

Multivariate Analysis - Homework 2

1. Consider the two covariance matrices Σ_1 and Σ_2 of two random vectors \mathbf{y}_1 and \mathbf{y}_2 , respectively, where

$$\Sigma_1 = \begin{pmatrix} 14 & 8 & 3 \\ 8 & 5 & 2 \\ 3 & 2 & 1 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 6 & 6 & 1 \\ 6 & 8 & 2 \\ 1 & 2 & 1 \end{pmatrix}$$

Compute their generalized variance and total variances, compare them, and explain why it is the case.

2. Suppose that $\mathbf{y} = (Y_1, Y_2, Y_3)' \sim N_3(\mathbf{0}, \Sigma)$, where

$$\Sigma = \begin{pmatrix} 1 & \rho & 0 \\ \rho & 1 & \rho \\ 0 & \rho & 1 \end{pmatrix}$$

Is there a value of ρ for which $Y_1 + Y_2 + Y_3$ and $Y_1 - Y_2 - Y_3$ are independent? Prove or disprove.

3. Suppose \mathbf{y} is $N_3(\boldsymbol{\mu}, \Sigma)$, where

$$\boldsymbol{\mu} = \begin{pmatrix} 3 \\ 2 \\ 4 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 6 & 1 & -2 \\ 1 & 13 & 4 \\ -2 & 4 & 4 \end{pmatrix}$$

- (a) Find the distribution of $Z = 2Y_1 - Y_2 + 3Y_3$.
 - (b) Find the joint distribution of $Z_1 = Y_1 + Y_2 + Y_3$ and $Z_2 = Y_1 - Y_2 + 2Y_3$.
 - (c) Find the distribution of Y_2 .
 - (d) Find the joint distribution of Y_1 and Y_3 .
 - (e) Find the joint distribution of Y_1 , Y_3 and $\frac{1}{2}(Y_1 + Y_3)$.
 - (f) Find the conditional distribution of (Y_1, Z_1) given (Y_2, Z_2) .
4. Let $Y_1 \sim N(0, 1)$, and let

$$Y_2 = \begin{cases} -Y_1 & \text{if } -1 \leq Y_1 \leq 1 \\ Y_1 & \text{otherwise.} \end{cases}$$

Show each of the following.

- (a) Y_2 also has an $N(0, 1)$ distribution. (Hint: Compute the CDF of Y_2 .)
- (b) Y_1 and Y_2 do not have a bivariate normal distribution.
5. Suppose $\mathbf{y} = (\mathbf{y}'_1, \mathbf{y}'_2)' \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $|\boldsymbol{\Sigma}| \neq 0$; $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are partitioned accordingly.
- (a) Check that
- $$\begin{aligned}
 (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) &= \{\mathbf{y}_1 - \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{y}_2 - \boldsymbol{\mu}_2)\}' \\
 &\quad \times (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21})^{-1} \{\mathbf{y}_1 - \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{y}_2 - \boldsymbol{\mu}_2)\} \\
 &\quad + (\mathbf{y}_2 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{y}_2 - \boldsymbol{\mu}_2)
 \end{aligned}$$
- (b) Derive the conditional density $f(\mathbf{y}_1 | \mathbf{y}_2)$ using the result in (a), rather than the decorrelation method, and verify that it is still normal.
6. For one sample case, prove that the likelihood ratio test leads to Hotelling's T^2 test for multivariate normal samples.
7. (R exercise.) The world's 10 largest companies (2005 database) yield the following data (also attached as .txt file):

The World's 10 Largest Companies ¹			
Company	x_1 = sales (billions)	x_2 = profits (billions)	x_3 = assets (billions)
Citigroup	108.28	17.05	1,484.10
General Electric	152.36	16.59	750.33
American Intl Group	95.04	10.91	766.42
Bank of America	65.45	14.14	1,110.46
HSBC Group	62.97	9.52	1,031.29
ExxonMobil	263.99	25.33	195.26
Royal Dutch/Shell	265.19	18.54	193.83
BP	285.06	15.73	191.11
ING Group	92.01	8.10	1,175.16
Toyota Motor	165.68	11.13	211.15

¹From www.Forbes.com partially based on *Forbes* The Forbes Global 2000, April 18, 2005.

For all the three variables:

- (a) Construct individual QQ plots to investigate univariate normality. Interpret the output.
- (b) Conduct formal statistical tests for the individual normality. Explain the results.

- (c) Check the multivariate normality of $(X_1, X_2, X_3)'$ using the pairwise scatter plot matrix and the χ^2 QQ plot.
8. (R exercise) Recall the relationship between the hypothesis testing and the confidence interval, i.e. the conclusion of a test can be directly obtained from the related confidence interval. For the multivariate case, the confidence interval becomes the “confidence region”.
- (a) Analogous to the definition of confidence interval, define the $1 - \alpha$ confidence region for the population mean vector $\boldsymbol{\mu}$. And derive the mathematical formula of this confidence region when the covariance matrix $\boldsymbol{\Sigma}$ is unknown. Suppose the sample is $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, each \mathbf{y}_i is p -variate, the sample mean vector is $\bar{\mathbf{y}}$ and the sample covariance matrix is \mathbf{S} .
- (b) For the sweat data (attached as sweat.dat), suppose we only have the information of the first two variables with mean μ_1 and μ_2 . Find the 95% confidence region for $\boldsymbol{\mu} = (\mu_1, \mu_2)'$.
- (c) Describe the confidence region geometrically using the eigenvalue and eigenvectors of the sample covariance matrix \mathbf{S} .
- (d) Construct the 95% univariate confidence interval for each variable.
- (e) Consider the test $H_0: \boldsymbol{\mu} = \boldsymbol{\mu}_0$. Give an example of $\boldsymbol{\mu}_0$ such that the multivariate test rejects H_0 but both univariate tests fail to do so. You should answer this question based on the confidence region and confidence intervals obtained in (b) and (d).