

# Module 8 Lab Submission

Nora Quick

For this exploration, we will simulate some data and then use the `regsubsets()` function to select models.

The data we simulate below will have some predictor variables (X1, X2, and X3) that contribute to the response, and some other predictor variables (W1, W2, and W3) that do not contribute to the response. The goal is to explore how well these methods of model selection do at identifying the correct variables.

```
n <- 30

set.seed(12345)
X1 <- rnorm(n)
X2 <- rnorm(n)
X3 <- 0.5*X1 + 0.5*X2 + rnorm(n, 0, 0.5)

W1 <- rnorm(n)
W2 <- rnorm(n)
W3 <- 0.4*X1 + 0.3*X2 + rnorm(n, 0, 0.4)

Y <- 1 + 0.3*X1 + 0.3*X2 + 0.5*X3 + rnorm(n)

lab8Data <- data.frame(Y, X1, X2, X3, W1, W2, W3)
```

Note that the true model is

$$Y_i = 1 + 0.3X_{1i} + 0.3X_{2i} + 0.5X_{3i} + \epsilon_i$$

so our hope is that we would select X1, X2, and X3, and not W1, W2, or W3.

1. Use the `regsubsets()` function to perform best subset selection, modeling Y as a function of all the other variables in the data set `lab8Data`. Name the resulting object `lab8.regfit.best`

```
lab8.regfit.best <- regsubsets(Y ~ ., data = lab8Data)
# lab8.regfit.best
```

2. Use the `summary()` function to summarize the `lab8.regfit.best` object that you got from the previous step, and store the summary object as `lab8.reg.best.summary`. What are the RSS values for the best subsets of each size (1 - 6)?

```
lab8.reg.best.summary <- summary(lab8.regfit.best)
# lab8.reg.best.summary

lab8.reg.best.summary$rss
```

```
## [1] 15.85931 15.01322 14.42506 14.30083 14.28906 14.28876
```

The RSS values for the best subsets of each size (1-6) are as follows; 15.86, 15.01, 14.43, 14.30, 14.29, and 14.29.

3. Use the `glm()` function to fit the best model with three predictors (in this case,  $X_1$ ,  $X_3$ , and  $W_2$ ) and then use the `cv.glm()` function to find the LOOCV error estimates for this model. Store the object resulting from `cv.glm()` as `lab8.cv.err`. What is the value of the first element of the `delta` component of this object?

```
glm.fit <- glm(Y ~ X1 + X3 + W2, data = lab8Data)
lab8.cv.err <- cv.glm(lab8Data, glm.fit)

lab8.cv.err$delta
```

```
## [1] 0.6057914 0.6035681
```

The first element of the 'delta' component of this object is 0.61.

4. Repeat the above step, but with the true predictors in the model instead ( $X_1$ ,  $X_2$ , and  $X_3$ ). How do the delta values compare for the 'best' 3-predictor model vs. the true model?

```
glm.fit2 <- glm(Y ~ X1 + X2 + X3, data = lab8Data)
lab8.cv.err2 <- cv.glm(lab8Data, glm.fit)

lab8.cv.err2$delta
```

```
## [1] 0.6057914 0.6035681
```

The delta values for the best 3-predictor model vs. the true model are the same. The delta values are 0.61 and 0.60.