

Name: Solutions

1. Suppose that we have observable random variables  $y_1, y_2, y_3$  and  $y_4$  satisfying  $E(y_1) = 2\beta_1 - \beta_2 + \beta_3 - \beta_4$ ,  $E(y_2) = 2\beta_1 + \beta_3$ , and  $E(y_3) = \beta_2$ ,  $E(y_4) = 2\beta_1 + \beta_2 + \beta_3$ . Let  $\mathbf{Y} = (y_1, y_2, y_3, y_4)^T$ , and  $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, \beta_4)^T$ . Answer part (a)–(e) in this scenario.

- (a) Find  $\mathbf{X}$  and  $\boldsymbol{\varepsilon}$  such that a model for  $\mathbf{Y}$  can be expressed in the form of  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ .

$$\mathbf{X} = \begin{pmatrix} 2 & -1 & 1 & -1 \\ 2 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \end{pmatrix}$$

$$\varepsilon_i = y_i - E(y_i) \quad (i=1,2,3,4).$$

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4)^T$$

- (b) Is  $\mathbf{X}$  in your model in (a) of full rank? Why or why not?

no.

It is easy to see that  $E(y_4) = E(y_2) + E(y_3)$ , thus  $\mathbf{X}$  is not of full row rank.

It is also easy to see that <sup>twice</sup> the third column of  $\mathbf{X}$  is exactly the first column of  $\mathbf{X}$ .

- (c) Clearly and precisely state minimal condition(s) under which your model in part (a) is a Gauss-Markov model.

$$\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I} \quad \text{for the } \boldsymbol{\varepsilon} \text{ defined in part (a).}$$

(Question 1, continued...)

(d) Is  $\beta_4$  estimable? If yes, find a linear unbiased estimator for  $\beta_4$ . If no, why?

Because  $\beta_4 = -(\bar{E}(y_1) - \bar{E}(y_2) + \bar{E}(y_3))$ , it is estimable.

Since  $\bar{E}(y)$  is always estimable, and so does any linear combination of its components.

$\hat{\beta}_4 = -(y_1 - y_2 + y_3)$  is a linear unbiased estimator for  $\beta_4$ .

(e) Let  $\theta_1 = \beta_1$ ,  $\theta_2 = 2\beta_2 - \beta_3$ ,  $\theta_3 = \beta_2 + 2\beta_3$ ,  $\theta_4 = \beta_4$ , and  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)^T$ . Find  $Z$  such that your model in (a) for  $Y$  can be written in the form of  $Y = Z\theta + \epsilon$ .

$$\left. \begin{array}{l} \theta_2 = 2\beta_2 - \beta_3 \\ \theta_3 = \beta_2 + 2\beta_3 \end{array} \right\} \Rightarrow \beta_2 = \frac{2}{5}\theta_2 + \frac{1}{5}\theta_3, \quad \beta_3 = -\frac{1}{5}\theta_2 + \frac{2}{5}\theta_3$$

$$\begin{aligned} \Rightarrow \bar{E}(y_1) &= 2\beta_1 - \beta_2 + \beta_3 - \beta_4 = 2\theta_1 - \left(\frac{2}{5}\theta_2 + \frac{1}{5}\theta_3\right) + \left(-\frac{1}{5}\theta_2 + \frac{2}{5}\theta_3\right) - \theta_4 \\ &= 2\theta_1 - \frac{3}{5}\theta_2 + \frac{1}{5}\theta_3 - \theta_4 \end{aligned}$$

$$\bar{E}(y_2) = 2\beta_1 + \beta_3 = 2\theta_1 - \frac{1}{5}\theta_2 + \frac{3}{5}\theta_3$$

$$\bar{E}(y_3) = \beta_2 = \frac{2}{5}\theta_2 + \frac{1}{5}\theta_3$$

$$\begin{aligned} \bar{E}(y_4) &= 2\beta_1 + \beta_2 + \beta_3 = 2\theta_1 + \frac{2}{5}\theta_2 + \frac{1}{5}\theta_3 + \frac{1}{5}\theta_2 + \frac{2}{5}\theta_3 \\ &= 2\theta_1 + \frac{3}{5}\theta_2 + \frac{3}{5}\theta_3 \end{aligned}$$

$$\Rightarrow Z = \begin{pmatrix} 2 & -\frac{3}{5} & \frac{1}{5} & -1 \\ 2 & -\frac{1}{5} & \frac{3}{5} & 0 \\ 0 & \frac{2}{5} & \frac{1}{5} & 0 \\ 2 & \frac{3}{5} & \frac{3}{5} & 0 \end{pmatrix}$$

2. Consider an experiment with two factors:  $A$  (with levels  $A_1$  and  $A_2$ ) and  $B$  (with levels  $B_1$  and  $B_2$ ). Let  $y_{ijk}$  be the outcome of the  $k$ th unit at the level of  $A_i$  factor and  $B_j$  factor ( $i, j = 1, 2$ ). Data are collected as in the following table:

Factor $A$	Factor $B$	Outcome
$A_1$	$B_1$	20
$A_1$	$B_1$	25
$A_1$	$B_2$	30
$A_1$	$B_2$	35
$A_2$	$B_1$	55
$A_2$	$B_2$	40
$A_2$	$B_2$	30

Consider a full rank Gauss-Markov model for the outcome data that takes the form

$$\mathbf{Y} = \begin{pmatrix} y_{111} \\ y_{112} \\ y_{121} \\ y_{122} \\ y_{211} \\ y_{221} \\ y_{222} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} + \boldsymbol{\varepsilon} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

In this question, assume that  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . Answer part (a)–(h) in this scenario.

