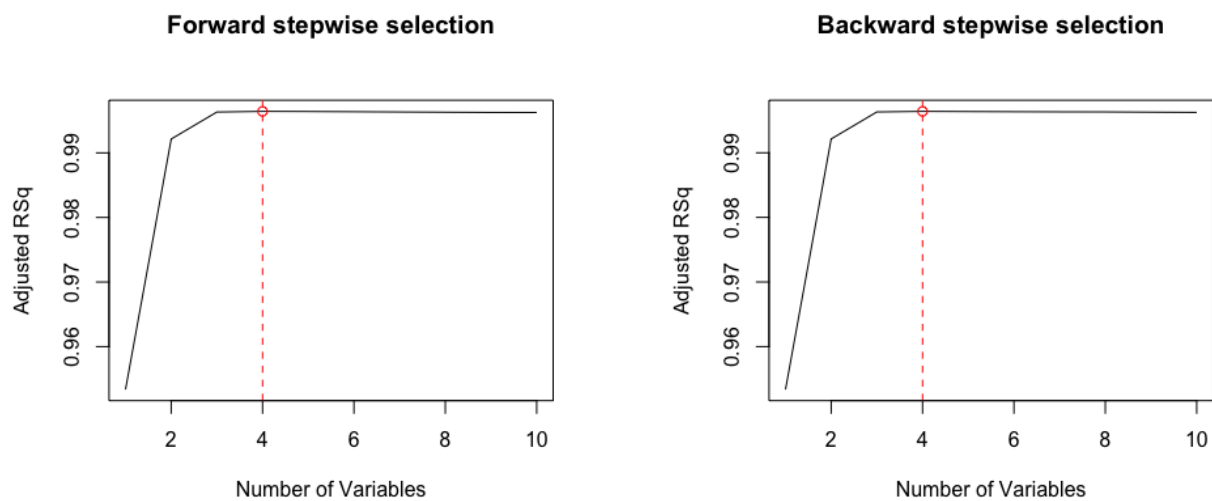
(b) Backward stepwise selection

Figure 10: Best model obtained according to C_p



(a) Forward stepwise selection

(b) Backward stepwise selection

Figure 11: Best model obtained according to BIC 

(a) Forward stepwise selection

(b) Backward stepwise selection

Figure 12: Best model obtained according to $Adjusted R^2$

Ans:

From Figure 10, the best model using forward and backward stepwise selection method according to C_p contains 4 predictors, which is the same number with best subset selection. The best model obtained by forward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_5 X^5 \\ &= 1.072 + 2.387X + 2.846X^2 + 4.558X^3 + 0.081X^5\end{aligned}$$

The best model obtained by backward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_9 X^9 \\ &= 1.079 + 2.232X + 2.833X^2 + 4.820X^3 + 0.001X^9\end{aligned}$$

From Figure 11, the best model using forward stepwise selection method according to BIC contains 3 predictors, which is the same number with best subset selection. The best model obtained by forward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 \\ &= 1.062 + 1.975X + 2.876X^2 + 5.018X^3\end{aligned}$$

The best model obtained by backward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 \\ &= 1.062 + 1.975X + 2.876X^2 + 5.018X^3\end{aligned}$$

From Figure 12, the best model using forward stepwise selection method according to $Adjusted R^2$ contains 4 predictors, which is the same number with best subset selection. The best model obtained by forward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_5 X^5 \\ &= 1.072 + 2.387X + 2.846X^2 + 4.558X^3 + 0.081X^5\end{aligned}$$

The best model obtained by backward stepwise selection is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_9 X^9 \\ &= 1.079 + 2.232X + 2.833X^2 + 4.820X^3 + 0.001X^9\end{aligned}$$

The model from C_p and $Adjusted R^2$ still have a noise variable, although whose coefficient is small. The estimated coefficients of signal variables are close to the true value.

In summary, in this problem setting, the forward stepwise selection gets the same models with best subset selection, no matter what criterion C_p or BIC or $Adjusted R^2$ is used. However, the backward stepwise selection gets the different models with best subset selection, when C_p or $Adjusted R^2$ is used. Meanwhile, the backward stepwise selection gets the same model with best subset selection when BIC is used.

(e)

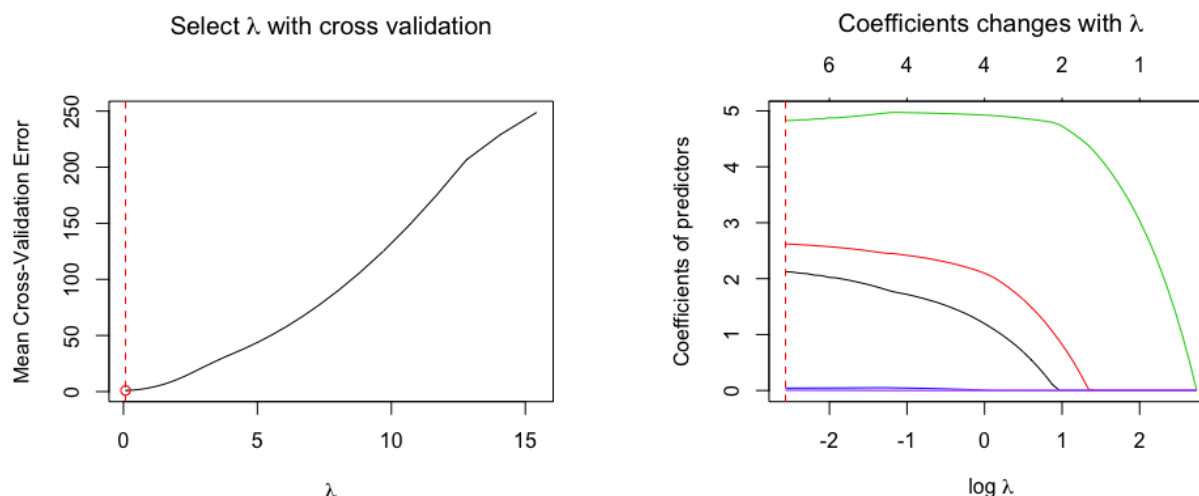
Code:

```

1 library("glmnet")
2 x_lasso = model.matrix(response~., simulatedData)[,-1]
3 y_lasso = simulatedData$response
4 cv.lasso <- cv.glmnet(x_lasso, y_lasso, nfolds = 10, parallel=TRUE,
5   standardize=TRUE, type.measure='mse')
6 # plot(cv.lasso)
7 plot(cv.lasso$lambda, cv.lasso$cvm, main=expression(paste("Select  $\lambda$ ",
8   lambda, "with cross validation")), xlab=expression(lambda), ylab="
9   Mean Cross-Validation Error", type = 'l')
10 abline(v=cv.lasso$lambda.min, lty=2, col="red")
11 points(cv.lasso$lambda.min, cv.lasso$cvm[which(cv.lasso$lambda==cv.
12   lasso$lambda.min)], col="red")
13 plot(cv.lasso$glmnet.fit, xvar="lambda", label=TRUE)
14 abline(v=log(cv.lasso$lambda.min), lty=2, col="red")
15 cv.lasso$lambda.min
16 # cv.lasso$lambda.1se
17 coef(cv.lasso, s=cv.lasso$lambda.min)

