# Module 4 Lab Submission

## Nora Quick

First, we will explore the Brain size data in the data set `case0902` from the `Sleuth3` library. You can read more about this data set by viewing the help file:

```
help(case0902)
head(case0902)
```

```
##                Species  Brain     Body Gestation Litter
## 1             Aardvark    9.6     2.20        31    5.0
## 2             Acouchis    9.9     0.78        98    1.2
## 3 African elephant 4480.0  2800.00       655    1.0
## 4              Agoutis   20.3     2.80       104    1.3
## 5            Axis deer  219.0    89.00       218    1.0
## 6               Badger   53.0     6.00        60    2.2
```

1. **Fit a linear model with `Brain` as the response variable, and `Body`, `Gestation`, and `Litter` as the predictor variables.**

```
head(case0902)
```

```
##                Species  Brain     Body Gestation Litter
## 1             Aardvark    9.6     2.20        31    5.0
## 2             Acouchis    9.9     0.78        98    1.2
## 3 African elephant 4480.0  2800.00       655    1.0
## 4              Agoutis   20.3     2.80       104    1.3
## 5            Axis deer  219.0    89.00       218    1.0
## 6               Badger   53.0     6.00        60    2.2
```

```
Brain_fit <- lm(Brain ~ Body + Gestation + Litter, data = case0902)
summary(Brain_fit)
```

```
##
## Call:
## lm(formula = Brain ~ Body + Gestation + Litter, data = case0902)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1026.68   -62.08    17.29    51.73   988.76
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -225.29213   83.05875  -2.712  0.00797 **
## Body           0.98588    0.09428  10.457  < 2e-16 ***
```

1

```
## Gestation      1.80874      0.35445    5.103 1.79e-06 ***
## Litter         27.64864     17.41429   1.588  0.11579
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 224.6 on 92 degrees of freedom
## Multiple R-squared:   0.81,  Adjusted R-squared:  0.8038
## F-statistic: 130.7 on 3 and 92 DF,  p-value: < 2.2e-16
```

2. **Calculate the case influence measures for this model using the `augment()` function from the package broom. Which species has the highest leverage for this model? Which species has the highest Cook's Distance?**

```
Brain_diag <- augment(Brain_fit)
head(Brain_diag)
```

```
## # A tibble: 6 x 10
##    Brain    Body Gestation Litter .fitted .resid   .hat .sigma .cooksd .std.resid
##    <dbl>   <dbl>     <int>  <dbl>   <dbl>  <dbl>  <dbl>  <dbl>   <dbl>      <dbl>
## 1    9.6     2.2        31      5   -28.8   38.4 0.0361   226. 2.85e-4      0.174
## 2    9.9    0.78        98    1.2   -14.1   24.0 0.0288   226. 8.70e-5      0.108
## 3 4480     2800        655      1  3748.   732.  0.719   173. 2.43e+1       6.16
## 4   20.3     2.8       104    1.3    1.52   18.8 0.0250   226. 4.60e-5     0.0847
## 5  219       89        218      1   284.   -65.4 0.0189   226. 4.16e-4     -0.294
## 6   53        6         60    2.2   -50.0  103.  0.0263   226. 1.46e-3      0.465
```

```
Brain_diag[Brain_diag$.hat > 0.4, ]
```

```
## # A tibble: 1 x 10
##   Brain  Body Gestation Litter .fitted .resid  .hat .sigma .cooksd .std.resid
##   <dbl> <dbl>     <int>  <dbl>   <dbl>  <dbl> <dbl>  <dbl>   <dbl>      <dbl>
## 1  4480  2800       655      1   3748.   732. 0.719   173.    24.3       6.16
```

Based on this R code we can see that litter 1 has the highest leverage in this model. In addition to that we can see from the .cooksd column that litter 1 also has the highest cook's distance.

Now we will continue investigating multicollinearity. Recall the simulated scenario considered in the `M4Lab-examples.Rmd` file, where we followed these steps:

1. Define $\beta_0 = 0.5$, $\beta_1 = 0.3$, and $\beta_2 = 0.7$
2. Define the mean of $X_1$ and $X_2$
3. Generate correlated/uncorrelated $X_1$ and $X_2$ data
4. Generate the response variable; use model equation and add N(0,1) noise
5. Fit a MLR model
6. Extract the coefficient estimate; $\hat{\beta}_0$, $\hat{\beta}_1$, or $\hat{\beta}_2$
7. Repeat steps (4) through (6) many times.

