MATH377201

This question paper consists of 4
printed pages, each of which is
identified by the reference MATH377201.

New Cambridge Elementary
Statistical Tables are provided.
Only approved basic scientific
calculators may be used.

# ©UNIVERSITY OF LEEDS

Examination for the Module MATH3772
(January 2003)

## MULTIVARIATE ANALYSIS

Time allowed: **2 hours**

Attempt not more than THREE questions.
All questions carry equal marks.

1. The $p \times 1$ random vector x has a multivariate normal distribution with probability density function

$$f(\mathbf{x}) = |2\pi\Sigma|^{-1/2} \exp\left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}, \qquad \mathbf{x} \in \mathbb{R}^p,$$

and moment generating function

$$M_{\mathbf{x}}(\mathbf{t}) = \exp\left\{ \mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2}\mathbf{t}^T \Sigma \mathbf{t} \right\}.$$

The matrix $\Sigma$ is non-singular. We write $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$.

(a) Let x be partitioned into $p_1$ and $p_2$ components, $p_1 + p_2 = p$, with corresponding partitions

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \qquad \text{and} \qquad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Using the moment generating function, show that the marginal distribution of $\mathbf{x}_1$ is $N_{p_1}(\boldsymbol{\mu}_1, \Sigma_{11})$.

(b) Consider $\mathbf{y} = \mathbf{x}_1 + M\mathbf{x}_2$, where $M$ is a $p_1 \times p_2$ matrix. Show that

$$\mathrm{Cov}(\mathbf{y}, \mathbf{x}_2) = \Sigma_{12} + M\Sigma_{22}.$$

What value of $M$ results in independence between y and $\mathbf{x}_2$?

(c) Suppose

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N_2\left( \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 4 \end{pmatrix} \right).$$

Construct a new random vector $\mathbf{z} = (z_1, z_2)^T$ with $z_1 = 2x_1 + x_2$, $z_2 = x_1 - x_2$. Find the variance matrix and correlation matrix of z.

**CONTINUED...**

**2.** **(a)** Let $x_1, \ldots, x_n$ be a random sample from $N_p(\boldsymbol{\mu}_x, \Sigma)$ and $y_1, \ldots, y_m$ be a random sample from $N_p(\boldsymbol{\mu}_y, \Sigma)$. The two samples are independent of one another. Show that the union intersection test of the hypothesis $H_0 : \boldsymbol{\mu}_x = \boldsymbol{\mu}_y$ vs. $H_1 : \boldsymbol{\mu}_x \neq \boldsymbol{\mu}_y$, where $\Sigma$ is unknown, leads to the test statistic

$$T^2 = \frac{nm}{n+m}(\bar{\mathbf{x}} - \bar{\mathbf{y}})^T S_p^{-1}(\bar{\mathbf{x}} - \bar{\mathbf{y}}),$$

where $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are sample mean vectors and $S_p$ is the pooled within groups estimate of $\Sigma$.

**(b)** If $H_0$ is rejected in the overall test, how can simultaneous confidence intervals be used to give insight into the reasons for rejection?

**(c)** In a study to compare male gorilla skulls with female skulls the following variables were measured:

variable 1 = braincase length, variable 2 = braincase height.

A sample of 11 males and 11 females yielded the following summary statistics:

$$\bar{\mathbf{x}} = \begin{pmatrix} 152 \\ 103 \end{pmatrix}, \qquad S_x = \begin{pmatrix} 40 & 10 \\ 10 & 25 \end{pmatrix},$$

$$\bar{\mathbf{y}} = \begin{pmatrix} 142 \\ 101 \end{pmatrix}, \qquad S_y = \begin{pmatrix} 32 & 4 \\ 4 & 17 \end{pmatrix}.$$

The pooled covariance matrix for the two samples and its inverse are given by

$$S_p = \begin{pmatrix} 36 & 7 \\ 7 & 21 \end{pmatrix}, \qquad S_p^{-1} = \begin{pmatrix} 0.03 & -0.01 \\ -0.01 & 0.05 \end{pmatrix}.$$

Explain how $S_p$ is calculated. Compare the sexes on the basis of the information provided.

