ISLR | Chapter 6 Exercises

 $Marshall\ McQuillen$ 8/3/2018

Conceptual

1

• A. For a model with k predictors, Best Subset Selection will always have the best training RSS. The reason for this is, given a fixed k, there are $\binom{p}{k}$ possible models, and Best Subset Selection considers all of those $\binom{p}{k}$ possibilities.

In Forward Stepwise Selection, of the total $\binom{p}{k}$ possible models, only the models that contain the (k-1) model produced by Forward Stepwise Selection will be considered for the "best" k-variable model.

In Backward Stepwise Selection, of the total $\binom{p}{k}$ possible models, the predictors in the "best" k-variable model must be a subset of the model with (k+1) predictors.

In short, Best Subset Selection will have the best (lowest) training RSS for a model with k predictors because it considers all the possible $\binom{p}{k}$ models, whereas Forward and Backward Stepwise Selection only consider a **subset** of all the possible $\binom{p}{k}$ models.

- B. There is no definitive answer for which subset selection method will have the lowest testing RSS (overfitting). If there is a large number of predictors, Best Subset Selection has the possibility of finding a model that has a low training RSS but a high testing RSS. Cross validation could be used to estimate the testing error of three models (one for Best Subset Selection, one for Forward Stepwise Selection and one for Backward Stepwise Selection) and a decision on which model has the lowest testing RSS could be made in consideration of the CV results.
- C.
 - i. True.
 - ii. True.
 - iii. False, the predictors in the k-variable model identified by Backward Subset Selection are **not** a subset of the predictors in the (k+1)-variable model identified by Forward Subset Selection.
 - iv. False, the predictors in the k-variable model identified by Forward Stepwise Selection are **not** a subset of the predictors in the (k + 1)-variable model identified by Backward Stepwise Selection.
 - v. False, the predictors in the k-variable model identified by Best Subset Selection are **not necessarily** a subset of the predictors in the (k + 1)-variable model identified by Best Subset Selection.

$\mathbf{2}$

- A. The lasso, relative to least squares is, *iii*, less flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.
- **B.** Ridge Regression, relative to least squares is, *iii*, less flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.
- C. Non-linear methods, relative to least squares are, *ii*, more flexible and hence will give improved prediction accuracy when their increase in variance is less than their decrease in bias.

In the (alternate) cost function for the Lasso...

$$\sum_{i=1}^{n} \left(y_k - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 subject \ to \ \sum_{j=1}^{p} |\beta_j| \le s$$

As we increase s from 0...

- A. ...the training RSS will, iv, steadily decrease. As the budget (s) for the sum of the regression coefficients increases from 0, each β_j will approach the value it would reach in ordinary least squares regression (no constraint). Therefore, without a constraint, we are letting the regression coefficients "roam freely" to reach their ordinary least squares values. Using the constraint, we are "lassoing" them in (pun intended).
- B. ...the testing RSS will, ii, decrease initially, and then eventually start increasing in a U shape. When s is 0, all the regression coefficients are 0, and the "model" is simply the intercept, β_0 , the mean of the response (highly biased, very low variance). As s increases from 0, the model becomes more flexible, allowing for an increasingly better fit to the data up to a point (bottom of the U). Once this point is reached, the model becomes **overly** flexible, overfitting the training data and leading to an increase in the test error (right side of U).
- C. ...variance will, *iii*, steadily increase. As s increases from 0, the model becomes more and more flexible, and as we know, more flexible models have a higher variance and lower bias. The model will be more influenced by the data it is trained on.
- **D**. ...(squared) bias will, *iv*, steadily decrease. As *s* increases from 0, the model becomes more and more flexible, and as we know, more flexible models have a higher variance and lower bias. The model will have a better chance of representing the *true* relationship between the predictors and the response.
- E. ...the irreducible error will, v, remain constant. The irreducible error is just that, irreducible.

4

In the cost function for Ridge Regression...

$$\sum_{i=1}^{n} \left(y_k - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

...as we increase λ from 0...

- A. ...the training RSS will, iii, steadily increase. As λ increases from 0, more "weight" is given to the second term in the cost function, thus penalizing large regression coefficients more and more. Increasing λ from 0 restricts the regression coefficients more and more.
- B. ...the testing RSS will, ii, decrease initially, and then eventually start increasing in a U shape. As λ increases from 0, increasingly strict restrictions are put on the magnitude that the regression coefficients can grow too. This has the effect of making the model less flexible and more generalizable, **up to a point**. The left side of the U represents the model's increase in bias being less than it's decrease in variance, and the right side of the U represents the decrease in variance no longer being worth the increase in bias.
- C. ...the variance will, iv, steadily decrease. As λ increases from 0, increasingly strict limits are placed on the regression coefficients, making it less flexible.

- **D**. ...the (squared) bias will, *iii*, steadily increase. As λ increases from 0, increasingly strict limits are placed on the regression coefficients, making it less flexible. This will make it increasingly harder for the model to estimate the true relationship between the response and the predictors.
- E. ...the irreducible error will, v, remain constant. The irreducible error is just that, **irreducible**.

5

• A. The optimization problem for Ridge Regression is:

$$\sum_{i=1}^{n} \left(y_k - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

Written out when n = p = 0, the equation comes to:

$$(y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2)$$

When $\beta_0 = 0$, a small simplification can be made such that the above equation becomes:

$$(y_1 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_1 x_{21} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2)$$

• **B.** Given that $x_{11} = x_{12}$ and $x_{21} = x_{22}$, I will refer to these as simply x_1 and x_2 respectively where applicable. A small rewrite of the preceding equation gives:

$$f(x) = (y_1 - \beta_1 x_1 - \beta_2 x_1)^2 + (y_2 - \beta_1 x_2 - \beta_2 x_2)^2 + \lambda (\beta_1^2 + \beta_2^2)$$

Given that the problem above is one of *minimization*, taking the derivative is the first step in showing that $\beta_1 = \beta_2$.

