1.) Chp 7 Exercise 1 (a) (b) (c):

(a) Identify the optimal model(s) based on $R^2_{adj}$ AIC & BIC from the approach based on all possible subsets.

| Criteria | Subset Size | Model |
|---|---|---|
| $R^2_{adj}$ | 2 or 3 | $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ or $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$ |
| AIC | 2 | $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ |
| BIC | 2 | $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ |

(b) Identify the optimal model(s) based on AIC & BIC from the forward selection approach

Optimal Model based on AIC & BIC using the forward selection approach is

$$Y = X_3.$$

(c) Carefully explain why different models are chosen in (a) & (b).

- In (b) we used forward selection to determine the best model.
  The forward selection algorithm is "greedy" in that it doesn't look ahead and only adds the most significant variable at a given step.

2) Four Treatments: A, B, C, D.   $n = 200$,  $n_i = 50$ $\forall i$

a model,  $y = X\beta + e$   was fit.

$$X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

$\hat\beta' = \begin{bmatrix} 37.5 & -11.5 & 1.0 & -27.7 \end{bmatrix}$

$SE(\hat\beta) = (2.75, 3.89, 3.89, 3.89)$ & $\hat\sigma_e = 19.45$

$$(X'X)^{-1} = \begin{bmatrix} 0.02 & -0.02 & -0.02 & -0.02 \\ -0.02 & 0.04 & 0.02 & 0.02 \\ -0.02 & 0.02 & 0.04 & 0.02 \\ -0.02 & 0.02 & 0.02 & 0.04 \end{bmatrix}$$

H.D.5 (6)  → (a) Interpret each of the four regression parameters.

$\beta_0$ = The mean response for treatment A   ($\mu_A$)

$\beta_1$ = The difference between the mean response for treatment B and Treatment A  ($\mu_B - \mu_A$)

$\beta_2$ = "     "     C and treatment A  ($\mu_C - \mu_A$)

$\beta_3$ = "     "     D and treatment A  ($\mu_D - \mu_A$)

H.D.5 (26) → (b) What is an approximate 95% CI for the mean difference in response between treatment groups B & A ($\mu_B - \mu_A$)?

[NOTE: This is just a CI for $\beta_1$.]

$$\boxed{\begin{array}{l} 95\% \ CI: \quad \hat\beta_1 \pm t^*_{200-3-1, 0.975} \ SE(\hat\beta) \\[2mm] \qquad -11.5 \pm 1.972141 \ (2.75) = (-19.17, -3.83) \end{array}}$$

(c) What is an appropriate 95% CI for the mean response in treatment group B?

$$\hat\beta_0 + \hat\beta_1 = \hat\mu_A + \hat\mu_B - \hat\mu_A = \hat\mu_B = 26.00 = \bar{y}_B$$

$$95\% \ CI: \quad \bar{y}_B \pm t_{196, 0.975} \ \frac{\hat\sigma_e}{\sqrt{n_i}} \quad \Rightarrow \quad 26 \pm 1.972141 \left( \frac{19.45}{\sqrt{50}} \right)$$

$$\boxed{(20.57534, 31.42466)}$$

3) Chp 6 Exercise 3:

n = 234. vars: y = Suggested Retail Price, $x_1$ = Engine Size, $x_2$ = Cylinders, $x_3$ = HP

$x_4$ = Highway MPG, $x_5$ = Weight, $x_6$ = Wheel Base.

$x_7$ = Hybrid (A dummy var = 1 for Hybrid cars)

- model: $Y = B_0 + B_1 x_1 + B_2 x_2 + B_3 x_3 + B_4 x_4 + B_5 x_5 + B_6 x_6 + B_7 x_7 + e$    (6.36)

(a) Decide if the proposed model is valid. Give reasons for your answers.

- The proposed model doesn't seem to be valid. We see in the plot of the sqrt(std. residuals) vs fitted values that there is clearly a positive trend.
- Additionally we can see in the QQ plot of the residuals that our residuals aren't normally distributed, they seem to be skewed right.
- We also have a few points that could be classified as bad leverage points.

