

STAT 608, Spring 2022 - Assignment 4  
SOLUTIONS

1. Suppose that for the model  $y_i = \alpha + e_i$ , the errors are independent with mean 0. Also suppose that measurements are taken using one device for the first  $n_1$  measurements, and then a more precise instrument was used for the next  $n_2$  measurements. Thus  $\text{Var}(e_i) = \sigma^2$ ,  $i = 1, 2, \dots, n_1$  and  $\text{Var}(e_i) = \sigma^2/2$ ,  $i = n_1 + 1, n_1 + 2, \dots, n$ .

- (a) Ignore the fact that the errors have different variances, and derive the least squares estimator for  $\hat{\alpha}$  using matrix notation and  $\hat{\alpha} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ .

NOTICE FIRST THAT THE DESIGN MATRIX IS SIMPLY A COLUMN OF ONES SINCE WE'RE FITTING AN INTERCEPT-ONLY MODEL. IGNORING ANY WEIGHTING, WE HAVE  $\hat{\alpha} = (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}'\mathbf{y} = \bar{y}$ , THE MEAN OF THE  $y$ -VALUES.

- (b) Derive the weighted least squares estimator for  $\alpha$ ,  $\hat{\alpha}_{\text{WLS}}$ .

TAKING  $w_i = 1/\sigma_i^2$ , WE HAVE A DIAGONAL MATRIX WITH  $n_1$  ELEMENTS EQUAL TO 1 AND  $n_2$  ELEMENTS EQUAL TO 2:

$$\mathbf{W} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 2 \end{bmatrix}$$

THEN

$$\hat{\alpha}_{\text{WLS}} = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}\mathbf{y} = \frac{1}{n_1 + 2n_2} \left( \sum_{i=1}^{n_1} y_i + 2 \sum_{i=n_1+1}^n y_i \right)$$

- (c) Suppose that  $n_1 = n_2$ . Compute the expected values and variances of the two estimators above. Which is a better estimator and why? (Use theoretical  $\text{MSE}(\hat{\alpha}) = \text{Bias}^2(\hat{\alpha}) + \text{Var}(\hat{\alpha})$  as your definition of "better.")

I'LL DEFINE  $n_1 = n_2 = k$  SO THAT  $n = 2k$ .

$$\begin{aligned} \text{E}(\hat{\alpha}) &= \text{E}(\bar{Y}) \\ &= \frac{1}{n} \sum \text{E}(Y_i) \\ &= \frac{1}{n} \sum \text{E}(\alpha + e_i) \\ &= \frac{1}{n} n\alpha \\ &= \alpha \end{aligned}$$

$$\begin{aligned}
\text{VAR}(\hat{\alpha}) &= \text{VAR}(\bar{Y}) \\
&= \text{VAR}\left(\frac{1}{n} \sum Y_i\right) \\
&= \frac{1}{n^2} \sum \text{VAR}(Y_i) \\
&= \frac{1}{(2k)^2} \sum \text{VAR}(\alpha + e_i) \\
&= \frac{1}{(2k)^2} (k\sigma^2 + k\sigma^2/2) \\
&= \frac{k\sigma^2}{(2k)^2} \left(1 + \frac{1}{2}\right) \\
&= \frac{\sigma^2}{4k} \left(\frac{3}{2}\right) \\
&= \frac{3\sigma^2}{8k}
\end{aligned}$$

$$\begin{aligned}
\text{E}(\hat{\alpha}_{\text{WLS}}) &= \text{E}\left[\frac{1}{k+2k} \left(\sum_{i=1}^k y_i + 2 \sum_{i=k+1}^{2k} y_i\right)\right] \\
&= \frac{1}{3k} \left(\sum_{i=1}^k \text{E}(\alpha + e_i) + 2 \sum_{i=k+1}^{2k} \text{E}(\alpha + e_i)\right) \\
&= \frac{1}{3k} (k\alpha + 2k\alpha) \\
&= \alpha
\end{aligned}$$

$$\begin{aligned}
\text{Var}(\hat{\alpha}_{\text{WLS}}) &= \text{Var}\left(\frac{1}{k+2k} \left(\sum_{i=1}^k y_i + 2 \sum_{i=k+1}^{2k} y_i\right)\right) \\
&= \frac{1}{(3k)^2} \left(\sum_{i=1}^k \text{Var}(\alpha + e_i) + 4 \sum_{i=k+1}^{2k} \text{Var}(\alpha + e_i)\right) \\
&= \frac{1}{(9k^2)} (k\sigma^2 + 4k\sigma^2/2) \\
&= \frac{k\sigma^2}{9k^2} (1 + 2) \\
&= \frac{\sigma^2}{3k} < \frac{3\sigma^2}{8k}
\end{aligned}$$

