

Exam 1

Advanced Linear Models (PHST 781)
Department of Biostatistics and Bioinformatics

12th October, 2018

Name:

- There are a total of 115 points in this Question Paper. Answer as much as you can. If your acquired score is greater than equal to 100, it will be counted as 100%.
- The Exam is scheduled for 3 hours. "Time Left" reminders will be posted in 1.5 hrs, 2:30 hrs, 2:45 hrs from the beginning of the Exam time.
- There are three ★ marked problems that are more involved than the rest. In case you are stuck in one of those, it might be a good idea to consider solving other problems first and then continue with the * marked problems.
- You may take help from the "Exam Assistance Note" containing a few required definitions, lemma and theorem statements.

Let Y be Response variable, X_1, X_2 denote the Explanatory variables, ε be unknown random errors and β_1, \dots, β_3 are unknown parameters of interest. Determine whether the following relationship equation is a linear model. Relationship Equation: $Y = \beta_0 + \beta_1 e^{X_1} + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$.

1. (a)

Score:
Total Score: 5

Ans:

☐ Not Linear Model

☐ Linear Model

Identify if the following matrix is a Orthogonal Projection matrix.

$$M = \frac{1}{9} \begin{bmatrix} 3 & 2 & 4 \\ 2 & 1 & 2 \\ 4 & 2 & 4 \end{bmatrix}$$

(b)

Score:
Total Score: 5

Ans:

☐ Orthogonal Projection

☐ Not Orthogonal Projection

Consider the matrices $A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$. Does the matrix A have a Generalized Inverse ? If yes, Construct a generalized inverse of A .

(c)

Score:
Total Score: 1+5

Ans: The matrix A have a generalized inverse.

The generalized inverse of A is $A^- = \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$.

Consider the matrix

$$D = \frac{1}{2} \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}.$$

(d)

- i. Represent D as a Kronecker product of two lower-dimensional matrices.
- ii. Prove that D is a **Orthogonal Projection** matrix.

Score:
Total Score: 3+7

Ans:

- (e) Consider the model $\underline{\mathbf{y}} = \mathbf{X}\underline{\boldsymbol{\beta}} + \underline{\boldsymbol{\varepsilon}}$, where $\underline{\boldsymbol{\beta}} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$, $\mathbf{X} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 0 & -1 \end{bmatrix}$ and $\underline{\boldsymbol{\varepsilon}}$ is a mean zero error vector with variance co-variance matrix $\sigma^2 I_{4 \times 4}$.
 Prove that the parameter β_1 is not estimable

Score:
 Total Score: 7

Ans:

Let $\mathbf{x}_1 = [1 \ 1 \ \dots \ 1]_{2n \times 1}^T$ and $\mathbf{x}_2 = [1 \ -1 \ 1 \ \dots \ -1]_{2n \times 1}^T$, i.e. the i^{th} entry of the vector \mathbf{x}_2 is $(-1)^{i-1}$ for $i = 1, \dots, 2n$. Consider the linear model

$$\mathbf{Y} = \mathbf{X}_{2n \times 2} \boldsymbol{\beta} + \boldsymbol{\varepsilon}_{2n \times 1},$$

2.

where the two columns of the design matrix \mathbf{X} are \mathbf{x}_1 and \mathbf{x}_2 , $\boldsymbol{\beta} = [\beta_1 \ \beta_2]^T$. Furthermore, assume that $\boldsymbol{\varepsilon}_{2n \times 1} \sim N(\mathbf{0}, \sigma^2 I_{2n \times 2n})$. \mathbf{Y} denotes the random vector corresponding to the **response variable**. Let a statistician is interested in inference for the parametric function $\theta = \beta_1 - \beta_2$

Is θ a linear parametric function of $\boldsymbol{\beta}$? If yes, find a vector $\boldsymbol{\lambda}_1$ such that $\theta = \boldsymbol{\lambda}_1^T \boldsymbol{\beta}$

(a)

Score:
Total Score: 1+ 4

Ans:

Is θ estimable parametric function? provide appropriate justification to your answer.

(b)

Score:
Total Score: 1+ 6

Ans:

Find the **Orthogonal Projection Matrix** for $\mathcal{C}(\mathbf{X})$, the column space of \mathbf{X} . (Justify your steps with the reference to the results/theorem/lemma you are using.)

(c)

Score:
Total Score: 8

Ans:

Find the **Best Linear Unbiased Estimator** for θ ? (Show your steps and justify your steps with the reference to the results/theorem/lemma you are using.)