Calculating the derivative with respect to β_1 can be broken down into the derivative of the three separate terms in the above equation:

$$\frac{\partial}{\partial \beta_1} f(x) = \frac{\partial}{\partial \beta_1} \left(y_1 - \beta_1 x_1 - \beta_2 x_1 \right)^2 + \frac{\partial}{\partial \beta_1} \left(y_2 - \beta_1 x_2 - \beta_2 x_2 \right)^2 + \frac{\partial}{\partial \beta_1} \left(\lambda \beta_1^2 + \lambda \beta_2^2 \right)$$

First Term Derivative

$$\frac{\partial}{\partial \beta_1} (y_1 - \beta_1 x_1 - \beta_2 x_1)^2 = 2 (y_1 - \beta_1 x_1 - \beta_2 x_1) \cdot (-x_1) = 2 (-y_1 x_1 + \beta_1 x_1^2 + \beta_2 x_1^2)$$

Second Term Derivative

$$\frac{\partial}{\partial \beta_1} \left(y_2 - \beta_1 x_2 - \beta_2 x_2 \right)^2 = 2 \left(y_2 - \beta_1 x_2 - \beta_2 x_2 \right) \cdot \left(-x_2 \right) = 2 \left(-y_2 x_2 + \beta_1 x_2^2 + \beta_2 x_2^2 \right)$$

Third Term Derivative

$$\frac{\partial}{\partial \beta_1} \left(\lambda \beta_1^2 + \lambda \beta_2^2 \right) = 2\lambda \beta_1$$

Bringing this all together, the derivative of the full equation comes out to:

$$\frac{\partial}{\partial \beta_1} f(x) = 2 \left(-y_1 x_1 + \beta_1 x_1^2 + \beta_2 x_1^2 \right) + 2 \left(-y_2 x_2 + \beta_1 x_2^2 + \beta_2 x_2^2 \right) + 2\lambda \beta_1$$

Setting the derivative equal to 0 and solving for β_1 :

$$2\left(-y_1x_1 + \beta_1x_1^2 + \beta_2x_1^2\right) + 2\left(-y_2x_2 + \beta_1x_2^2 + \beta_2x_2^2\right) + 2\lambda\beta_1 = 0$$

Divide by 2:

$$(-y_1x_1 + \beta_1x_1^2 + \beta_2x_1^2) + (-y_2x_2 + \beta_1x_2^2 + \beta_2x_2^2) + \lambda\beta_1 = 0$$

Group β_1 terms and β_2 terms together:

$$(\beta_1 x_1^2 + \beta_1 x_2^2 + \lambda \beta_1) + \beta_2 x_1^2 + \beta_2 x_2^2 - y_1 x_1 - y_2 x_2 = 0$$

Factor out β_1 and β_2 :

$$\beta_1 (x_1^2 + x_2^2 + \lambda) + \beta_2 (x_1^2 + x_2^2) - y_1 x_1 - y_2 x_2 = 0$$

Add y_1x_1 and y_2x_2 to both sides of the equation:

$$\beta_1 (x_1^2 + x_2^2 + \lambda) + \beta_2 (x_1^2 + x_2^2) = y_1 x_1 + y_2 x_2$$

Repeating the same process, this time taking the derivative with respect to β_2 , would yield:

$$\beta_2 (x_1^2 + x_2^2 + \lambda) + \beta_1 (x_1^2 + x_2^2) = y_1 x_1 + y_2 x_2$$

Using substitution, we can set the two equations equal to each other:

$$\beta_1 (x_1^2 + x_2^2 + \lambda) + \beta_2 (x_1^2 + x_2^2) = \beta_2 (x_1^2 + x_2^2 + \lambda) + \beta_1 (x_1^2 + x_2^2)$$

Factoring λ out into it's own term:

$$\beta_1 \left(x_1^2 + x_2^2 \right) + \beta_1 \lambda + \beta_2 \left(x_1^2 + x_2^2 \right) = \beta_2 \left(x_1^2 + x_2^2 \right) + \beta_2 \lambda + \beta_1 \left(x_1^2 + x_2^2 \right)$$

Subract $\beta_1(x_1^2 + x_2^2)$ and $\beta_2(x_1^2 + x_2^2)$ then divide by λ :

$$\beta_1 \lambda = \beta_2 \lambda$$
 thus $\beta_1 = \beta_2$

• C. The lasso optimization problem is similar to that of Ridge Regression, with a small change to the third term in the equation:

$$f(x) = (y_1 - \beta_1 x_1 - \beta_2 x_1)^2 + (y_2 - \beta_1 x_2 - \beta_2 x_2)^2 + \lambda (|\beta_1| + |\beta_2|)$$

• D. NEED TO COME BACK TOO

• A. Considering the Ridge Regression loss function where p=1, the equation...

