

STAT 608, Spring 2022 - Assignment 5
SOLUTIONS

1. Chapter 7, Question 1 parts (a), (b), and (c).

- (a) USING TABLE 7.4, THE BEST MODEL BASED ON ALL THREE CRITERIA IS THE SECOND MODEL, WITH PREDICTORS X_1 AND X_2 . WE NOTICE THAT R^2_{ADJ} IS HIGHEST FOR THIS MODEL (IT IS THE SAME VALUE FOR THE THIRD MODEL, BUT BY DEFAULT WE'LL CHOOSE THE SIMPLER MODEL), AND AIC AND BIC ARE SMALLEST FOR THE SECOND MODEL.
- (b) BOTH AIC AND BIC TAKE A STEP FORWARD TO ADD X_3 , AND THEN STOP.
- (c) IT APPEARS THAT X_3 IS A PRETTY GOOD VARIABLE FOR EXPLAINING Y AND DOES A BETTER JOB BY ITSELF THAN EITHER X_1 OR X_2 ALONE, WHILE THE LINEAR COMBINATION OF X_1 AND X_2 TOGETHER EXPLAIN Y PERFECTLY. BUT SINCE X_3 WORKS PRETTY WELL, NEITHER X_1 NOR X_2 ALONE HELPS X_3 ENOUGH TO BE INCLUDED IN THE MODEL IN THE NEXT STEP FORWARD.

2. (Old Qualifying Exam Question) A randomized trial was conducted to investigate the relationship between a continuous response y and four treatments A, B, C, and D. The sample size was $n = 200$, with 50 observations in each of the four treatment groups. Let y be the 200×1 vector of response values, ordered so that the first 50 entries are for treatment group A, the next 50 for B, then C, and finally D. The regression model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ was fit, where \mathbf{X} is the 200×4 design matrix given by

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

and where each entry is a column vector of length 50. The estimated regression coefficients were $\hat{\boldsymbol{\beta}}' = [37.5, -11.5, 1.0, -27.7]$, with standard errors 2.75, 3.89, 3.89, 3.89, and residual standard deviation $\hat{\sigma} = 19.45$. Also:

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

- (a) Interpret each of the four regression parameters.
 β_0 IS THE MEAN RESPONSE FOR GROUP A. (NOTICE THAT THIS IS THE PARAMETER, NOT THE STATISTIC, SO I DON'T HAVE TO SAY IT'S THE "APPROXIMATE" MEAN OF GROUP A.) β_1 IS THE MEAN DIFFERENCE IN RESPONSE BETWEEN GROUPS B AND A. β_2 IS THE MEAN DIFFERENCE IN RESPONSE BETWEEN GROUPS C AND A. β_3 IS THE MEAN DIFFERENCE IN RESPONSE BETWEEN GROUPS D AND A.
- (b) What is an approximate 95% confidence interval for the mean difference in response between treatment groups B and A (so, the difference $\mu_B - \mu_A$)?

FIRST, NOTE THAT $\beta_1 = \mu_B - \mu_A$. THE CONFIDENCE INTERVAL FOR β_1 IS

$$\begin{aligned}\hat{\beta}_1 \pm t_{n-p-1}^* \text{SE}(\hat{\beta}_1) &= -11.5 \pm t_{196,0.025} \sqrt{19.45^2 \times 0.04} \\ &= -11.5 \pm 1.97 \times 3.89 \\ &= (-19.1633, -3.8367)\end{aligned}$$

- (c) What is an approximate 95% confidence interval for the mean response in treatment group B?

THE PARAMETER OF INTEREST HERE IS $\beta_0 + \beta_1$. NOTE THAT WE CAN WRITE $\mathbf{A}\hat{\boldsymbol{\beta}} = \hat{\beta}_0 + \hat{\beta}_1$ WITH $\mathbf{A} = [1 \ 1 \ 0 \ 0]$. WE KNOW THAT FOR A MATRIX OF CONSTANTS \mathbf{A} , $\text{VAR}(\mathbf{A}\hat{\boldsymbol{\beta}}) = \mathbf{A}\text{VAR}(\hat{\boldsymbol{\beta}})\mathbf{A}' = \sigma^2 \mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}' = 0.02\hat{\sigma}^2$.

