

Homework 2

Yi Yang (yy3089)

October 24, 2021

1 P1

1.1 a

In this question, the RSS for the linear regression is:

$$\begin{aligned} RSS &= \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij} \right)^2 \\ &= (\hat{\beta}_0 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{12})^2 + (\hat{\beta}_0 + \hat{\beta}_1 x_{21} + \hat{\beta}_2 x_{22})^2 \end{aligned} \quad (1.1)$$

So the Ridge Regression Optimization is:

$$\begin{aligned} \min RSS + \lambda \sum_{j=1}^p \beta_j^2 &= \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2 \\ &= (\hat{\beta}_0 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{12})^2 + (\hat{\beta}_0 + \hat{\beta}_1 x_{21} + \hat{\beta}_2 x_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) \end{aligned} \quad (1.2)$$

1.2 b

For the optimization result:

$$\begin{aligned} \frac{\partial f}{\partial \beta_1} &= -2x_{11}(y_1 - \beta_1 x_{11} - \beta_2 x_{12}) - 2x_{21}(y_2 - \beta_1 x_{21} - \beta_2 x_{22}) + 2\lambda\beta_1 = 0 \\ \frac{\partial f}{\partial \beta_2} &= -2x_{12}(y_1 - \beta_1 x_{11} - \beta_2 x_{12}) - 2x_{22}(y_2 - \beta_1 x_{21} - \beta_2 x_{22}) + 2\lambda\beta_2 = 0 \end{aligned} \quad (1.3)$$

Denote that $x_1 = x_{11} = x_{21}$, $x_2 = x_{12} = x_{22}$. To solve the two equations, we need to make some transform:

$$\begin{aligned} (x_1^2 + x_2^2 + \lambda) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 &= x_1 y_1 + x_2 y_2 \\ (x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2 + \lambda) \hat{\beta}_2 &= x_1 y_1 + x_2 y_2 \\ (x_1^2 + x_2^2 + \lambda) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 &= (x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2 + \lambda) \hat{\beta}_2 \\ (x_1^2 + x_2^2 + \lambda) (\hat{\beta}_1 - \hat{\beta}_2) + (x_1^2 + x_2^2) (\hat{\beta}_2 - \hat{\beta}_1) &= 0 \\ \lambda \hat{\beta}_1 &= \lambda \hat{\beta}_2 \end{aligned} \quad (1.4)$$

Since $\lambda \neq 0$, $\hat{\beta}_1 = \hat{\beta}_2$

1.3 c

LASSO Regression Optimization is:

$$\begin{aligned} \min RSS + \lambda \sum_{j=1}^p |\beta_j| &= \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\hat{\beta}_j| \\ &= (\hat{\beta}_0 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{12})^2 + (\hat{\beta}_0 + \hat{\beta}_1 x_{21} + \hat{\beta}_2 x_{22})^2 + \lambda (|\hat{\beta}_1| + |\hat{\beta}_2|) \end{aligned} \quad (1.5)$$

1.4 d

For the optimization result:

$$\begin{aligned} \frac{\partial f}{\partial \beta_1} &= -2x_{11}(y_1 - \beta_1 x_{11} - \beta_2 x_{12}) - 2x_{21}(y_2 - \beta_1 x_{21} - \beta_2 x_{22}) \pm \lambda = 0 \\ \frac{\partial f}{\partial \beta_2} &= -2x_{12}(y_1 - \beta_1 x_{11} - \beta_2 x_{12}) - 2x_{22}(y_2 - \beta_1 x_{21} - \beta_2 x_{22}) \pm \lambda = 0 \end{aligned} \quad (1.6)$$

Denote that $x_1 = x_{11} = x_{21}$, $x_2 = x_{12} = x_{22}$. To solve the two equations, we need to make some transform:

$$\begin{aligned} (x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 &= x_1 y_1 + x_2 y_2 \pm \lambda \\ (x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 &= x_1 y_1 + x_2 y_2 \pm \lambda \end{aligned} \quad (1.7)$$

Only when $\hat{\beta}_1 > 0$ and $\hat{\beta}_2 > 0$ or $\hat{\beta}_1 < 0$ and $\hat{\beta}_2 < 0$, there will be solutions for this problem. Due to the two equations are same in form, there will be more than one solution for this problem. if $\hat{\beta}_1 > 0$ and $\hat{\beta}_2 > 0$, the equation will be:

$$(x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 = x_1 y_1 + x_2 y_2 - \lambda \quad (1.8)$$

Otherwise, the equation will be:

$$(x_1^2 + x_2^2) \hat{\beta}_1 + (x_1^2 + x_2^2) \hat{\beta}_2 = x_1 y_1 + x_2 y_2 + \lambda \quad (1.9)$$

2 P2

2.1 a

When $p = 1$, the optimization function is $(y - \beta_1)^2 + \lambda\beta_1^2 = (1 + \lambda)\beta_1^2 - 2y\beta_1 + y^2$. To find the β_1 that minimize the y , we need to calculate $2(1 + \lambda)\beta_1 - 2y = 0$

We can plot the image:

```
1 y = 2
2 lambda = 1
3 beta = seq(-5,5,0.01)
4
5 plot(beta,(y-beta)^2 + lambda*beta^2,type="l",xlab = "beta"
6 ,ylab = "Optimization function",xlim=c(-5,5))
7
8 min_x = y/(1+lambda)
9 min_y = (y-min_x)^2 + lambda*min_x^2
10 points(min_x,min_y,col="red")
11
12 title("Optimization Result for Ridge Regression")
```

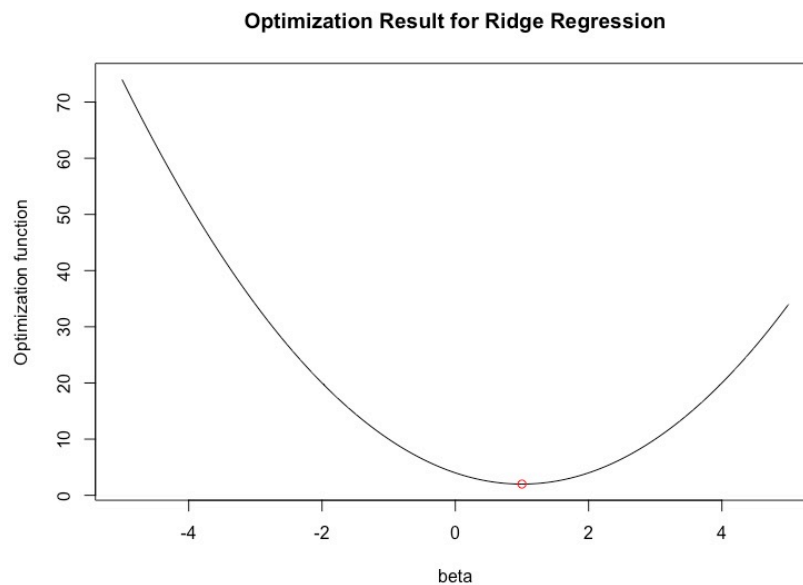


Figure 2.1: Screenshot for 'Ridge Regression'

2.2 b

When $p = 1$, the optimization function is $(y - \beta_1)^2 + \lambda|\beta_1| = \beta_1^2 - (2y \pm \lambda)\beta_1 + y^2$. To find the β_1 that minimize the y , we need to calculate $2\beta_1 - (2y \pm \lambda) = 0$. If $\beta_1 \geq 0$, that's $2\beta_1 - (2y - \lambda) = 0$. Otherwise, that's $2\beta_1 - (2y + \lambda) = 0$.

We can plot the image for all three conditions:

```
1 y = 2
2 lambda = 3
3 beta = seq(-5,5,0.01)
```

```

4
5 plot(beta,(y-beta)^2 + lambda*abs(beta),type="l",xlab = "beta"
6 ,ylab = "Optimization_function",xlim=c(-5,5))
7
8 if(y > lambda/2) {
9     min_x1 = y - lambda/2
10 } else {
11     if(y < -lambda/2){
12         min_x1 = y + lambda/2
13     } else {
14         min_x1 = 0
15     }
16 }
17
18 min_y1 = (y-min_x1)^2 + lambda*abs(min_x1)
19 points(min_x1,min_y1,col="red")
20
21 title("Optimization_Result_for_LASSO_Regression")

```

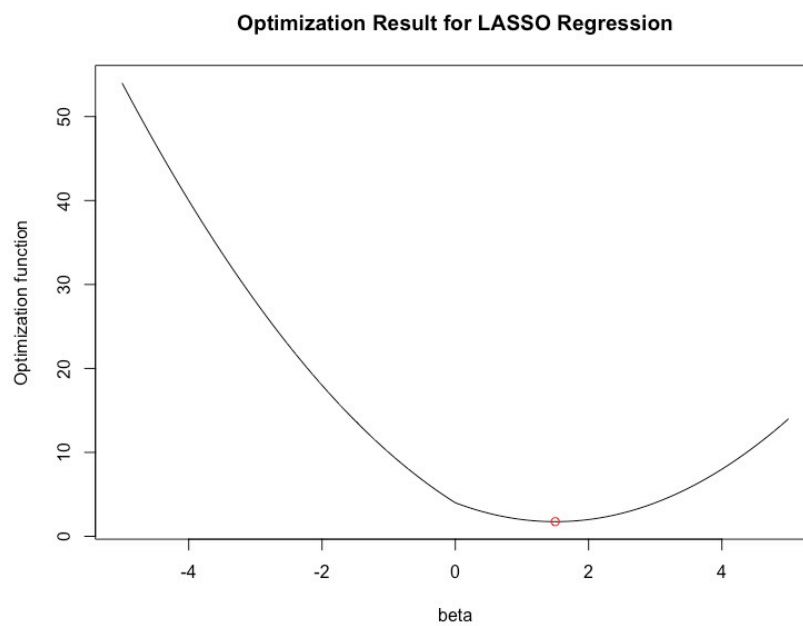


Figure 2.2: Screenshot for 'LASSO Regression'