- (a) Is this data set a balanced one?

No. combination  $A_2, B_1$  has only one outcome, while others have two.

- (b) Express the mean outcomes  $\mu_{ij} = E(y_{ijk})$  ( $i, j = 1, 2$ ) corresponding to all possible combinations of the factors  $A$  and  $B$  as functions of  $\beta_1, \beta_2, \beta_3, \beta_4$ .

$$\begin{aligned} \mu_{11} &= \beta_1 + \beta_2 + \beta_3 + \beta_4 & \mu_{12} &= \beta_1 + \beta_2 - \beta_3 - \beta_4 \\ \mu_{21} &= \beta_1 - \beta_2 + \beta_3 - \beta_4 & \mu_{22} &= \beta_1 - \beta_2 - \beta_3 + \beta_4 \end{aligned}$$

- (c) Express the overall mean of the outcomes as a function of  $\beta_1, \beta_2, \beta_3, \beta_4$ .

$$\begin{aligned} \mu_{\cdot} &= \frac{1}{4} (\mu_{11} + \mu_{12} + \mu_{21} + \mu_{22}) \\ &= \frac{1}{4} (4\beta_1) = \beta_1 \end{aligned}$$

(Question 2, continued ...)

Let  $\mathbf{P}$  be the projection matrix associated with  $\mathbf{X}$ , and  $\mathbf{I}$  be the identity matrix. Now you are given the following:

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \end{pmatrix} = \begin{pmatrix} 36.25 \\ -8.75 \\ 2.5 \\ -7.5 \end{pmatrix}, (\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{32} \cdot \begin{pmatrix} 5 & -1 & 1 & -1 \\ -1 & 5 & -1 & 1 \\ 1 & -1 & 5 & -1 \\ -1 & 1 & -1 & 5 \end{pmatrix} \text{ and } \mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{Y} = 75$$

where  $\hat{\beta}_i$  ( $i = 1, \dots, 4$ ) are the ordinary least squares estimates for the parameters. You don't have to simplify your expression in answering the following parts, but you need to clearly define any notation used and exactly specify which distribution and what quantile you are using.

- (d) Find a 95% level confidence interval for  $\mu_{21} - \mu_{11}$ .

$$\mu_{21} - \mu_{11} = -2(\beta_2 + \beta_4)$$

~~different~~  
~~function~~

use the formula  $\mathbf{c}^T \hat{\beta} \pm t_{\alpha/2} \sqrt{\text{MSE}} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}}$

$$\begin{aligned} \hat{\mathbf{c}}^T \hat{\beta} &= -2(-8.75 - 7.5) \\ &= 32.5 \end{aligned}$$

$t_{\alpha/2}$  is the upper 0.05 quantile of  $t$  distribution with 3 d.f.

$$\text{MSE} = 75/3 = 25$$

$$\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} = (0, -2, 0, -2)^T \begin{pmatrix} 5 & -1 & 1 & -1 \\ -1 & 5 & -1 & 1 \\ 1 & -1 & 5 & -1 \\ -1 & 1 & -1 & 5 \end{pmatrix} \begin{pmatrix} 0 \\ -2 \\ 0 \\ -2 \end{pmatrix}$$

- (e) Find the  $F$  statistic for  $H_0: \mu_{12} = \mu_{22} = 35$  vs  $H_a: \text{not } H_0$ , and give its degrees of freedom.

use  $H_0: \mathbf{C}\beta = \mathbf{d}$ .

$$\mathbf{C} = \begin{pmatrix} 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}, \mathbf{d} = \begin{pmatrix} 35 \\ 35 \end{pmatrix}$$

$$\begin{aligned} F &= \frac{(\hat{\mathbf{C}}\hat{\beta} - \mathbf{d})^T (\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T)^{-1} (\hat{\mathbf{C}}\hat{\beta} - \mathbf{d}) / 2}{\text{MSE}} \\ &= \frac{(32.5, 35)^T (\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T)^{-1} \begin{pmatrix} 32.5 \\ 35 \end{pmatrix} / 2}{25} \end{aligned}$$

$$\text{d.f.} = \underline{2, 3}$$



(Question 2, continued ...)