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

...can be re-written as the following.

$$f(\beta_j) = (y_j - \beta_j)^2 + \lambda \beta_j^2$$

A quick derivation shows that equation 6.12 is solved by 6.14:

$$\frac{\partial f(\beta_j)}{\partial \beta_j} = 2(y_j - \beta_j)(-1) + 2\lambda \beta_j$$

$$2(y_j - \beta_j)(-1) + 2\lambda \beta_j = 0$$

$$-2(y_j - \beta_j) + 2\lambda \beta_j = 0$$

$$-2y_j + 2\beta_j + 2\lambda \beta_j = 0$$

$$-y_j + \beta_j + \lambda \beta_j = 0$$

$$-y_j + \beta_j(1 + \lambda) = 0$$

$$\beta_j(1 + \lambda) = y_j$$

$$\beta_j = \frac{y_j}{(1 + \lambda)}$$

The following code holds $y_j = \lambda = 1$ and plots a sequence of β_j estimates in blue. The minimum of those estimates is circled in red, and the minimum estimated by the equation above is circled in green, both encompassing 0.5

```
ridge <- function(y, lambda, beta) {
    return((y - beta)^2 + lambda*beta^2)
}

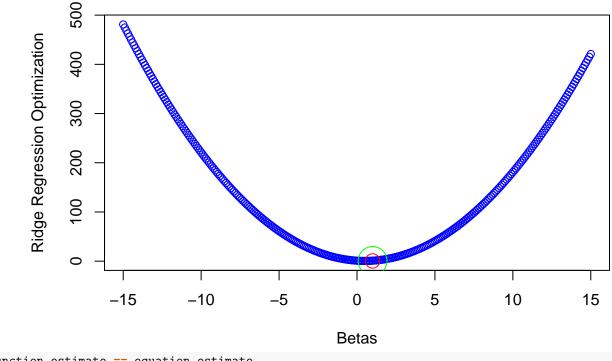
y = 1
lambda = 1
betas <- seq(-15, 15, 0.1)

f.x <- ridge(y, lambda, betas)

function_estimate <- f.x[which.min(f.x)]</pre>
```

```
equation_estimate <- y/(lambda + 1)
plot(x = betas,
     y = f.x,
     col = "blue",
     xlab = 'Betas',
     ylab = "Ridge Regression Optimization",
     main = "Ridge Regression Minimization")
points(function_estimate, col = 'red', cex = 2)
points(equation_estimate, col = 'green', cex = 4)
```

Ridge Regression Minimization



function_estimate == equation_estimate

[1] TRUE

• B. Considering the Lasso Regression loss function where p=1, the equation...

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

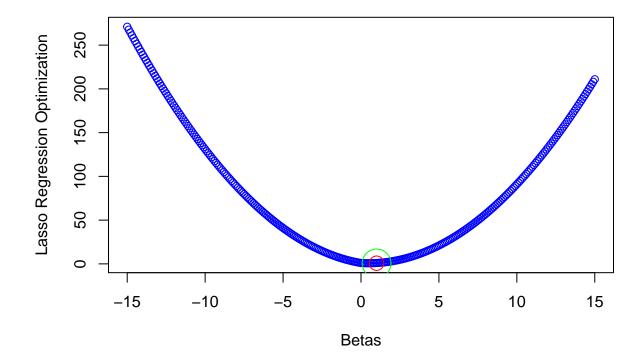
...can be re-written as the following.

$$f(\beta_j) = (y_j - \beta_j)^2 + \lambda |\beta_j|$$

The following code holds $y_j = \lambda = 1$ and plots a sequence of β_j estimates in blue. The minimum of those estimates is circled in red, and the minimum estimated by the equation (where $y_j > \frac{\lambda}{2}$) above is circled in green.

```
lasso <- function(beta, y=y, lambda=lambda) {</pre>
    return((y - beta)^2 + lambda*abs(beta))
}
y = 1
lambda = 1
betas \leftarrow seq(-15, 15, 0.1)
f.x <- lasso(betas, y, lambda)</pre>
function_estimate <- f.x[which.min(f.x)]</pre>
equation_estimate <- y - (lambda/2)
plot(x = betas,
     y = f.x,
     col = "blue",
     xlab = 'Betas',
     ylab = "Lasso Regression Optimization",
     main = "Minimization of Lasso where y > lambda/2")
points(function_estimate, col = 'red', cex = 2)
points(equation_estimate, col = 'green', cex = 4)
```

Minimization of Lasso where y > lambda/2



7 NEED TO COME BACK TOO

Applied

8

• A. set.seed(5)

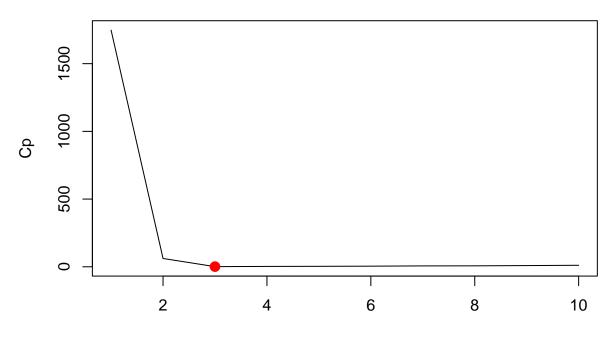
```
X <- rnorm(n=100)
noise <- rnorm(n=100)

• B.
beta_0 <- 1
beta_1 <- -2
beta_2 <- 3.75
beta_3 <- 5
Y <- beta_0 + beta_1*X + beta_2*(X^2) + beta_3*(X^3) + noise</pre>
```

• C. Using Best Subset Selection below, C_p , BIC and adjusted R^2 all select the model with three variables. The coefficient estiminates come out to $Y = 1.072 - 1.583X + 3.695X^2 + 4.890X^3$