[Hints:

**1.** You may use the fact that the Hotelling $T^2$ and $F$ distribution are related by $T^2(p, \nu) = \{\nu p / (\nu - p + 1)\} F(p, \nu - p + 1)$.

**2.** Simultaneous $100\alpha$ percent confidence intervals for this problem can be written in the form

$$(\mathbf{a}^T(\bar{\mathbf{x}} - \bar{\mathbf{y}}) - c, \mathbf{a}^T(\bar{\mathbf{x}} - \bar{\mathbf{y}}) + c)$$

where $c = \left\{ T_\alpha^2(p, \nu)\frac{n+m}{nm}\mathbf{a}^T S_p \mathbf{a} \right\}^{\frac{1}{2}}$ and $T_\alpha^2(p, \nu)$ is the $100\alpha$ percentage point of the $T^2(p, \nu)$ distribution.]

**3.** **(a)** Let x be a $p$-dimensional random vector with mean vector $\mu$ and covariance matrix $\Sigma$.

    **(i)** Define the principal components y of x in terms of the standardized eigenvectors of $\Sigma$.

    **(ii)** Obtain the variance-covariance matrix of the principal components y.

    **(iii)** If $\Sigma = \alpha\alpha^T$ for some vector $\alpha$, find the first principal component. What can you say about the other principal components?

**(b)** Data were collected on 50 irises from the species *Iris setosa*. The variables are: $x_1$=sepal length; $x_2$=sepal width; $x_3$=petal length; $x_4$=petal width. The sample mean vector and correlation matrix were

$$x = \begin{pmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{pmatrix}, \qquad R = \begin{pmatrix} 1 & 0.74 & 0.27 & 0.28 \\ 0.74 & 1 & 0.18 & 0.23 \\ 0.27 & 0.18 & 1 & 0.33 \\ 0.28 & 0.23 & 0.33 & 1 \end{pmatrix}.$$

Principal component analysis gave eigenvalues 2.06, 1.02, 0.67, 0.25. The corresponding eigenvectors were the columns of

$$\begin{pmatrix} 0.60 & -0.33 & 0.07 & 0.72 \\ 0.58 & -0.44 & 0.00 & -0.69 \\ 0.38 & 0.63 & 0.68 & -0.09 \\ 0.40 & 0.55 & -0.73 & -0.01 \end{pmatrix}.$$

Interpret these principal components briefly. Assess their relative contribution to total variation. What methods can be used to decide on the number of components we should retain?

4. (a) Let $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^p$, denote two probability density functions for populations $\Pi_1$ and $\Pi_2$ respectively, with prior probabilities $\pi_1$ and $\pi_2$, where $\pi_1 + \pi_2 = 1$. Consider the allocation rule which assigns $\mathbf{x}$ to $\Pi_1$ if

$$\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{\pi_2}{\pi_1}$$

and to $\Pi_2$ otherwise. Show that this rule is admissible.

(b) If the two populations have $N_p(\boldsymbol{\mu}_1, \Sigma)$ and $N_p(\boldsymbol{\mu}_2, \Sigma)$ distributions and $\pi_1 = 2\pi_2$, show that this rule classifies $\mathbf{x}$ as coming from the first population if

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\Sigma^{-1}\mathbf{x} \geq c$$

and to the second population otherwise, for a suitable constant $c$. What is the value of $c$?

(c) Let

$$\Sigma = \begin{pmatrix} 2 & 1 \\ 1 & 5 \end{pmatrix}, \qquad \boldsymbol{\mu}_1 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \qquad \boldsymbol{\mu}_2 = \begin{pmatrix} 5 \\ 1 \end{pmatrix}.$$

Calculate and sketch the boundary between the two classification regions for $\pi_1 = 2\pi_2$ and $\pi_1 = \pi_2$ respectively. Compare these boundaries. How would you classify $\mathbf{x} = (3,2)^T$ under each rule?

END