```

Result:



(a) Ten-fold cross-validation MSE for the lasso

(b) The corresponding lasso coefficient estimates

Figure 13: Lasso model with cross-validation to select optimal λ

```

1 > cv.lasso$lambda.min
2 [1] 0.07676824
3 > coef(cv.lasso, s=cv.lasso$lambda.min)

```

```

4 | 11 x 1 sparse Matrix of class "dgCMatrix"
5 |                               1
6 | (Intercept) 1.183697416
7 | X_1         2.123302063
8 | X_2         2.621261843
9 | X_3         4.826107793
10 | X_4         0.042987976
11 | X_5         0.010370603
12 | X_6         .
13 | X_7         0.003844961
14 | X_8         .
15 | X_9         .
16 | X_10        .

```

Ans:

Using ten-fold cross-validation, the λ in lasso is selected as $\lambda = 0.07676824$. The small λ means the lasso will have similar performance with least squared estimation. The model obtained by lasso is:

$$\begin{aligned}\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_4 X^4 + \hat{\beta}_5 X^5 + \hat{\beta}_7 X^7 \\ &= 1.184 + 2.123X + 2.621X^2 + 4.826X^3 + 0.043X^4 + 0.010X^5 + 0.004X^7\end{aligned}$$

This model shows variable selection but still includes some noise variables, although whose coefficients are small. The estimated coefficients of signal variables are close to their true value.

(f)**Code:**

```

1 | # (f)
2 | beta0 = 1
3 | beta7 = 13
4 | Y = beta0 + beta7*X^7 + e
5 | # best subset selection
6 | # create a single data set containing both X and Y
7 | simulatedData = data.frame(
8 |   response = Y,
9 |   X_1 = X,
10 |   X_2 = X^2,
11 |   X_3 = X^3,
12 |   X_4 = X^4,
13 |   X_5 = X^5,
14 |   X_6 = X^6,
15 |   X_7 = X^7,
16 |   X_8 = X^8,
17 |   X_9 = X^9,
18 |   X_10 = X^10,
19 |   stringsAsFactors = FALSE

```

```

20 )
21 str(simulatedData) # Get the structure of the data frame.
22 library(leaps)
23 # perform best subset selection, using 10 predictors
24 regfit.full = regsubsets(response~., simulatedData, nvmax = 10)
25 reg.summary = summary(regfit.full)
26 # RSS
27 plot(reg.summary$rss, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="RSS", type = "l")
28 # Adjusted RSq
29 plot(reg.summary$adjr2, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="Adjusted_RSq", type = "l")
30 abline(v=which.max(reg.summary$adjr2), lty=2, col="red")
31 points(which.max(reg.summary$adjr2), reg.summary$adjr2[which.max(reg.
    summary$adjr2)], col="red")
32 coef(regfit.full, which.max(reg.summary$adjr2))
33 # Cp
34 plot(reg.summary$cp, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="Cp", type = "l")
35 abline(v=which.min(reg.summary$cp), lty=2, col="red")
36 points(which.min(reg.summary$cp), reg.summary$cp[which.min(reg.summary$
    cp)], col="red")
37 coef(regfit.full, which.min(reg.summary$cp))
38 # BIC
39 plot(reg.summary$bic, main="Best_subset_selection", xlab="Number_of_
    Variables", ylab="BIC", type = "l")
40 abline(v=which.min(reg.summary$bic), lty=2, col="red")
41 points(which.min(reg.summary$bic), reg.summary$bic[which.min(reg.
    summary$bic)], col="red")
42 coef(regfit.full, which.min(reg.summary$bic))
43
44 # lasso
45 library("glmnet")
46 x_lasso = model.matrix(response~., simulatedData)[, -1]
47 y_lasso = simulatedData$response
48 cv.lasso <- cv.glmnet(x_lasso, y_lasso, nfolds = 10, parallel=TRUE,
    standardize=TRUE, type.measure='mse')
49 # plot(cv.lasso)
50 plot(cv.lasso$lambda, cv.lasso$cv, main=expression(paste("Select_",
    lambda, "_with_cross_validation")), xlab=expression(lambda), ylab="
    Mean_Cross-Validation_Error", type = 'l')
51 abline(v=cv.lasso$lambda.min, lty=2, col="red")
52 points(cv.lasso$lambda.min, cv.lasso$cv[which(cv.lasso$lambda==cv.
    lasso$lambda.min)], col="red")
53 plot(cv.lasso$glmnet.fit, xvar="lambda", main=expression(paste("
    Coefficients_changes_with_", lambda)), xlab=expression(paste('log_',
    lambda)), ylab="Coefficients_of_predictors")

```

```

54 abline(v=log(cv.lasso$lambda.min), lty=2, col="red")
55 cv.lasso$lambda.min
56 # cv.lasso$lambda.1se
57 coef(cv.lasso, s=cv.lasso$lambda.min)

```

Result:

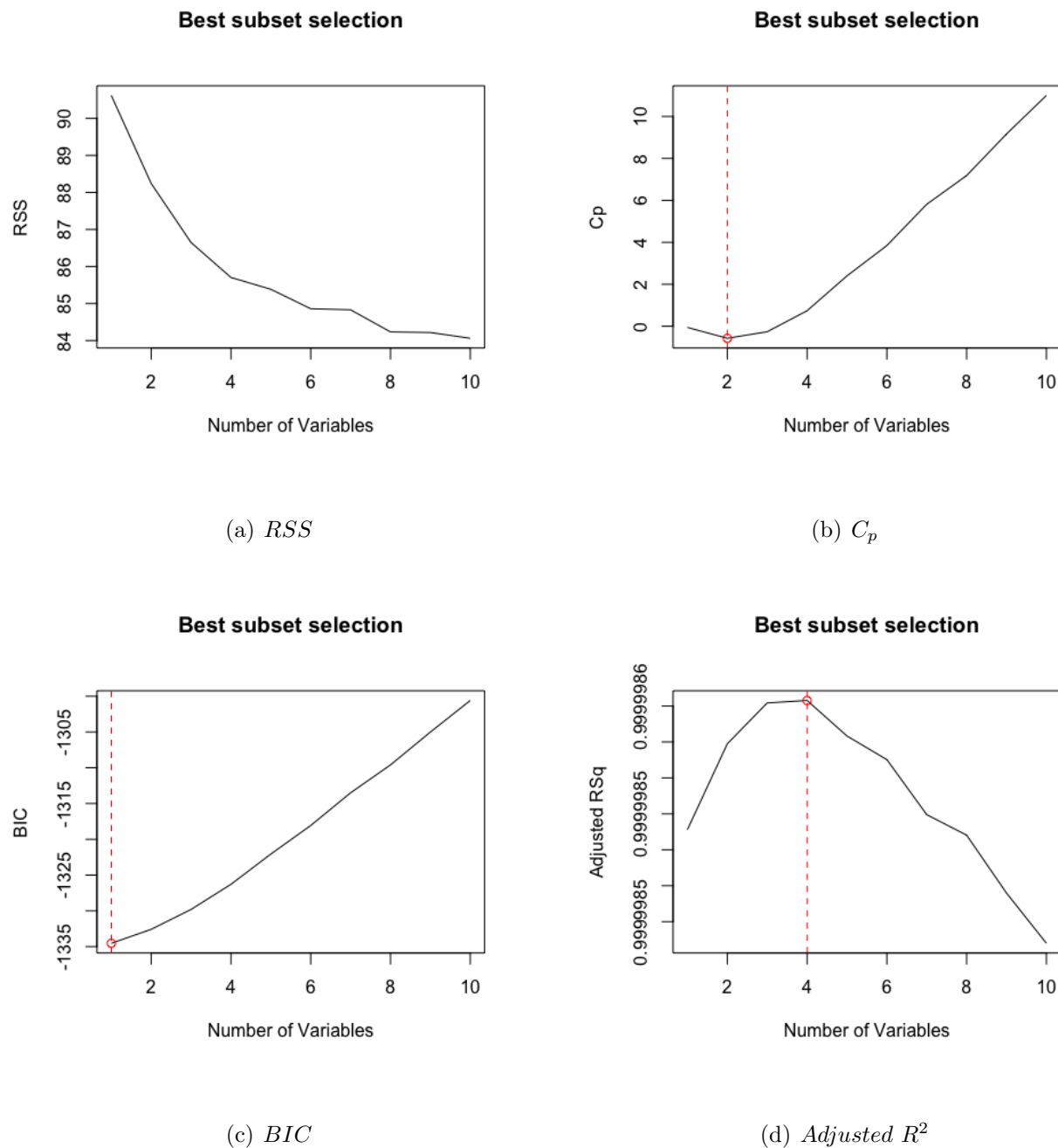


Figure 14: Best subset selection


```

1 > str(simulatedData) # Get the structure of the data frame.
2 'data.frame':  100 obs. of  11 variables:
3 $ response: num  -0.113  1.042  -3.61  342.983  0.351  ...
4 $ X_1      : num  -0.626  0.184  -0.836  1.595  0.33  ...
5 $ X_2      : num  0.3924  0.0337  0.6983  2.5449  0.1086  ...
6 $ X_3      : num  -0.24585  0.00619  -0.5835  4.05986  0.03578  ...
7 $ X_4      : num  0.15401  0.00114  0.48759  6.47662  0.01179  ...
8 $ X_5      : num  -0.096482  0.000209  -0.407443  10.332031  0.003884  ...
9 $ X_6      : num  6.04e-02  3.84e-05  3.40e-01  1.65e+01  1.28e-03  ...
10 $ X_7      : num  -3.79e-02  7.04e-06  -2.85e-01  2.63e+01  4.22e-04  ...
11 $ X_8      : num  2.37e-02  1.29e-06  2.38e-01  4.19e+01  1.39e-04  ...
12 $ X_9      : num  -1.49e-02  2.38e-07  -1.99e-01  6.69e+01  4.58e-05  ...
13 $ X_10     : num  9.31e-03  4.36e-08  1.66e-01  1.07e+02  1.51e-05  ...

```

```

1 > coef(regfit.full, which.max(reg.summary$adjr2))
2 (Intercept)      X_1      X_2      X_3      X_7
3  1.0762524    0.2914016   -0.1617671   -0.2526527   13.0091338

```

```

1 > coef(regfit.full, which.min(reg.summary$cp))
2 (Intercept)      X_2      X_7
3  1.0704904   -0.1417084   13.0015552

```

```

1 > coef(regfit.full, which.min(reg.summary$bic))
2 (Intercept)      X_7
3  0.9589402   13.0007705

```

```

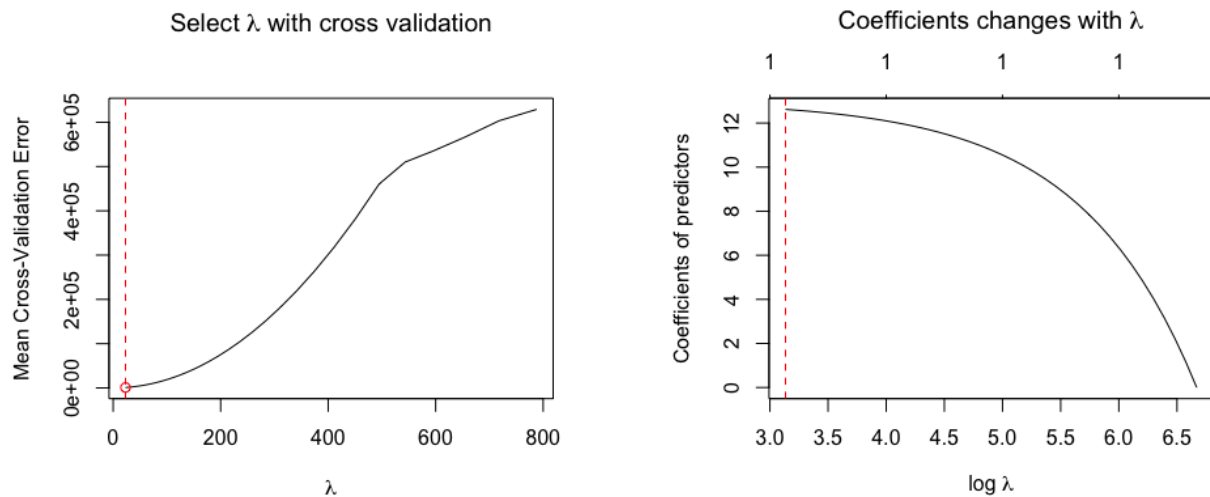
1 > cv.lasso$lambda.min
2 [1] 22.96953

```

```

1 > coef(cv.lasso, s=cv.lasso$lambda.min)
2 11 x 1 sparse Matrix of class "dgCMatrix"
3
4 (Intercept)  2.558369
5 X_1          .
6 X_2          .
7 X_3          .
8 X_4          .
9 X_5          .
10 X_6          .
11 X_7         12.621791
12 X_8          .
13 X_9          .
14 X_10         .

