

Applied Statistics Preliminary Examination

Theory of Linear Models

May 2025

Instructions:

- Do all 3 Problems. Neither calculators nor electronic devices of any kind are allowed. Show all your work, clearly stating any theorem or fact that you use. Each of the 12 parts carries an equal weight of 10 points.
- Abbreviations/Acronyms.
 - IID (independent and identically distributed).
 - LSE (least squares estimator); BLUE (best linear unbiased estimator). Sometimes the LSE may be designated OLS (ordinary least squares) estimator, in order to differentiate it from the GLS (generalized least squares) estimator.
- Notation.
 - \mathbf{x}^T or \mathbf{A}^T : indicates transpose of vector \mathbf{x} or matrix \mathbf{A} .
 - $\text{tr}(\mathbf{A})$ and $|\mathbf{A}|$: denotes the trace and determinant, respectively, of matrix \mathbf{A} .
 - \mathbf{I}_n : the $n \times n$ identity matrix.
 - $\mathbf{j}_n = (1, \dots, 1)^T$ is an n -vector of ones, and $\mathbf{J}_{m,n}$ is an $m \times n$ matrix of ones.
 - $\mathbb{E}(\mathbf{x})$ and $\mathbb{V}(\mathbf{x})$: expectation and variance of random vector \mathbf{x} .
 - $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$: the m -dimensional random vector \mathbf{x} has a normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
 - $X \sim t(n, \lambda)$: a t distribution with n degrees of freedom and noncentrality parameter λ . If $\lambda = 0$ we write simply: $X \sim t(n)$.
 - $X \sim F(n_1, n_2, \lambda)$: an F distribution with n_1 and n_2 numerator and denominator degrees of freedom respectively, and noncentrality parameter λ . If $\lambda = 0$ we write simply: $X \sim F(n_1, n_2)$.
- Possibly useful results.
 - For a generic linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where the design matrix \mathbf{X} is of dimensions $n \times k$, and $\widehat{\boldsymbol{\beta}}$ is the LSE of $\boldsymbol{\beta}$, the R^2 is defined as:

$$R^2 = \left(\widehat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} - n\bar{y}^2 \right) / S_{yy}, \quad \bar{y} = \sum_{i=1}^n y_i / n, \quad S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2.$$

We also establish the following notation: if x_{ij} is the (i, j) -th element of the design matrix \mathbf{X} , let $\bar{x}_j = \sum_{i=1}^n x_{ij} / n$ be the mean of the j -th column.

Problems:

1. Let \mathbf{X}_k be a full-rank $n \times k$ matrix whose column space $C(\mathbf{X}_k)$ consists of the linear span of the vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, i.e., $C(\mathbf{X}_k) = \mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_k)$, and let \mathbf{X}_{k+1} be a full-rank $n \times (k+1)$ matrix whose column space consists of the linear span of the vectors that constitute $C(\mathbf{X}_k)$ and the extra vector \mathbf{x}_{k+1} , so that $C(\mathbf{X}_{k+1}) = \mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_k, \mathbf{x}_{k+1})$. In addition, let $\mathbf{G}_k = (\mathbf{X}_k^T \mathbf{X}_k)^{-1}$, $\mathbf{G}_{k+1} = (\mathbf{X}_{k+1}^T \mathbf{X}_{k+1})^{-1}$, $\mathbf{P}_k = \mathbf{X}_k \mathbf{G}_k \mathbf{X}_k^T$, and $\mathbf{P}_{k+1} = \mathbf{X}_{k+1} \mathbf{G}_{k+1} \mathbf{X}_{k+1}^T$.

(a) Show that \mathbf{P}_k is the projection matrix onto $C(\mathbf{P}_k)$.

(b) Show that $C(\mathbf{P}_k) = C(\mathbf{X}_k)$.

(c) Show that $\mathbf{P}_{k+1} - \mathbf{P}_k$ is a projection matrix.

(d) Is $C(\mathbf{P}_k)$ orthogonal to $C(\mathbf{P}_{k+1} - \mathbf{P}_k)$? Justify.

2. Consider the two full-rank linear regression models for the same response vector $\mathbf{y} = (y_1, \dots, y_n)^T$ with sample mean \bar{y} and sums of squares about the mean S_{yy} (see these definitions in the “Possibly useful results” of page 1):

$$\mathbf{y} = \mathbf{X}_k \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (1)$$

and

$$\mathbf{y} = \mathbf{X}_{k+1} \mathbf{b} + \mathbf{e}, \quad (2)$$

where the design matrices are as defined in Problem 1, so that $C(\mathbf{X}_k) \subset C(\mathbf{X}_{k+1})$. As the notation suggests, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)^T$ and $\mathbf{b} = (b_1, \dots, b_{k+1})^T$ are the corresponding coefficient vectors of appropriate dimension, with associated LSEs $\hat{\boldsymbol{\beta}}$ and $\hat{\mathbf{b}}$. The ultimate goal of this problem is to prove that, if R_k^2 and R_{k+1}^2 are the R^2 values for models (1) and (2) respectively, then R_{k+1}^2 cannot be smaller than R_k^2 . (To simplify the problem we also assume that the first column of each design matrix consists of a vector of 1's, i.e., $\mathbf{x}_1 = \mathbf{j}_n$.) In what follows, x_{ij} is the (i, j) -th element of the design matrix \mathbf{X}_k , and $\bar{x}_j = \sum_{i=1}^n x_{ij}/n$ is the mean of the j -th column (see “Possibly useful results” of page 1).