(d)

Score:
Total Score: 10

Ans:

Consider the data set

Response:	15.8	20.1	19.1	16.7	14.3	19.1
Explanatory Variable:	7	6	4	6.5	3	2

3. A statistician believes that the data is being generated from a piece-wise linear model specified as,

$$Y = \begin{cases} \alpha_0 + \alpha_1 X + \varepsilon & \text{if } X \leq 5 \\ \gamma_0 + \gamma_1 X + \varepsilon & \text{if } X > 5. \end{cases} \quad (1)$$

Represent the above model in terms of the matrix form of the standard linear model

$$\underline{\mathbf{y}} = \mathbf{X}\underline{\boldsymbol{\beta}} + \underline{\boldsymbol{\varepsilon}},$$

- (a) by explicitly specifying (In terms of numbers when known and in terms of appropriate symbols when unknown) the following quantities a) \mathbf{X} , the design Matrix, b) $\underline{\mathbf{y}}$, the response vector and Vector of regression coefficient, $\underline{\boldsymbol{\beta}}$.

Score:
Total Score: 6+1+1

Ans:

(b)

The model defined in Equation (1) is not guaranteed to be continuous. Show that the fitted model will be continuous if the constraint $\alpha_0 - \gamma_0 = 5\gamma_1 - 5\alpha_1$ is imposed on the regression coefficients. Is it a linear constraint on the parameters?

Score:
Total Score: 2+1

Ans:

(c) *

Construct a model

$$\underline{\mathbf{y}}_{\star} = \mathbf{X}_{\star} \underline{\beta}_{\star} + \underline{\varepsilon},$$

incorporating the constraint in part(b). Specify (In terms of numbers when known and in terms of appropriate symbols when unknown) the following quantities a) \mathbf{X}_{\star} , the design Matrix, b) $\underline{\mathbf{y}}_{\star}$, the response vector and c) Vector of regression coefficient, $\underline{\beta}_{\star}$.

Score:
Total Score: 7+1+3

Ans: (You have more space in page 10 to write answer to this question.)

Consider a linear model given by

$$\mathbf{Y} = \mathbf{X}\tilde{\beta} + \tilde{\varepsilon}$$

4. where $\mathbf{Y} \in \mathbb{R}^n, \tilde{\beta} \in \mathbb{R}^p$ and $\tilde{\varepsilon}$ is a Normally distributed random vector with mean $\mathbf{0}$ and variance $\sigma^2 I_{n \times n}$ and \mathbf{X} is an $n \times p$ matrix with rank $r < p < n$ (i.e. The design matrix \mathbf{X} does not have full column rank). The parameter σ^2 is an unknown positive number. Assume that $\tilde{\varepsilon} \sim N(\mathbf{0}, \sigma^2 I_{n \times n})$. Let $\mathcal{C}(\mathbf{X})$ denotes the column space of \mathbf{X} , the vector space containing all possible linear combination of the column vectors of the matrix \mathbf{X} . It can be shown that $P = \mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is the **Orthogonal Projection** matrix for the $\mathcal{C}(\mathbf{X})$, where $(\mathbf{X}^T \mathbf{X})^-$ is any Generalized inverse of $(\mathbf{X}^T \mathbf{X})$.

Show that $E(\mathbf{Y}^T (I_{n \times n} - P) \mathbf{Y}) = (n - r) \sigma^2$.

(a)

Score:
Total Score: 5

Ans:

Prove that the statistics $\mathbf{Y}^T \mathbf{P} \mathbf{Y}$ and $\mathbf{Y}^T (\mathbf{I}_{n \times n} - \mathbf{P}) \mathbf{Y}$ are independent (Mention if you are using any result).

(b)

Score:
Total Score:5

Ans:

Derive the distribution of $\frac{(\mathbf{Y}^T P \mathbf{Y})/r}{(\mathbf{Y}^T (I_{n \times n} - P) \mathbf{Y})/(n-r)}$? (Show your steps)

(c)

Score:
Total Score: 10

Ans:

Let $\hat{\beta}_{LSE}$ be the **Least Square Estimator** for the parameter β . Show that

$$(\mathbf{Y} - \mathbf{X}\hat{\beta}_{LSE})^T P(\mathbf{Y} - \mathbf{X}\hat{\beta}_{LSE}) = (\hat{\beta}_{LSE} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta}_{LSE} - \beta).$$

(d) ★

Score:
Total Score: 5

Ans:

(e) ★

Derive the distribution of $\frac{\|X(\hat{\beta}_{LSE} - \beta)\|^2}{\sigma^2}$. (You may use the relation in part(d) to get your answer.)