$$\begin{aligned}\hat{\beta}_0 + \hat{\beta}_1 \pm t_{n-p-1}^* \text{SE}(\hat{\beta}_0 + \hat{\beta}_1) &= 37.5 - 11.5 \pm t_{196,0.025} \sqrt{19.45^2 \times 0.02} \\ &= 26 \pm 1.97 \times 2.75 \\ &= (20.58, 31.42)\end{aligned}$$

3. Chapter 6, Question 3

- (a) NO, THIS IS NOT A VALID MODEL, FOR THE FOLLOWING REASONS:

- THE PLOT OF THE RESIDUALS VERSUS THE FITTED VALUES DOES NOT HAVE A NICE RANDOM SCATTER.
 - THE SMOOTH LINE THROUGH THE SCALE-LOCATION PLOT HAS INCREASING SLOPE.
 - THERE ARE SEVERAL HIGH LEVERAGE POINTS WHICH ARE ALSO HIGH RESIDUAL POINTS, MEANING THEY STRONGLY INFLUENCE THE MODEL.
 - THE SCATTERPLOT MATRIX SUGGESTS STRONG CORRELATIONS AMONG PREDICTORS, SO I ANTICIPATE HIGH VIF VALUES.
- (b) BECAUSE THE RELATIONSHIPS BETWEEN THE PREDICTORS ARE NOT ALL LINEAR, WE AREN'T ABLE TO SAY WHAT IS WRONG WITH THE MODEL. WE NEED TO USE TRANSFORMATIONS BEFORE SEARCHING FOR MODEL IMPROVEMENTS.
- (c) POINT 223 IS IDENTIFIED BY THE CUTOFF IN R AS HAVING A COOK'S DISTANCE TOO LARGE; I AM WILLING TO GUESS THAT POINTS 222 AND 67 WOULD ALSO HAVE COOK'S DISTANCE VALUES THAT OUR CUTOFF WOULD FLAG. BUT IT APPEARS THAT POINT 229 SHOULD ALSO BE INVESTIGATED. IT HAS A VERY LARGE RESIDUAL AND LARGE FITTED VALUE, SO IT STANDS APART FROM MANY OTHER POINTS, LIKE POINTS 222 AND 223, SO ALL THESE POINTS PROBABLY NEED TO BE PULLED IN BY TRANSFORMATIONS.
- (d) THIS MODEL IS CERTAINLY AN IMPROVEMENT, IN THAT THE SCALE-LOCATION PLOT LINE LOOKS MUCH MORE HORIZONTAL, AND THE RESIDUALS VS. FITTED VALUES PLOT LOOKS MUCH MORE RANDOM. HOWEVER, THE MODEL IS STILL NOT VALID, BECAUSE
- THE VIF VALUES STILL EXCEED 5, SO THE PREDICTORS ARE EXPLAINING MUCH OF THE SAME INFORMATION.

- THERE ARE STILL SEVERAL INFLUENTIAL POINTS: 66, 67, AND 88 ARE BAD OUTLIERS FOR THIS MODEL.
 - I'D PREFER THE RESIDUALS VS. FITTED VALUES PLOT TO HAVE A SMOOTH LINE WITH A MORE HORIZONTAL SLOPE.
- (e) OUR HYPOTHESES ARE $H_0 : \beta_4 = \beta_6 = 0$ VS. H_a : AT LEAST ONE OF THE TWO PARAMETERS IS NOT ZERO. IF THE FULL MODEL WERE VALID, THE FOLLOWING P-VALUE WOULD BE VALID:

$$F = \frac{(0.1781^2 \times 228 - 0.1724^2 \times 226) / 2}{0.1724^2} \quad (\text{P-VALUE} = 0.0002)$$

SO, WE WOULD HAVE VERY STRONG EVIDENCE THAT THE TWO MODELS WERE SIGNIFICANTLY DIFFERENT, EVEN THOUGH EACH INDIVIDUAL PREDICTORS IS NOT SIGNIFICANT.