```



(a) Ten-fold cross-validation MSE

(b) The corresponding lasso coefficient estimates

Figure 15: Lasso

Ans:

Best subset selection and lasso both do the variable selection.

Best subset selection will get different models when different criteria are used:

$$\begin{aligned}
 \text{Adjusted } R^2 : \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_7 X^7 \\
 &= 1.076 + 0.291X - 0.162X^2 - 0.253X^3 + 13.009X^7 \\
 C_p : \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_2 X^2 + \hat{\beta}_7 X^7 \\
 &= 1.070 - 0.142X^2 + 13.002X^7 \\
 BIC : \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_7 X^7 \\
 &= 0.959 + 13.001X^7
 \end{aligned}$$

Adjusted R² and *C_p* still include some noise variables, whose coefficients are very small. *BIC* doesn't include noise variable.

Lasso gets a model without noise variables as:

$$\begin{aligned}
 \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_7 X^7 \\
 &= 2.558 + 12.622X^7
 \end{aligned}$$

However the estimated coefficients of the model are not as close to their true value as best subset selection, because the regularization modification $\lambda \sum_{j=1}^p |\beta_j|$ will affect the estimation of coefficients.

P5

(a)

Code:

```

1 # (a)
2 rm(list=ls())
3 set.seed(1)
4 College = read.csv("/Users/yangchenye/Downloads/College.csv", header=T,
5   na.strings="?")
6 dim(College)
7 College=na.omit(College) # remove incomplete cases
8 dim(College)
9 names(College)
10 # 75% of the sample size
11 smp_size = floor(0.75 * nrow(College))
12 train_ind = sample(seq_len(nrow(College)), size = smp_size)
13 # select train data and test data, drop the name 'X' of College
14 train = subset(College[train_ind, ], select = -X)
15 test = subset(College[-train_ind, ], select = -X)

```

Result:

```

1 > dim(College)
2 [1] 777  19
3 > names(College)
4 [1] "X"           "Private"      "Apps"         "Accept"
5 [5] "Enroll"      "Top10perc"    "Top25perc"    "F.Undergrad"
6 [9] "P.Undergrad" "Outstate"     "Room.Board"   "Books"
7 [13] "Personal"    "PhD"          "Terminal"     "S.F.Ratio"
8 [17] "perc.alumni" "Expend"       "Grad.Rate"
9 > dim(test)
10 [1] 195  18
11 > dim(train)
12 [1] 582  18

```

(b)

Code:

```

1 # (b)
2 train_fit = lm(train$Grad.Rate~., data = train)
3 summary(train_fit)
4 test_lm_predict = predict(train_fit, test[-18]) # predict
5 test_lm_MSE = mean((test$Grad.Rate - test_lm_predict) ^ 2) # test MSE

```

Result:

```

1 > summary(train_fit)
2
3 Call:
4 lm(formula = train$Grad.Rate ~ ., data = train)
5
6 Residuals:
7      Min       1Q   Median       3Q      Max
8 -50.657  -7.049  -0.529   7.117  51.438
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) 26.8066318   5.8737195   4.564 6.17e-06 ***
13 PrivateYes   4.7305052   1.9991510   2.366 0.018306 *
14 Apps        0.0010352   0.0005326   1.944 0.052435 .
15 Accept     -0.0007207   0.0010828  -0.666 0.505961
16 Enroll      0.0024664   0.0028023   0.880 0.379161
17 Top10perc   0.0922866   0.0830377   1.111 0.266877
18 Top25perc   0.0811958   0.0621955   1.305 0.192257
19 F.Undergrad -0.0002828   0.0004785  -0.591 0.554695
20 P.Undergrad -0.0016878   0.0004434  -3.806 0.000157 ***
21 Outstate    0.0008544   0.0002648   3.227 0.001324 **
22 Room.Board  0.0021872   0.0006795   3.219 0.001362 **
23 Books       -0.0015732   0.0033449  -0.470 0.638309
24 Personal    -0.0016231   0.0008667  -1.873 0.061621 .
25 PhD         0.0964858   0.0654935   1.473 0.141252
26 Terminal    -0.0477951   0.0724298  -0.660 0.509599
27 S.F.Ratio    0.2824858   0.1955257   1.445 0.149083
28 perc.alumni  0.3194714   0.0572827   5.577 3.80e-08 ***
29 Expend      -0.0003060   0.0001671  -1.831 0.067599 .
30
31 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
32
33 Residual standard error: 12.8 on 564 degrees of freedom
34 Multiple R-squared:  0.4723,    Adjusted R-squared:  0.4564
35 F-statistic: 29.7 on 17 and 564 DF,  p-value: < 2.2e-16

```

```

1 > test_lm_MSE
2 [1] 163.6218

```

Ans:

The test mean squared error obtained is 163.6218.

(c)

Code:

```
1 # (c) ridge
```

```

2 library("glmnet")
3 x_ridge = model.matrix(train$Grad.Rate~., data = train)[, -1]
4 y_ridge = train$Grad.Rate
5 # Ridge: Alpha = 0
6 cv.ridge = cv.glmnet(x_ridge, y_ridge, alpha=0, nfolds = 10, parallel=
  TRUE, standardize=TRUE, type.measure='mse')
7 # CVE~lambda
8 plot(cv.ridge$lambda, cv.ridge$cv, main=expression(paste("Select_",
  lambda, "_with_cross_validation")), xlab=expression(lambda), ylab="
  Mean_Cross-Validation_Error", type = 'l')
9 abline(v=cv.ridge$lambda.min, lty=2, col="red")
10 points(cv.ridge$lambda.min, cv.ridge$cv[which(cv.ridge$lambda==cv.
  ridge$lambda.min)], col="red")
11 plot(cv.ridge$glmnet.fit, xvar="lambda", main=expression(paste("
  Coefficients_changes_with_", lambda)), xlab=expression(paste('log_',
  lambda)), ylab="Coefficients_of_predictors")
12 abline(v=log(cv.ridge$lambda.min), lty=2, col="red")
13 # best lambda
14 cv.ridge$lambda.min
15 coef(cv.ridge, s=cv.ridge$lambda.min)
16 # best ridge regression with best lambda
17 ridge_best = glmnet(x_ridge, y_ridge, alpha=0, lambda = cv.ridge$
  lambda.min, standardize=TRUE)
18 test_ridge_predict = predict(ridge_best, model.matrix(test$Grad.Rate~
  ., data = test)[, -1], s="lambda.min") # predict
19 test_ridge_MSE = mean((test$Grad.Rate - test_ridge_predict) ^ 2) #
  test MSE

