

**Applied Statistics Preliminary Examination**  
**Theory of Linear Models**  
**May 2022**

**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  $\lambda$ . 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  $\lambda$ . If  $\lambda = 0$  we write simply:  $X \sim F(n_1, n_2)$ .
- Possibly useful results.
  - If  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is given in partitioned form as

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix},$$

with  $m_1 = \dim(\mathbf{x}_1)$ , then the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2$  is

$$\mathbf{x}_1 | \mathbf{x}_2 \sim N_{m_1}(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}).$$

- Let the  $n$ -dimensional vector  $\mathbf{y} \sim N(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n)$ , where  $\boldsymbol{\mu} = \mathbf{X}\mathbf{b}$  for an  $(n \times k)$  full-rank matrix of constants  $\mathbf{X}$  with  $k < n$ . Define the matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , and the quadratic forms  $Q_1 = \mathbf{y}^T \mathbf{H} \mathbf{y}$  and  $Q_2 = \mathbf{y}^T (\mathbf{I}_n - \mathbf{H}) \mathbf{y}$ .
  - Show that both of the matrices  $\mathbf{H}$  and  $\mathbf{I}_n - \mathbf{H}$  are symmetric and idempotent.
  - Find all the eigenvalues and ranks of each of  $\mathbf{H}$  and  $\mathbf{I}_n - \mathbf{H}$ .
  - Find a constant  $c$  such that  $cQ_1$  and  $cQ_2$  have familiar distributions, and specify what these resulting distributions are.
  - Find the distribution of

$$W = \frac{(n-k)Q_1}{kQ_2}.$$

- Consider fitting the linear model where the full-rank  $(n \times k)$  model matrix is decomposed as  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ , with  $\mathbf{X}_1$   $(n \times k_1)$ ,  $\mathbf{X}_2$   $(n \times k_2)$ ,  $k = k_1 + k_2$ ,  $\boldsymbol{\beta}^T = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)$ , and the normal spherical errors assumption,  $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$ :
 
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}.$$

Suppose however that the vector of observations  $\mathbf{y}$  comes from the reduced model  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}$ , and it is known that the column spaces of  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are orthogonal, i.e.,  $\mathcal{C}(\mathbf{X}_1) \perp \mathcal{C}(\mathbf{X}_2)$ . Thus, by fitting the model with model matrix  $\mathbf{X}$  we are in fact *overfitting*. Let  $\widehat{\boldsymbol{\beta}}^T = (\widehat{\boldsymbol{\beta}}_1^T, \widehat{\boldsymbol{\beta}}_2^T)$ , denote the resulting LSE of  $\boldsymbol{\beta}$ , and  $s^2 = (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}})^T(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}})/(n-k)$  the usual estimator of  $\sigma^2$ .

- Find  $\mathbb{E}(\widehat{\boldsymbol{\beta}}_1)$  and  $\mathbb{E}(\widehat{\boldsymbol{\beta}}_2)$ . Is  $\widehat{\boldsymbol{\beta}}_1$  unbiased for  $\boldsymbol{\beta}_1$ ?
  - Compute  $\mathbb{V}(\widehat{\boldsymbol{\beta}}_1)$ ,  $\mathbb{V}(\widehat{\boldsymbol{\beta}}_2)$ , and  $\text{Cov}(\widehat{\boldsymbol{\beta}}_1, \widehat{\boldsymbol{\beta}}_2)$ .
  - If  $\mathcal{P}_{\mathcal{C}(\mathbf{X})}$ ,  $\mathcal{P}_{\mathcal{C}(\mathbf{X}_1)}$ , and  $\mathcal{P}_{\mathcal{C}(\mathbf{X}_2)}$  denote the respective projection matrices onto the column spaces  $\mathcal{C}(\mathbf{X})$ ,  $\mathcal{C}(\mathbf{X}_1)$ , and  $\mathcal{C}(\mathbf{X}_2)$ , show that  $\mathcal{P}_{\mathcal{C}(\mathbf{X})} = \mathcal{P}_{\mathcal{C}(\mathbf{X}_1)} + \mathcal{P}_{\mathcal{C}(\mathbf{X}_2)}$ .
  - Show that even though the model is overfit,  $s^2$  is still unbiased for  $\sigma^2$ .
- Consider the linear model  $y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \epsilon_i$ , where  $i = 1, \dots, n$ ,  $\{\epsilon_i\} \sim \text{IID } N(0, \sigma^2)$ , and the  $x_{i,j}$  are elements of the  $(n \text{ by } p+1)$  design matrix  $\mathbf{X}$  (of rank  $k < p+1 < n$ ). Let  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\bar{\mathbf{y}} = \mathbf{j}_n^T \mathbf{y} / n$  be the sample mean of the  $y_i$ , and let  $\widehat{\boldsymbol{\beta}} = \mathbf{G} \mathbf{X}^T \mathbf{y}$  denote the LSE of  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T$  corresponding to the generalized inverse  $\mathbf{G}$  of  $\mathbf{X}^T \mathbf{X}$ . Additionally, let  $R^2 = (\widehat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} - n\bar{y}^2) / (\mathbf{y}^T \mathbf{y} - n\bar{y}^2)$  be the usual *coefficient of determination*, and define the quadratic forms:
 