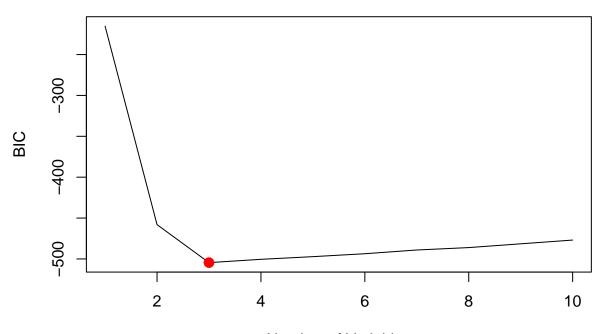
```
library(leaps)
df <- data.frame(X, Y)</pre>
# Best Subset Selection
# in poly(), setting raw = TRUE return raw polynomials, as opposed to
# orthogonal
regfit.full <- regsubsets(Y ~ poly(X, 10, raw = TRUE),</pre>
                           data = df,
                           nvmax = 10,
                           method = 'exhaustive')
reg.summary <- summary(regfit.full)</pre>
# Cp
plot(reg.summary$cp,
     xlab = 'Number of Variables',
     ylab = 'Cp',
     main = "BSS Cp Chooses the 3 Variable Model",
     type = '1')
points(which.min(reg.summary$cp),
       reg.summary$cp[which.min(reg.summary$cp)],
       col = 'red',
       cex = 2,
       pch = 20)
```

BSS Cp Chooses the 3 Variable Model



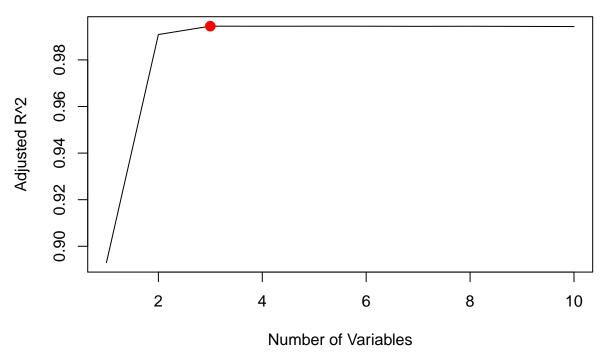
Number of Variables

BSS BIC Chooses the 3 Variable Model



Number of Variables

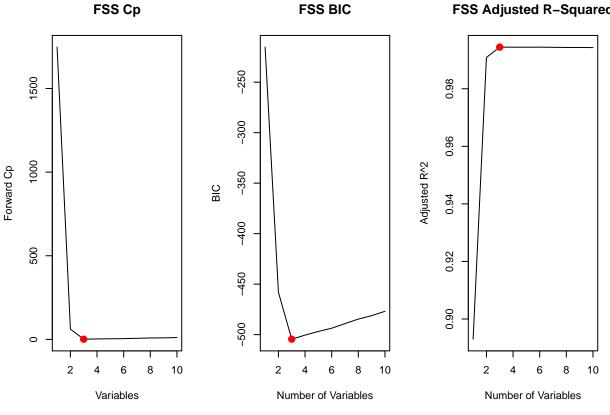
BSS Adjusted R-Squared Chooses the 3 Variable Model



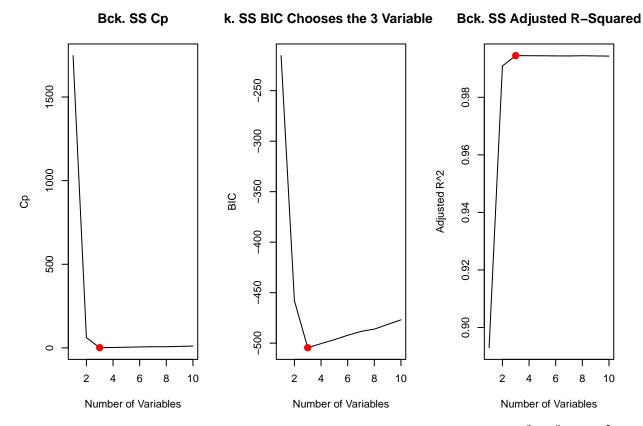
• D. C_p , BIC and adjusted \mathbb{R}^2 from both Forward and Backward Stepwise Selection all select the model with three variables as well.

```
# Forward and Stepwise Selection
regfit.forward <- regsubsets(Y ~ poly(X, 10, raw = TRUE),</pre>
                           data = df,
                           nvmax = 10,
                           method = 'forward')
reg.summary.forward <- summary(regfit.forward)</pre>
par(mfrow = c(1,3))
# Cp
plot(reg.summary.forward$cp,
     xlab = 'Variables',
     ylab = 'Forward Cp',
     main = "FSS Cp",
     type = '1')
points(which.min(reg.summary.forward$cp),
       reg.summary.forward$cp[which.min(reg.summary.forward$cp)],
       col = 'red',
       cex = 2,
       pch = 20)
```

```
# BIC
plot(reg.summary.forward$bic,
     xlab = 'Number of Variables',
     ylab = 'BIC',
     main = "FSS BIC",
     type = '1')
points(which.min(reg.summary.forward$bic),
       reg.summary.forward$bic[which.min(reg.summary.forward$bic)],
       col = 'red',
       cex = 2,
       pch = 20)
# Adj R62
plot(reg.summary.forward$adjr2,
     xlab = 'Number of Variables',
     ylab = 'Adjusted R^2',
     main = "FSS Adjusted R-Squared",
     type = '1')
points(which.max(reg.summary.forward$adjr2),
      reg.summary.forward$adjr2[which.max(reg.summary.forward$adjr2)],
      col = 'red',
      cex = 2,
      pch = 20)
             FSS Cp
                                           FSS BIC
                                                                 FSS Adjusted R-Squared
```

