

A CONSTRUCTION METHOD OF CERTAIN MATRICES REQUIRED IN THE MULTIVARIATE HETEROSCEDASTIC METHOD

HIROTO HYAKUTAKE

(Received Mar. 4, 1985; revised Apr. 25, 1985)

Summary

For statistical inference about several normal means, the heteroscedastic method was proposed by Dudewicz and Bishop (1979, *Optimizing Methods in Statistics*, Academic Press, 183-203). However, the practical application in the multivariate case was not possible because it had not been known how to construct the certain matrices required in the method. In this paper, a construction method of the matrices is given.

1. Introduction

We consider a statistical inference on means μ_1, \dots, μ_k of k p -variate normal populations $N_p(\mu_i, \Sigma_i)$, $\Sigma_i > 0$, $i=1, \dots, k$, where μ_i 's and Σ_i 's are unknown, and Σ_i 's are different. When $k=1$, Stein [7] gave the two-stage sampling scheme for the univariate case in which the power for a test, a confidence coefficient for a confidence interval, and etc. are completely free from the population variance. This sampling scheme was generalized to the multivariate case by Chatterjee [1]. Using this procedure, it is possible to give the heteroscedastic method proposed by Dudewicz and Bishop [2] to overcome the complexity arisen in many cases of comparison of several (k) normal populations with different variances or covariance matrices. This is available to multiple comparisons, construction of simultaneous confidence intervals, ranking and selection problems, and etc. However, in the multivariate case, it was not possible to use the heteroscedastic method in practice up to the present because although the existence of certain matrices was known, an algorithm for calculation of the matrices was not known. In this paper, we give the matrices required in the multivariate sampling procedure.

Keywords and phrases: Multivariate normal populations, heteroscedastic method, two-stage sampling scheme, matrix quadratic equation.

In the sampling procedure of the multivariate heteroscedastic method, we first take samples of size N_0 from each of k populations and compute

$$(1.1) \quad \bar{\mathbf{x}}^{(i)} = \frac{1}{N_0} \sum_{r=1}^{N_0} \mathbf{x}_r^{(i)}, \quad S_i = \frac{1}{\nu} \sum_{r=1}^{N_0} (\mathbf{x}_r^{(i)} - \bar{\mathbf{x}}^{(i)})(\mathbf{x}_r^{(i)} - \bar{\mathbf{x}}^{(i)})',$$

for $i=1, \dots, k$, where $\nu = N_0 - 1 \geq p$. Then for each i , we define N_i by

$$(1.2) \quad N_i = \max \{N_0 + p^2, [c \cdot \text{tr}(TS_i)] + 1\},$$

where c is a positive constant, T is a $p \times p$ given positive definite (p.d.) matrix, and $[y]$ denotes the greatest integer not greater than y . Note that we use the same c and T for all $i=1, \dots, k$. Next we take $N_i - N_0$ additional observations $\mathbf{x}_{N_0+1}^{(i)}, \dots, \mathbf{x}_{N_i}^{(i)}$ and construct the basic random vector \mathbf{z}_i ($i=1, \dots, k$) in the following way:

We choose p matrices $A_i^{(i)}: p \times N_i = [\mathbf{a}_{i1}^{(i)}, \dots, \mathbf{a}_{iN_0}^{(i)}, \mathbf{a}_{iN_0+1}^{(i)}, \dots, \mathbf{a}_{iN_i}^{(i)}]$, $i=1, \dots, k$, in such a way that

$$(I) \quad \mathbf{a}_{i1}^{(i)} = \dots = \mathbf{a}_{iN_0}^{(i)},$$

$$(II) \quad A_i^{(i)} \mathbf{j}_{N_i} = \mathbf{e}_i, \text{ where } \mathbf{j}_{N_i}: N_i \times 1 = (1, \dots, 1)' \text{ and } \mathbf{e}_i: p \times 1 = (0, \dots, 0, 1, 0, \dots, 0)',$$

$$(III) \quad A^{(i)} A^{(i)'} = (1/c) T^{-1} \otimes S_i^{-1}, \text{ where } A^{(i)}: p^2 \times N_i = [A_1^{(i)'}, A_2^{(i)'}, \dots, A_p^{(i)'}]'$$

Then we define

$$(1.3) \quad \mathbf{z}_i: p \times 1 = [\text{tr}(A_1^{(i)} X^{(i)}), \text{tr}(A_2^{(i)} X^{(i)}), \dots, \text{tr}(A_p^{(i)} X^{(i)})]',$$

where $X^{(i)}: p \times N_i = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{N_0}^{(i)}, \mathbf{x}_{N_0+1}^{(i)}, \dots, \mathbf{x}_{N_i}^{(i)}]$.