(b) The plot of residuals vs fitted values produces a curved pattern. Describe what, if anything can be learned about the model from this plot.

A curved pattern suggests that we may need to adjust the variables by transforming them. It also suggests we might need to add a quadratic variable to our model.

(c) Identify any bad leverage pts.

- Looking at the plot of Cook's-D it seems like pt 223 is a bad leverage pt.

- The multivariate box-cox was used to transform the predictors while a log transform was used for the response.

$\log(y) = B_0 + B_1 x_1^{0.33} + B_2 \log(x_2) + B_3 \log(x_3) + B_4(\frac{1}{x_4}) + B_5 x_5 + B_6 \log(x_6) + B_7 x_7 + e$    (6.37)

(d) Decide if the transformed model is valid.

- The proposed model still doesn't seem valid. There still seems to be a trend in the plot of residuals vs fitted values and we can also see in the QQ plot of the residuals that our residuals aren't normally distributed; we now have more pts that seem like bad leverage pts.

3.) (contd)

(e) To obtain a final model, the analyst wants to simply remove the two insignificant predictors $\left(\frac{1}{x_4}\right)$ and $\log(x_6)$. Perform a partial F-test to see if this is a sensible strategy.

$$F = \frac{RSS_{reduced} - RSS_{Full} / k}{RSS_{Full} / n-p-1} = \frac{7.232671 - 6.717118 \mid 2}{6.717118 / 226} = 8.662901$$

- $P(F_{2,226} > 8.046067) = 0.0008371326$
- We have significant evidence that the variables in our model were significant, so the analysts strategy does not seem reasonable.

(5) Describe how model could be expanded in order to evaluate the effect of manufacturer on suggested retail price.

- We could create dummy variables for each of the manufactures we would like to include in the model and then we could do a number of things. We could start with a simple EDA by making added variable plots to see if the manufacturers might explain some of the residuals in our model. If they seem to be important, we could add them to the model and then conduct a partial F-test to see if the added variables are actually significant. We would need to watch for multicolinearity to make sure the added variables are not affecting the variance of of coefficient estimates.

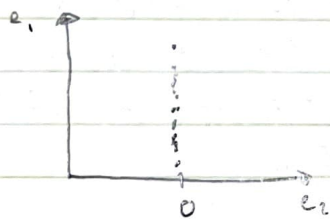4.) We are interested in the linear Model $Y = \beta_0 + \beta_1 x_1 + \beta_2 V_2 + e$

    (a) Fit model to data for which $x_1 = 2.2 x_2$ (w/ no error).
    Describe the appearance of the AVP for $x_2$ after $x_1$ has been
    added to the model. Explain why. Assume $Y$ has a correlation w/
    the predictors that is either 0 or 1.

(H 0.6 pg 37) → • AVP: On y-axis → plot residuals from model $Y = \beta_0 + \beta_1 x_1 + e_2$
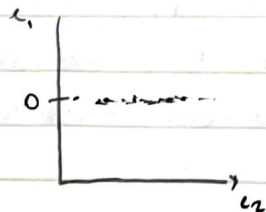                    On x-axis → plot residuals from model $x_2 = 2.2 x_1 + e_2$

    • Our AVP would look like the following. B/c the residuals
       from the model $x_2 = 2.2 x_1 + e_2$ are all 0, we would
       just stack the residuals from the first model



    (b) Again referring to the above models, assume $y = 3x_1$ w/o any error.
    Describe the appearance of the added variable plot for $x_2$ after
    $x_1$ had been added to the model. Explain assume the true that the
    correlation btwn $x_1$ & $x_2$ is between 0 & 1.

    • Our AVP would look like a straight horozontal line at $y = 0$. This is b/c
    $Y$ is perfectly explained by $x_1$, so all the residuals from $y = 3x_1$ are 0.