- (f) WE COULD SIMPLY ADD A SET OF INDICATOR (DUMMY) VARIABLES, WITH $k - 1$ INDICATORS FOR THE k MANUFACTURERS. IF, HOWEVER, SOME MANUFACTURERS MAKE ONLY MODELS WITH ONE OR TWO DIFFERENT NUMBERS OF CYLINDERS, WE MAY FIND THAT THE INFORMATION CARRIED IN NUMBER OF CYLINDERS IS TOO SIMILAR TO THE INFORMATION CARRIED IN THE SET OF MANUFACTURER INDICATOR VARIABLES, AND NEED TO DROP ONE OF THOSE VARIABLES.

4. (From Weisberg, 2005) We are interested in the linear model $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$.

- (a) Suppose we fit the model above to data for which $x_1 = 2.2x_2$ with no error (that is, all residuals = 0). For example, X_1 could be a weight in pounds, and X_2 the weight of the same object in kg. Describe the appearance of the added-variable plot for x_2 after x_1 had been added to the above model. Explain why. Assume that Y has a correlation with the predictors that is neither 0 nor 1. Hint: Think about what goes on the x -axis and the y -axis of the added variable plot. You should notice something interesting about one of those residual vectors.

IN THIS CASE, X_2 IS PERFECTLY EXPLAINED BY X_1 . THIS MEANS THAT THE RESIDUALS OF X_2 REGRESSED ON X_1 ARE ALL EXACTLY 0. SO, THE PLOTTED POINTS IN THE ADDED VARIABLE PLOT WILL ALL HAVE HORIZONTAL AXIS VALUE 0. THIS MEANS THE PLOT WILL BE A STRAIGHT VERTICAL LINE OF POINTS SCATTERED IN THE VERTICAL DIRECTION.

- (b) Again referring to the model above, this time suppose that Y and X_1 are perfectly correlated, so $Y = 3X_1$, without any error. Describe the appearance of the added-variable plot for x_2 after x_1 had been added to the model. Explain. Assume this time that the correlation between the predictors is between 0 and 1.

IN THIS CASE, Y IS PERFECTLY EXPLAINED BY X_1 . THIS MEANS THAT THE RESIDUALS OF Y REGRESSED ON X_1 ARE ALL EXACTLY 0. SO, THE PLOTTED POINTS IN THE ADDED VARIABLE PLOT WILL ALL HAVE VERTICAL AXIS VALUE 0. THIS MEANS THAT PLOT WILL BE A STRAIGHT HORIZONTAL LINE OF POINTS SCATTERED IN THE HORIZONTAL DIRECTION.

5. Suppose we are interested in the linear model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + e_i$. Also suppose the columns \mathbf{x}_1 and \mathbf{x}_2 of the design matrix for this model have mean 0 and *length* 1. (That is, $\mathbf{x}_1' \mathbf{x}_1 = 1$ and $\mathbf{x}_2' \mathbf{x}_2 = 1$. This is a very particular situation that is unlikely to happen in practice; it just makes our arithmetic easier for a moment.). Then if r is the correlation between \mathbf{x}_1 and \mathbf{x}_2 , we have the following:

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

- (a) In our setup where the predictors have mean 0 and length 1, explain why $\text{SXX} = 1$. Use that to show that the VIF formula on page 203 matches $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ (above).

WE HAVE THAT

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

SO $\text{VAR}(\hat{\beta}_1|\mathbf{X}) = \text{VAR}(\hat{\beta}_2|\mathbf{X}) = \sigma^2/(1-\rho^2)$. IF $\rho \rightarrow 0$, THEN $\text{VAR}(\hat{\beta}_1|\mathbf{X}) = \text{VAR}(\hat{\beta}_2|\mathbf{X}) \rightarrow \sigma^2$. IF $|\rho| \rightarrow 1$, THEN $\text{VAR}(\hat{\beta}_1|\mathbf{X}) = \text{VAR}(\hat{\beta}_2|\mathbf{X}) \rightarrow \infty$.

AS $|\rho| \rightarrow 1$, THE VARIABLES x_1 AND x_2 BECOME MORE CORRELATED, WHICH INDICATES BOTH VARIABLES CONTAIN MOSTLY THE SAME INFORMATION. THIS IS A MULTICOLLINEARITY PROBLEM WHICH WILL LEAD TO ESTIMATION OF β BEING LESS ACCURATE.

