

**Assignment 3** (3 pages)  
Statistics 32950-24620 (Spring 2023)  
Due 9 am Tuesday, April 11th.

Requirements

- Your answers should be typed or clearly written, started with your name, Assignment 3, STAT 24620 or 32950; saved as LastnameFirstnamePset3.pdf. Make sure to upload to Gradescope under the correct section: 246Pset3 or 329Pset3.
- When you use R (or others) to solve problems, select only relevant parts of the output, edit, then insert in your writing.
- You may discuss approaches with others. However the assignment should be devised and written by yourself. Capturing contents from other sources then pasting as your answers are not allowed.

**Reference:** Chapters 5, 6, and 10 of the text by Johnson and Wichern.

**Problem assignments:**

1. (*Hands on Hotelling's  $T^2$ , small data set*)

- (a) The data consist of four observations  $\mathbf{x}'_j = (x_{j1}, x_{j2})$ : (2, 12), (8, 9), (6, 9), (8, 10). Let  $\boldsymbol{\mu}_o = \begin{bmatrix} 7 \\ 11 \end{bmatrix}$ . In the following, leave your results in integer or fraction form (instead of approximated by decimals).
- Calculate the sample mean vector  $\bar{\mathbf{x}}$ .
  - Find sample covariance matrix  $\mathbf{S}$ .
  - Obtain  $\mathbf{S}^{-1}$ .
  - Evaluate Hotelling's  $T^2$ .
  - Specify the distribution of  $T^2$  under  $H_o : \boldsymbol{\mu} = \boldsymbol{\mu}_o$ .
- (b) Let  $C = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$ . Transform the data  $\{\mathbf{x}_j\}$  in (a) into  $\mathbf{y}_j = C\mathbf{x}_j$ . Let  $\boldsymbol{\mu}_o^* = C\boldsymbol{\mu}_o = \begin{bmatrix} 18 \\ 4 \end{bmatrix}$ .
- Calculate the sample mean vector  $\bar{\mathbf{y}}$ .
  - Derive the new sample covariance matrix  $S_y = C\mathbf{S}C'$ .
  - Evaluate Hotelling's  $T^2$  for  $\{\mathbf{y}_j\}$  under  $H_o : \boldsymbol{\mu}_y = \boldsymbol{\mu}_o^*$ .
- (c) Prove that, in general, if  $C$  is  $p \times p$  invertible matrix, then the transformed data  $\mathbf{y}_j = C\mathbf{x}_j$  has the same Hotelling's  $T^2$  statistic (under  $H_o : \boldsymbol{\mu}_y = C\boldsymbol{\mu}_o$ ).

2. (*Hands on canonical correlation analysis*)

(Note: "Hands on" means providing step details, using R for computing is allowed.)

The  $2 \times 1$  random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  have joint mean vector  $\boldsymbol{\mu}$  and joint covariance matrix  $\boldsymbol{\Sigma}$ ,

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \dots \\ \mu_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 2 \\ \dots \\ 0 \\ 1 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \text{ with } \Sigma_{11} = \begin{bmatrix} 8 & 2 \\ 2 & 5 \end{bmatrix}, \Sigma_{22} = \begin{bmatrix} 6 & -2 \\ -2 & 7 \end{bmatrix}, \Sigma'_{21} = \Sigma_{12} = \begin{bmatrix} 3 & 1 \\ -1 & 3 \end{bmatrix}.$$

- (a) Let  $\rho_1^*$  be the largest canonical correlation between  $\mathbf{X}$  and  $\mathbf{Y}$ . Then  $\rho_1^{*2}$  is a common eigenvalue of several matrices. Write out three such matrices in terms of the  $\Sigma_{ij}$ 's.

- (b) Calculate the canonical correlation  $\rho_1^*$  (the largest),  $\rho_2^*$  (the second largest).

(Hint: Choose an easy one for the calculation based on previous question.)

- (c) Determine the canonical variate pairs  $(U_1, V_1)$  and  $(U_2, V_2)$  corresponding to  $\rho_1^*$  and  $\rho_2^*$ .
- (d) Let  $\mathbf{U} = [U_1, U_2]'$ ,  $\mathbf{V} = [V_1, V_2]'$ . Evaluate

$$E\left(\begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}\right), \quad Cov\left(\begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}\right) = \begin{bmatrix} \Sigma_{UU} & \Sigma_{UV} \\ \Sigma_{VU} & \Sigma_{VV} \end{bmatrix}$$

- (e) Comment on the correlation structure between and within  $\mathbf{U}$ ,  $\mathbf{V}$ .

### 3. (Practice canonical correlation analysis, low dimensions)

Download the data from <https://www.stat.uchicago.edu/~meiwan/courses/s23-mva/stiffness.DAT>

(R command: `stiff = read.table("stiffness.DAT")`)

The data were obtained by taking four different measures of stiffness,  $x_1, x_2, x_3$  and  $x_4$  of each of  $n = 30$  boards. The first measurement involves sending a shock wave down the board, the second measurement is determined while vibrating the board, and the last two measurements are obtained from static tests. The squared distances  $d_j^2 = (\mathbf{x}_j - \bar{\mathbf{x}})' \mathbf{S}^{-1} (\mathbf{x}_j - \bar{\mathbf{x}})$  are also included as the last column in the data. (ref. Table 4.3 in J&W)

Let  $\mathbf{X} = [X_1, X_2]'$  be the vector of variables representing the dynamic measures of stiffness, and let  $\mathbf{Y} = [X_3, X_4]'$  be the vector of variables representing the static measures of stiffness.

