

STAT 8260 Exam 2 – Tuesday, April 10
SHOW ALL WORK

Name: Answer Key

1. Suppose we take a random sample of $n = 6$ adult father-son pairs and we measure the height of each man. Suppose that 3 of the sons in this sample were first-born children in their family and the other 3 happened to be second-born. Let (y_{ij}, x_{ij}) by the (son's height, father height) pair of measurements corresponding to the j th son of the i th birth order, $i = 1, 2, j = 1, 2, 3$. Consider the following five models for such data.

- (i) $y_{ij} = \alpha_i + e_{ij}$
- (ii) $y_{ij} = \alpha + \beta x_{ij} + e_{ij}$
- (iii) $y_{ij} = \alpha_i + \beta_i x_{ij} + e_{ij}$
- (iv) $y_{ij} = \mu + \alpha I(i = 2) + \beta x_{ij} + e_{ij}$
- (v) $y_{ij} = \mu + \alpha I(i = 2) + \beta x_{ij} I(i = 1) + \gamma x_{ij} I(i = 2) + e_{ij}$

where $I(\cdot)$ denotes an indicator variable taking the value 1 when the condition inside the parentheses is true and the value 0 otherwise. All models have the same assumptions on the e_{ij} s: $e_{11}, \dots, e_{23} \stackrel{iid}{\sim} N(0, \sigma^2)$.

- a. (8 pts.) Which, if any, of these models are equivalent.

(iii) and (v) are equivalent, the others are all different.

You could write down the model matrices for all 5 models and check which ones have the same column space, but I intended this to be much easier than that. Model (i) specifies 2 means, α_1 and α_2 for the 2 groups, respectively. Model (ii) specifies a single regression line (common intercept and slope β) for all subjects. Model (iii) allows two regression lines, $\alpha_1 + \beta_1 x$ for group 1, and $\alpha_2 + \beta_2 x$ for group 2. Model (iv) specifies 2 different intercepts: μ and $\mu + \alpha$ for groups 1 and 2, respectively, but a single common slope β . Finally, Model (v), like model (iii) specifies distinct intercepts, μ and $\mu + \alpha$, for the 2 groups, and distinct slopes, β and γ , for the 2 groups.

b. (9 pts.) Write down model (iii) in vector/matrix notation.

$$\begin{pmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{21} \\ y_{22} \\ y_{23} \end{pmatrix} = \begin{pmatrix} 1 & 0 & X_{11} & 0 \\ 1 & 0 & X_{12} & 0 \\ 1 & 0 & X_{13} & 0 \\ 0 & 1 & 0 & X_{21} \\ 0 & 1 & 0 & X_{22} \\ 0 & 1 & 0 & X_{23} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} e_{11} \\ e_{12} \\ e_{13} \\ e_{21} \\ e_{22} \\ e_{23} \end{pmatrix}$$

$$\underline{y} = \underline{X} \underline{\beta} + \underline{e}$$

c. (9 pts.) Under the maintained hypothesis that model (v) holds, express the hypothesis that $H_0 : \{\text{model (ii) holds}\}$, in the form of the general linear hypothesis on the parameters of model (v).

$$\underline{C} \underline{\beta} = \underline{t} \quad \text{where} \quad C = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}, \quad \underline{\beta} = \begin{pmatrix} \mu \\ \alpha \\ \beta \\ \gamma \end{pmatrix}$$

$$\text{and } \underline{t} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

- d. (9 pts.) Based upon the results in the table below, compute an appropriate test statistic (give me its numeric value for these data) and give its reference distribution for testing the hypothesis in part (c).

Model	MSE	R^2
(ii)	9.9268	.0432
(v)	3.2262	.8445

(Several way to compute F)

$$\bar{F} = \frac{(R_{FM}^2 - R_{RM}^2)/h}{(1 - R_{FM}^2)/(n-k-1)} = \frac{(.8445 - .0432)/2}{(1 - .8445)/(6-3-1)} = 5.15$$

$$\sim F(h, n-k-1) = F(2, 2) \text{ under } H_0$$

2. Consider the regression model

$$\begin{pmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \end{pmatrix} = \beta \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} + \begin{pmatrix} e_{11} \\ e_{12} \\ e_{21} \\ e_{22} \end{pmatrix} \\ = \beta \mathbf{j}_4 + \mathbf{e}$$

where $\begin{pmatrix} e_{11} \\ e_{12} \end{pmatrix} \sim N\left(0, \sigma_1^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$ independent of $\begin{pmatrix} e_{21} \\ e_{22} \end{pmatrix} \sim N(0, \sigma_2^2 \mathbf{I})$ where $\sigma_2^2 = 2\sigma_1^2$ and $\rho = .5$.

