

1. Suppose you are designing a new study of the yield of a chemical process like the one partially analyzed in problems 4 through 6 on assignment 3. Suppose the engineers assigned to your project wish to run the process at the same five temperature values for each of two new catalysts. Call them catalyst A and catalyst B . The proposed model for the observed yield when the process is run with the i -th catalyst at the j -th temperature level is

$$Y_{ijk} = \mu + \alpha_i + \beta(T_{ij} - 100) + \epsilon_{ijk}, \quad i = 1, 2, \quad j = 1, 2, \dots, 5, \quad \text{and} \quad k = 1, \dots, n,$$

where T_{ij} is the temperature at which the process was run, and $\epsilon_{ijk} \sim NID(0, \sigma^2)$. When the runs are made, we will have n replicates for each the ten temperature/catalyst combinations. The engineers want to test the null hypothesis $H_0 : \alpha_1 = \alpha_2$ against the alternative $H_A : \alpha_1 \neq \alpha_2$ using a type I error level of $\alpha = 0.05$. Relative to the value of the error variance, σ^2 , they wish to make the number of replicates (n) large enough to have probability of at least 0.90 of rejecting the null hypothesis if $\alpha_1 - \alpha_2 = 0.5\sigma$. What is the smallest value of n that satisfies these conditions?

2. Marcuse(1949, Biometrics, 5) recorded moisture content for three types of cheese made by two different methods. Two pieces of cheese were measure for each type and each method. The data are shown below.

Treatment	Moisture Content Measurements	
Type A made with Method 1	$Y_{11} = 39.02$	$Y_{12} = 38.79$
Type B made with Method 1	$Y_{21} = 35.74$	$Y_{22} = 35.41$
Type C made with Method 1	$Y_{31} = 37.02$	$Y_{32} = 36.00$
Type A made with Method 2	$Y_{41} = 38.96$	$Y_{42} = 39.01$
Type B made with Method 2	$Y_{51} = 35.58$	$Y_{52} = 35.52$
Type C made with Method 2	$Y_{61} = 35.70$	$Y_{62} = 36.04$

Consider the model

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij},$$

where $\epsilon_{ijk} \sim NID(0, \sigma^2)$, $i = 1, 2, 3, 4, 5, 6$, and $j = 1, 2$. This model can be expressed in matrix form as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

equivalently,

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \\ Y_{41} \\ Y_{42} \\ Y_{51} \\ Y_{52} \\ Y_{61} \\ Y_{62} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \alpha_4 \\ \alpha_5 \\ \alpha_6 \end{bmatrix} + \begin{bmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{31} \\ \epsilon_{32} \\ \epsilon_{41} \\ \epsilon_{42} \\ \epsilon_{51} \\ \epsilon_{52} \\ \epsilon_{61} \\ \epsilon_{62} \end{bmatrix}$$

- (a) What is the distribution of $\mathbf{Y} = (Y_{11}, Y_{12}, Y_{21}, Y_{22}, Y_{31}, Y_{32}, Y_{41}, Y_{42}, Y_{51}, Y_{52}, Y_{61}, Y_{62})^T$?
- (b) Determine which of the following are testable hypotheses. You only need to state if the hypothesis is testable or not testable.

- i. $H_0 : \alpha_1 = \alpha_2 = \alpha_3$
- ii. $H_0 : \alpha_1 = \alpha_2 = \alpha_3 = 0$
- iii. $H_0 : \mu + \alpha_1 = 39$ and $\mu + \alpha_4 = 39$
- iv.

$$H_0 : \begin{bmatrix} 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 & -1 & 0 & 1 \end{bmatrix} \boldsymbol{\beta} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

v.

$$H_0 : \begin{bmatrix} 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 & -1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \boldsymbol{\beta} = \begin{bmatrix} 0 \\ 0 \\ 38 \end{bmatrix}$$

vi.

$$H_0 : \begin{bmatrix} 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 & -1 & 0 & 1 \\ 0 & 0 & 1 & -1 & 0 & -1 & 1 \end{bmatrix} \boldsymbol{\beta} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

- (c) Express each of the following hypotheses in the form $H_0 : C\boldsymbol{\beta} = \mathbf{0}$. If the hypothesis is testable, compute the value of the corresponding F -statistic, report the degrees of freedom and the p -value for the test, and state your conclusion.
- i. After averaging across the two methods of making cheese, the average moisture content is the same for all three types of cheese.
 - ii. For each type of cheese, the average moisture content is not affected by the method for making cheese. (This hypothesis allows the average moisture content to vary across types of cheese).