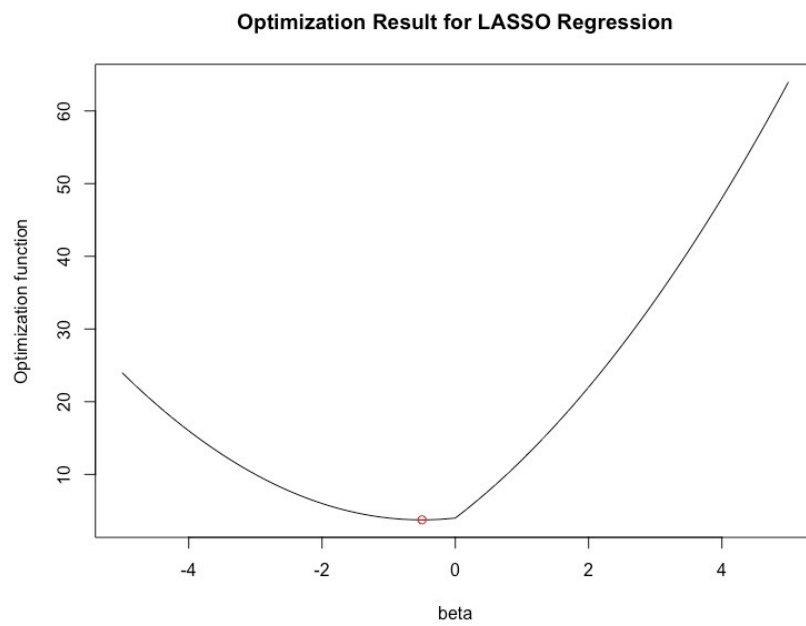


Figure 2.3: Screenshot for '*LASSO Regression*'

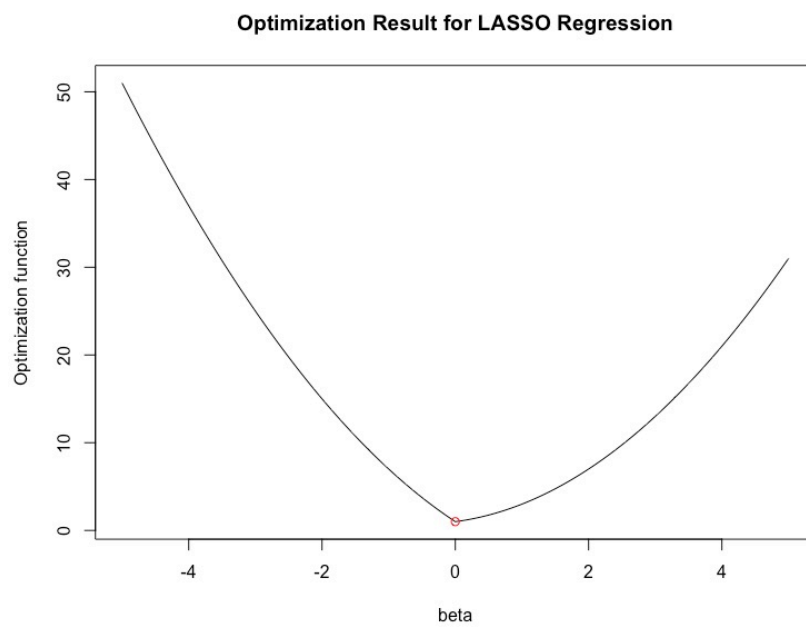


Figure 2.4: Screenshot for '*LASSO Regression*'

3 P3

3.1 a

We know that $y_i \sim \mathcal{N}(\beta_0 + \sum_{j=1}^p x_{ij}\beta_j, \sigma^2)$, so the likelihood function is:

$$\begin{aligned}\mathcal{L}(X, \beta, \sigma | Y) &= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2}{2\sigma^2}} \\ &= \frac{1}{(\sigma\sqrt{2\pi})^n} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2}\end{aligned}\quad (3.1)$$

3.2 b

We know that posterior $\propto \mathcal{L} \times \text{prior}$, and $\text{prior}(\beta) = \frac{1}{2b} e^{(-\frac{\sum_{j=1}^p |\beta_j|}{b})}$, thus:

$$\begin{aligned}\text{posterior} &\propto \frac{1}{(\sigma\sqrt{2\pi})^n} e^{(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2)} * \frac{1}{2b} e^{(-\frac{\sum_{j=1}^p |\beta_j|}{b})} \\ &\propto \frac{1}{2b} * \frac{1}{(\sigma\sqrt{2\pi})^n} e^{(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 - \frac{\sum_{j=1}^p |\beta_j|}{b})}\end{aligned}\quad (3.2)$$

3.3 c

To prove the argument, we need to prove that the maximum of the posterior is given by the LASSO regression with λ :

$$\log(\text{posterior}) = \log\left(\frac{1}{2b} * \frac{1}{(\sigma\sqrt{2\pi})^n}\right) + \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 - \frac{\sum_{j=1}^p |\beta_j|}{b}\right) \quad (3.3)$$

To find the maximum for the posterior, we need to find the minimum for the

$$\left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{\sum_{j=1}^p |\beta_j|}{b}\right):$$

$$\begin{aligned}\arg\min_{\beta} \left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{\sum_{j=1}^p |\beta_j|}{b}\right) &= \arg\min_{\beta} \left(\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{2\sigma^2 \sum_{j=1}^p |\beta_j|}{b}\right) \\ &= \arg\min_{\beta} \left(\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{2\sigma^2}{b} \sum_{j=1}^p |\beta_j|\right)\end{aligned}\quad (3.4)$$

Let $\lambda = \frac{2\sigma^2}{b}$, that's the same form with LASSO Regression.

3.4 d

$$\begin{aligned}\text{prior}(\beta) &= \prod_{j=1}^p \mathbf{P}(\beta_j) \\ &= \prod_{j=1}^p \frac{1}{\sqrt{2\pi c}} e^{\left(-\frac{\beta_j^2}{2c}\right)} \\ &= \frac{1}{(\sqrt{2\pi c})^p} e^{\left(-\frac{1}{2c} \sum_{j=1}^p \beta_j^2\right)}\end{aligned}\quad (3.5)$$

posterior $\propto \mathcal{L} \times \text{prior}$

$$\begin{aligned} &\propto \frac{1}{(\sigma\sqrt{2\pi})^n} e^{\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2\right)} * \frac{1}{(\sqrt{2\pi c})^p} e^{\left(-\frac{1}{2c} \sum_{j=1}^p \beta_j^2\right)} \\ &\propto \frac{1}{(\sigma\sqrt{2\pi})^n} * \frac{1}{(\sqrt{2\pi c})^p} e^{\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 - \frac{1}{2c} \sum_{j=1}^p \beta_j^2\right)} \end{aligned} \quad (3.6)$$

3.5 e

To prove the argument, we need to prove that the maximum of the posterior is given by the Ridge Regression with λ :

$$\log(\text{posterior}) = \log\left(\frac{1}{2b} * \frac{1}{(\sigma\sqrt{2\pi})^n}\right) + \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 - \frac{1}{2c} \sum_{j=1}^p \beta_j^2\right) \quad (3.7)$$

To find the maximum for the posterior, we need to find the minimum for the

$\left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{1}{2c} \sum_{j=1}^p \beta_j^2\right)$:

$$\begin{aligned} \arg\min_{\beta} \left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{1}{2c} \sum_{j=1}^p \beta_j^2 \right) &= \arg\min_{\beta} \left(\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{2\sigma^2 \sum_{j=1}^p |\beta_j|}{2c} \right) \\ &= \arg\min_{\beta} \left(\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij}\beta_j - \beta_0)^2 + \frac{\sigma^2}{c} \sum_{j=1}^p |\beta_j| \right) \end{aligned} \quad (3.8)$$

Let $\lambda = \frac{\sigma^2}{c}$, that's the same form with Ridge Regression.