```

Result:

```

1 > cv.ridge$lambda.min
2 [1] 3.717836
3 > coef(cv.ridge, s=cv.ridge$lambda.min)
4 18 x 1 sparse Matrix of class "dgCMatrix"
5
6 (Intercept) 3.046609e+01
7 PrivateYes 4.222649e+00
8 Apps 3.903074e-04
9 Accept 3.282553e-04
10 Enroll 4.136841e-04
11 Top10perc 9.958251e-02
12 Top25perc 8.206982e-02
13 F. Undergrad -6.921736e-05
14 P. Undergrad -1.385642e-03
15 Outstate 6.635847e-04
16 Room.Board 1.878762e-03
17 Books -1.775428e-03
18 Personal -1.760206e-03

```

```

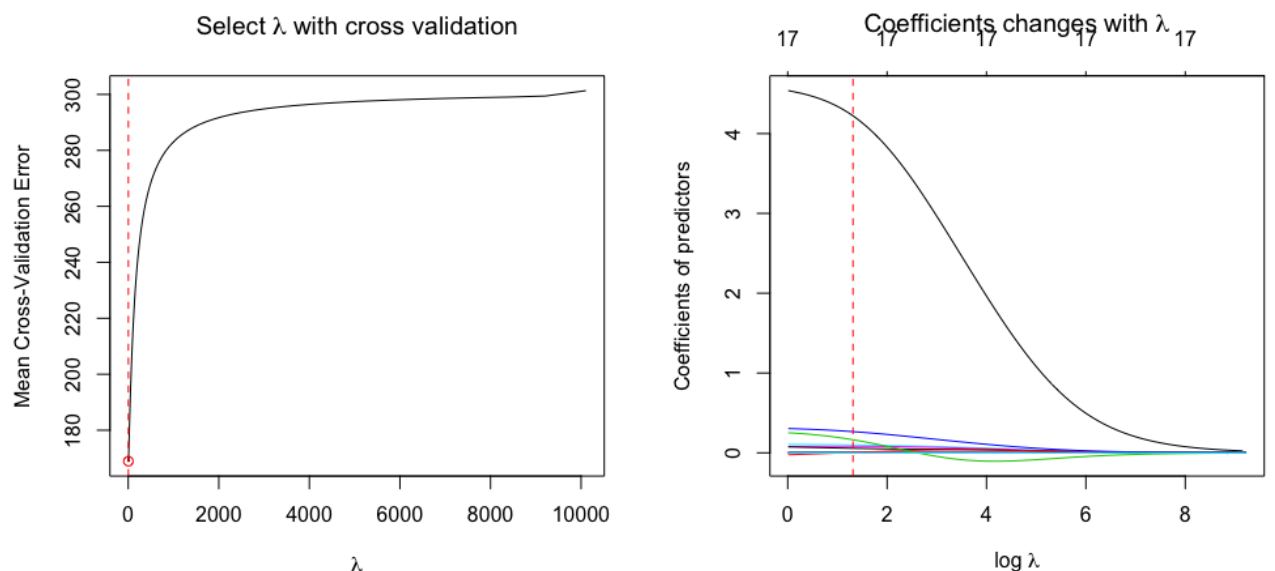
19 PhD                5.689299e-02
20 Terminal           5.003101e-03
21 S.F.Ratio          1.641541e-01
22 perc.alumni        2.657793e-01
23 Expend             -9.817824e-05

```

```

1 > test_ridge_MSE
2 [1] 164.5001

```



(a) Ten-fold cross-validation MSE

(b) The corresponding ridge regression coefficients

Figure 16: Ridge regression model

Ans:

The λ chosen in ridge regression is 3.717836. The test mean squared error obtained is 164.5001.

(d)**Code:**

```

1 # (d) lasso
2 library("glmnet")
3 x_lasso = model.matrix(train$Grad.Rate~., data = train)[,-1]
4 y_lasso = train$Grad.Rate
5 # Lasso: Alpha = 1
6 cv.lasso = cv.glmnet(x_lasso, y_lasso, alpha=1, nfolds = 10, parallel=
  TRUE, standardize=TRUE, type.measure='mse')
7 # CVE~lambda
8 plot(cv.lasso$lambda, cv.lasso$cvm, main=expression(paste("Select ",
  lambda, " with cross validation")), xlab=expression(lambda), ylab="
  Mean Cross-Validation Error", type = 'l')

```

```

9 abline(v=cv.lasso$lambda.min, lty=2, col="red")
10 points(cv.lasso$lambda.min, cv.lasso$cvm[which(cv.lasso$lambda==cv.
    lasso$lambda.min)], col="red")
11 plot(cv.lasso$glmnet.fit, xvar="lambda", main=expression(paste("
    Coefficients changes with ", lambda)), xlab=expression(paste('log ',
    lambda)), ylab="Coefficients of predictors")
12 abline(v=log(cv.lasso$lambda.min), lty=2, col="red")
13 # best lambda
14 cv.lasso$lambda.min
15 cv.lasso$nzzero[which(cv.lasso$lambda==cv.lasso$lambda.min)]
16 coef(cv.lasso, s=cv.lasso$lambda.min)
17 # best ridge regression with best lambda
18 lasso_best = glmnet(x_lasso, y_lasso, alpha=0, lambda = cv.lasso$
    lambda.min, standardize=TRUE)
19 test_lasso_predict = predict(lasso_best, model.matrix(test$Grad.Rate~
    ., data = test)[-1], s="lambda.min") # predict
20 test_lasso_MSE = mean((test$Grad.Rate - test_lasso_predict) ^ 2) #
    test MSE

```

Result:

```

1 > cv.lasso$lambda.min
2 [1] 0.2446088
3 > cv.lasso$nzzero[which(cv.lasso$lambda==cv.lasso$lambda.min)]
4 s40
5 13
6 > coef(cv.lasso, s=cv.lasso$lambda.min)
7 18 x 1 sparse Matrix of class "dgCMatrix"
8
9 (Intercept) 29.0167980959
10 PrivateYes 3.6706143690
11 Apps 0.0006514975
12 Accept .
13 Enroll .
14 Top10perc 0.0833713627
15 Top25perc 0.0821703158
16 F.Undergrad .
17 P.Undergrad -0.0015119149
18 Outstate 0.0007874207
19 Room.Board 0.0019193219
20 Books -0.0006307748
21 Personal -0.0015104017
22 PhD 0.0465084115
23 Terminal .
24 S.F.Ratio 0.1911923409
25 perc.alumni 0.3104190913
26 Expend -0.0001341114