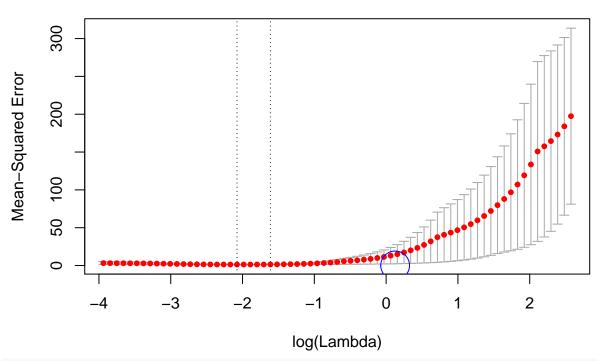


```
method = 'backward')
reg.summary.backward <- summary(regfit.backward)</pre>
par(mfrow = c(1,3))
# Cp
plot(reg.summary.backward$cp,
     xlab = 'Number of Variables',
     ylab = 'Cp',
     main = "Bck. SS Cp",
     type = '1')
points(which.min(reg.summary.backward$cp),
       reg.summary.backward$cp[which.min(reg.summary.backward$cp)],
      col = 'red',
       cex = 2,
       pch = 20)
# BIC
plot(reg.summary.backward$bic,
     xlab = 'Number of Variables',
     ylab = 'BIC',
     main = "Bck. SS BIC Chooses the 3 Variable Model",
     type = '1')
points(which.min(reg.summary.backward$bic),
       reg.summary.backward$bic[which.min(reg.summary.backward$bic)],
       col = 'red',
       cex = 2,
       pch = 20)
# Adj R62
plot(reg.summary.backward$adjr2,
     xlab = 'Number of Variables',
     ylab = 'Adjusted R^2',
     main = "Bck. SS Adjusted R-Squared",
     type = '1')
points(which.max(reg.summary.backward$adjr2),
      reg.summary.backward$adjr2[which.max(reg.summary.backward$adjr2)],
      col = 'red',
      cex = 2,
      pch = 20)
```



• E. Using the λ value selected using cross validation, Lasso Regression selects X^2 , X^3 and X^5 with coefficients of 3.5, 3.6 and 0.2 respectively (along with an intercept of 1.31). This is different from the true coefficients of $\beta_0 = 1$, $\beta_1(X) = -2$, $\beta_2(X^2) = 3.75$, $\beta_3(X^3) = 5$.

```
# Lasso Regression
set.seed(5)
library(glmnet)
## Loading required package: Matrix
## Loading required package: foreach
## Loaded glmnet 2.0-16
model.x <- model.matrix(Y ~ poly(X, 10, raw = TRUE), data = df)[, -1]</pre>
par(mfrow = c(1,1))
set.seed(1)
cv.out <- cv.glmnet(model.x, Y, alpha = 1)</pre>
plot(cv.out)
best_lam <- cv.out$lambda.min</pre>
points(best_lam,
      best_lam,
      col = 'blue',
      cex = 4)
```



best_lam

```
## [1] 0.1253245
```

```
lasso.mod <- glmnet(model.x, Y, alpha = 1, lambda = best_lam)
coef(lasso.mod)</pre>
```

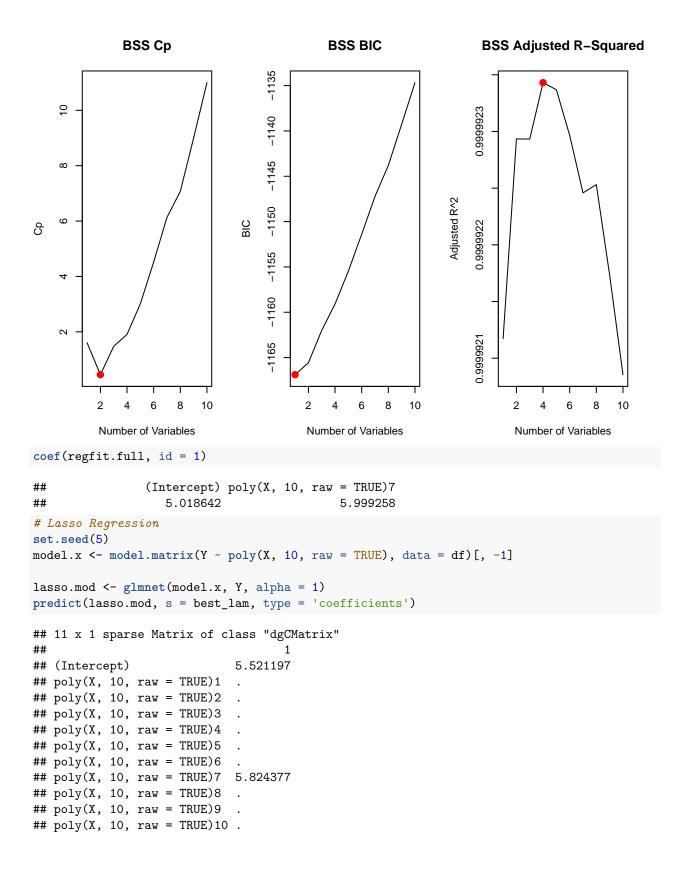
```
## 11 x 1 sparse Matrix of class "dgCMatrix"
## (Intercept)
                             1.3135583
## poly(X, 10, raw = TRUE)1
## poly(X, 10, raw = TRUE)2
                             3.5019559
## poly(X, 10, raw = TRUE)3
                             3.5942028
## poly(X, 10, raw = TRUE)4
## poly(X, 10, raw = TRUE)5
                             0.1930986
## poly(X, 10, raw = TRUE)6
## poly(X, 10, raw = TRUE)7
## poly(X, 10, raw = TRUE)8
## poly(X, 10, raw = TRUE)9
## poly(X, 10, raw = TRUE)10
```