a. (13 pts.) Derive a simple, non-matrix formula for $\hat{\beta}_{\text{BLUE}}$, the best linear unbiased estimator of β in this model.

$$\underline{\mathbf{e}} \sim N(0, \sigma_1^2 \mathbf{V}) \text{ where } \mathbf{V} = \begin{pmatrix} 1 & .5 & 0 & 0 \\ .5 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix}$$

$\hat{\beta}_{\text{BLUE}}$ is the GLS estimator $(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$ where $\mathbf{X} = \mathbf{j}_4$

$$\mathbf{V}^{-1} = \begin{pmatrix} \frac{4}{3} \begin{pmatrix} 1 & -.5 \\ -.5 & 1 \end{pmatrix} & \underline{\mathbf{0}} \\ \underline{\mathbf{0}} & \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \end{pmatrix} \quad \mathbf{j}^T \mathbf{V}^{-1} = \begin{pmatrix} \frac{2}{3} & \frac{2}{3} & \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad (\text{sum of columns of } \mathbf{V}^{-1}) \\ (\mathbf{j}^T \mathbf{V}^{-1} \mathbf{j})^{-1} = \frac{3}{7}$$

$$(\mathbf{j}^T \mathbf{V}^{-1} \mathbf{j})^{-1} \mathbf{j}^T \mathbf{V}^{-1} \mathbf{y} = \frac{3}{7} \begin{pmatrix} \frac{2}{3} & \frac{2}{3} & \frac{1}{2} & \frac{1}{2} \end{pmatrix} \mathbf{y} = \begin{pmatrix} \frac{2}{7} & \frac{2}{7} & \frac{3}{14} & \frac{3}{14} \end{pmatrix} \mathbf{y}$$

$$= \frac{2}{7} y_1 + \frac{3}{14} y_2.$$

b. (13 pts.) Assuming the model is correct, compute the ratio

$$\frac{\text{var}(\bar{y}_{..})}{\text{var}(\hat{\beta}_{\text{BLUE}})},$$

where $\bar{y}_{..} = \frac{1}{4} \sum_{i=1}^2 \sum_{j=1}^2 y_{ij}$ is the sample mean.

$\bar{y}_{..}$ is the OLS estimator, with variance

$$\begin{aligned} \text{var}(\bar{y}_{..}) &= (X^T X)^{-1} X^T V X (X^T X)^{-1} = (\underline{j}^T \underline{j})^{-1} \underline{j}^T \sigma_1^2 V \underline{j} (\underline{j}^T \underline{j})^{-1} \\ &= \sigma_1^2 \frac{1}{16} \underline{j}^T V \underline{j} \quad , \quad \underline{j}^T V \underline{j} = \left(\frac{3}{2}, \frac{3}{2}, 2, 2 \right) \underline{j} = 7 \\ &= \sigma_1^2 7/16 \end{aligned}$$

$$\text{var}(\hat{\beta}_{\text{BLUE}}) = \sigma_1^2 (X^T V^{-1} X)^{-1} = \sigma_1^2 (\underline{j}^T V^{-1} X)^{-1} = \sigma_1^2 3/7$$

$$\Rightarrow \frac{\text{var}(\bar{y}_{..})}{\text{var}(\hat{\beta}_{\text{BLUE}})} = \frac{\sigma_1^2 \frac{7}{16}}{\sigma_1^2 \frac{3}{7}} = \frac{49}{48}$$

↑
notice > 1

3. Consider the classical linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

where \mathbf{X} is $n \times (k+1)$ with $\text{rank}(\mathbf{X}) = k+1 < n$. Suppose this model is fit to data \mathbf{y} yielding the ordinary least squares regression parameter estimate $\hat{\boldsymbol{\beta}}$ and MSE, s^2 . We wish to form a prediction interval for an additional observation y_0 conforming to this model. That is, $y_0 = \mathbf{x}_0^T \boldsymbol{\beta} + e_0$ where \mathbf{x}_0 is a vector of explanatory variables corresponding to y_0 (\mathbf{x}_0^T can be considered as an additional row of \mathbf{X}) and $e_0 \sim N(0, \sigma^2)$ independent of \mathbf{e} .