4 P4

4.1 a

```
1 set.seed(1)
2 x = rnorm(100)
3 eps = rnorm(100)
```

4.2 b

```
1 beta = c(1, -1, 2, 0.5)
2 y = beta[1] + beta[2]*x + beta[3]*x^2 + beta[4]*x^3 + eps
```

4.3 c

```
1 library(leaps)
2 ans.df = data.frame(y = y, x = x)
3 ans.model = regsubsets(y ~ poly(x, 10, raw = T),
4                       data = ans.df, nvmax = 10)
5 ans.summary = summary(ans.model)
6
7 min_p = which.min(ans.summary$cp)
8 plot(ans.summary$cp, xlab = "Subset_size", ylab = "Cp",
9      col = "blue", pch = 20, type = "o")
10 points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
11 title("Best_Subset_with_Cp")
12 coefficients(ans.model, id = which.min(ans.summary$cp))
13
14 min_p = which.min(ans.summary$bic)
15 plot(ans.summary$bic, xlab = "Subset_size", ylab = "BIC",
16      col = "blue", pch = 20, type = "o")
17 points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
18 title("Best_Subset_with_BIC")
19 coefficients(ans.model, id = which.min(ans.summary$bic))
20
21 max_p = which.max(ans.summary$adjr2)
22 plot(ans.summary$adjr2, xlab = "Subset_size", ylab = "AdjR2",
23      col = "blue", pch = 20, type = "o")
24 points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
25 title("Best_Subset_with_AdjR2")
26 coefficients(ans.model, id = which.max(ans.summary$adjr2))
```

The best model uses a X^5 term. $|\beta_0 - \hat{\beta}_0| = 0.072$, $|\beta_1 - \hat{\beta}_1| = 0.445$, $|\beta_2 - \hat{\beta}_2| = 0.157$

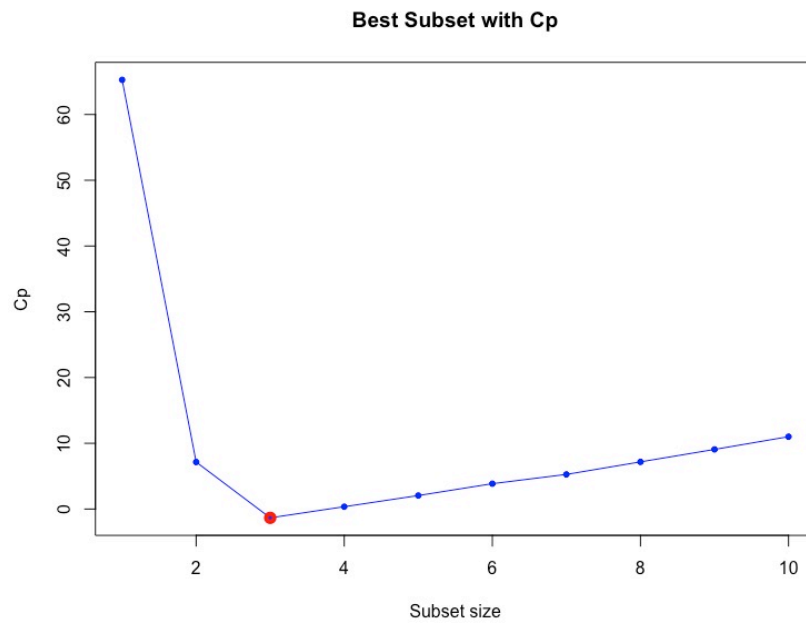


Figure 4.1: Screenshot for 'Best Subset with C_p '

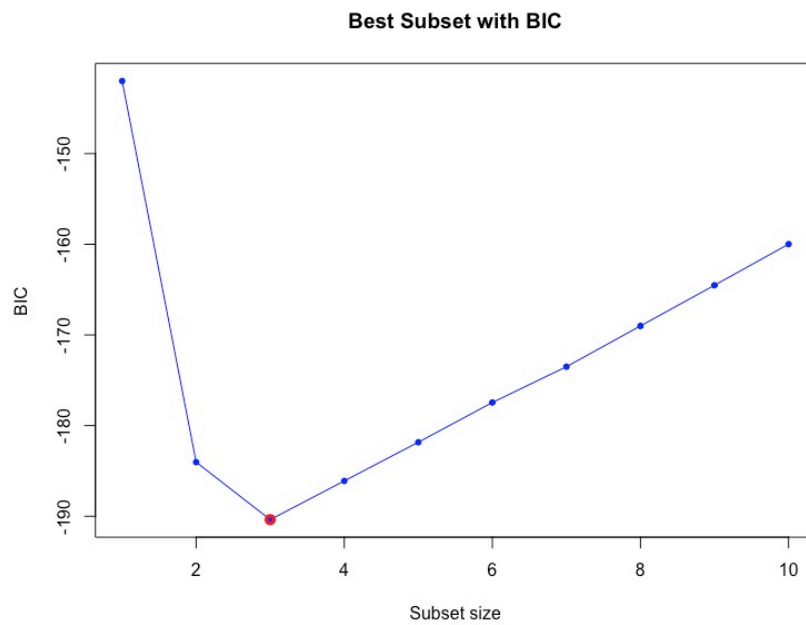


Figure 4.2: Screenshot for 'Best Subset with BIC'

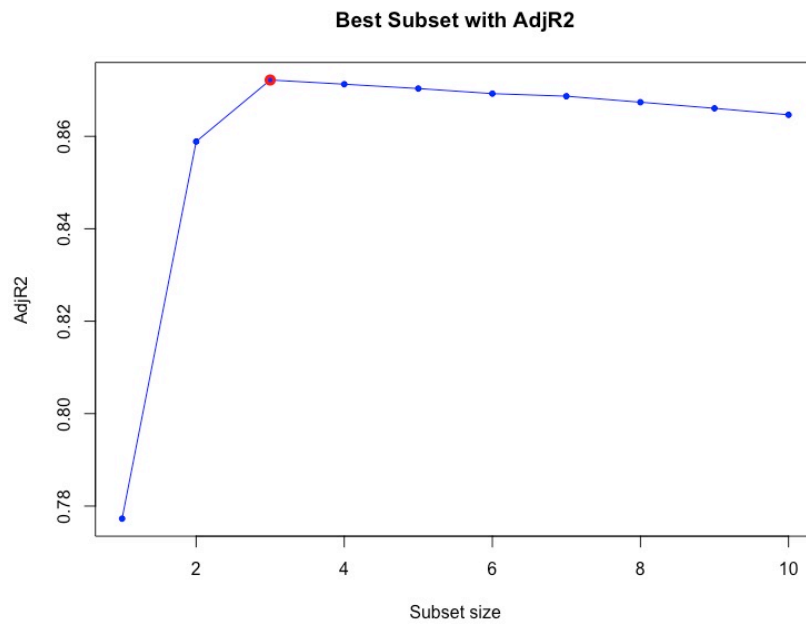


Figure 4.3: Screenshot for 'Best Subset with Adjust R^2 '

```

R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
> ans.model = regsubsets(y ~ poly(x, 10, row = T), data = ans.df, nvmax = 10)
> ans.summary = summary(ans.model)
> min_p = which.min(ans.summary$cp)
> plot(ans.summary$cp, xlab = "Subset size", ylab = "Cp", col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
> title("Best Subset with Cp")
> coefficients(ans.model, id = which.min(ans.summary$cp))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)5
1.07219472 -0.55485280 1.84323764 0.09022577
> min_p = which.min(ans.summary$bic)
> plot(ans.summary$bic, xlab = "Subset size", ylab = "BIC", col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
> title("Best Subset with BIC")
> coefficients(ans.model, id = which.min(ans.summary$bic))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)5
1.07219472 -0.55485280 1.84323764 0.09022577
> max_p = which.max(ans.summary$adjr2)
> plot(ans.summary$adjr2, xlab = "Subset size", ylab = "AdjR2", col = "blue", pch = 20, type = "o")
> points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
> title("Best Subset with AdjR2")
> coefficients(ans.model, id = which.max(ans.summary$adjr2))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)5
1.07219472 -0.55485280 1.84323764 0.09022577
>

```

Figure 4.4: Screenshot for 'Coefficients'

4.4 d

If we use forward selection:

```
1 library(leaps)
2 ans.df = data.frame(y = y, x = x)
3 ans.model = regsubsets(y ~ poly(x, 10, raw = T),
4                       data = ans.df, nvmax = 10, method = "forward")
5 ans.summary = summary(ans.model)
6
7 min_p = which.min(ans.summary$cp)
8 plot(ans.summary$cp, xlab = "Subset_size", ylab = "Cp",
9      col = "blue", pch = 20, type = "o")
10 points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
11 title("Best Subset with Cp using 'forward'")
12 coefficients(ans.model, id = which.min(ans.summary$cp))
13
14 min_p = which.min(ans.summary$bic)
15 plot(ans.summary$bic, xlab = "Subset_size", ylab = "BIC",
16      col = "blue", pch = 20, type = "o")
17 points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
18 title("Best Subset with BIC using 'forward'")
19 coefficients(ans.model, id = which.min(ans.summary$bic))
20
21 max_p = which.max(ans.summary$adjr2)
22 plot(ans.summary$adjr2, xlab = "Subset_size", ylab = "AdjR2",
23      col = "blue", pch = 20, type = "o")
24 points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
25 title("Best Subset with AdjR2 using 'forward'")
26 coefficients(ans.model, id = which.max(ans.summary$adjr2))
```

Forward stepwise uses a X^5 and a X^7 term. $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$ do not change a lot.