(a) Show that $\bar{y} = \sum_{j=1}^k \hat{\beta}_j \bar{x}_j$.

(b) If x_{ij} is the (i, j) -th element of \mathbf{X}_k , and $\bar{x}_j = \sum_{i=1}^n x_{ij}/n$ is the mean of the j -th column, show that $S_{yy} R_k^2 = \sum_{j=1}^k \sum_{i=1}^n \hat{\beta}_j (x_{ij} - \bar{x}_j) y_i$.

(c) By considering the difference $S_{yy}(R_{k+1}^2 - R_k^2)$ and noting the results of Problem 1, or otherwise, show that $R_k^2 \leq R_{k+1}^2$.

3. Consider the observations $\{y_{ij}\}$ from the following linear model on 6 distinct subjects:

$$y_{ij} = \alpha_i + (\mu + \gamma_i)x_{ij} + \epsilon_{ij}, \quad i = 1, 2, \quad j = 1, 2, 3,$$

where x_{ij} is the age of subject j whose gender is indicated by the index i . Let $\boldsymbol{\beta}^T = (\alpha_1, \alpha_2, \mu, \gamma_1, \gamma_2)$, assume that all subjects have different ages, and that $\{\epsilon_{ij}\} \sim \text{IID } N(0, \sigma^2)$.

(a) Express the model in the usual vector-matrix form, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, write down the design matrix \mathbf{X} , and determine its rank.

(b) Is the hypothesis $H_0 : \alpha_1 = \alpha_2 = \mu = 0$ testable? Carefully justify your answer.

- (c) It is desired to test if the model can be reduced to common intercepts ($\alpha_1 = \alpha_2$) and slopes ($\gamma_1 = \gamma_2$) for the two genders. Carefully express this as an appropriate null hypothesis, and show that it is testable.
- (d) Propose a test statistic for the test in (c) and specify its distribution under H_0 . Give detailed explanations and provide numerical values for all the quantities that can actually be calculated.

Applied Statistics Preliminary Examination

Theory of Linear Models

August 2025

Instructions:

- Do all 3 Problems. Neither calculators nor electronic devices of any kind are allowed. Show all your work, clearly stating any theorem or fact that you use. Each of the 14 parts carries an equal weight of 10 points.
- Abbreviations/Acronyms.
 - IID (independent and identically distributed).
 - SSE (sum of squared errors). Also called *residual sum of squares*.
 - LSE (least squares estimator); BLUE (best linear unbiased estimator). Sometimes the LSE may be designated OLS (ordinary least squares) estimator, in order to differentiate it from the GLS (generalized least squares) estimator.
- Notation.
 - \mathbf{x}^T or \mathbf{A}^T : indicates transpose of vector \mathbf{x} or matrix \mathbf{A} .
 - $\text{tr}(\mathbf{A})$ and $|\mathbf{A}|$: denotes the trace and determinant, respectively, of matrix \mathbf{A} .
 - \mathbf{I}_n : the $n \times n$ identity matrix.
 - $\mathbf{j}_n = (1, \dots, 1)^T$ is an n -vector of ones, and $\mathbf{J}_{m,n}$ is an $m \times n$ matrix of ones.
 - $\mathbb{E}(\mathbf{x})$ and $\mathbb{V}(\mathbf{x})$: expectation and variance of random vector \mathbf{x} .
 - $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$: the m -dimensional random vector \mathbf{x} has a normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
 - $X \sim t(n, \lambda)$: a t distribution with n degrees of freedom and noncentrality parameter λ . If $\lambda = 0$ we write simply: $X \sim t(n)$.
 - $X \sim F(n_1, n_2, \lambda)$: an F distribution with n_1 and n_2 numerator and denominator degrees of freedom respectively, and noncentrality parameter λ . If $\lambda = 0$ we write simply: $X \sim F(n_1, n_2)$.
- Possibly useful results.
 - Note the inverse for the patterned matrix \mathbf{A} :

$$\mathbf{A} = \begin{bmatrix} 3 & 0 & -2 \\ 0 & 2 & 1 \\ -2 & 1 & 3 \end{bmatrix} \implies \mathbf{A}^{-1} = \frac{1}{7} \begin{bmatrix} 5 & -2 & 4 \\ -2 & 5 & -3 \\ 4 & -3 & 6 \end{bmatrix}.$$

- If \mathbf{a} and \mathbf{b} are vectors in a vector space with inner-product $\langle \mathbf{a}, \mathbf{b} \rangle$, then the *Cauchy-Schwarz Inequality* states that:

$$|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \sqrt{\langle \mathbf{a}, \mathbf{a} \rangle \langle \mathbf{b}, \mathbf{b} \rangle}.$$