- (f) Suppose ~~two~~ <sup>three</sup> new outcomes are observed at the condition with respectively level  $A_1$  and  $B_2$ , find a 95% prediction interval for the average of the three new observations.

use 
$$\hat{c}^T \hat{\beta} \pm t_{\alpha/2} \sqrt{MSE} \sqrt{\gamma + c^T \text{cov}(c)}$$

with  $\gamma = 1/3$ ,  $c = (1, 1, -1, -1)$

- (g) Suppose that the fifth row (outcome=55) in the data table is now removed. Does it change the estimability of any of the parameters  $\beta_j$  ( $j = 1, \dots, 4$ ) in the model? Why or why not?

- Since the fifth row is the only outcome associated with  $A_2, B_1$ , removing it will reduce the rank of  $X$  from 4 to 3.
- In the cell-mean model  $y_{ijk} = \mu_{ij} + \epsilon_{ijk}$ ,  $\mu_{21}$  will become not estimable, any contrast involving  $\mu_{21}$  will also be not estimable.
- $\beta_1$  will become not estimable

- (h) Does the change in the previous part (g) have any impact on the least square estimation of  $\mu_{11}$  and  $\mu_{12}$ ? Explain your answer.

No. the OLS estimator for  $\mu_{11}$  &  $\mu_{12}$  remains the same.

by verifying the OLS estimators are

always  $\frac{20+25}{2}$  &  $\frac{30+35}{2}$ .

3. Consider the model  $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$  where  $\epsilon_i$  independently follow standard normal distribution for  $i = 1, \dots, n$ . Let  $x_0$  be the value at which the function  $E(y) = \beta_0 + \beta_1 x + \beta_2 x^2$  is maximized (or minimized). Answer parts (a) and (b) in this scenario.

(a) Find the maximum likelihood estimator for  $x_0$ .

- This is a full rank model unless in extreme case.

- The OLS estimator  $\hat{\beta} = (X^T X)^{-1} X^T Y$  is the MLE for  $\beta$ .

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

-  $E(y)$  is minimized or maximized at  $x_0 = -\frac{\beta_1}{2\beta_2}$

depending on the sign of  $\beta_2$

-  $\Rightarrow$  MLE of  $x_0$ ,  $\hat{x}_0 = -\frac{\hat{\beta}_1}{2\hat{\beta}_2}$  for  $\hat{\beta}_2 \neq 0$ , since the ratio is a continuous function.

(b) Find a  $(1 - \alpha)$  level confidence interval for  $x_0$ .

One solution:

$$\text{let } \Sigma = \text{var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}, \quad \hat{\Sigma} = \text{MSE}(X^T X)^{-1}$$

using the delta method.

$$\hat{x}_0 \approx -\frac{\beta_1}{2\beta_2} + \frac{\partial x_0}{\partial \beta_1} (\hat{\beta}_1 - \beta_1) + \frac{\partial x_0}{\partial \beta_2} (\hat{\beta}_2 - \beta_2)$$

$\Rightarrow$  An approximate variance of  $\hat{x}_0$  is

$$\hat{V}^2 = \left( \frac{\partial x_0}{\partial \beta_1} \right)^2 \text{var}(\hat{\beta}_1) + \left( \frac{\partial x_0}{\partial \beta_2} \right)^2 \text{var}(\hat{\beta}_2)$$

$$+ 2 \left( \frac{\partial x_0}{\partial \beta_1} \frac{\partial x_0}{\partial \beta_2} \right) \text{cov}(\hat{\beta}_1, \hat{\beta}_2) \quad \text{evaluated at } \hat{\beta}_1 \text{ and } \hat{\beta}_2$$

$$\text{with } \frac{\partial x_0}{\partial \beta_1} = -\frac{1}{2\beta_2} \quad \frac{\partial x_0}{\partial \beta_2} = \frac{\beta_1}{2\beta_2^2}$$

Then an approximate CI, is

$$\hat{x}_0 \pm z_{\alpha/2} \hat{V}$$

Q3 part (b) for the extra point without using large sample assumption

$$\text{let } \beta = (\beta_0, \beta_1, \beta_2) \quad C^{(x_0)} = (0, 1, 2x_0)$$

$$\text{Then } C^T(x_0) \beta = \beta_1 + 2\beta_2 x_0$$

$$\text{Note that At the truth } x_0 = -\frac{\beta_1}{2\beta_2} \quad C^T(x_0) \beta = 0.$$

Then, the F statistic

$$F^{(x_0)} = \frac{C^T(x_0) \hat{\beta} (C^T(x_0) (X^T X)^{-1} C(x_0))^{-1} C^T(x_0) \hat{\beta}}{SSE/n-3}$$

where 3 is the d.f. of  $X$  in this case

$F(x_0)$  follows  $F(1, n-3)$  distribution under the null hypothesis  $x_0 = -\beta_1/2\beta_2$

Thus, a  $1-\alpha$  level confidence region for  $x_0$

$$\text{is } \{ x_0 : F_{1, n-3, \alpha/2} < F(x_0) < F_{1, n-3, 1-\alpha/2} \}$$