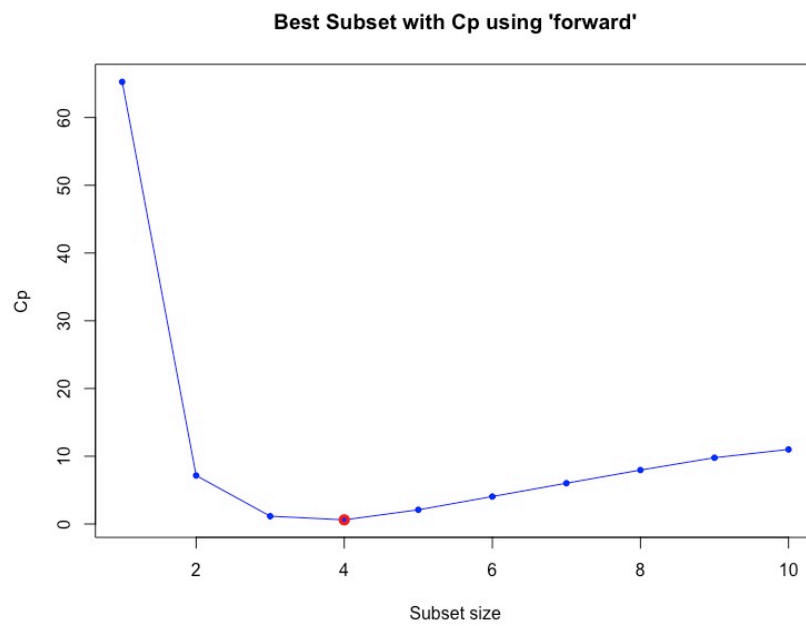


Figure 4.5: Screenshot for 'Best Subset with C_p using 'forward' '

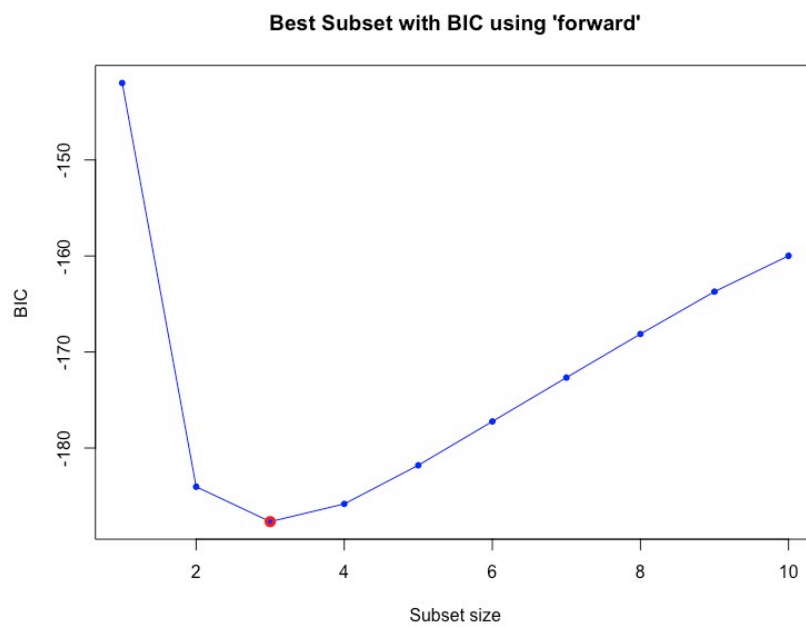


Figure 4.6: Screenshot for 'Best Subset with BIC using 'forward' '

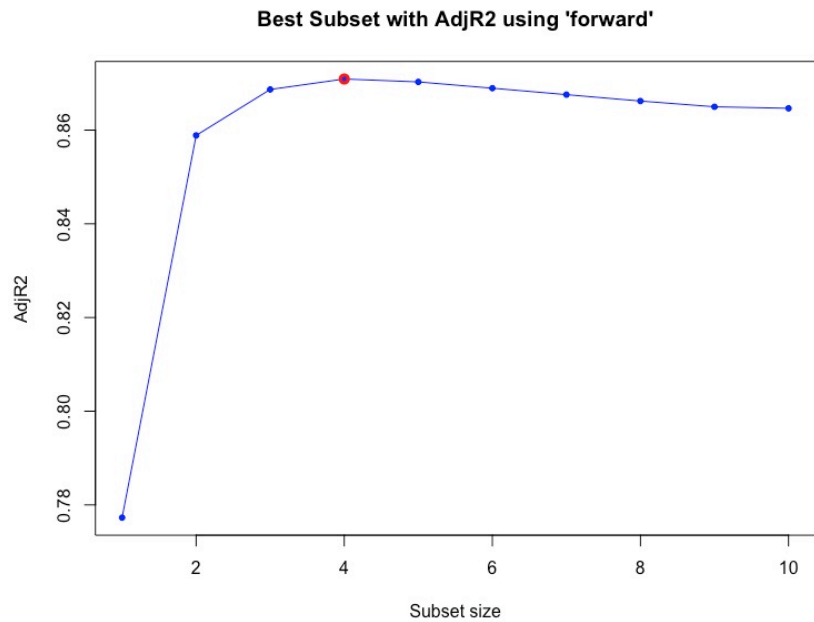


Figure 4.7: Screenshot for 'Best Subset with Adjust R^2 using forward'

```

R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
> plot(ans.summary$cp, xlab = "Subset size", ylab = "Cp",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
> title("Best Subset with Cp using 'forward'")
> coefficients(ans.model, id = which.min(ans.summary$cp))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)5 poly(x, 10, row = T)7
1.070684876 -0.581441777 1.847950163 0.104710260 -0.002672665
> min_p = which.min(ans.summary$bic)
> plot(ans.summary$bic, xlab = "Subset size", ylab = "BIC",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
> title("Best Subset with BIC using 'forward'")
> coefficients(ans.model, id = which.min(ans.summary$bic))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)7
1.07636566 -0.36989702 1.82056631 0.01610623
> max_p = which.max(ans.summary$adjr2)
> plot(ans.summary$adjr2, xlab = "Subset size", ylab = "AdjR2",
+      col = "blue", pch = 20, type = "o")
> points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
> title("Best Subset with AdjR2 using 'forward'")
> coefficients(ans.model, id = which.max(ans.summary$adjr2))
(Intercept) poly(x, 10, row = T)1 poly(x, 10, row = T)2 poly(x, 10, row = T)5 poly(x, 10, row = T)7
1.070684876 -0.581441777 1.847950163 0.104710260 -0.002672665
>

```

Figure 4.8: Screenshot for 'Coefficients using forward'

If we use backward selection:

```

1 library(leaps)
2 ans.df = data.frame(y = y, x = x)
3 ans.model = regsubsets(y ~ poly(x, 10, raw = T),
4                       data = ans.df, nvmax = 10, method = "backward")
5 ans.summary = summary(ans.model)
6
7 min_p = which.min(ans.summary$cp)
8 plot(ans.summary$cp, xlab = "Subset_size", ylab = "Cp",
9      col = "blue", pch = 20, type = "o")
10 points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
11 title("Best Subset with Cp using 'backward'")
12 coefficients(ans.model, id = which.min(ans.summary$cp))
13
14 min_p = which.min(ans.summary$bic)
15 plot(ans.summary$bic, xlab = "Subset_size", ylab = "BIC",
16      col = "blue", pch = 20, type = "o")
17 points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
18 title("Best Subset with BIC using 'backward'")
19 coefficients(ans.model, id = which.min(ans.summary$bic))
20
21 max_p = which.max(ans.summary$adjr2)
22 plot(ans.summary$adjr2, xlab = "Subset_size", ylab = "AdjR2",
23      col = "blue", pch = 20, type = "o")
24 points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
25 title("Best Subset with AdjR2 using 'backward'")
26 coefficients(ans.model, id = which.max(ans.summary$adjr2))

```

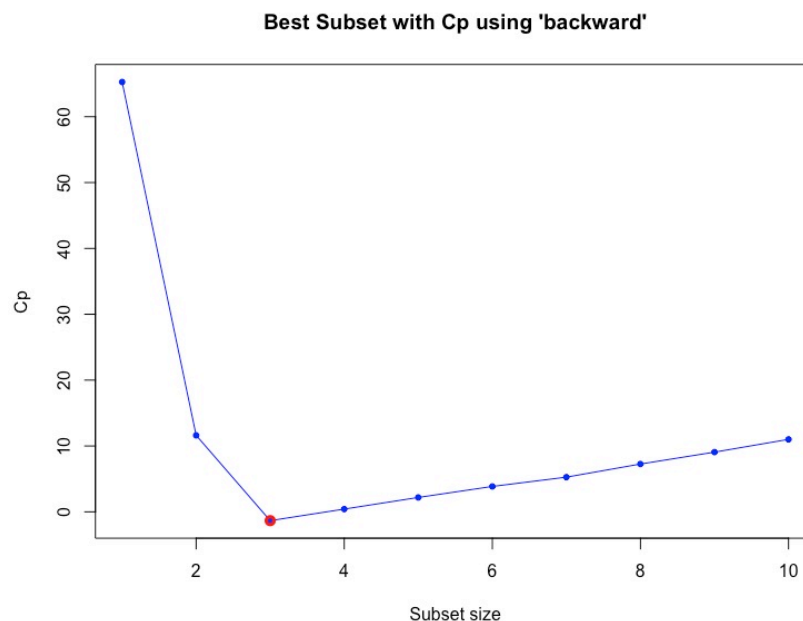


Figure 4.9: Screenshot for 'Best Subset with C_p using 'backward''

Backward stepwise uses a X^5 term. $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$ do not change a lot.