• F. Setting a new response variable such that $Y = 5 + 6X^7 + \epsilon$, C_p , BIC and Adjusted R^2 suggest the 2, 1 and 4 variable model respectively. BIC did correctly identify X^7 , as well as correctly identify the coefficients as 5.02 and 6 for β_0 and β_1 respectively.

The Lasso also selected the correct coefficients, β_0 and β_7 , as well as chose coefficients that are relatively close to the true values, although not as close as those obtained using Best Subset Selection. This goes to show that when the response Y is a function of multiple predictors, the Lasso will produced a biased estimate of the model. However, when the response is a function of only a couple coefficients (in this case 1), the Lasso will approximate the true relationship relatively well.

```
set.seed(5)
X = rnorm(100)
```

```
noise <- rnorm(100)</pre>
beta_0 <- 5
beta_7 <- 6
Y \leftarrow beta_0 + beta_7*X^7 + noise
df <- data.frame(X, Y)</pre>
# Best Subset Selection
regfit.full <- regsubsets(Y ~ poly(X, 10, raw = TRUE),</pre>
                           data = df,
                           nvmax = 10,
                           method = 'exhaustive')
reg.summary <- summary(regfit.full)</pre>
par(mfrow = c(1,3))
# Cp
plot(reg.summary$cp,
     xlab = 'Number of Variables',
     ylab = 'Cp',
     main = "BSS Cp",
     type = '1')
points(which.min(reg.summary$cp),
       reg.summary$cp[which.min(reg.summary$cp)],
       col = 'red',
       cex = 2,
       pch = 20)
# BIC
plot(reg.summary$bic,
     xlab = 'Number of Variables',
     ylab = 'BIC',
     main = "BSS BIC",
     type = '1')
points(which.min(reg.summary$bic),
       reg.summary$bic[which.min(reg.summary$bic)],
       col = 'red',
       cex = 2,
       pch = 20)
# Adj R62
plot(reg.summary$adjr2,
     xlab = 'Number of Variables',
     ylab = 'Adjusted R^2',
     main = "BSS Adjusted R-Squared",
     type = '1')
points(which.max(reg.summary$adjr2),
      reg.summary$adjr2[which.max(reg.summary$adjr2)],
      col = 'red',
      cex = 2,
      pch = 20)
```



• A. Train/Test Split

```
library(ISLR)
library(Metrics)

##
## Attaching package: 'Metrics'

## The following object is masked from 'package:glmnet':
##
## auc

set.seed(5)
train_idx <- sample(c(TRUE, FALSE), nrow(College), rep = TRUE)
test_idx <- (!train_idx)</pre>
```

• B. Ordinary Least Squares returns a testing RMSE of 998.93.

[1] 1143.314

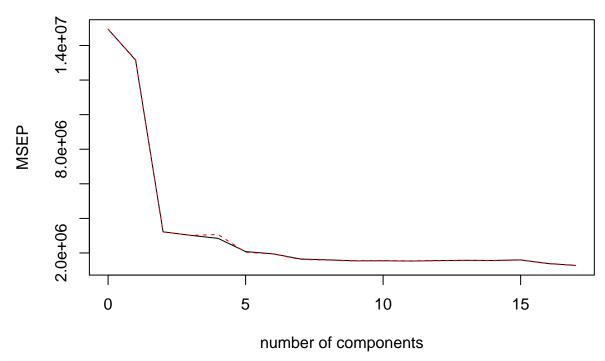
• C. Ridge Regression produces a slightly higher testing RMSE than OLS, coming in at 1003.

[1] 1002.997

• D. Similar to Ridge Regression, the Lasso has a slightly higher testing RMSE than the OLS, at 993.41. The coefficients are shown below.

```
best.lambda <- cv.out$lambda.min</pre>
lasso.mod <- glmnet(x, y, alpha = 0, lambda = best.lambda)</pre>
lasso.pred <- predict(lasso.mod,</pre>
                      s = best.lambda,
                      newx = x[test_idx,])
lasso.test.mse <- mse(y[test_idx], lasso.pred)</pre>
sqrt(lasso.test.mse)
## [1] 993.4065
predict(lasso.mod, s = best.lambda, type = 'coefficients')
## 19 x 1 sparse Matrix of class "dgCMatrix"
##
## (Intercept) -565.44436277
## (Intercept)
## PrivateYes -501.14562863
## Accept 1.52059509
## Enroll
               -0.64842120
## Top10perc 46.78469349
## Top25perc -12.20864720
## F.Undergrad 0.04354884
## P.Undergrad
                  0.04275676
## Outstate
                 -0.07954578
## Room.Board 0.15882475
## Books 0.02993927
## Personal 0.02665051
## PhD
                 -8.24595036
## Terminal -3.63832744
## S.F.Ratio 15.43376088
## perc.alumni -0.96459601
## Expend
                  0.07843444
## Grad.Rate
                  8.91838396
  • E. Using 5 principal components, PCR returns a testing RMSE of 1885.01, double the error of OLS.
library(pls)
##
## Attaching package: 'pls'
## The following object is masked from 'package:stats':
##
##
       loadings
set.seed(5)
pcr.fit <- pcr(Apps ~ .,</pre>
               data = College[train_idx,],
               scale = TRUE,
               validation = 'CV')
validationplot(pcr.fit, val.type = 'MSEP')
```