The heteroscedastic inference method is based on the vectors $\mathbf{z}_1, \dots, \mathbf{z}_k$. The distributions of \mathbf{z}_i and the statistics based on these random vectors are fully independent of the population covariance matrices. Hyakutake and Siotani [4] and Hyakutake, Siotani, Li and Mustafid [5] gave the asymptotic approximations for their distributions.

Our purpose is to find the matrices $A^{(i)} = [A_1^{(i)'}, \dots, A_p^{(i)'}]'$, $i=1, \dots, k$, which satisfy the conditions (I), (II) and (III), for practical application.

2. Main result

In this section, we give a method of constructing the matrices $A^{(i)} = [A_1^{(i)'}, \dots, A_p^{(i)'}]'$ which satisfy the conditions (I), (II) and (III). This construction is done independently for each population in the heteroscedastic method. Since a result from one population does not affect results from other populations, we can drop the suffix i without loss of generality.

First of all, we put $A = [A_1', \dots, A_p']' = [A_0; B]$, where $A_0: p^2 \times (N-m) = [\mathbf{a}_0, \dots, \mathbf{a}_0] = \mathbf{a}_0 \mathbf{j}_{N-m}'$, B is a $p^2 \times m$ matrix, and $p^2 \leq m \leq N - N_0$ by (I). From (II), $A \mathbf{j}_N = [A_0; B] \mathbf{j}_N = (N-m) \mathbf{a}_0 + B \mathbf{j}_m = \mathbf{e}$, hence

$$(2.1) \quad \mathbf{a}_0 = \frac{1}{N-m}(\mathbf{e} - B\mathbf{j}_m),$$

where $\mathbf{e}: p^2 \times 1 = (\mathbf{e}'_1, \dots, \mathbf{e}'_p)'$. Substituting (2.1) into (III), the condition (III) is written in the term of B as

$$(2.2) \quad \frac{1}{c}T^{-1} \otimes S^{-1} = B \left[I_m + \frac{1}{N-m} \mathbf{j}_m \mathbf{j}'_m \right] B' \\ - \frac{1}{N-m} \mathbf{e} \mathbf{j}'_m B' - \frac{1}{N-m} B \mathbf{j}_m \mathbf{e}' + \frac{1}{N-m} \mathbf{e} \mathbf{e}',$$

which is a matrix quadratic equation. Let G be an $m \times m$ matrix satisfying

$$(2.3) \quad GG' = I_m + \frac{1}{N-m} \mathbf{j}_m \mathbf{j}'_m.$$

In general, the existence of such a decomposition of a p.d. or positive semi-definite (p.s.d.) matrix and its algorithm are well known (see e.g. appendix in Muirhead [6]). The matrix G may be a lower triangular matrix or a symmetric matrix. A solution for symmetric matrix is given by $G = I_m - (1/m) \mathbf{j}_m \mathbf{j}'_m \pm (\sqrt{N}/m\sqrt{N-m}) \mathbf{j}_m \mathbf{j}'_m$. Letting $D = BG$, the equation (2.2) can be written in term of D as

$$(2.4) \quad (D - E)(D - E)' = \frac{1}{c}T^{-1} \otimes S^{-1} - \frac{1}{N} \mathbf{e} \mathbf{e}',$$

where

$$E = \frac{1}{N-m} \mathbf{e} \mathbf{j}'_m (GG')^{-1} G = \frac{1}{N-m} \mathbf{e} \mathbf{j}'_m \left(I_m - \frac{1}{N} \mathbf{j}_m \mathbf{j}'_m \right) G = \frac{1}{N} \mathbf{e} \mathbf{j}'_m G.$$

To show the existence of the solution in the equation (2.4), it is necessary that $(1/c)T^{-1} \otimes S^{-1} - (1/N)\mathbf{e}\mathbf{e}'$ is a p.d. or p.s.d. matrix. This is shown by considering the inverse matrix of $(1/c)T^{-1} \otimes S^{-1} - (1/N)\mathbf{e}\mathbf{e}'$, which is given by

$$(2.5) \quad cT \otimes S + (cT \otimes S)\mathbf{e}\mathbf{e}'(cT \otimes S)/(N - \mathbf{e}'(cT \otimes S)\mathbf{e}).$$

Since $c > 0$, T is a p.d. matrix, and S is a p.d. matrix with probability one, the first term of (2.5) is a p.d. matrix. The second term of (2.5) is a p.s.d. matrix clearly because of $\mathbf{e}'(cT \otimes S)\mathbf{e} = c \cdot \text{tr}(TS) < N$ by (1.2). Therefore the matrix in (2.5) is a p.d. matrix, and hence $(1/c)T^{-1} \otimes S^{-1} - (1/N)\mathbf{e}\mathbf{e}'$ is a p.d. matrix.