CLEARLY, BOTH ESTIMATORS ARE UNBIASED. THE VARIANCE OF  $\hat{\alpha}_{\text{WLS}}$  IS SMALLER, SO WE PREFER THE WEIGHTED LEAST SQUARES MODEL; ON AVERAGE, ITS ESTIMATE WILL BE CLOSER TO THE TRUE PARAMETER VALUE.

2. Question 2, Chapter 4

FIRST, WE FIND THAT  $w_i = 1/x_i^2$ , SO THAT  $\text{VAR}(\sqrt{w_i}e_i) = x_i^2\sigma^2/x_i^2 = \sigma^2$  IS CONSTANT. THUS OUR WEIGHT MATRIX  $\mathbf{W}$  IS  $\text{DIAG}(1/x_i^2)$ . THE DESIGN MATRIX HERE IS SIMPLY THE VECTOR OF VALUES  $x_i$ , SO WE SEE, USING  $\hat{\beta} = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{Y}$ , THAT:

$$\mathbf{X}'\mathbf{W}\mathbf{X} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} 1/x_1^2 & 0 & \cdots & 0 \\ 0 & 1/x_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/x_n^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum 1 = n$$

AND

$$\mathbf{X}'\mathbf{W}\mathbf{Y} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} 1/x_1^2 & 0 & \cdots & 0 \\ 0 & 1/x_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/x_n^2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum \frac{y_i}{x_i}$$

So,  $\hat{\beta} = \frac{\sum y_i/x_i}{n}$ .

3. Question 3, Chapter 4

- (a) FIRST, WE SEE FROM THE PLOT THAT THE VARIANCES OF THE RESPONSE VARIABLE ARE NOT EQUAL AT DIFFERENT VALUES OF EITHER PREDICTOR VARIABLE. THAT IS ONE ASSUMPTION OF A VALID MODEL, SO ANY MODEL ASSUMING EQUAL VARIANCES IS NOT VALID. SINCE THE RESPONSE VARIABLE IS A MEDIAN PRICE PER SQUARE FOOT, AND WE KNOW THAT THE VARIANCE OF A MEDIAN IS PROPORTIONAL TO  $1/n_i$ , WHERE  $n_i$  IS THE NUMBER OF OBSERVATIONS FROM WHICH THE MEDIAN IS CALCULATED,  $n_i$  IS THE APPROPRIATE CHOICE FOR THE WEIGHTS.
- (b) WE STILL DO NOT HAVE A VALID MODEL (EVEN AFTER USING THESE WEIGHTS) BECAUSE THE RESIDUALS DO NOT HAVE CONSTANT VARIANCE (INDICATED BY THE INCREASING SLOPE IN THE LAST PLOT) ACROSS THE DIFFERENT FITTED VALUES FOR THE MODEL. VARIANCES OF THE RESIDUALS SHOULD BE CONSTANT ACROSS ANY LINEAR COMBINATION OF THE  $x$  VALUES, INCLUDE THE  $y$ -HAT VALUES. ALSO, THE RELATIONSHIPS BETWEEN  $y$  AND THE  $x$  VARIABLES DO NOT APPEAR TO BE LINEAR.
- (c) FIRST, I WOULD CONSIDER ADDING PREDICTORS, SUCH AS INDICATOR VARIABLES FOR WHEN EACH OF THE CURRENT PREDICTORS EQUAL ZERO, OR A CATEGORICAL VARIABLE FOR BEING LOCATED IN AN URBAN OR SUBURBAN AREA OF HOUSTON, OR BEING LOCATED INSIDE OR OUTSIDE THE LOOP OR THE BELTWAY. SOME OF THE LARGEST HOME VALUES MAY BE EXPLAINED BETTER USING ONE OF THOSE VARIABLES. PERHAPS THE AVERAGE INCOME OF THE NEIGHBORHOOD (POSSIBLY AVAILABLE FROM CENSUS DATA) COULD ALSO BE A PREDICTOR. ALSO, IT IS OCCASIONALLY THE CASE THAT A 0 IS USED TO INDICATE THAT THE INFORMATION IS ACTUALLY MISSING. NEXT, I WOULD TRY AGAIN TO USE THE WEIGHTS DISCUSSED HERE, AND CONSIDER TRANSFORMATIONS OF THE PREDICTORS AND RESPONSE UNTIL I HAD THE FUNCTIONAL FORM OF THE MODEL CORRECT AND CONSTANT VARIANCE OF THE RESIDUALS. IF I STILL HAD AN OUTLIER, I MIGHT DOUBLE-CHECK THE VALUES OF THE OBSERVATION FOR ACCURACY AND THEN CONSIDER AN INDICATOR FOR THAT OUTLIER TO PREVENT IT FROM BEING A BAD LEVERAGE POINT IN THE MODEL.

