

Preliminary Examination: Linear Models

Q1 counts for 60 points and Q2, 40 points.

Answer questions by showing all of your work.

1. (60 points) For a linear model $y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \varepsilon_i$, $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$, consider the centered model $y_i = \alpha + \beta_1(x_{i1} - \bar{x}_1) + \cdots + \beta_k(x_{ik} - \bar{x}_k) + \epsilon_i$ where $\bar{x}_j = \sum_{i=1}^n x_{ij}/n$ for $j=1, \dots, k$, which is also written by $\mathbf{y} = [\mathbf{1} \ \mathbf{X}_c] \begin{pmatrix} \alpha \\ \boldsymbol{\beta}_1 \end{pmatrix} + \boldsymbol{\varepsilon}$ where

$$\mathbf{X}_c = \begin{pmatrix} x_{11} - \bar{x}_1 & \cdots & x_{1k} - \bar{x}_k \\ \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & \cdots & x_{nk} - \bar{x}_k \end{pmatrix}.$$

- (a) Prove that the (ordinary) least squares estimators for α and $\boldsymbol{\beta}_1$ are \bar{y} and $(\mathbf{X}_c' \mathbf{X}_c)^{-1} \mathbf{X}_c' \mathbf{y}$, respectively.
- (b) Find the distribution of SSR_c/σ^2 where $\text{SSR}_c = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ and $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_k x_{ik}$. Check all the conditions if you apply a theorem you memorize.
- (c) Prove that $\hat{\boldsymbol{\varepsilon}} \sim N(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{P}))$ where $\hat{\boldsymbol{\varepsilon}} = (\hat{\varepsilon}_1 \ \dots \ \hat{\varepsilon}_n)^\top$ and $\hat{\varepsilon}_i = y_i - \hat{y}_i$ after defining \mathbf{P} appropriately.
- (d) Find the distribution of SSE/σ^2 where $\text{SSE} = \sum_{i=1}^n \hat{\varepsilon}_i^2$. Check all the conditions if you apply a theorem you memorize.
- (e) The response variables can be written as $\mathbf{y} = \hat{\mathbf{y}} + \hat{\boldsymbol{\varepsilon}}$ from the original model where $\hat{\mathbf{y}} = \mathbf{P}\mathbf{y}$ and $\mathbf{y} = \hat{\alpha}\mathbf{1} + \hat{\mathbf{y}}_c + \hat{\boldsymbol{\varepsilon}}$ from the centered form where $\hat{\mathbf{y}}_c = \mathbf{P}_c\mathbf{y}$.
 - i. Show that $\mathbf{P}_c\mathbf{y} = (\mathbf{P} - \frac{1}{n}\mathbf{J})\mathbf{y}$.
 - ii. For SSE defined in part (d), show that $\text{SSE} = \mathbf{y}'(\mathbf{I} - \mathbf{P}_c - \frac{1}{n}\mathbf{J})\mathbf{y}$.
 - iii. Show that SSR_c defined in part (b) and SSE defined in part (d) are independent.

2. (40 points) We consider the linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where

$$\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$$

for some $\mathbf{X}_1 \in \mathbb{R}^{n \times p_1}$ and $\mathbf{X}_2 \in \mathbb{R}^{n \times p_2}$. To accommodate for this block matrix form, we write $\boldsymbol{\beta} \in \mathbb{R}^{p_1+p_2}$ as,

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix},$$

where $\boldsymbol{\beta}_1 \in \mathbb{R}^{p_1}$ and $\boldsymbol{\beta}_2 \in \mathbb{R}^{p_2}$. Throughout this question, assume that the design matrix \mathbf{X} is deterministic with full rank and $n > p_1 + p_2$. The error vector $\boldsymbol{\varepsilon}$ has multivariate normal distribution with zero mean and a diagonal covariance matrix with common diagonal element σ^2 .

Let

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{p_1+p_2}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|, \quad \hat{\boldsymbol{\beta}}_1 = \arg \min_{\boldsymbol{\beta}_1 \in \mathbb{R}^{p_1}} \|\mathbf{Y} - \mathbf{X}_1\boldsymbol{\beta}_1\|, \quad \hat{\boldsymbol{\beta}}_2 = \arg \min_{\boldsymbol{\beta}_2 \in \mathbb{R}^{p_2}} \|\mathbf{Y} - \mathbf{X}_2\boldsymbol{\beta}_2\|,$$

where $\|\cdot\|$ denotes the Euclidean norm of a vector.

- (a) Construct an example of (\mathbf{Y}, \mathbf{X}) where $\hat{\boldsymbol{\beta}}^T \neq (\hat{\boldsymbol{\beta}}_1^T, \hat{\boldsymbol{\beta}}_2^T)$.
- (b) Prove that $\hat{\boldsymbol{\beta}}^T = (\hat{\boldsymbol{\beta}}_1^T, \hat{\boldsymbol{\beta}}_2^T)$ if $\mathbf{X}_1^T \mathbf{X}_2$ is a zero matrix.
- (c) Prove that $\hat{\boldsymbol{\beta}}_1^T$ and $\hat{\boldsymbol{\beta}}_2^T$ are independent if $\mathbf{X}_1^T \mathbf{X}_2$ is a zero matrix.
- (d) Prove that $\|\mathbf{Y} - \mathbf{X}_1\boldsymbol{\beta}_1\|^2 + \|\mathbf{Y} - \mathbf{X}_2\boldsymbol{\beta}_2\|^2 \geq \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$ if $\mathbf{X}_1^T \mathbf{X}_2$ is a zero matrix.
- (e) Construct an example of (\mathbf{Y}, \mathbf{X}) where

$$\|\mathbf{Y} - \mathbf{X}_1\boldsymbol{\beta}_1\|^2 + \|\mathbf{Y} - \mathbf{X}_2\boldsymbol{\beta}_2\|^2 = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$