- (a) Perform a canonical correlation analysis of these data. (R command: `cancor(X,Y)`)
- (b) Write the first canonical variates  $U_1$  and  $V_1$  as linear combinations of the components of  $\mathbf{X}$  and  $\mathbf{Y}$  respectively.
- (c) Produce two scatterplots of the data: one in the coordinate plane of the first canonical variate pair  $(U_1, V_1)$ , one in the plane of the second pair  $(U_2, V_2)$ .
- (d) Based on the two plots and the values of the canonical correlations  $\{\rho_1^*, \rho_2^*\}$ , comment on the correlation structure “captured” by each canonical pair.

### 4. (Multivariate vs univariate inference)

Input the dataset <https://www.stat.uchicago.edu/~meiwan/courses/s23-mva/fly.dat>

of 15 observations on  $X_1 = \text{antenna length}$  (mm) and  $X_2 = \text{wing length}$  (mm) of two species of flies.

R command for data input:

`fly = read.table("fly.dat")`

Define two univariate variables  $Y_1 = \text{antenna length} + \text{wing length}$ ,  $Y_2 = \text{wing length}$ .

Treat the data as bivariate samples from two populations (species **Af** and **Apf**) with equal covariance.

- (a) (Property of Hotelling's  $T^2$ )
- Compute a Hotelling's  $T^2$ -statistic for the hypothesis of equality of the mean vectors in the two species based on  $(Y_1, Y_2)$ . Is the hypothesis of equality of the means accepted?
  - Should you get the same results if you use the original variables  $(X_1, X_2)$ ? Why?
- (b) If you conduct (univariate) two-sample  $t$ -tests at a test level  $\alpha = 0.05$  performed on each of the variables  $Y_1$  and  $Y_2$  separately (i.e. assuming independence of  $Y_1$  and  $Y_2$ ), would the hypothesis of equality of species means be accepted? What if the test level  $\alpha = 0.01$ ?

- (c) Draw a scatterplot of  $Y_1$  vs.  $Y_2$  for the data in both species groups, marking the data points of the two species groups with different symbols, and explain how it can happen that (a) and (b) have different conclusions.
- (d) (*Comparison: Confidence region vs simultaneous confidence intervals*)
- Draw a 98% confidence ellipse for the species mean differences of  $Y_1$  and  $Y_2$  based on Hotelling's  $T^2$ .
  - In the same graph, draw a rectangle corresponding to univariate (marginal) 99% confidence interval for the mean differences of  $Y_1$  and  $Y_2$ .
  - Explain that the rectangle is a 98% confidence region by Bonferroni method.
  - Is the zero vector  $\mathbf{0} = (0, 0)$  in any of the two regions (ellipse by i, rectangle by ii)? Compare and comment on the goodness of the two regions. Which one is better?

5. (*Derivation in canonical correlation analysis*)

The  $2 \times 1$  random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  have joint covariance matrix  $\Sigma$ ,

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \text{ with } \Sigma_{11} = \begin{bmatrix} 1 & \rho_x \\ \rho_x & 1 \end{bmatrix}, \Sigma_{22} = \begin{bmatrix} 1 & \rho_y \\ \rho_y & 1 \end{bmatrix}, \Sigma_{21} = \Sigma_{12} = \begin{bmatrix} r & r \\ r & r \end{bmatrix}.$$

where  $\rho_x, \rho_y, r \in (0, 1)$ .

- Derive  $\rho_1^*$ , the largest canonical correlation between  $\mathbf{X}$  and  $\mathbf{Y}$ . Show your work.
- Derive the canonical variate pairs  $(U_1, V_1) = (\mathbf{a}_1' \mathbf{X}, \mathbf{b}_1' \mathbf{Y})$  corresponding to  $\rho_1^*$ , with normalization  $\mathbf{a}_1' \Sigma_{11} \mathbf{a}_1 = 1, \mathbf{b}_1' \Sigma_{22} \mathbf{b}_1 = 1$ .

6. (*Derivations for multivariate data and random variables*)

- Let  $\mathbf{X} = \begin{bmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_n \end{bmatrix}$  be an  $n \times p$  data matrix of  $n$  observations, the  $j$ th observation is  $\mathbf{x}_j \in \mathbb{R}^p, j = 1, \dots, n$ .

The sample covariance matrix can be expressed as  $\mathbf{S} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$ . Show that

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}' \mathbf{H} \mathbf{X}, \quad \text{with } \mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n'$$

where  $\mathbf{I}_n$  is the identity matrix of dimension  $n \times n$ , and  $\mathbf{1}_n$  is the  $n$ -vector with all elements = 1.

- Let  $\mathbf{W} = \mathbf{A} \mathbf{Y} + \mathbf{c}$ , where  $\mathbf{Y}$  is a  $p$ -dimensional random vector,  $\mathbf{A}$  is a fixed, scalar matrix of dimension  $k \times p$ , and  $\mathbf{c} \in \mathbb{R}^k$  is a fixed vector. Show that  $\text{Cov}(\mathbf{W}) = \mathbf{A} \text{Cov}(\mathbf{Y}) \mathbf{A}'$ .
- (**For 32950 only.** Optional for 24620.) Let  $\mathbf{Y}$  be a  $p$ -dimensional random vector and  $\mathbf{W}$  be a  $q$ -dimensional random vector,  $\mathbf{a}, \mathbf{b}$  are fixed vectors with dimensions  $p$  and  $q$  respectively. Show that

$$\text{Cov}(\mathbf{a}' \mathbf{Y}, \mathbf{b}' \mathbf{W}) = \mathbf{a}' \text{Cov}(\mathbf{Y}, \mathbf{W}) \mathbf{b} = \mathbf{b}' \text{Cov}(\mathbf{W}, \mathbf{Y}) \mathbf{a}$$