We used a function, included here, to perform steps 4. through 6., and then repeated that function many times (step 7.)

```
fitmodel <- function(X1, X2, beta0, beta1, beta2){
  n <- length(X1)
  Y <- beta0 + beta1*X1 + beta2*X2 + rnorm(n, 0, 1) # Generate/calculate response
  fit <- lm(Y ~ X1 + X2) # Fit the model
  fit$coefficients # Return estimated coefficient values
}
```

To run this function, we have to define the coefficient values (Step 1.), and set the mean and covariance matrix to generate predictor variables (Steps 2. and 3.).

```
# Step 1
beta0 <- 0.5 # define beta_0
beta1 <- 0.3 # define beta_1,
beta2 <- 0.7 # define beta_2

# Step 2
mu <- matrix(c(0,0)) # Set means for X_1, X_2
sigma1 <- matrix(c(1, 0, 0, 1), ncol = 2) # Cov Matrix: Cov(X_1, X_2) = 0

# Step 3
set.seed(1822) # Francis Galton born, invented regression concept

n <- 250
X <- mvrnorm(250, mu=c(0,0), Sigma=sigma1)
X1 <- X[,1]
X2 <- X[,2]

# Step 7
beta_estimates <- replicate(10000, fitmodel(X1, X2, beta0, beta1, beta2))
```

Finally, we calculated the standard deviation of the estimates of $\beta_0$ that resulted from these simulated datasets:

```
sd(beta_estimates[1,])
```

```
## [1] 0.0634904
```

3. **Now it is your turn to calculate the standard deviation of the estimates of $\beta_1$ and $\beta_2$ in the uncorrelated case; and $\beta_0$, $\beta_1$, and $\beta_2$ in the correlated case. As you run the simulations, fill in the standard errors in the table below. Note: In the correlated case, use `sigma2 <- matrix(c(1, 0.9, 0.9, 1), ncol = 2)` to define the covariance matrix.**

| Parameter | $SE(\hat{\beta}_i)$ |
|---|---|
| *Uncorrelated* | |
| $\beta_0$ | 0.063 |
| $\beta_2$ | 0.061 |
| $\beta_3$ | 0.062 |
| *Correlated* | |
| $\beta_0$ | 0.063 |
| $\beta_2$ | 0.068 |

| Parameter | $SE(\hat{\beta}_i)$ |
|:---:|:---:|
| $\beta_3$ | 0.065 |

```r
# Step 1
beta0 <- 0.5 # define beta_0
beta1 <- 0.3 # define beta_1,
beta2 <- 0.7 # define beta_2

# Step 2
mu <- matrix(c(0,0)) # Set means for X_1, X_2
sigma1 <- matrix(c(1, 0, 0, 1), ncol = 2) # Cov Matrix: Cov(X_1, X_2) = 0

# Step 3
# set.seed(1822) # Francis Galton born, invented regression concept

n <- 250
X <- mvrnorm(250, mu=c(0,0), Sigma=sigma1)
X1 <- X[,1]
X2 <- X[,2]

# Step 7
beta_estimates <- replicate(10000, fitmodel(X1, X2, beta0, beta1, beta2))

# SD
sd(beta_estimates[1,])
```

```
## [1] 0.06253699
```

4. **The variances (and therefore standard deviations) of $\hat{\beta}_1$ and $\hat{\beta}_2$ are much larger when $X_1$ and $X_2$ are correlated than when they are uncorrelated. Does it make sense that $\hat{\beta}_0$ is unaffected? Explain your reasoning.**

Yes, because we have to have one base situation that no matter if they're correlated or uncorrelated the data follows the same pattern.

5. **Recall the sample VIFs calculated (in M4Lab-examples.Rmd) for some simulated data in the correlated case:**

```
      X1        X2
5.304359 5.304359
```

**Compare the variances (*squared standard deviations*) in the table above for the correlated predictor setting to the variances for the uncorrelated predictor setting: what is the ratio of the variance of $\hat{\beta}_1$ in the correlated predictor setting to the variance of $\hat{\beta}_1$ in the uncorrelated predictor setting? Similarly, what is the variance of $\hat{\beta}_2$ in the correlated predictor setting to the variance of $\hat{\beta}_2$ in the uncorrelated predictor setting? Do these ratios seem close to the VIFs that we calculated?**

The correlated beta_1 is about 0.007 larger than the uncorrelated and the correlaed bata_2 is 0.003 larger than the uncorrelated one. It does not appear that these ratios seem close to the VIFs that we calculated in the lab documentation.