3. The following are part of the data reported by Ryan, et.al. (1976, Jour. of Atmo. Sci. 33) on the formation of ice crystals. The ice crystals were formed in a growth chamber maintained at a fixed temperature ($-5^{\circ}C$) and a fixed level of saturation of air with water. Ice crystals were harvested at various times (seconds) and the axial length (micrometers) of each ice crystal was measured. The objective was to model how mean length of ice crystals increases with time. The data file is posted as *crystals.txt*. It contains measurements on 16 ice crystals. The first row of the file contains variable names *time* and *length*. The data are also shown below.

time	length
60	18
60	21
80	25
80	28
100	30
100	29
100	33
120	36
120	34
120	28
140	32
140	35
160	38
160	30
160	37
180	37

Note that more than one ice crystal was measured at some time points. A file containing *R* code for assisting you in answering some the following questions is posted in the file *crystals.r*. A corresponding file with SAS code is posted as *crystals.sas*. SAS users should read the data from the file posted as *crystals.dat*.

- (a) Compute least squares estimates for the parameters in the model $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$, where $\epsilon_i \sim NID(0, \sigma^2)$. This notation means that the random errors (and the observations) have normal distributions and satisfy the Gauss-Markov property. Report the estimates and their standard errors.

- (b) Define

$$X = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \quad \text{and} \quad \mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \text{and} \quad P_X = X(X^T X)^{-1} X^T \quad \text{and} \quad P_1 = \mathbf{1}(\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T.$$

Here X is the model matrix for the model in part (a). What is the distribution of the quadratic form

$$R(\beta_1 | \beta_0) = \frac{1}{\sigma^2} \mathbf{Y}^T (P_X - P_1) \mathbf{Y}?$$

- (c) For the model in part (a),

$$SS_{residuals} = \frac{1}{\sigma^2} \mathbf{Y}^T (I - P_X) \mathbf{Y}?$$

has a central chi-square distribution with $(n - 2) = 14$ degrees of freedom. Define

$$MS_{residuals} = \frac{SS_{residuals}}{n - 2}$$

and use the results from part (b) to derive the distribution of

$$F = \frac{R(\beta_1 | \beta_0)}{MS_{residuals}}.$$

Report degrees of freedom and a formula for the noncentrality parameter.

- (d) What is the null hypothesis associated with the F statistic in part (c)? Justify your answer by showing that the noncentrality parameter in part (c) is zero if and only if the null hypothesis is true.
- (e) Report the value of the test statistics in part (c) and state your conclusion.
- (f) Examine the plot of the estimated line and observations and the residual plots provided by the code posted. What do these plots suggest?

4. Suppose the model proposed in part (a) of problem 3 is incorrect. In particular, suppose that the correct model is

$$Y_i = \lambda_1 + \lambda_2 X_i + \lambda_3 X_i^2 + \eta_i,$$

where $\eta_i \sim NID(0, \omega^2)$. This model can be expressed in matrix notation as $\mathbf{Y} = Z\boldsymbol{\lambda} + \boldsymbol{\eta}$, where

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad Z = \begin{bmatrix} 1 & X_1 & X_1^2 \\ 1 & X_2 & X_2^2 \\ \vdots & \vdots & \vdots \\ 1 & X_n & X_n^2 \end{bmatrix} = \begin{bmatrix} X & \mathbf{d} \end{bmatrix}, \quad \mathbf{d} = \begin{bmatrix} X_1^2 \\ X_2^2 \\ \vdots \\ X_n^2 \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\lambda} = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix}.$$