$$Q_1 = (\bar{\mathbf{y}} \mathbf{j}_n - \mathbf{y})^T (\bar{\mathbf{y}} \mathbf{j}_n - \mathbf{y}), \quad Q_2 = (\bar{\mathbf{y}} \mathbf{j}_n - \mathbf{X} \widehat{\boldsymbol{\beta}})^T (\bar{\mathbf{y}} \mathbf{j}_n - \mathbf{X} \widehat{\boldsymbol{\beta}}), \quad Q_3 = (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}).$$

- Without using calculus, prove that  $\mathbf{X}^T (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}) = \mathbf{0}$ . (Your argument should use the form of  $\widehat{\boldsymbol{\beta}}$  and derive any properties of generalized inverses that are deemed necessary in the proof.)
- Prove that  $Q_1 = Q_2 + Q_3$ .
- If  $\mathcal{C}(\mathbf{X})$  denotes the column space of  $\mathbf{X}$ , prove or disprove the statement:  $R^2 = 1$  if and only if  $\mathbf{y} \in \mathcal{C}(\mathbf{X})$ .
- Assuming that  $\eta = 3\beta_1 - \beta_0$  is *estimable*, show how to construct a test of  $H_0 : \eta = 2$  vs.  $H_1 : \eta \neq 2$ . Clearly define the test statistic, its distribution under both  $H_0$  and  $H_1$ , and the rejection rule.

## Design of Experiment: Prelim Problems

May 2022

Please Do All Problems. Each of the 16 parts carries an equal weight of 10 points.

1. In an experiment on the preparation of chocolate cakes, 3 recipes (factor A) for the batter were compared. Recipes I and II differed in that the chocolate was added at 40 degrees Celsius and 60 degrees Celsius, respectively, while recipe III contained extra sugar. Six different baking temperatures (factor B), ranging from 175 to 225 degrees Celsius were tested with each recipe. Assume factor A has fixed effect and factor B has random effect. Interaction effects between factors A and B are the only interaction effects. The moisture level of each cake is measured and used as response variable. Please write down the statistical models and assumptions for each of the following designs.
  - a) Suppose for each of the 18 treatment combinations, three cakes are made.
  - b) Suppose 2 bags of flour (factor D with random effect) were used. Each bag of flour is used to make two cakes for each of the 18 treatment combinations.
  - c) Suppose 3 bags of flour (factor D with random effect) are randomly assigned to each of 18 treatment combinations. In total, 54 bags of flour are used. Two cakes are made from each bag.
  - d) There were 10 replications (factor C with random effect) conducted over time. Each time that a mix was made by any recipe, enough batter was prepared for 6 cakes, each of which was baked at a different temperature.

2. In a study, the research interest is the total concentration of a specific alkaloid for Herb A as measured by infusion extraction (steeping the tea bag). The process of herb drying, mixing, and bagging was found to be highly reliable (accurate and precise), but the method of growing and harvesting Herb A was found to affect alkaloid concentration. The focus of this problem is to analyze the experimental growing data for Herb A to assess the resulting alkaloid concentration measured in controlled lab conditions.