4. Return to Question 4 from Homework 2, about coins being put on a scale. Now suppose that the variance in  $Y$  is proportional to the number of coins put on the scale. I recommend double-checking using both linear algebra and (if you're working in R) the linear model function `lm( $y \sim x$ , weight = w)`, where  $w$  is a vector of weights (the diagonal of the weight matrix) and you invent your own  $y$ .

- (a) Design an appropriate matrix of weights  $\mathbf{W}$ .

SINCE THE FIRST TWO MEASUREMENTS HAVE ONE COIN AND THE LAST TWO HAVE TWO COINS, THE VARIANCES CAN BE WRITTEN AS  $\sigma^2, \sigma^2, 2\sigma^2, 2\sigma^2$ . THE APPROPRIATE WEIGHT MATRIX IS THUS:

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 1/2 \end{bmatrix}$$

- (b) Calculate the new least-squares estimates of the weights of the coins using weighted least squares.

THIS TIME, WE HAVE

$$\begin{aligned} \hat{\mu}_{\text{WLS}} &= (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1} \mathbf{X}'\mathbf{W}\mathbf{y} \\ &= \left( \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 1/2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 1/2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} y_1 + y_3/2 + y_4/2 \\ y_2 + y_3/2 + y_4/2 \end{bmatrix} \\ &= \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} y_1 + y_3/2 + y_4/2 \\ y_2 + y_3/2 + y_4/2 \end{bmatrix} \\ &= \frac{1}{3} \begin{bmatrix} 2y_1 - y_2 + y_3/2 + y_4/2 \\ -y_1 + 2y_2 + y_3/2 + y_4/2 \end{bmatrix} \end{aligned}$$

5. For the model  $y_i = \beta_0 + \beta_1 x_i + e_i$ , the errors are iid with mean 0. The four observed values of  $x_i$  are  $\mathbf{x} = [1 \ 2 \ 3 \ 4]'$ . The estimator of  $\beta_1$  is  $\tilde{\beta}_1 = (y_4 + 2y_3 - 2y_2 - y_1) / 5$ . For this model, do the following:

- (a) Is  $\tilde{\beta}_1$  unbiased? Show why or why not.

YES, IT IS UNBIASED, SINCE

$$\begin{aligned} E(\tilde{\beta}_1) &= E[(y_4 + 2y_3 - 2y_2 - y_1) / 5] \\ &= (1/5) [(\beta_0 + 4\beta_1) + 2(\beta_0 + 3\beta_1) - 2(\beta_0 + 2\beta_1) - (\beta_0 + \beta_1)] \\ &= (1/5) [4\beta_1 + 6\beta_1 - 4\beta_1 - \beta_1] \\ &= \beta_1 \end{aligned}$$

- (b) What is the sampling variance of  $\tilde{\beta}_1$ ?

$$\text{Var}((y_4 + 2y_3 - 2y_2 - y_1)/5) = (1/25)(\sigma^2 + 4\sigma^2 + 4\sigma^2 + \sigma^2) = (10/25)\sigma^2 = (2/5)\sigma^2$$

- (c) Get the usual least squares estimator of  $\beta_1$  and calculate its sampling variance.