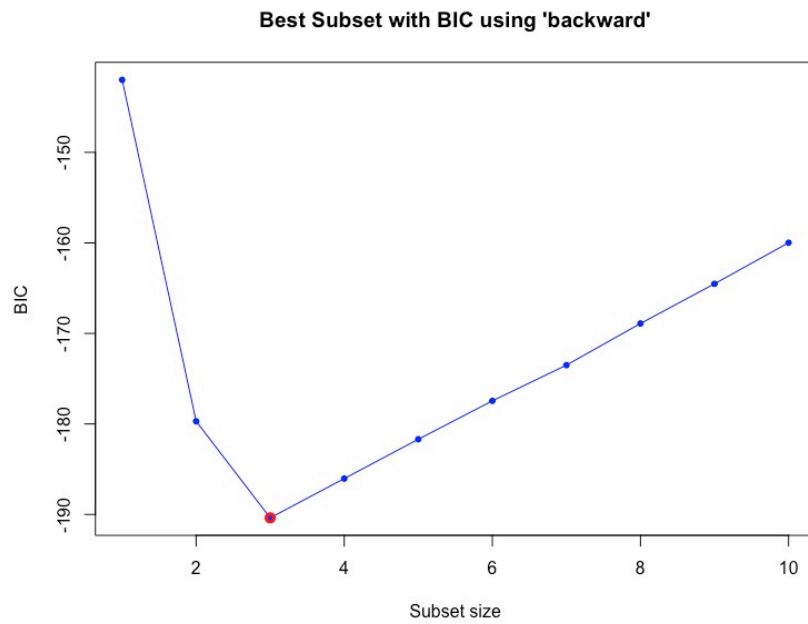


Figure 4.10: Screenshot for 'Best Subset with BIC using 'backward''

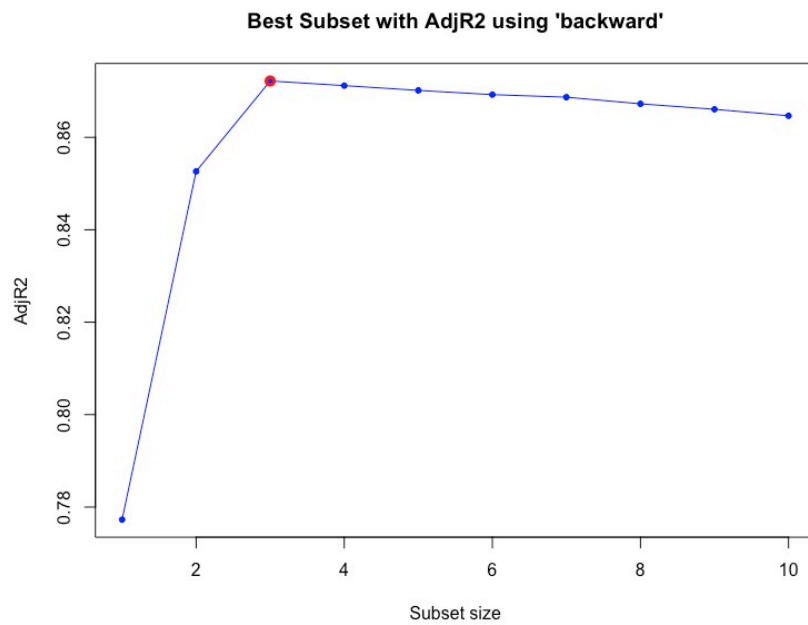


Figure 4.11: Screenshot for 'Best Subset with Adjust R^2 using 'backward''

```
R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
> plot(ans.summary$cp, xlab = "Subset size", ylab = "Cp",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
> title("Best Subset with Cp using 'backward'")
> coefficients(ans.model, id = which.min(ans.summary$cp))
      (Intercept) poly(x, 10, raw = T)1 poly(x, 10, raw = T)2 poly(x, 10, raw = T)5
      1.07219472      -0.55485280      1.84323764      0.09022577
> min_p = which.min(ans.summary$bic)
> plot(ans.summary$bic, xlab = "Subset size", ylab = "BIC",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
> title("Best Subset with BIC using 'backward'")
> coefficients(ans.model, id = which.min(ans.summary$bic))
      (Intercept) poly(x, 10, raw = T)1 poly(x, 10, raw = T)2 poly(x, 10, raw = T)5
      1.07219472      -0.55485280      1.84323764      0.09022577
> max_p = which.max(ans.summary$adjr2)
> plot(ans.summary$adjr2, xlab = "Subset size", ylab = "AdjR2",
+      col = "blue", pch = 20, type = "o")
> points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
> title("Best Subset with AdjR2 using 'backward'")
> coefficients(ans.model, id = which.max(ans.summary$adjr2))
      (Intercept) poly(x, 10, raw = T)1 poly(x, 10, raw = T)2 poly(x, 10, raw = T)5
      1.07219472      -0.55485280      1.84323764      0.09022577
> |
```

Figure 4.12: Screenshot for 'Coefficients using 'backward''

4.5 e

```

1 library(glmnet)
2 xmat = model.matrix(y ~ poly(x, 10, raw = T), data = ans.df)[, -1]
3 ans.lasso = cv.glmnet(xmat, y, alpha = 1)
4 ans.lambda = ans.lasso$lambda.min
5 ans.lambda
6 plot(ans.lasso)
7
8 ans.lasso = glmnet(xmat, y, alpha = 1)
9 predict(ans.lasso, s = ans.lambda, type = "coefficients")

```

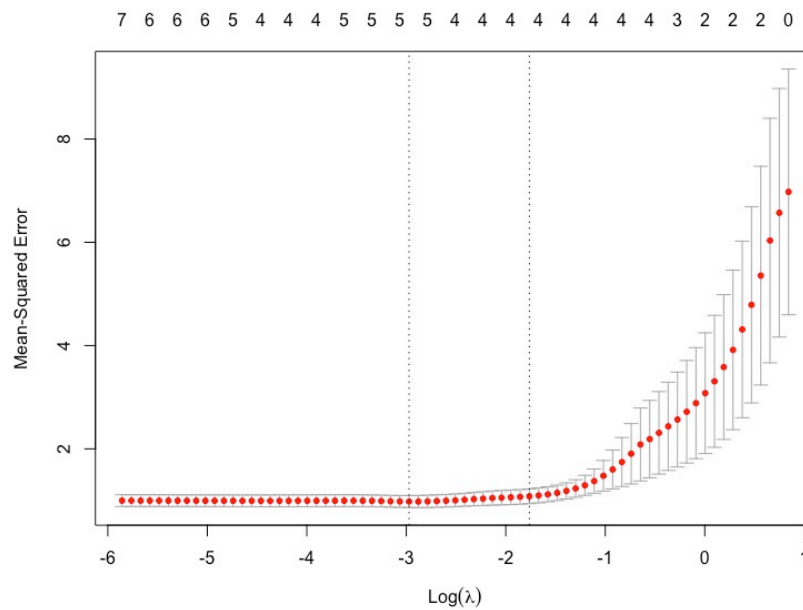


Figure 4.13: Screenshot for 'LASSO Cross-validation'



Figure 4.14: Screenshot for 'Coefficients'

LASSO uses a X^4 and X^5 term. $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$ do not change a lot.

4.6 f

```
1 beta_7 = 3
2 yy = beta[1] * beta_7*x^7 + eps
3 ans.df = data.frame(y = yy, x = x)
4 ans.model = regsubsets(y ~ poly(x, 10, raw = T),
5                       data = ans.df, nvmax = 10)
6 ans.summary = summary(ans.model)
7
8 min_p = which.min(ans.summary$cp)
9 plot(ans.summary$cp, xlab = "Subset_size", ylab = "Cp",
10      col = "blue", pch = 20, type = "o")
11 points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
12 title("Best_Subset_with_Cp")
13 coefficients(ans.model, id = which.min(ans.summary$cp))
14
15 min_p = which.min(ans.summary$bic)
16 plot(ans.summary$bic, xlab = "Subset_size", ylab = "BIC",
17      col = "blue", pch = 20, type = "o")
18 points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
19 title("Best_Subset_with_BIC")
20 coefficients(ans.model, id = which.min(ans.summary$bic))
21
22 max_p = which.max(ans.summary$adjr2)
23 plot(ans.summary$adjr2, xlab = "Subset_size", ylab = "AdjR2",
24      col = "blue", pch = 20, type = "o")
25 points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
26 title("Best_Subset_with_AdjR2")
27 coefficients(ans.model, id = which.max(ans.summary$adjr2))
```

Among all three criterion, BIC picks up the best estimation of the true model.

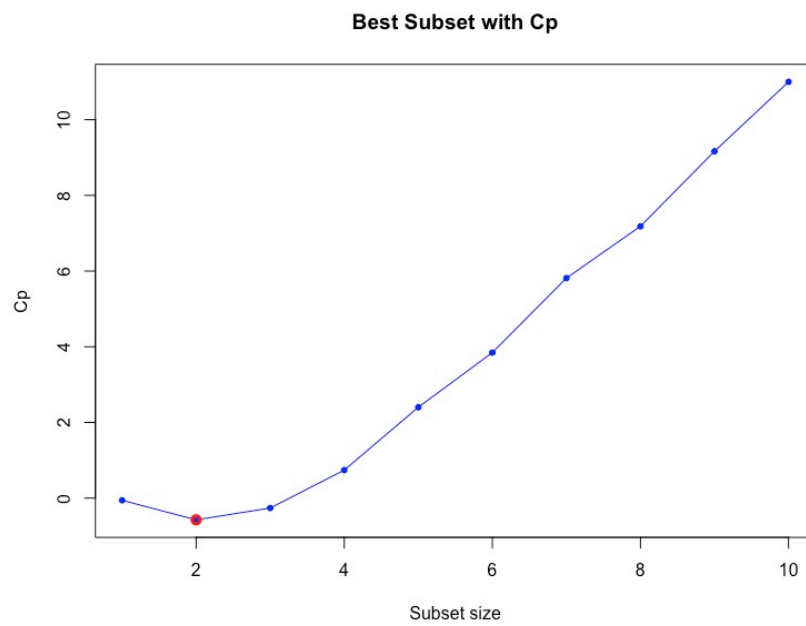


Figure 4.15: Screenshot for 'Best Subset with C_p '

