

Multivariate Analysis - Homework 1

1. Prove that the sample covariance matrix \mathbf{S} can be written as

$$\mathbf{S} = \frac{1}{n-1} \mathbf{Y}'(\mathbf{I} - \frac{1}{n} \mathbf{J})\mathbf{Y} = \frac{1}{n-1} (\mathbf{Y} - \bar{\mathbf{Y}})'(\mathbf{Y} - \bar{\mathbf{Y}}),$$

where \mathbf{Y} is the data matrix, \mathbf{I} is the n -dimensional identity matrix, and \mathbf{J} is the $n \times n$ matrix with all elements 1's, and $\bar{\mathbf{Y}}$ is the sample mean matrix with all rows to be the sample mean vector $\bar{\mathbf{y}}$.

2. Let matrix \mathbf{A} be

$$\mathbf{A} = \begin{bmatrix} 4 & 8 & 8 \\ 3 & 6 & -9 \end{bmatrix}$$

Compute

- (a) $\mathbf{A}\mathbf{A}'$, and its eigenvalues and eigenvectors.
 - (b) $\mathbf{A}'\mathbf{A}$, and its eigenvalues and eigenvectors. Compare the nonzero eigenvalues between (a) and (b).
 - (c) Obtain the spectral decomposition of $\mathbf{A}\mathbf{A}'$.
 - (d) Self-study the definition of singular value decomposition of a matrix. Obtain the singular value decomposition of \mathbf{A} .
3. Suppose the random vector $\mathbf{y} = (Y_1, Y_2, Y_3, Y_4)$ with mean vector $\boldsymbol{\mu} = (4, 3, 2, 1)'$ and covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} 3 & 0 & 2 & 2 \\ 0 & 1 & 1 & 0 \\ 2 & 1 & 9 & -2 \\ 2 & 0 & -2 & 4 \end{bmatrix}$$

Partition \mathbf{y} as $\mathbf{y}^{(1)} = (Y_1, Y_2)'$ and $\mathbf{y}^{(2)} = (Y_3, Y_4)'$. Furthermore, let matrix

$$\mathbf{A} = [1, 2] \text{ and } \mathbf{B} = \begin{bmatrix} 1 & -2 \\ 2 & -1 \end{bmatrix}$$

and consider the linear combinations $\mathbf{A}\mathbf{y}^{(1)}$ and $\mathbf{B}\mathbf{y}^{(2)}$. Find

- (a) $E(\mathbf{y}^{(1)})$
- (b) $E(\mathbf{A}\mathbf{y}^{(1)})$
- (c) $COV(\mathbf{y}^{(1)})$

(d) $COV(\mathbf{A}\mathbf{y}^{(1)})$

(e) $E(\mathbf{B}\mathbf{y}^{(2)})$

(f) $COV(\mathbf{B}\mathbf{y}^{(2)})$

(g) $COV(\mathbf{y}^{(1)}, \mathbf{y}^{(2)})$

(h) $COV(\mathbf{A}\mathbf{y}^{(1)}, \mathbf{B}\mathbf{y}^{(2)})$

4. (R exercise.) The following table (data attached) gives partial data from three variables measured in milliequivalents per 100g:

y_1 = available soil calcium,

y_2 = exchangeable soil calcium,

y_3 = turnip green calcium.

Table. Calcium in Soil and Turnip Greens

Location Number	y_1	y_2	y_3
1	35	3.5	2.80
2	35	4.9	2.70
3	40	30.0	4.38
4	10	2.8	3.21
5	6	2.7	2.73
6	20	2.8	2.81
7	35	4.6	2.88
8	35	10.9	2.90
9	35	8.0	3.28
10	30	1.6	3.20

Define

$$z_1 = y_1 + y_2 + y_3,$$

$$z_2 = 2y_1 - 3y_2 + 2y_3,$$

$$z_3 = -y_1 - 2y_2 - 3y_3.$$

- (a) Find the sample mean vector $\bar{\mathbf{z}}$, sample covariance matrix \mathbf{S}_z of $\mathbf{z} = (Z_1, Z_2, Z_3)'$
- (b) Find the sample correlation matrix \mathbf{R}_z from \mathbf{S}_z .
- (c) Find the generalized variance and total variance of \mathbf{z} .

- (d) Realize the spectral decomposition and Cholesky decomposition of both \mathbf{S}_z and \mathbf{R}_z , and get the square root matrix of them.
5. (R exercise.) The attached data are 42 measurements on air-pollution variables recorded at 12:00 noon in the Los Angeles area on different days.
- (a) Plot the pairwise scatter plot matrix for all the variables in R. And comment on the output.
 - (b) Construct the sample mean vector, sample covariance matrix and sample correlation matrix. Interpret the entries in the sample correlation matrix.
 - (c) Compute the Euclidean distance matrix and the Mahalanobis/statistical distance matrix among the first five days. Explain the advantage of the Mahalanobis distance.
 - (d) Describe the overall variability of the data.
 - (e) Get the Spectral decomposition and Cholesky decomposition of the sample covariance matrix. Observe the difference between the two decompositions.
 - (f) Obtain a 3-D scatter plot for any three variables that you think make sense. Use any package/command in R **except for the one given in the slides**.
6. (R exercise.) Generate 100 random pairs of numbers for the bivariate random vector $(X, Y)'$, and plot the scatterplots between X and Y under the following respective settings:
- (a) X and Y are positively correlated;
 - (b) X and Y are negatively correlated;
 - (c) X and Y are perfectly positive-correlated;
 - (d) X and Y are uncorrelated;
 - (e) X and Y are nonlinearly correlated.