FIRST, WE CALCULATE  $\bar{x} = 2.5$ ,  $\sum x_i^2 = 30$ , AND  $\text{SXX} = 5$ . THEN WE HAVE

$$\begin{aligned}\hat{\beta}_1 &= (1/5) \left( \sum x_i y_i - n\bar{x}\bar{y} \right) \\ &= (1/5) (y_1 + 2y_2 + 3y_3 + 4y_4 - 10\bar{y}) \\ &= (1/5) \left( y_1 + 2y_2 + 3y_3 + 4y_4 - 10 \times \frac{y_1 + y_2 + y_3 + y_4}{4} \right) \\ &= (1/5) \left( -\frac{3}{2}y_1 - \frac{1}{2}y_2 + \frac{1}{2}y_3 + \frac{3}{2}y_4 \right)\end{aligned}$$

THUS

$$\begin{aligned}\text{VAR}(\hat{\beta}_1) &= (1/25) \left( \frac{9}{4}\sigma^2 + \frac{1}{4}\sigma^2 + \frac{1}{4}\sigma^2 + \frac{9}{4}\sigma^2 \right) \\ &= (1/5)\sigma^2\end{aligned}$$

- (d) Compare the sampling variance of  $\tilde{\beta}_1$  with the sampling variance of the least squares estimator of  $\beta_1$ .

THE VARIANCE OF  $\hat{\beta}_1$  IS SMALLER THAN THAT OF  $\tilde{\beta}_1$ . SO WE FIND ANOTHER LINEAR UNBIASED ESTIMATOR ( $\hat{\beta}_1$ ) WHOSE VARIANCE IS SMALLER THAN  $\tilde{\beta}_1$ . THUS  $\tilde{\beta}_1$  CANNOT BE A BLUE. (ACTUALLY  $\hat{\beta}_1$  IS A BLUE AND HAS THE SMALLEST VARIANCE AMONG ALL LINEAR UNBIASED ESTIMATORS.)

6. A food manufacturing company is interested in modeling whether people prefer  $x_1 =$  Type A or Type B hotdog buns with their hot dogs. They also want to control for  $x_2 =$  different amounts of sodium in the hot dogs themselves and are testing the hot dog buns at a variety of sodium contents, giving each taster both a hot dog and a bun with no condiments. The response variable is  $y =$  perceived taste of the bun, on a scale of 1 to 10.

- (a) In order to find out whether Type A or Type B is preferred, is it necessary to have an interaction term? Why or why not?

THE FULL ANSWER IS THAT IT DEPENDS. AS A FIRST PASS, WE GENERALLY THINK OF AN INTERACTION-FREE SITUATION: THAT IS, THERE ARE PARALLEL LINES, AND AS SODIUM INCREASES, PERCEIVED TASTE MIGHT INCREASE LINEARLY, AND WE'LL JUST BE INTERESTED IN WHETHER TYPE A OR TYPE B HAS A HIGHER INTERCEPT. IN THAT CASE, NO, WE DON'T NEED AN INTERACTION. TO VISUALLY CHECK WHETHER AN INTERACTION IS NECESSARY, WE COULD PLOT  $y$  VS.  $x_2$  AND COLOR CODE THE POINTS BY  $x_1$ . IF IT APPEARS THAT THE DIFFERENT COLORED POINTS HAVE DIFFERENT SLOPES, WE MIGHT USE AN INTERACTION TERM IN THE MODEL.

- (b) Develop a linear model for this study, interpreting all parameters in the context of the problem. You can assume that there is no interaction in the model. Write down your hypothesis to be testing in terms of your model parameters. (You don't have any data to conduct the test; just write down the hypotheses.)

LET'S ASSUME THERE IS NO INTERACTION FOR THIS MODEL:  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + e_i$ . IN THIS CASE, IF  $x_1$  IS AN INDICATOR OF BEING TYPE A, THEN

- $\beta_0$  IS THE MEAN PERCEIVED TASTE OF BUN TYPE B WHEN THE HOT DOG HAS NO SODIUM.
- $\beta_1$  IS THE MEAN DIFFERENCE IN TASTE BETWEEN TYPES A AND B, **holding sodium constant**.
- $\beta_2$  IS THE SLOPE: THE AMOUNT OF INCREASE OR DECREASE IN MEAN TASTE WHEN SODIUM INCREASES BY ONE UNIT (PROBABLY MG), **holding the bun type constant**.