(13 pts.) For $\hat{y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}}$, show that

$$\frac{y_0 - \hat{y}_0}{s \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}} \sim t(n-k-1).$$

Since $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$, $\hat{y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}} \sim N(\mathbf{x}_0^T \boldsymbol{\beta}, \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0)$

Since \hat{y}_0 is computed ~~based on~~ ^{from} $\hat{\boldsymbol{\beta}}$ which is based on \mathbf{y} a sample independent of y_0 , $\text{Cov}(\hat{y}_0, y_0) = 0$. In addition $y_0 \sim N(\mathbf{x}_0^T \boldsymbol{\beta}, \sigma^2)$. So,

$$y_0 - \hat{y}_0 \sim N(0, \sigma^2 + \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0)$$

We know that $\frac{SSE}{\sigma^2} \sim \chi^2(n-k-1)$ and $s = \sqrt{\frac{SSE}{dfe}}$

and SSE independent of $\hat{\boldsymbol{\beta}}$ hence indep of $\hat{y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}}$
also, SSE indep of y_0 (different samples).

$$\text{Then since } t = \frac{(y_0 - \hat{y}_0) / \left[\sigma \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \right]}{\sqrt{\frac{SSE/\sigma^2}{dfe}}} = \frac{\text{num}}{\text{den}}$$

when num $\sim N(0, 1)$ \leftarrow independent

$$\text{den} = \sqrt{\chi^2/df} \leftarrow \Rightarrow t \sim t(dfe) = t(n-k-1)$$

4. Consider the classical linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

where \mathbf{X} is $n \times p$ with $\text{rank}(\mathbf{X}) = p < n$. Define $\tilde{\sigma}_\ell^2 = SSE/(n-\ell)$, an estimator of σ^2 . Recall the mean square error of an estimator T of a parameter θ is $E\{(T-\theta)^2\}$ but can also be represented as variance + bias².

a. (13 pts.) Find the mean square error of $\tilde{\sigma}_\ell^2$, for an arbitrary value of the constant ℓ .

$$\begin{aligned} \text{Var}(\tilde{\sigma}_\ell^2) &= \frac{1}{(n-\ell)^2} \text{Var}(\mathbf{y}^T \mathbf{P}_{C(\mathbf{X})} \mathbf{y}) \\ &= \frac{1}{(n-\ell)^2} \left[2 + \text{tr}(\mathbf{P}_{C(\mathbf{X})} \sigma^2 \mathbf{I} \mathbf{P}_{C(\mathbf{X})} + \sigma^2 \mathbf{I}) + 4(\mathbf{X}\boldsymbol{\beta})^T \mathbf{P}_{C(\mathbf{X})} \sigma^2 \mathbf{I} \underbrace{\mathbf{P}_{C(\mathbf{X})} \mathbf{X}\boldsymbol{\beta}}_{=0} \right] \\ &= \frac{1}{(n-\ell)^2} 2\sigma^4 \underbrace{\text{tr}(\mathbf{P}_{C(\mathbf{X})})}_{=n-p} = \frac{n-p}{(n-\ell)^2} 2\sigma^4 \\ \text{Bias}(\tilde{\sigma}_\ell^2) &= \text{Bias}\left(\frac{n-p}{n-\ell} \underbrace{\frac{SSE}{n-p}}_{=s^2}\right) = E\left(\frac{n-p}{n-\ell} s^2\right) - \sigma^2 \\ &= \frac{n-p}{n-\ell} \underbrace{E(s^2)}_{=\sigma^2} - \sigma^2 = \sigma^2 \left(\frac{n-p}{n-\ell} - 1\right) \\ \text{MSE}(\tilde{\sigma}_\ell^2) &= \frac{n-p}{(n-\ell)^2} 2\sigma^4 + \sigma^4 \left(\frac{n-p}{n-\ell} - 1\right)^2 \end{aligned}$$

b. (13 pts.) Find the value of l that yields the estimator $\hat{\sigma}_l^2$ with the smallest mean square error.

From part (a) $MSE(\hat{\sigma}_l^2) = \frac{4\sigma^4(n-p)}{(n-l)^3} + \frac{2\sigma^4(n-p)}{(n-l)^2} \left(\frac{n-p}{n-l} - 1\right)^2 \equiv M$

$$\frac{\partial M}{\partial l} = \frac{4\sigma^4(n-p)}{(n-l)^3} + 2\sigma^4 \left(\frac{n-p}{n-l} - 1\right) \left(-\frac{n-p}{(n-l)^2}(-1)\right)$$

$$= \frac{4\sigma^4(n-p)}{(n-l)^3} + \frac{2\sigma^4(n-p)}{(n-l)^2} \left(\frac{n-p}{n-l} - 1\right) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow 2 + (n-l) \left(\frac{n-p}{n-l} - 1\right) = 0 \quad \left(\text{By multiplying through by } \frac{(n-l)^3}{2\sigma^4(n-p)}\right)$$

$$\Rightarrow 2 + (n-p - n + l) = 0 \Rightarrow l = p - 2$$

$\Rightarrow \frac{SSE}{n-p+2}$ is the minimum MSE estimator of σ^2