5.)  $\beta$: $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + e_i$. Also, the columns of the design matrix have mean 0 and length 1. That is $x_1'x_1 = 1$ ; $x_2'x_2 = 1$. Then if $r$ is the correlation between $x_1$ & $x_2$ we have the following:

$$(X'X) = \begin{bmatrix} n & 0 & 0 \\ 0 & 1 & r \\ 0 & r & 1 \end{bmatrix} \quad ; \quad (X'X)^{-1} = \begin{bmatrix} 1/n & 0 & 0 \\ 0 & 1/(1-r^2) & -r/(1-r^2) \\ 0 & -r/(1-r^2) & 1/(1-r^2) \end{bmatrix}$$

H.O.6 (40) →
Textbook (203)

(a) Explain why $S_{xx} = 1$. Use that to show that the VIF formula on pg. 203 matches $\sigma^2 (X'X)^{-1} = \text{var}(\hat\beta)$

$$\left[ \text{NOTE: } S_{XX} = \sum (x_i - \bar{x})^2 \right]$$

• Since we're given $\bar{x}_1$ & $\bar{x}_2$   $SSX_1 = \sum x_{1i}^2 = x_1'x_1 = 1$

$$SSX_2 = \sum x_{2i}^2 = x_2'x_2 = 1$$

$$\text{Var}(\hat\beta_j) = \frac{1}{1-r_{12}^2} * \frac{\sigma^2}{(n-1)S^2_{x_j}} = \frac{1}{1-r_{12}^2} * \frac{\sigma^2}{(n-1)\frac{SSX_j}{(n-1)}} = \frac{1}{1-r^2} * \frac{\sigma^2}{1} = \frac{\sigma^2}{1-r^2}$$

which is the formula for VIF = $\frac{\sigma^2}{1-r^2}$

(b) Determine which values of $r$ will make the variances of $\hat\beta_1$ & $\hat\beta_2$ large. Explain why using what you know about the variance of the vector $\hat\beta$.

$$\text{var}(\hat\beta|X) = \sigma^2(X'X)^{-1} = \begin{bmatrix} \sigma^2/n & 0 & 0 \\ 0 & \sigma^2/(1-r^2) & -\sigma^2/(1-r^2) \\ 0 & -\sigma^2/(1-r^2) & \sigma^2/(1-r^2) \end{bmatrix}$$

$\Rightarrow \text{var}(\hat\beta_1|X) = \text{var}(\hat\beta_2|X) = \frac{\sigma^2}{(1-r^2)}$

As $|r| \to 1$, $\text{var}(\hat\beta_1|X) = \text{var}(\hat\beta_2|X) \to \infty$.

• This is why multicollinearity is a problem. As the correlation of our predictors increases, the variances of the coefficient estimates become highly unstable.

6.) In a study on weight gain in rabbits, researchers randomly assigned rabbits to 1,2 or 3 mg of Dickory supplements A or B. (one rabbit to each level of each supplement). Consider the linear model $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$, where $X_1$ is dose level & $X_2$ is an indicator variable indicating the type of supplement used.

(a) The VIF for both $\hat{\beta_1}$ & $\hat{\beta_2}$ are 0. This is b/c they are orthogonal vectors (i.e. $X_1$ & $X_2$ are independent) Thus the $cor(X_1, X_2) = r = 0 \Rightarrow VIF = \frac{1}{1-r^2} = 1$. VIF is the same for all values of $y$ b/c $y$ doesn't show up in the calculation of VIF.

(b) $VIF(V_1) = 1.375$. It is larger than the above because now the $cor(X_1, X_2) = 0.522233$. This is b/c the way it is coded

$$X_2 = \begin{cases} 0 & \text{if } A \\ 1 & \text{if } B \end{cases}$$

and the dosage values are higher for B. The correlation btwn $(X_1, X_2)$ is positive. If we had $X_2$ coded the opposite way, the VIF would be the same. b/c the corr would just flip its sign.

(c) Consider a ln model, $y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + e$. Under what circumstances would the VIF for all $p$ features be 1?

• The VIF for all the variables would be 1 if all the variables were orthogonal to eachother (i.e all the variables were independent).