THE RELEVANT NULL HYPOTHESIS IS  $H_0 : \beta_1 = 0$ , AND WE WOULD COMPARE THIS TO THE ALTERNATIVE HYPOTHESIS  $H_1 : \beta_1 \neq 0$ .

- (c) Assume now that an interaction term was added and it is statistically significant. How should we interpret this interaction in context?

IF THERE IS AN INTERACTION, IT MEANS THE EFFECT OF SODIUM IN THE HOT DOG IS DIFFERENT DEPENDING ON WHETHER TYPE A OR B BUNS ARE USED; TO PUT IT A DIFFERENT WAY, IT MAY BE THAT THE TYPE OF BUN THAT IS PREFERRED DIFFERS DEPENDING ON THE AMOUNT OF SODIUM IN THE HOT DOG (CROSSED LINES INTERACTION).

7. In a one-way ANOVA model with  $k = 3$  groups and 4 observations per group:

- (a) Use the F-statistic in Model Reduction Method 2 to derive a statistic for testing whether the average of the means of the first two groups is the same as the mean of the third group. That is, create the F-statistic for testing  $H_0 : (\mu_1 + \mu_2) / 2 = \mu_3$ . (Hint: Don't fit a model with a  $y$ -intercept. It makes everything easier.)

OUR NULL HYPOTHESIS IS  $H_0 : \frac{\mu_1}{2} + \frac{\mu_2}{2} - \mu_3 = 0$ , OR THAT  $\mathbf{A}\boldsymbol{\mu} = 0$  WHERE  $\mathbf{A} = [0.5 \ 0.5 \ -1]_{r \times (p+1)}$ . THE DESIGN MATRIX IS

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

ALSO,

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}_{(p+1) \times 1}$$

So,  $r = 1$  AND  $p = 2$ . WE CAN THEREFORE WRITE THE F-STATISTIC AS FOLLOWS:

$$\begin{aligned} F &= \frac{(\mathbf{A}\hat{\boldsymbol{\beta}} - h)' (\mathbf{A} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}')^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - h) / r}{\text{SSE}/(n - p - 1)} \\ &= \frac{\left( \mathbf{A} \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \hat{\mu}_3 \end{bmatrix} \right)' \left( \mathbf{A} \begin{bmatrix} 0.25 & 0 & 0 \\ 0 & 0.25 & 0 \\ 0 & 0 & 0.25 \end{bmatrix} \mathbf{A}' \right)^{-1} \left( \mathbf{A} \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \hat{\mu}_3 \end{bmatrix} \right)}{\text{SSE}/(12 - 2 - 1)} \\ &= \frac{24 (0.5\hat{\mu}_1 + 0.5\hat{\mu}_2 - \hat{\mu}_3)^2}{\text{SSE}} \end{aligned}$$

- (b) Where  $\hat{\mu}_1 = 5.6$ ,  $\hat{\mu}_2 = 7.9$ ,  $\hat{\mu}_3 = 6.1$ , and  $\text{SSE} = 12.8$ , test your hypothesis. Use  $\alpha = 0.05$ . Note that the degrees of freedom for the F-statistic are  $r$  (the number of rows of your  $\mathbf{A}$  matrix) and  $n - p - 1$ .

PLUGGING IN THE PROVIDED VALUES, WE HAVE

$$F = \frac{24(0.5 \times 5.6 + 0.5 \times 7.9 - 6.1)^2}{12.8} = 0.792 < F_{1,9}(0.05) = 5.12$$

BECAUSE THE F-STATISTIC FALLS BELOW THE CRITICAL VALUE, WE FAIL TO REJECT THE NULL HYPOTHESIS. THERE IS NOT ENOUGH EVIDENCE TO SUGGEST THAT THE AVERAGE OF THE FIRST TWO MEANS IS DIFFERENT FROM THE THIRD.