Three varieties of Herb A (A1, A2, and A3), were grown at research greenhouses with either 0, 23, 45, or 68  $g/m^2$  of fertilizer (primarily nitrogen supplied as soy bean meal). Eighteen temporary fields were created by placing standard top soil in previously unplanted areas, with small experimental green houses placed over each one. Each of the 3 varieties were randomly assigned to 6 of the 18 green houses, so that each green house had only one variety (to protect from confusion at harvest). The four levels of fertilizer were randomly assigned to equal-sized quarters of each greenhouse. At harvest time, 3 leaves were collected from each of several center plants in each quarter and combined as the sample. Each sample was uniformly processed (dried, crushed, and blended), then portions shipped to three labs for alkaloid concentration measurements. We will analyze the data from one of the labs. Of interest are differences in variety and fertilizer levels. Assume there are no interaction effects.

Denote  $y_{iu_j}$  as the alkaloid concentration measurement for the sample from  $i^{th}$  Herb A variety (factor A, assumed with fixed effect),  $u^{th}$  assigned field, and  $j^{th}$  fertilizer level (factor B, assumed with random effect). We have some summary information about the data as follows:  $\sum_j (\bar{y}_{..j} - \bar{y}_{...})^2 = 0.119$ ;  $\bar{y}_{1..} = 1.667$ ;  $\bar{y}_{2..} = 1.725$ ;  $\bar{y}_{3..} = 1.588$ .

- a) Please write down a model for this experiment. Clearly define all the terms and state all relevant assumptions.
- b) Fill in the ANOVA table below.

Source	d.f.	MS	E(MS)
variety	_____	_____	_____
Error(w)	_____	0.128	_____
Total(w)	_____	0.126	_____
fertilizer	_____	_____	_____
Error(S)	_____	0.0173	_____
Total(S)	_____	_____	_____

- c) State and conduct statistical tests for Variety effect at 0.05 significance level, and interpret the test results.
- d) State and conduct statistical tests for Fertilizer effect at 0.05 significance level, and interpret the test results.
- e) Calculate E(Mean Square for Variety) and E(MSE(W)).

- f) Use Tukey's method to construct simultaneous 95% confidence intervals for the pairwise effect difference of the three varieties.
- g) Between Variety effect and Fertilizer effect, which one is more accurately estimated? What adjustment do you want to make about the experimental design, if we want to estimate the other effect more accurately?

3. Suppose that a furniture manufacturer has developed a new method for assembling a particular type of rocking chair at its assembly plant and wishes to determine if the new method (method 2) produces a smaller mean time to assemble the rocking chair than does the current method (method 1). Because industrial research shows that different training techniques are sometimes needed to effectively train different age groups to use a new assembly method, the plant supervisor also wishes to assess whether the plant's training techniques for the new assembly method will effectively train the different age groups of workers in the plant. To this end, the plant supervisor defines three age groups for the plant's workers, where group 1 represents younger workers, group 2 represents intermediate age workers, and group 3 represents older workers. Three workers are randomly selected from the workers in each age group and are trained to use the new assembly method, where a training technique that is thought to be appropriate for a particular age group is used to train each randomly selected worker in that age group. Each trainee then assembles two rocking chairs using method 1 (the current method) and two rocking chairs using method 2 (the new method). (The order in which the methods are used is randomly selected for each trainee.)

$y_{ijkl}$  represents the  $l$  th assembly time used by trainee  $k$  nested inside age group  $j$  and using assembly method  $i$  to assemble the rocking chair. Assume assembly method (factor A) and age group (factor B) both have fixed effect. Trainee (factor C) has random effect. Assume there are interaction effects between assembly method and age group. We know that  $\bar{y}_{1...} - \bar{y}_{2...} = 0.85$ .

- a) Write a linear model equation appropriate for analyzing the responses  $y_{ijkl}$  from this experiment, clearly state the assumptions and conditions for this model.
- b) Fill in the ANOVA table below.

Source	d.f.	MS	E(MS)
A			
B		8.03	
		6.543	
AB		0.59	
Error		2.31	
Total			

- c) Test if the new method (method 2) produces a smaller mean time to assemble the rocking chair than does the current method (method 1) at 0.05 significance level.
- d) State and conduct statistical test for age group effects at 0.05 significance level, and interpret the test results.
- e) State and conduct statistical test for whether the variance for trainee effect is greater than the variance of random error.