Score:
Total Score: 5

Ans:

5. Consider the linear model when the observed response variable

$$Y_{i,j} = \mu + \tau_i + \varepsilon_{i,j} \text{ for } i = 1, 2, \dots, 5; j = 1, 2, \dots, 3,$$

where $\varepsilon_{i,j} \stackrel{iid}{\sim} N(0, \sigma^2)$ and $\sigma^2, \mu, \tau_1, \dots, \tau_5$ are unknown parameters of the model. μ is called the ‘baseline effect’ or ‘mean effect’ while τ_1, \dots, τ_5 are called ‘treatment effects’. Consider the notations,

$$\mathbf{Y} = \left[\begin{array}{c|c|c|c} \text{1st Treatment} & & & \\ \hline Y_{1,1}, \dots, Y_{1,3} & \dots & Y_{i,1}, \dots, Y_{i,3} & \dots \\ \hline & & \text{5th Treatment} & \\ \hline & & Y_{5,1}, \dots, Y_{5,3} & \end{array} \right]^T,$$

$$\beta = [\mu, \tau_1, \dots, \tau_5]^T = [\mu \quad \tau^T]^T, \text{ with } \tau^T = [\tau_1, \dots, \tau_5],$$

$$\varepsilon = \left[\begin{array}{c|c|c|c} \text{1st Treatment} & & & \\ \hline \varepsilon_{1,1}, \dots, \varepsilon_{1,3} & \dots & \varepsilon_{i,1}, \dots, \varepsilon_{i,3} & \dots \\ \hline & & \text{5th Treatment} & \\ \hline & & \varepsilon_{5,1}, \dots, \varepsilon_{5,3} & \end{array} \right]^T,$$

$$\mathbf{1}_5 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}_{5 \times 1}, \quad \mathbf{1}_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}_{3 \times 1}, \quad \mathbf{0} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} \quad \text{and}$$

$$\mathbf{X} = \left[\begin{array}{c|c|c|c|c} \mathbf{1}_3 & \mathbf{1}_3 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{1}_3 & \mathbf{0} & \mathbf{1}_3 & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_3 & \mathbf{0} & \dots & \dots & \mathbf{1}_3 \end{array} \right]_{15 \times 6} = \left[\underbrace{\mathbf{1}_5 \otimes \mathbf{1}_3}_{X_\mu} \mid \underbrace{I_5 \otimes \mathbf{1}_3}_{X_\tau} \right] = [X_\mu \mid X_\tau],$$

where $X_\mu = \mathbf{1}_5 \otimes \mathbf{1}_3$, $X_\tau = I_5 \otimes \mathbf{1}_3$, then we can represent the model as

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon.$$

- Consider a linear parametric function $\alpha^T \tau$ where $\alpha \in \mathbb{R}^5$. **Prove that, if $\alpha^T \tau$ is estimable then it must be a contrast between the treatments.**
- Let a linear parametric function $\alpha^T \tau$ is a contrast between the treatment effects for some $\alpha^T \in \mathbb{R}^5$. **Prove that $\alpha^T \tau$ is estimable.**
- Note that $\tau_1 - \tau_2$ estimable as $\tau_1 - \tau_2$ is a contrast between the treatments.

Derive $\widehat{\tau_1 - \tau_2}$ the Best Linear Unbiased Estimator for $\tau_1 - \tau_2$ and express your answer in terms of $\bar{Y}_{i,\bullet}$ for $i = 1, \dots, 5$.

$$\text{It is given that } P_X \mathbf{Y} = \begin{bmatrix} \bar{Y}_{1,\bullet} \\ \vdots \\ \bar{Y}_{5,\bullet} \end{bmatrix} \otimes \mathbf{1}_3 \text{ where } \bar{Y}_{i,\bullet} = \frac{1}{3} \sum_{j=1}^3 Y_{i,j} \text{ for } i = 1, \dots, 5.$$

Here P_X denotes the Orthogonal Projection matrix for the column space of \mathbf{X}

HINT: You may have already constructed a vector $\mathbf{a} \in \mathbb{R}^{15}$ such that $\begin{bmatrix} 0 & \boldsymbol{\alpha}^T \end{bmatrix} = \mathbf{a}^T \mathbf{X}$.

- (d) What is P_X (mention if you are using any nontrivial result to derive P_X)?

Derive the variance of $\widehat{\tau_1 - \tau_2}$?

- (e) What is the value of the constant \mathbf{v} so that

$$\hat{\sigma}^2 = \frac{\mathbf{Y}^T (I - P_X) \mathbf{Y}}{\mathbf{v}} = \frac{\sum_{i=1}^5 \sum_{j=1}^3 Y_{i,j}^2 - 3 \sum_{i=1}^5 \bar{Y}_{i,\bullet}^2}{\mathbf{v}}$$

is unbiased estimator of σ^2 . (You do not need to prove)

- (f) How many unique primary contrasts of the treatments τ_1, \dots, τ_5 are there ?

(Note: primary contrasts are the contrasts of the form $\tau_i - \tau_j$ for $1 \leq i \neq j \leq 5$. Note that , $\tau_i - \tau_j$ and $\tau_j - \tau_i$ are assumed to be 'same' primary contrast for a pair of distinct indices $i \neq j$.)

- (g) Write down the definition of the **simultaneous confidence intervals** for all the unique primary contrasts.
- (h) Construct a 95% Bonferroni's simultaneous confidence intervals for all the unique primary contrasts above.