Apps

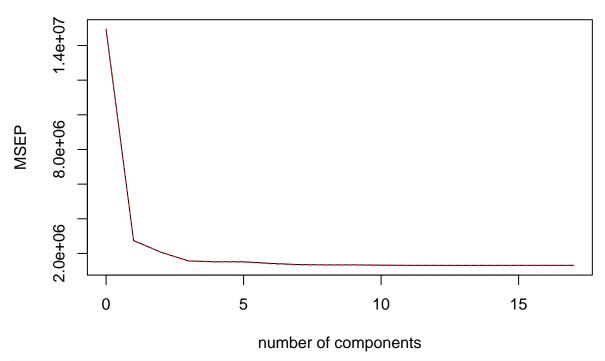


```
pcr.pred <- predict(pcr.fit, College[test_idx, ], ncomp = 5)
pcr.test.mse <- mse(College[test_idx, 'Apps'], pcr.pred)
sqrt(pcr.test.mse)</pre>
```

[1] 1885.047

• F. Partial Least Squares with 3 components returns a testing RMSE of 1665.24, also well above the RMSE of OLS.

Apps

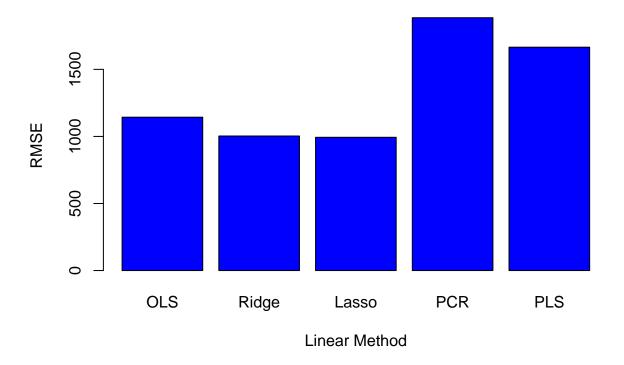


```
pls.pred <- predict(pls.fit, College[test_idx,], ncomp = 3)
pls.test.mse <- mse(College[test_idx, 'Apps'], pls.pred)
sqrt(pls.test.mse)</pre>
```

[1] 1665.244

• G. Plotting the testing RMSE's of the five methods tested, it is clear that OLS, Ridge and Lasso are far superior to Principal Component Regression and Partial Least Squares. Since OLS produced a testing RMSE of just under 1000, this means that, on average, the prediction for the number of applications will be off by 983 using OLS.

Principal Component Methods Fall Behind



10

• A. Creating a 20 attributes randomly generated from a Normal Distribution and a response variable Y such that $Y = 567 - X_1 - X_3 + 2X_5 + 3X_7 + 5X_9 + 8X_{11} + 13X_{13} + 21X_{15} + \epsilon$. That is to say that X_2 , X_4 , X_6 , X_8 , X_{10} , X_{12} , X_{14} , X_{16} , X_{17} , X_{18} , X_{19} and X_{20} have no affect on the response.

```
set.seed(5)
X <- matrix(rnorm(1000* 20), nrow = 1000, ncol = 20)
betas <- sample(seq(-50, 50, 0.25), 20)
noise <- rnorm(20)
for (i in 1:20) {
    if (i %% 2 == 0) {
        betas[i] <- 0
    }
}
y <- X %*% betas + noise
df <- data.frame(x = X, y = y)</pre>
```

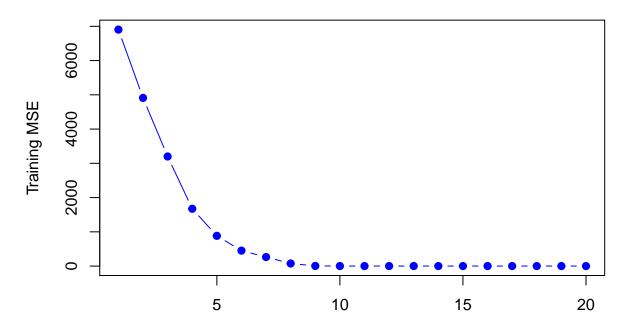
• B.

```
train_idx <- sample(1:nrow(df), 100)</pre>
```

• C. As one would expect, the training error approaches zero.

```
par(mfrow = c(1,1))
plot(best.subset.summary$rss / 100,
    pch = 19,
    col = 'blue',
    type = 'b',
    xlab = 'Number of Predictors',
    ylab = 'Training MSE',
    main = 'Training MSE Approaches 0 as No. Predictors Increases')
```