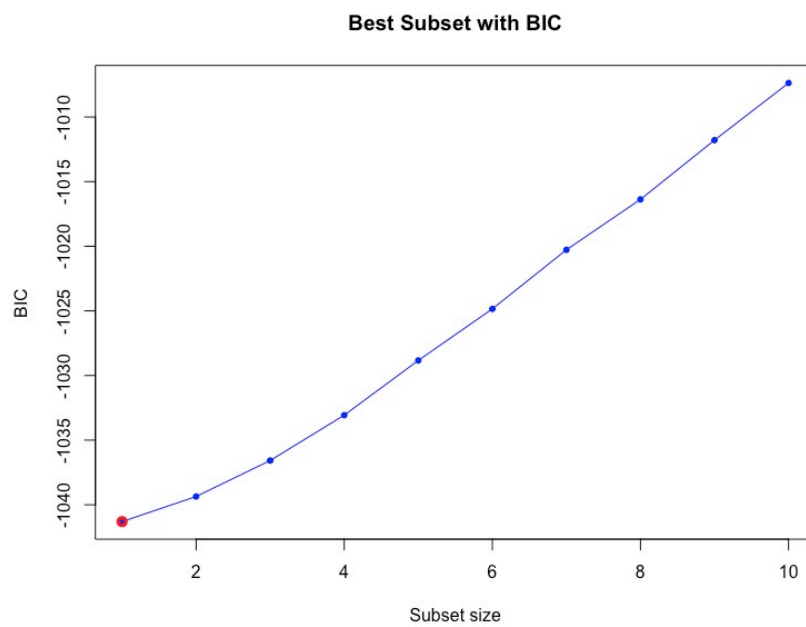


Figure 4.16: Screenshot for 'Best Subset with BIC'

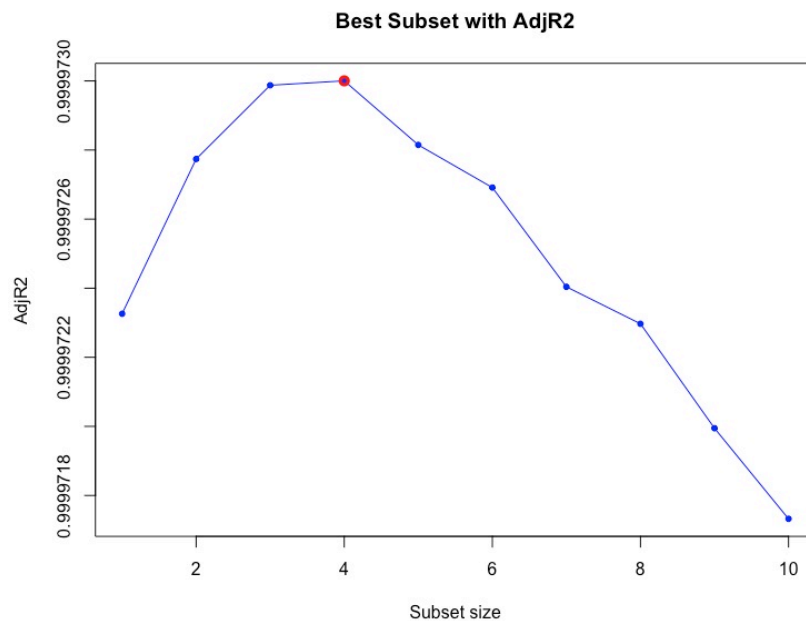


Figure 4.17: Screenshot for 'Best Subset with Adjust R^2 '

```

R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
> plot(ans.summary$cp, xlab = "Subset size", ylab = "Cp",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$cp[min_p], col = "red", lwd = 3)
> title("Best Subset with Cp")
> coefficients(ans.model, id = which.min(ans.summary$cp))
      (Intercept) poly(x, 10, raw = T)2 poly(x, 10, raw = T)7
      0.07049037      -0.14170843      3.00155519
> min_p = which.min(ans.summary$bic)
> plot(ans.summary$bic, xlab = "Subset size", ylab = "BIC",
+      col = "blue", pch = 20, type = "o")
> points(min_p, ans.summary$bic[min_p], col = "red", lwd = 3)
> title("Best Subset with BIC")
> coefficients(ans.model, id = which.min(ans.summary$bic))
      (Intercept) poly(x, 10, raw = T)7
      -0.04105975      3.00077047
> max_p = which.max(ans.summary$adjr2)
> plot(ans.summary$adjr2, xlab = "Subset size", ylab = "AdjR2",
+      col = "blue", pch = 20, type = "o")
> points(max_p, ans.summary$adjr2[max_p], col = "red", lwd = 3)
> title("Best Subset with AdjR2")
> coefficients(ans.model, id = which.max(ans.summary$adjr2))
      (Intercept) poly(x, 10, raw = T)1 poly(x, 10, raw = T)2 poly(x, 10, raw = T)3 poly(x, 10, raw = T)7
      0.07625245      0.29140161      -0.16176713      -0.25265268      3.00913375
>

```

Figure 4.18: Screenshot for 'Coefficients'

```

1 xmat = model.matrix(y ~ poly(x, 10, raw = T), data = ans.df)[, -1]
2 ans.lasso = cv.glmnet(xmat, y, alpha = 1)
3 ans.lambda = ans.lasso$lambda.min
4 ans.lambda
5 plot(ans.lasso)
6
7 ans.lasso = glmnet(xmat, y, alpha = 1)
8 predict(ans.lasso, s = ans.lambda, type = "coefficients")

```

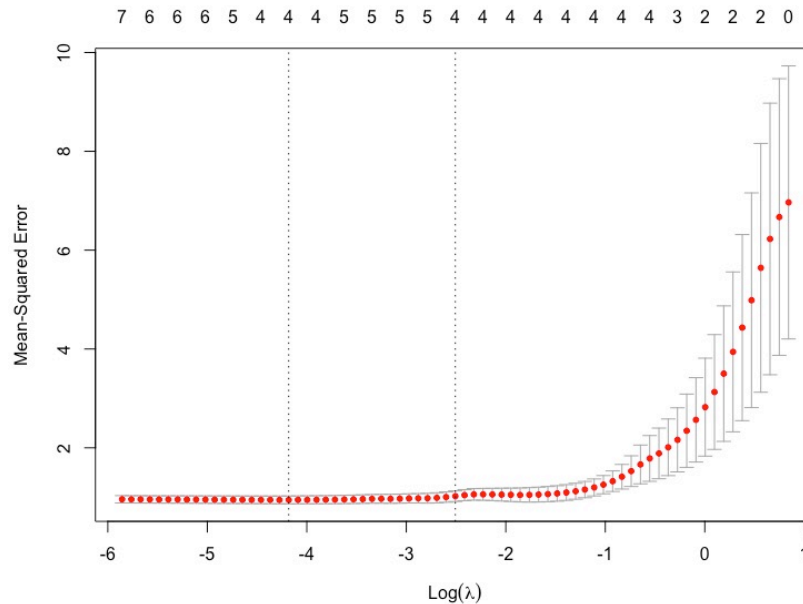


Figure 4.19: Screenshot for 'LASSO Cross-validation'

```

R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
(Intercept) poly(x, 10, raw = T)1 poly(x, 10, raw = T)2 poly(x, 10, raw = T)3 poly(x, 10, raw = T)7
0.07625245 0.29140161 -0.16176713 -0.25265268 3.00913375
> xmat = model.matrix(y ~ poly(x, 10, raw = T), data = ans.df)[, -1]
> ans.lasso = cv.glmnet(xmat, y, alpha = 1)
> ans.lambda = ans.lasso$lambda.min
> ans.lambda
[1] 0.0152686
> plot(ans.lasso)
> ans.lasso = glmnet(xmat, y, alpha = 1)
> predict(ans.lasso, s = ans.lambda, type = "coefficients")
11 x 1 sparse Matrix of class "dgMatrix"
      s1
(Intercept) 1.12529096
poly(x, 10, raw = T)1 -0.48567905
poly(x, 10, raw = T)2 1.68745208
poly(x, 10, raw = T)3 .
poly(x, 10, raw = T)4 0.03628426
poly(x, 10, raw = T)5 0.08384854
poly(x, 10, raw = T)6 .
poly(x, 10, raw = T)7 .
poly(x, 10, raw = T)8 .
poly(x, 10, raw = T)9 .
poly(x, 10, raw = T)10 .
>

```

Figure 4.20: Screenshot for 'Coefficients'

LASSO gives a best estimation of the true model. Only one variable besides the Intercept is given in the estimation.