Problems

1. Let $\mathbf{x} \sim N_3(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_3)$, with $\mathbf{x}^T = (x_1, x_2, x_3)$, $\boldsymbol{\mu}^T = (3, -2, 1)$, and define the following derived random variables and vectors:

$$3y = 2x_1^2 + 2x_2^2 + 2x_3^2 - 2x_1x_2 - 2x_1x_3 - 2x_2x_3, \quad \mathbf{w} = \begin{bmatrix} x_1 + x_2 + x_3 \\ x_1 - x_3 \end{bmatrix}, \quad z = x_1 + x_2 + x_3.$$

- (a) Find the distribution of \mathbf{w} , and hence show that its two components are independent.
- (b) Show that y and z are independent.
- (c) Are y and \mathbf{w} independent? Justify.
- (d) Find the distribution of y/σ^2 .
2. Consider the linear model in the usual vector-matrix form: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where the $\boldsymbol{\epsilon} \sim N(0, \sigma^2)$, and \mathbf{X} is $n \times p$ of full-rank, i.e., $\text{rk}(\mathbf{X}) = p$. Let $\widehat{\boldsymbol{\beta}}$ and s^2 denote respectively, the LSE of $\boldsymbol{\beta}$ and the usual unbiased estimate of σ^2 . Letting $\mathbf{G} = (\mathbf{X}^T \mathbf{X})^{-1}$, recall the following basic results:

$$\widehat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 \mathbf{G}), \quad \frac{(n-p)s^2}{\sigma^2} \sim \chi^2(n-p), \quad \widehat{\boldsymbol{\beta}} \text{ is independent of } s^2.$$

Consider now a single new model point (\mathbf{x}_0, y_0) , so that $y_0 = \mathbf{x}_0^T \boldsymbol{\beta} + \epsilon_0$. Given \mathbf{x}_0 , we know that the best point prediction of $\mathbb{E}(y_0) = \mathbf{x}_0^T \boldsymbol{\beta} \equiv \mu_0$ is $\widehat{y}_0 = \mathbf{x}_0^T \widehat{\boldsymbol{\beta}}$. Using the distribution of the standardized \widehat{y}_0 , one can then construct a (pointwise) confidence interval for the regression “line” μ_0 . The purpose of this Problem is to construct a confidence band of coverage $(1 - \alpha)100\%$ for the entire regression “line” that will hold simultaneously for all design points \mathbf{x}_0 .

- (a) Show that:

$$\frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{G}^{-1} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{p s^2} \sim F(p, n-p).$$

- (b) By defining an inner-product for vectors \mathbf{a} and \mathbf{b} in \mathbb{R}^p as $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{G} \mathbf{b}$, and taking $\mathbf{a} = \mathbf{x}_0$ and $\mathbf{b} = \mathbf{G}^{-1} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})$, show that:

$$(\widehat{y}_0 - \mu_0)^2 \leq (\mathbf{x}_0^T \mathbf{G} \mathbf{x}_0) (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{G}^{-1} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

- (c) Use the above results to derive the fact that a $(1 - \alpha)100\%$ simultaneous confidence band for the entire regression “line” is given by

$$\widehat{y}_0 \pm w,$$

and determine the expression for w as a function of \mathbf{x}_0 , \mathbf{G} , s^2 , and an appropriate F distribution quantile.

3. Consider the observations $\mathbf{y} = (y_{11}, \dots, y_{22})^T$ from the linear model:

$$y_{ij} = \mu + \alpha_i + \gamma_j + \epsilon_{ij}, \quad i = 1, 2, \text{ and } j = 1, 2,$$

where the $\{\epsilon_{ij}\} \sim \text{IID } N(0, \sigma^2)$, and $\boldsymbol{\beta} = (\mu, \alpha_1, \alpha_2, \gamma_1, \gamma_2)^T$ are unknown parameters to be estimated. Note that the model can be written in the usual vector-matrix form: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$.

- (a) Find the rank of the design matrix \mathbf{X} , and show that it is rank-deficient (not of full-rank). (Note: if $\text{rk}(\mathbf{X}) = k$, and the number of columns in \mathbf{X} is p , the rank-deficiency is still defined to be $q = p - k$, even though $n = 4 < 5 = p$ in this overparametrized model.)
- (b) Determine the general form of all the *estimable* functions: $\eta = \boldsymbol{\lambda}^T \boldsymbol{\beta}$.
- (c) Show (with justification) that the set of *side conditions* $\alpha_1 + \alpha_2 = \gamma_1 - \gamma_2 = 0$ is not a valid and complete set that can be used to remove the overparametrization.
- (d) In light of the previous part, propose a valid and complete set of side conditions that will successfully reparametrize the model to full rank.
- (e) Write down the new (full-rank) model matrix $\widetilde{\mathbf{X}}$ and new parameter vector $\widetilde{\boldsymbol{\beta}}$ that would result from using the valid and complete set of side conditions proposed in the previous part (whence the model can be reformulated as $\mathbf{y} = \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}} + \boldsymbol{\epsilon}$).