Training MSE Approaches 0 as No. Predictors Increases



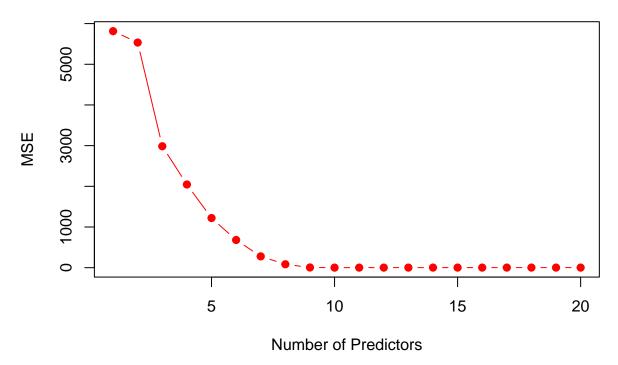
Number of Predictors

• D.

```
df.test.matrix <- model.matrix(y ~ ., data = df[-train_idx,])
test.errors <- NULL
for (i in 1:20) {
    coefs <- coef(best.subset, id = i)
        predictions <- df.test.matrix[, names(coefs)] %*% coefs
        test.errors[i] <- mse(y[-train_idx], predictions)
}

plot(test.errors,
    pch = 19,
    col = 'red',
    type = 'b',
    xlab = 'Number of Predictors',
    ylab = 'MSE',
    main = 'Testing MSE')</pre>
```

Testing MSE



• E. As one would expect, the testing MSE is minimized at a model with an intermediate amount of predictors. This resembles the classic "U" shape that the testing error will take when moving from a simple model to a more complex model (overfitting). (This is imperceptible to the eye in the above plot, however if I plot models 1 through 9 on one graph and 10 through 20 on another it becomes clear; these plots are shown below)

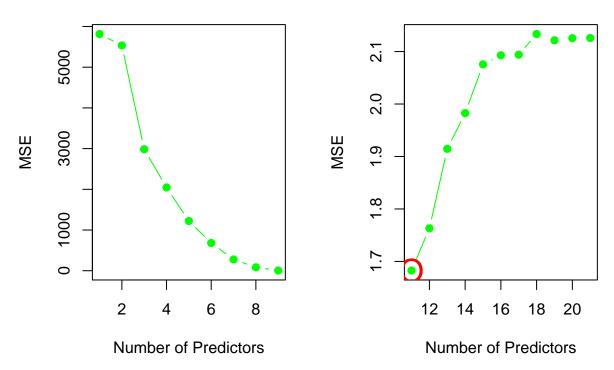
```
which.min(test.errors)
```

```
## [1] 10
```

```
par(mfrow = c(1, 2))
plot(test.errors[1:9],
     pch = 19,
     col = 'green',
     type = 'b',
     xlab = 'Number of Predictors',
     ylab = 'MSE',
     main = 'Testing MSE Approaching Minimum')
plot(test.errors[10:20],
     pch = 19,
     col = 'green',
     type = 'b',
     xlab = 'Number of Predictors',
     ylab = 'MSE',
     main = 'Testing MSE Increasing from Minimum',
     xaxt = 'n')
axis(1, at = seq(0, 10, 2), labels = c('10', '12', '14', '16', '18', '20'))
points(which.min(test.errors) - 9,
       test.errors[which.min(test.errors)],
       col = 'red',
       pch = '0',
```

cex = 2)

Testing MSE Approaching Minimu Testing MSE Increasing from Minim



• **F**. Illustrated in the table below, the coefficient estimates are very close to the true values (note that the coefficients excluded from the table were correctly set to 0 in the model).

```
estimates <- rep(0, 20)
estimates[1] <- coef(best.subset, id = 10)[2]
estimates[3] <- coef(best.subset, id = 10)[3]
estimates[5] <- coef(best.subset, id = 10)[4]
estimates[7] <- coef(best.subset, id = 10)[5]
estimates[9] <- coef(best.subset, id = 10)[6]</pre>
estimates[11] <- coef(best.subset, id = 10)[7]</pre>
estimates[13] <- coef(best.subset, id = 10)[8]
estimates[15] <- coef(best.subset, id = 10)[9]
estimates[17] <- coef(best.subset, id = 10)[10]
estimates[19] <- coef(best.subset, id = 10)[11]</pre>
coef_df <- data.frame(betas,estimates)</pre>
coef_df \leftarrow t(coef_df)[, c(1,3,5,7,9,11,13,15,17,19)]
rownames(coef_df) <- c("True_value", "Estimate")</pre>
colnames(coef_df) <- c('X1','X3','X5','X7','X9','X11','X13','X15','X17','X19')</pre>
knitr::kable(coef_df, digits = 2, caption = 'Coefficients')
```

Table 1: Coefficients

	X1	Х3	X5	X7	X9	X11	X13	X15	X17	X19
True_value	-8.25	-33.25	34.50	-22.50	-44.25	43.50	1.75	-19.25	12.50	-27.50
Estimate	-8.32	-33.27	34.52	-22.63	-44.17	43.57	1.86	-18.97	12.51	-27.54

• **G**. Plotting the error in the coefficient estimates as $\sqrt{\sum_{j=1}^{p}(\beta_j-\hat{\beta}_k^2)^2}$ on the y axis and the number of

predictors in the model as that number increases on the x axis, it is clear that this error is minimized at 10, the same value at which the testing MSE was minimized.

```
par(mfrow = c(1,1))
r \leftarrow seq(1, 20)
y \leftarrow rep(NA, 20)
col_names <- colnames(df)[1:20]</pre>
names(betas) <- col_names</pre>
for (i in 1:20) {
    cfs <- coef(best.subset, id = i)</pre>
    cfs <- cfs[2:length(cfs)]</pre>
    y[i] <- sqrt(sum((betas[names(cfs)] - cfs)^2))</pre>
plot(x = r,
     y = y,
     type = 'b',
     col = 'blue',
     pch = 19,
     xlab = 'Number of Predictors',
     ylab = 'Error',
     main = 'Error in Coefficient Estimates')
points(which.min(y),
       y[which.min(y)],
       col = 'red',
       pch = '0',
        cex = 2)
```

Error in Coefficient Estimates