5 P5

5.1 a

```
1 library(ISLR)
2 set.seed(1)
3 train.size = nrow(College) / 2
4 train = sample(1: nrow(College), train.size)
5 test = -train
6 College.train = College[train, ]
7 College.test = College[test, ]
```

5.2 b

```
1 College.lm = lm(Apps~., data = College.train)
2 summary(College.lm)
3
4 College.pred = predict(College.lm, College.test)
5 College.error = mean((College.test$Apps - College.pred)^2)
```

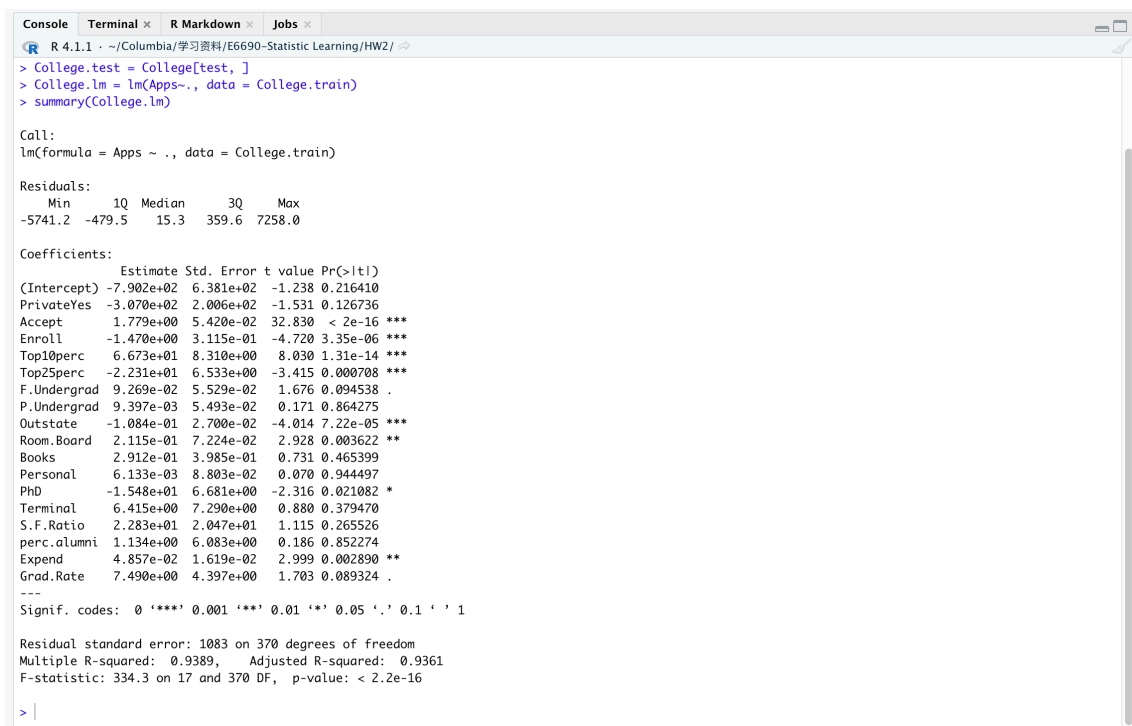


Figure 5.1: Screenshot for 'Model summary'

The error for test dataset is 1135758.

5.3 c

```
1 train_mat = model.matrix(Apps ~ ., data = College.train)
2 test_mat = model.matrix(Apps ~ ., data = College.test)
3 College.ridge = cv.glmnet(train_mat, College.train$Apps, alpha = 0)
4 plot(College.ridge)
```

```

5 College.lambda = College.ridge$lambda.min
6 College.pred = predict(College.ridge, newx=test_mat, s=College.lambda)
7 College.error = mean((College.test$Apps - College.pred)^2)

```

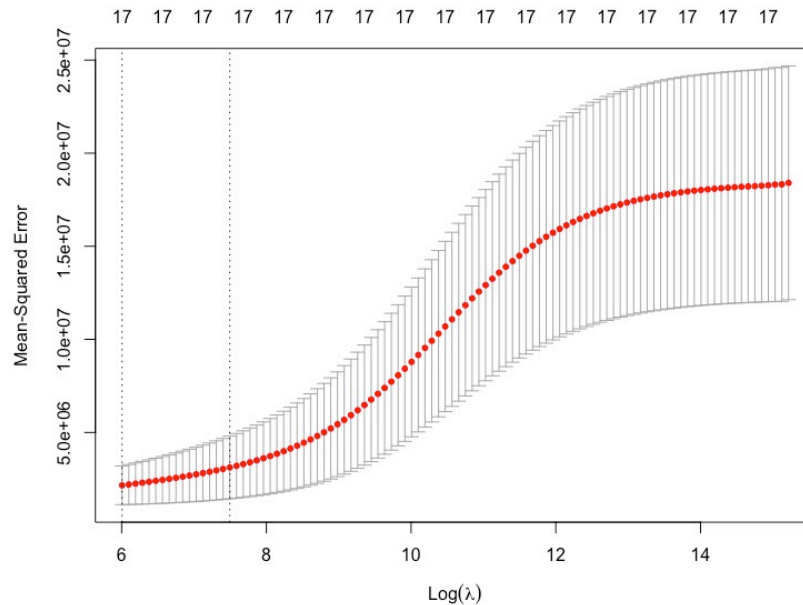


Figure 5.2: Screenshot for 'Ridge Regression'

The error for test dataset is 976261.5. Ridge Regression get a better result than Linear Regression.

5.4 d

```

1 College.lasso = cv.glmnet(train_mat, College.train$Apps, alpha = 1)
2 plot(College.lasso)
3 College.lambda = College.lasso$lambda.min
4 College.pred = predict(College.lasso, newx=test_mat, s=College.lambda)
5 College.error = mean((College.test$Apps - College.pred)^2)

```

The error for test dataset is 1115901. LASSO Regression get a better result than Linear Regression but not as good as Ridge Regression.

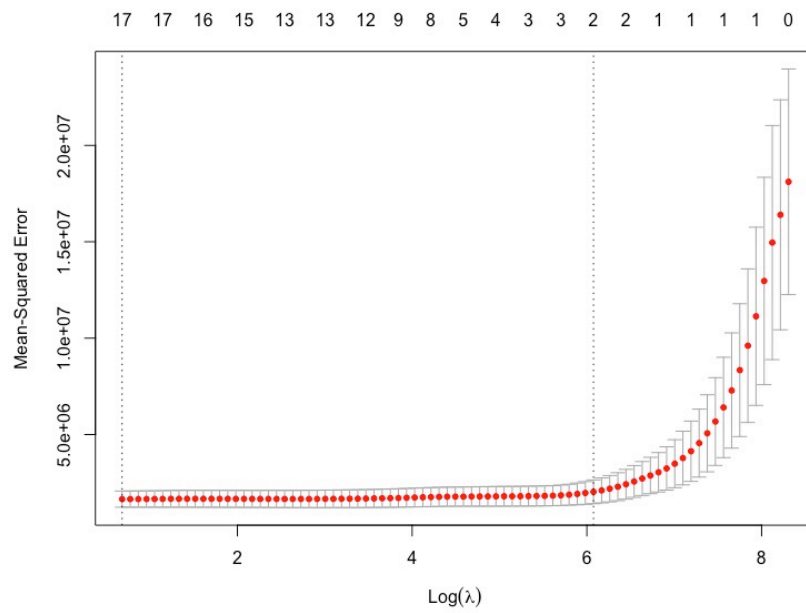


Figure 5.3: Screenshot for 'LASSO Regression'

6 P6

6.1 a

```
1 set.seed(1)
2 n = 1000
3 p = 20
4
5 xmat = matrix(rnorm(n*p), n, p)
6 eps = rnorm(n)
7
8 beta = rnorm(p)
9 beta[3] = 0
10 beta[6] = 0
11 beta[9] = 0
12 beta[12] = 0
13 beta[15] = 0
14 beta[18] = 0
15
16
17 y = xmat %*% beta + eps
```

I set $\beta_3, \beta_6, \beta_9, \beta_{12}, \beta_{15}, \beta_{18}$ to be zero.

6.2 b

```
1 train = sample(1: n, 100)
2 test = -train
3 xmat.train = xmat[train, ]
4 xmat.test = xmat[test, ]
5 y.train = y[train, ]
6 y.test = y[test, ]
```

6.3 c

```
1 ans.df = data.frame(x = xmat.train, y = y.train)
2 ans.model = regsubsets(y ~ ., data = ans.df, nvmax = p)
3 ans.summary = summary(ans.model)
4 ans.summary
5
6 ans.train_mse = rep(NA, p)
7 ans.xcols = colnames(xmat, do.NULL = F, prefix = "x.")
8 for(i in 1:p){
9   c_i = coef(ans.model, id = i)
10  if(i > 1){
11    y.train.pred = as.matrix(xmat.train[, ans.xcols %in% names(c_i)] %*%
12                          c_i[names(c_i) %in% ans.xcols])
13  }
14  else
15  {
16    y.train.pred = as.matrix(xmat.train[, ans.xcols %in% names(c_i)] *
17                          c_i[names(c_i) %in% ans.xcols])
18  }
```

```

19   ans.train_mse[i] = mean((y.train-y.train.pred)^2)
20 }
21 plot(ans.train_mse, ylab = "Training_MSE",xlab = "Subset_Size",
22       pch = 20, type = "o", col = "blue")
23 title("MSE-Subset_size_Curve")

```

```

R 4.1.1 - ~/Columbia/学习资料/E6690-Statistic Learning/HW2/
> ans.summary
Subset selection object
Call: regsubsets.formula(y ~ ., data = ans.df, nvmax = p)
20 Variables (and intercept)
Forced in Forced out
x.1 FALSE FALSE
x.2 FALSE FALSE
x.3 FALSE FALSE
x.4 FALSE FALSE
x.5 FALSE FALSE
x.6 FALSE FALSE
x.7 FALSE FALSE
x.8 FALSE FALSE
x.9 FALSE FALSE
x.10 FALSE FALSE
x.11 FALSE FALSE
x.12 FALSE FALSE
x.13 FALSE FALSE
x.14 FALSE FALSE
x.15 FALSE FALSE
x.16 FALSE FALSE
x.17 FALSE FALSE
x.18 FALSE FALSE
x.19 FALSE FALSE
x.20 FALSE FALSE
1 subsets of each size up to 20
Selection Algorithm: exhaustive
x.1 x.2 x.3 x.4 x.5 x.6 x.7 x.8 x.9 x.10 x.11 x.12 x.13 x.14 x.15 x.16 x.17 x.18 x.19 x.20
1 ( 1 )
2 ( 1 )
3 ( 1 )
4 ( 1 )
5 ( 1 )
6 ( 1 )
7 ( 1 )
8 ( 1 )
9 ( 1 )
10 ( 1 )
11 ( 1 )
12 ( 1 )
13 ( 1 )
14 ( 1 )
15 ( 1 )
16 ( 1 )
17 ( 1 )
18 ( 1 )
19 ( 1 )
20 ( 1 )
>