```

```

1 > test_lasso_MSE
2 [1] 163.8735

```

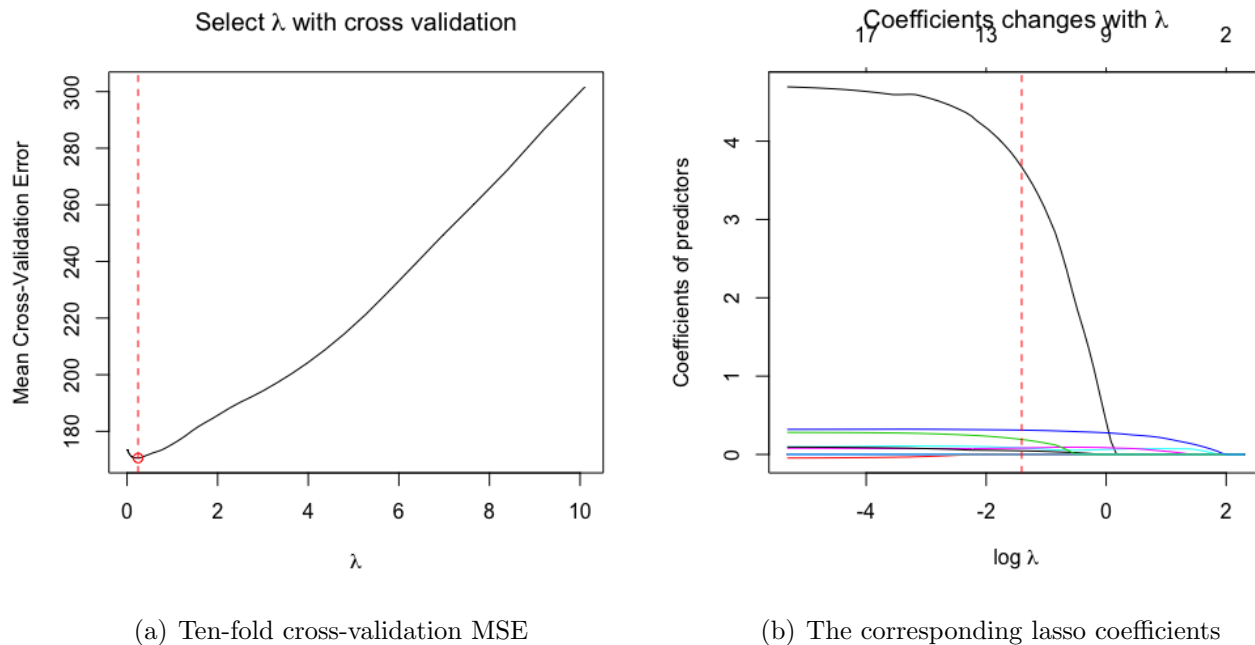


Figure 17: Lasso model

Ans:

The λ chosen in lasso is 0.2446088. The number of non-zero coefficient estimates is 13. The test mean squared error obtained is 163.8735.

P6

(a)

Code:

```
1 # (a)
2 rm(list=ls())
3 set.seed(1)
4 p = 20
5 n = 1000
6 X = rnorm(n, mean = 0, sd = 1)
7 e = rnorm(n, mean = 0, sd = 1)
8 c = seq(from = 1, to = 20, length.out = p)
9 beta = rep(0, p)
10 beta0 = 1
11 beta[1] = 2
12 beta[2] = 3
13 beta[3] = 5
14 Y = beta0 + beta[1]*X + beta[2]*X^2 + beta[3]*X^3 + e
```

(b)

Code:

```
1 # (b)
2 simulatedData = data.frame(
3   response = Y,
4   X_1 = X^c[1],
5   X_2 = X^c[2],
6   X_3 = X^c[3],
7   X_4 = X^c[4],
8   X_5 = X^c[5],
9   X_6 = X^c[6],
10  X_7 = X^c[7],
11  X_8 = X^c[8],
12  X_9 = X^c[9],
13  X_10 = X^c[10],
14  X_11 = X^c[11],
15  X_12 = X^c[12],
16  X_13 = X^c[13],
17  X_14 = X^c[14],
18  X_15 = X^c[15],
19  X_16 = X^c[16],
20  X_17 = X^c[17],
21  X_18 = X^c[18],
22  X_19 = X^c[19],
23  X_20 = X^c[20],
```

```

24 stringsAsFactors = FALSE
25 )
26 str(simulatedData) # Get the structure of the data frame.
27 # training set containing 100 observations, 1/10
28 smp_size = floor(0.1 * nrow(simulatedData))
29 train_ind = sample(seq_len(nrow(simulatedData)), size = smp_size)
30 # select train data and test data
31 train = simulatedData[train_ind, ]
32 # test = simulatedData[-train_ind, ]
33 test = model.matrix(~., data=simulatedData[-train_ind, ])

```

(c)

Code:

```

1 # (c)
2 library(leaps)
3 # perform best subset selection, using 20 predictors
4 regfit.full = regsubsets(response~., simulatedData, nvmax = p)
5 reg.summary = summary(regfit.full)
6 # RSS
7 plot(reg.summary$rss, main="Best_subset_selection", xlab="Number_of_
  Variables", ylab="RSS", type = "l")
8 # Adjusted RSq
9 plot(reg.summary$adjr2, main="Best_subset_selection", xlab="Number_of_
  Variables", ylab="Adjusted_RSq", type = "l")
10 abline(v=which.max(reg.summary$adjr2), lty=2, col="red")
11 points(which.max(reg.summary$adjr2), reg.summary$adjr2[which.max(reg.
  summary$adjr2)], col="red")
12 coef(regfit.full, which.max(reg.summary$adjr2))
13 # Cp
14 plot(reg.summary$cp, main="Best_subset_selection", xlab="Number_of_
  Variables", ylab="Cp", type = "l")
15 abline(v=which.min(reg.summary$cp), lty=2, col="red")
16 points(which.min(reg.summary$cp), reg.summary$cp[which.min(reg.summary$
  cp)], col="red")
17 coef(regfit.full, which.min(reg.summary$cp))
18 # BIC
19 plot(reg.summary$bic, main="Best_subset_selection", xlab="Number_of_
  Variables", ylab="BIC", type = "l")
20 abline(v=which.min(reg.summary$bic), lty=2, col="red")
21 points(which.min(reg.summary$bic), reg.summary$bic[which.min(reg.
  summary$bic)], col="red")
22 coef(regfit.full, which.min(reg.summary$bic))
23
24 # the training set MSE associated with the best model of each size

```

```

25 plot(regfit.full$rss[-1]/n, main="Training MSE", xlab="Number of
    predictors", ylab="MSE", col="blue", type="l", lty=2)
26 legend("topright", legend="Training", col="blue", lty=2)
27 which.min(regfit.full$rss[-1]/n)

```

Result:

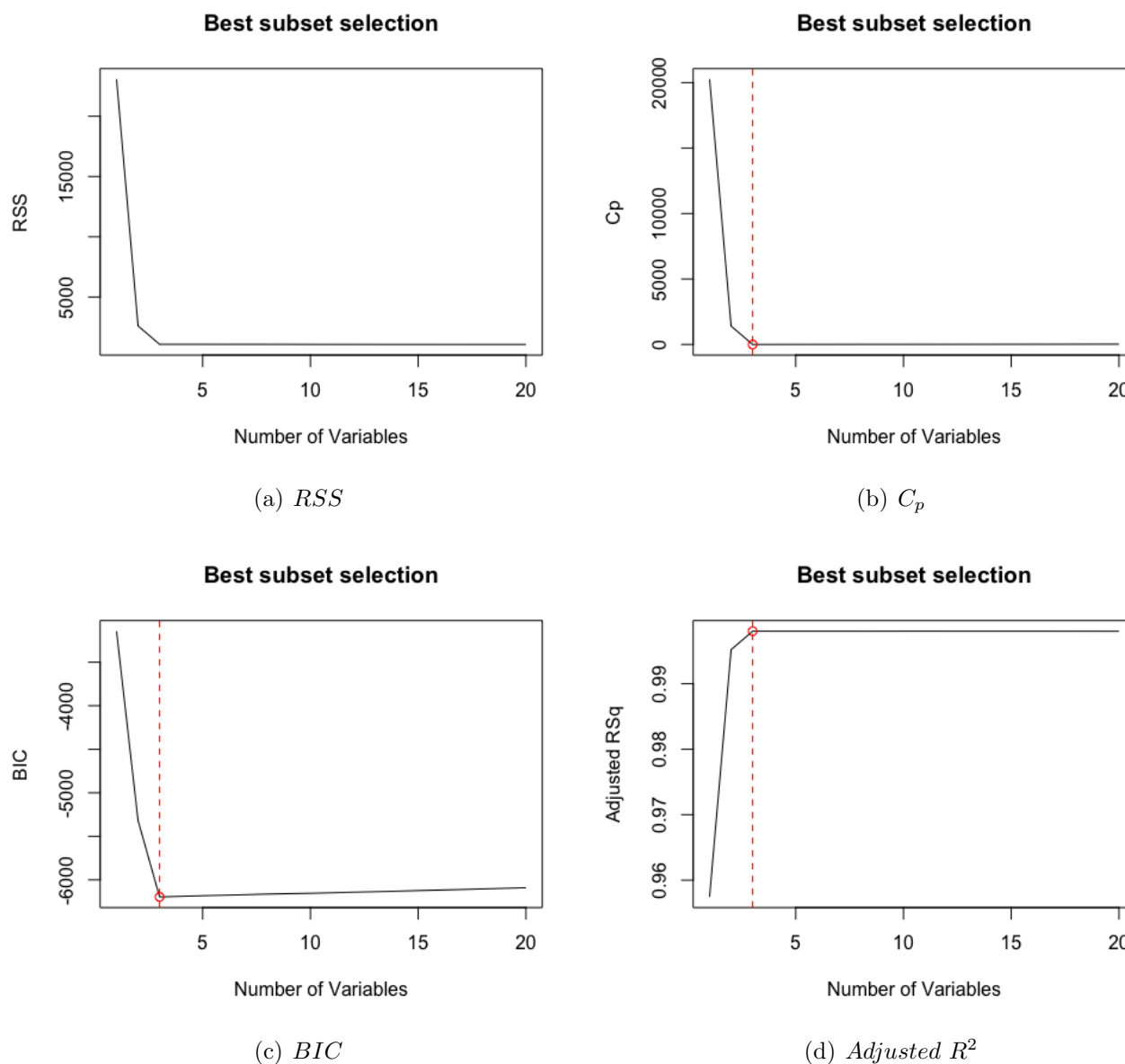


Figure 18: Best subset selection

```

1 > coef(regfit.full, which.max(reg.summary$adjr2))
2 (Intercept)      X_1      X_2      X_3
3  0.9648892  1.9270333  3.0181880  5.0249860
4 > coef(regfit.full, which.min(reg.summary$cp))
5 (Intercept)      X_1      X_2      X_3
6  0.9648892  1.9270333  3.0181880  5.0249860

```

```

7 > coef(regfit.full, which.min(reg.summary$bic))
8 (Intercept)      X_1      X_2      X_3
9   0.9648892   1.9270333   3.0181880   5.0249860

```

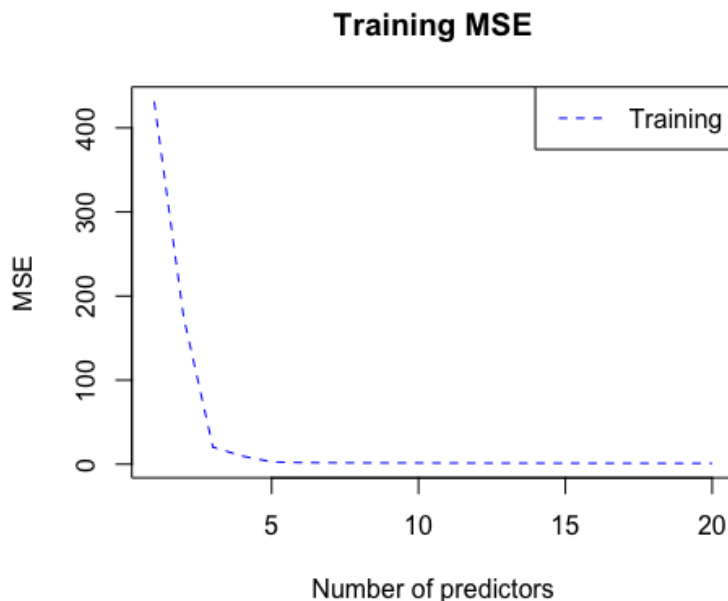


Figure 19: The training set MSE

```

1 > which.min(regfit.full$RSS[-1]/n)
2 [1] 20

```

(d)

Code:

```

1 # (d)
2 # Plot the test set MSE associated with the best model of each size
3 test.full.MSE = rep(NA, p) # test set MSE
4 for (i in 1:p) {
5   coefi = coef(regfit.full, id = i)
6   pred = test[, names(coefi)] %*% coefi
7   test.full.MSE[i] = mean((test[, "response"] - pred)^2)
8 }
9 plot(test.full.MSE, main="Testing MSE", xlab="Number of predictors",
10      ylab = "MSE", col = "red", type = "l", lty=1)
11 legend("topright", legend = "Testing", col = "red", lty=1)

```

Result:

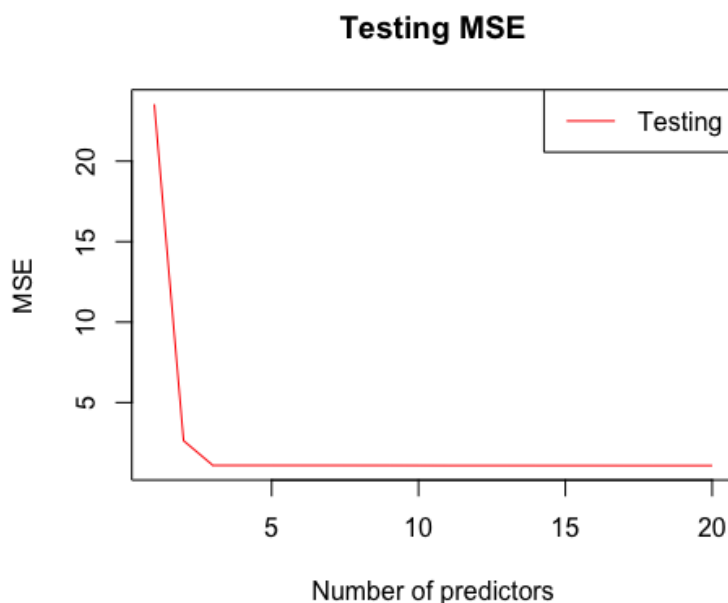


Figure 20: The testing set MSE

(e)

Code:

```
1 # (e)
2 which.min(test.full.MSE)
3 coef(regfit.full, id = which.min(test.full.MSE))
```

Result:

```
1 > which.min(test.full.MSE)
2 [1] 17
3 > coef(regfit.full, id = which.min(test.full.MSE))
4 (Intercept)      X_1      X_2      X_3      X_6
5 9.941718e-01  1.967462e+00  2.850958e+00  5.099776e+00  1.992884e-01
6      X_7      X_9      X_10      X_11      X_12
7 -4.737564e-01  5.198375e-01 -1.351559e-01 -2.331000e-01  7.858723e-02
8      X_13      X_14      X_15      X_16      X_17
9 5.409708e-02 -1.980249e-02 -6.854977e-03  2.563927e-03  4.507476e-04
10      X_18      X_19      X_20
11 -1.670087e-04 -1.204494e-05  4.337136e-06
```

Ans:

The test MSE takes its minimum value for the model with 17 predictors.

(f)

Ans:

The model at which test MSE minimized:

$$\begin{aligned}
\hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2 + \hat{\beta}_3 X^3 + \hat{\beta}_6 X^6 + \hat{\beta}_7 X^7 + \hat{\beta}_9 X^9 + \hat{\beta}_{10} X^{10} + \\
&\quad \hat{\beta}_{11} X^{11} + \hat{\beta}_{12} X^{12} + \hat{\beta}_{13} X^{13} + \hat{\beta}_{14} X^{14} + \hat{\beta}_{15} X^{15} + \hat{\beta}_{16} X^{16} + \\
&\quad \hat{\beta}_{17} X^{17} + \hat{\beta}_{18} X^{18} + \hat{\beta}_{19} X^{19} + \hat{\beta}_{20} X^{20} \\
&= 0.994 + 1.967X + 2.851X^2 + 5.100X^3 + 0.199X^6 - 0.473X^7 + \\
&\quad 0.520X^9 - 0.135X^{10} - 0.233X^{11} + 0.079X^{12} + \\
&\quad 0.054X^{13} - 0.020X^{14} - 0.007X^{15} + 0.003X^{16} + \\
&\quad 5 \times 10^{-4}X^{17} - 2 \times 10^{-4}X^{18} - 1 \times 10^{-5}X^{19} + 4 \times 10^{-6}X^{20}
\end{aligned}$$

The true model used to generate the data:

$$\begin{aligned}
Y &= \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 \\
&= 1 + 2X + 3X^2 + 5X^3
\end{aligned}$$

The estimated model contains a lot of noise variables, comparing to the true model. The estimated coefficients of noise variables are small. The estimated coefficients of signal variables are close to their true value.

The estimated model with minimal test MSE has 17 predictors may because the train data and test data both have lots of noise, introduced by ϵ , and the model tries to fit the noise in training and predict the same noise in validation. Moreover, the true model only has first three non-zero coefficients. Therefore, the test MSE increases little when the degree of freedom increases to large, like Figure 20.

(g)

Code:

```

1 # (g)
2 y = rep(NA, p)
3 for (r in 1:p){
4   temp = 0
5   coefi = coef(regfit.full, id=r) # intercept, X_1, X_2, ...
6   for (j in names(coefi)[-1]){
7     num = type.convert(substr(j, 3, 4), 'int')
8     temp = temp + (beta[num] - coefi[[j]])^2
9   }
10  y[r] = temp
11 }
12 y = sqrt(y)
13
14 plot(y, xlab="Number_of_predictors", ylab = "y", col = "black", type =
    "l", lty=1)

```

Result:

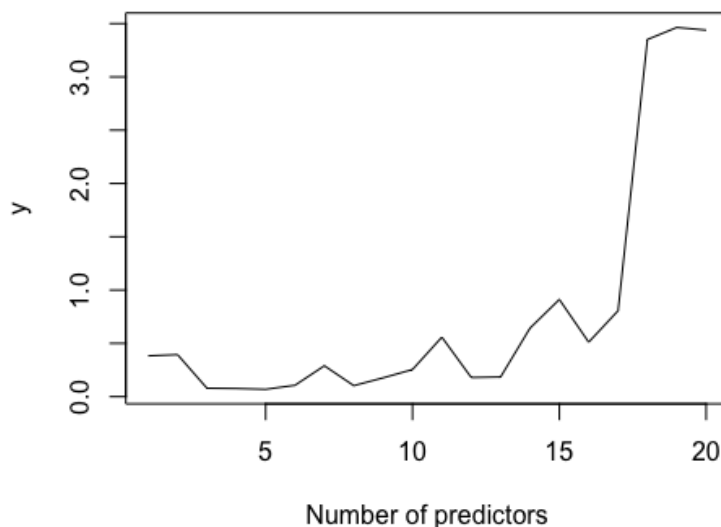


Figure 21: $y = \sqrt{\sum_{j=1}^p (\beta_j - \hat{\beta}_j^r)^2}$ for a range of values of r

Ans:

The plot displaying $\sqrt{\sum_{j=1}^p (\beta_j - \hat{\beta}_j^r)^2}$ for a range of values of r , number of coefficients, has a similar look with the regular test MSE plot. They both decrease first and then increase with the increasing of degree of freedom. But they are not exactly same because the $\sqrt{\sum_{j=1}^p (\beta_j - \hat{\beta}_j^r)^2}$ doesn't contain any form of X .

However, for this problem setting, they look different. Reasons:

The estimated model with minimal test MSE has 17 predictors may because the train data and test data both have lots of noise, introduced by ϵ , and the model tries to fit the noise in training and predict the same noise in validation. Moreover, the true model only has first three non-zero coefficients. Therefore, the test MSE increases little when the degree of freedom increases to large, like Figure 20.