BECAUSE THE MEAN OF X IS ZERO,

$$\text{SXX} = \sum (x_i - \bar{x})^2 = \sum x_i^2 = 1$$

AND

$$\text{VAR}(\hat{\beta}_j) = \frac{1}{1-\rho^2} \times \frac{\sigma^2}{(n-1)S_{x_j}^2} = \frac{1}{1-\rho^2} \times \frac{\sigma^2}{(n-1)\frac{\text{SXX}}{n-1}} = \frac{\sigma^2}{1-\rho^2}$$

WHICH MATCHES THE VIF SHOW ON PAGE 203.

- (b) Determine what values of r will make the variance of $\hat{\beta}_1$ and $\hat{\beta}_2$ large. Explain why, using what you know about the variance of the vector $\hat{\beta}$.

THE VARIANCE INFLATION FACTORS WOULD ALL EQUAL 1 IF THE PREDICTORS ARE ALL ORTHOGONAL TO EACH OTHER. IN THIS CASE, $R^2 = 0$ WHICH WILL RESULT IN $\text{VIF} = 1$. IF r IS CLOSE TO -1 OR 1, THE VIF BLOWS UP, MAKING THE VARIANCE OF THE $\hat{\beta}$ 'S LARGE AND UNSTABLE.

6. In a study on weight gain in rabbits, researchers randomly assigned rabbits to 1, 2, or 3 mg of one of dietary supplements A or B (one rabbit to each level of each supplement, which is not enough, of course). Consider the linear model $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$, where x_1 is the dosage level of the supplement, and x_2 is a dummy variable indicating the type of supplement used.

- (a) Compute the variance inflation factor for variable x_1 . You should be able to do this completely without the use of statistical software. Explain, using the word “orthogonal,” why the variance inflation factor is the value computed.

(Hint: To get started, you might go ahead and use R and the `vif()` function. You’ll have to invent a response vector y ; try

```
y <- c(1, 2, 3, 4, 5, 6)
```

to get you started. Notice that the VIF is the same no matter what values you use for y . Why? Then you might look at the formulas for VIF and notice that correlation is part of that formula. Calculate correlations between vectors to see what happens. Then you’ll see what is orthogonal to what.)

THE DESIGN MATRIX IS

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 0 \\ 1 & 3 & 1 \end{bmatrix}$$

THE VECTOR $\mathbf{x}_1 - \bar{x}_1 = (-1, -1, 0, 0, 1, 1)'$ IS ORTHOGONAL TO VECTOR $\mathbf{x}_2 - \bar{x}_2 = (-0.5, 0.5, -0.5, 0.5, -0.5, 0.5)'$, BECAUSE THEIR DOT PRODUCT IS 0. SINCE R^2 IS THE SQUARE OF THE CORRELATION COEFFICIENT, WE HAVE THAT

$$\text{VIF} = \frac{1}{1 - R_{12}^2} = \frac{1}{1 - r_{12}^2}$$

WHERE

$$r_{12} = \text{CORR}(x_1, x_2) = \frac{\text{COV}(x_1, x_2)}{\sigma_{x_1} \sigma_{x_2}}$$

WHERE

$$\text{COV}(x_1, x_2) = \text{E}[(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)] = \text{E}[(x_1 - 2)(x_2 - 0.5)] = 0$$

THEREFORE, $\text{VIF} = 1$.

- (b) Now suppose that the researcher used levels 1, 2, and 3 for supplement A, and levels 2, 3, and 4 for supplement B. Use software if desired. What is the variance inflation factor for variance x_1 in this case? Is it larger or smaller than in part (a) above? Why?

THE DESIGN MATRIX IS

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 3 & 0 \\ 1 & 2 & 1 \\ 1 & 3 & 1 \\ 1 & 4 & 1 \end{bmatrix}$$

WE HAVE $\text{COV}(x_1, x_2) = \text{E}[(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)] = 0.3$, $\sigma_{x_1} = 1.4088$, AND $\sigma_{x_2} = 0.5477$, SO $\text{VIF} = \frac{1}{1 - 0.273} = 1.375$. THE VARIANCE INFLATION FACTOR WILL BE LARGER THAN IN PART (A) BECAUSE THE TWO VARIABLES ARE NOW CORRELATED.