```

Figure 6.1: Screenshot for 'Model summary'

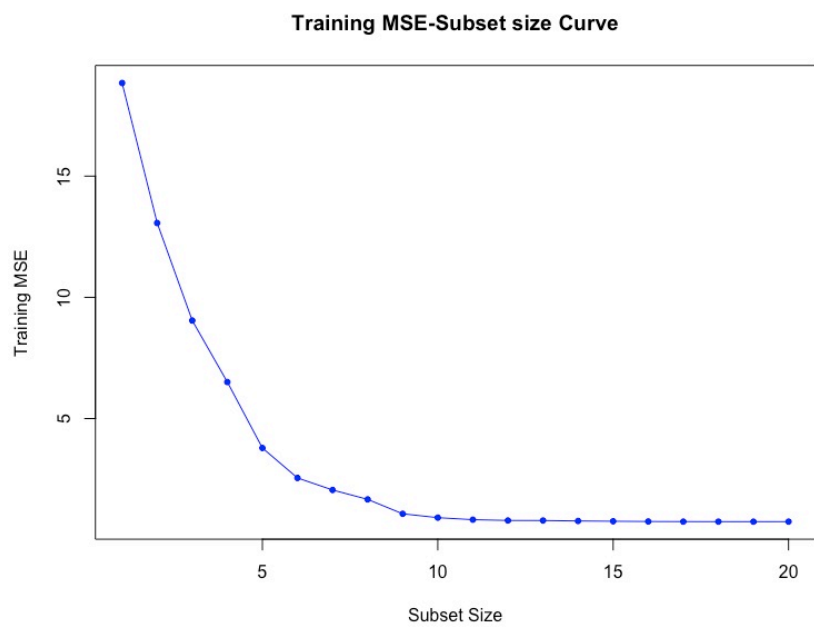


Figure 6.2: Screenshot for '*Training MSE-Subset size Curve*'

6.4 d

```
1 ans.test_mse = rep(NA, p)
2 ans.xcols = colnames(xmat, do.NULL = F, prefix = "x.")
3 for(i in 1:p){
4   c_i = coef(ans.model, id = i)
5   if(i > 1){
6     y.test.pred = as.matrix(xmat.test[, ans.xcols %in% names(c_i)] %*%
7                           c_i[names(c_i) %in% ans.xcols])
8   }
9   else
10  {
11    y.test.pred = as.matrix(xmat.test[, ans.xcols %in% names(c_i)] *
12                          c_i[names(c_i) %in% ans.xcols])
13  }
14  ans.test_mse[i] = mean((y.test-y.test.pred)^2)
15 }
16 plot(ans.test_mse, ylab = "Testing_MSE",xlab = "Subset_Size",
17      pch = 20, type = "o", col = "blue")
18 title("Testing_MSE-Subset_size_Curve")
```

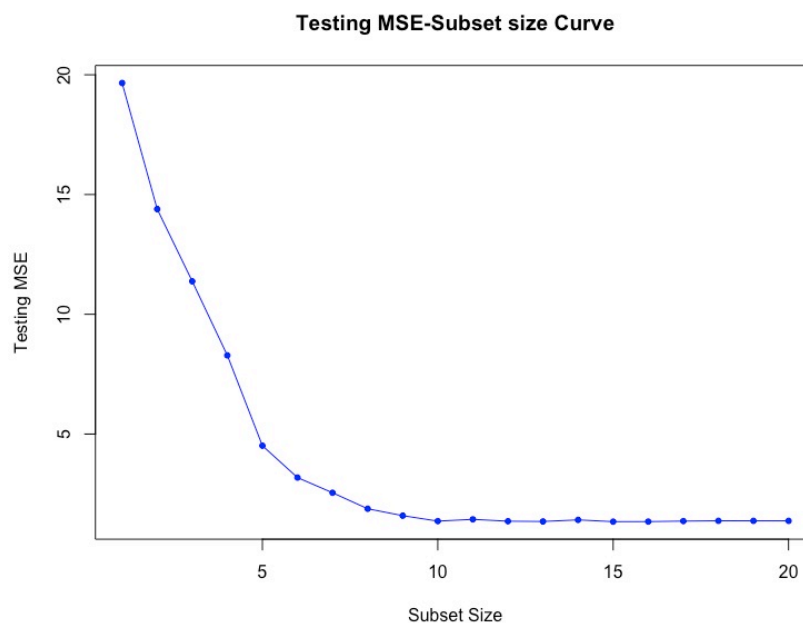


Figure 6.3: Screenshot for 'Testing MSE-Subset size Curve'

6.5 e

```
1 which.min(ans.test_mse)
```

The return of this function is 15, which means when there are 15 predictors in the estimated model, the model reaches the smallest test error.

6.6 f

```
1 coef(ans.model, id = which.min(ans.test_mse))
2 beta
```

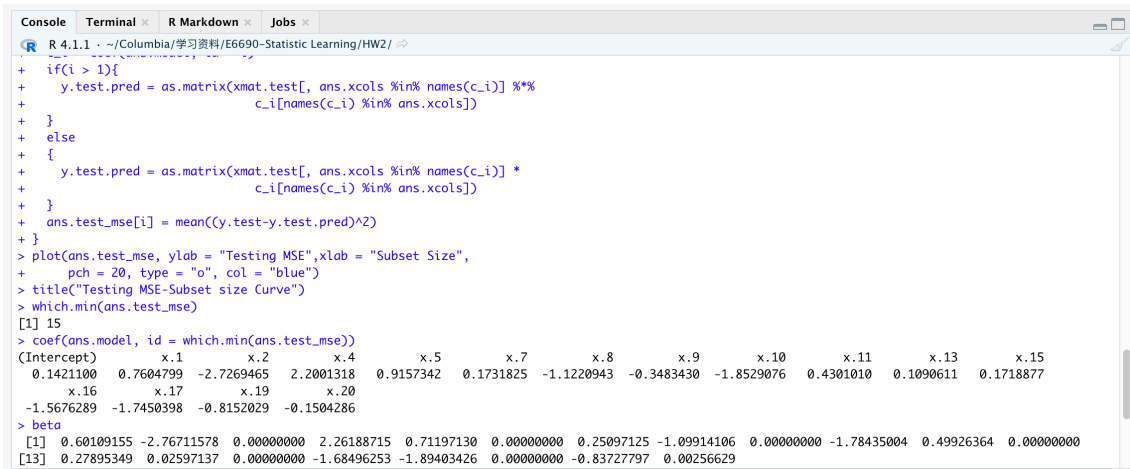


Figure 6.4: Screenshot for ' β value'

The model finds that $\beta_3, \beta_6, \beta_{12}, \beta_{14}, \beta_{18}$ are zeros.

6.7 g

```
1 ans.beta_error = rep(NA, p)
2 ans.xcols = colnames(xmat, do.NULL = F, prefix = "x.")
3 for(i in 1:p){
4   c_i = coef(ans.model, id = i)
5   ans.beta_error[i] = sqrt(sum((beta[ans.xcols %in% names(c_i)] -
6   c_i[names(c_i) %in% ans.xcols])^2) +
7   sum(beta[!(ans.xcols %in% names(c_i))])^2))
8 }
9 plot(x = 1:p, ans.beta_error, xlab = "Coefficient_#",
10      ylab = "Beta_error")
11 title("Beta_error_on_each_dimesion")
```

Compared with the result in (d), the model reaches the minimum beta error when there are 15 predictors. That's close to the result of 12 in (d). And when the number of predictors is 15, the β is also very small.

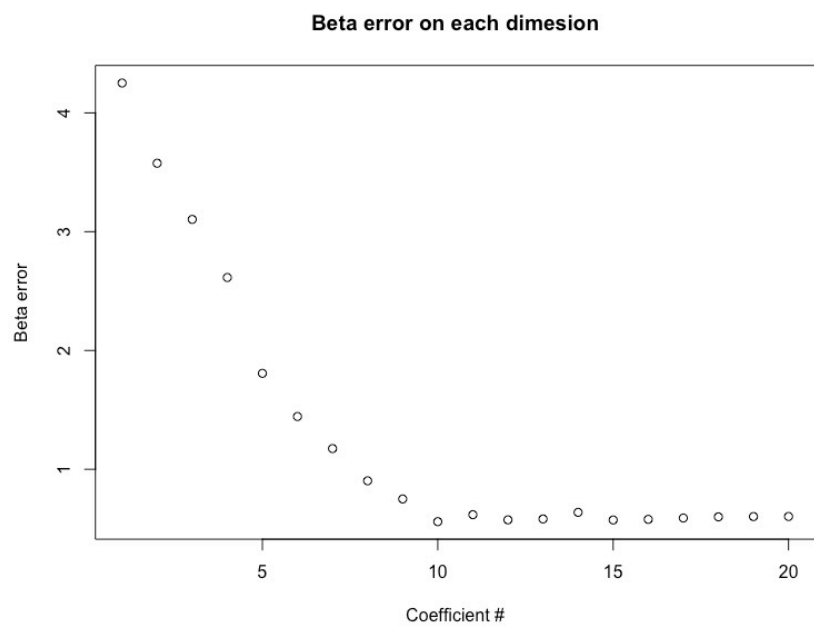


Figure 6.5: Screenshot for ' β error'