Assuming that this is the correct model for ice crystal growth, find

- (a) The distribution of $\mathbf{Y}^T (I - P_X) \mathbf{Y}$, where X is the model matrix from problem 3.
- (b) The distribution of $\mathbf{Y}^T (P_X - P_1) \mathbf{Y}$.
- (c) Does

$$\frac{\mathbf{Y}^T (P_X - P_1) \mathbf{Y}}{\mathbf{Y}^T (I - P_X) \mathbf{Y} / (n - 2)}$$

have an F -distribution? Explain.

5. Now, suppose that the model in problem 3 is correct, i.e. $\lambda_3 = 0$ and $\eta_i \sim NID(0, \sigma^2)$ for the model in problem 4.

- (a) Find the distribution of $\mathbf{Y}^T (I - P_Z) \mathbf{Y}$, where $P_Z = Z(Z^T Z)^{-1} Z^T$ and Z is the model matrix from problem 4.

(b) Does

$$\frac{\mathbf{Y}^T(P_X - P_1)\mathbf{Y}}{\mathbf{Y}^T(I - P_Z)\mathbf{Y}/(n - 3)}$$

have an F -distribution when the model in problem 3 is correct? Explain.

6. Problems 4 and 5 illustrate some of the consequences of incorrectly specifying the model. When you have replication at some sets of values of the explanatory variables, as we do for the ice crystal data, you can construct a lack-of-fit test for a proposed model. We will apply a lack-of-fit test to the quadratic model from problem 3. Consider the larger model

$$Y_{ij} = \gamma_0 + \gamma_1 X_i + \gamma_2 X_i^2 + \alpha_j + \epsilon_{ij},$$

where $\epsilon_{ij} \sim NID(0, \tau^2)$ and Y_{ij} denotes the observed axial length for the j -th ice crystal measure at time X_i . This model can be expressed in matrix notation as $\mathbf{Y} = W\boldsymbol{\gamma} + \boldsymbol{\epsilon}$, where

$$\mathbf{Y} = \begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \\ Y_{33} \\ Y_{41} \\ Y_{42} \\ Y_{43} \\ Y_{51} \\ Y_{52} \\ Y_{61} \\ Y_{62} \\ Y_{63} \\ Y_{71} \end{bmatrix}, \quad W = \begin{bmatrix} 1 & X_1 & X_1^2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & X_1 & X_1^2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & X_2 & X_2^2 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & X_2 & X_2^2 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & X_3 & X_3^2 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & X_3 & X_3^2 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & X_3 & X_3^2 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & X_4 & X_4^2 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & X_4 & X_4^2 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & X_4 & X_4^2 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & X_5 & X_5^2 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & X_5 & X_5^2 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & X_6 & X_6^2 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & X_6 & X_6^2 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & X_6 & X_6^2 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & X_7 & X_7^2 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad \boldsymbol{\gamma} = \begin{bmatrix} \gamma_0 \\ \gamma_1 \\ \gamma_2 \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \alpha_4 \\ \alpha_5 \\ \alpha_6 \\ \alpha_7 \end{bmatrix}$$

Note that the first three columns of W comprise the Z matrix from problem 4. Define $P_Z = Z(Z^T Z)^{-1} Z^T$ and $P_W = W(W^T W)^{-1} W^T$.

(a) Consider the test statistic

$$\frac{\mathbf{Y}^T(P_W - P_Z)\mathbf{Y}/(7 - 3)}{\mathbf{Y}^T(I - P_W)\mathbf{Y}/(n - 7)}.$$

Report a formula for the noncentrality parameter for the distribution of this statistic and use it to show that this is an appropriate lack-of-fit test.

(b) Would it be better to use $\mathbf{Y}^T(I - P_W)\mathbf{Y}/(n - 7)$ in the denominator of this test statistic instead of $\mathbf{Y}^T(I - P_Z)\mathbf{Y}/(n - 3)$? Explain.