Now we can choose a matrix $K: p^2 \times m$ such that

$$(2.6) \quad KK' = \frac{1}{c}T^{-1} \otimes S^{-1} - \frac{1}{N} \mathbf{e} \mathbf{e}'.$$

Hence

$$(2.7) \quad B = EG^{-1} + KG^{-1} = \frac{1}{N} \mathbf{e} \mathbf{j}'_m + KG^{-1}$$

and from (2.1),

$$(2.8) \quad \begin{aligned} A_0 = \mathbf{a}_0 \mathbf{j}'_{N-m} &= \frac{1}{N-m} \left\{ \mathbf{e} - \left(\frac{1}{N} \mathbf{e} \mathbf{j}'_m + KG^{-1} \right) \mathbf{j}_m \right\} \mathbf{j}'_{N-m} \\ &= \frac{1}{N} \mathbf{e} \mathbf{j}'_{N-m} - \frac{1}{N-m} KG^{-1} \mathbf{j}_m \mathbf{j}'_{N-m}. \end{aligned}$$

Thus we obtain

$$(2.9) \quad A = \frac{1}{N} \mathbf{e} \mathbf{j}'_N + KG^{-1} \left[-\frac{1}{N-m} \mathbf{j}_m \mathbf{j}'_{N-m}; I_m \right].$$

We summarize the above argument in the following.

THEOREM. *A matrix A satisfying the conditions (I), (II) and (III) is given by (2.9), where $G: m \times m$ and $K: p^2 \times m$ are satisfying (2.3) and (2.6), respectively, and $p^2 \leq m \leq N - N_0$.*

If we take $G = I_m - (1/m) \mathbf{j}_m \mathbf{j}'_m + (\sqrt{N}/m\sqrt{N-m}) \mathbf{j}_m \mathbf{j}'_m$, the matrix A in (2.9) is written as

$$(2.10) \quad \begin{aligned} A = \frac{1}{N} \mathbf{e} \mathbf{j}'_N + K \left[-\frac{1}{\sqrt{N(N-m)}} \mathbf{j}_m \mathbf{j}'_{N-m}; I_m - \frac{1}{m} \mathbf{j}_m \mathbf{j}'_m \right. \\ \left. + \frac{\sqrt{N-m}}{m\sqrt{N}} \mathbf{j}_m \mathbf{j}'_m \right]. \end{aligned}$$

It is noted that the matrix A is not unique since K and G are not unique.

3. Numerical examples

To give numerical examples of the previous result in the bivariate case, i.e. $p=2$, we use a part of the Iris Data discussed by Fisher [3]. We take the first-stage samples of size $N_0=10$, then the sample covariance matrix is

$$S = \begin{pmatrix} .5289 & .1933 \\ .1933 & .1157 \end{pmatrix}.$$

Now, if we take $c=1$ and $T=I_2$, the total sample size N defined in (1.2) is 14. Hence the additional (second-stage) sample size $N-N_0$ is 4 and $m=4$. Then the total observation matrix is given by

$$X = \begin{pmatrix} 7.0 & 6.4 & 6.9 & 5.5 & 6.5 & 5.7 & 6.3 & 4.9 & 6.6 & 5.2 & 5.0 & 5.9 & 6.0 & 6.1 \\ 3.2 & 3.2 & 3.1 & 2.3 & 2.8 & 2.8 & 3.3 & 2.4 & 2.9 & 2.7 & 2.0 & 3.0 & 2.2 & 2.9 \end{pmatrix}.$$

Next if we take the lower triangular matrices K and G with positive diagonal satisfying (2.6) with $e = (1, 0, 0, 1)'$ and (2.3), respectively, the matrices satisfying (I), (II) and (III) are

$$A_1 = \begin{pmatrix} -.13723 \dots -.13723 & 2.1580 & .07143 & .07143 & .07143 \\ .10105 \dots .10105 & -3.7925 & 2.7820 & 0 & 0 \end{pmatrix},$$

$$A_2 = \begin{pmatrix} -.17652 \dots -.17652 & -.17651 & -.17651 & 2.1182 & 0 \\ .15619 \dots .15619 & .12196 & .10953 & -3.6860 & 2.8926 \end{pmatrix}.$$

Then the basic random vector is

$$z = (7.366, 5.352)'.$$

On the other hand, if we use the formula (2.10) with the same lower triangular matrix K as before, we obtain

$$A_1 = \begin{pmatrix} -.11276 \dots -.11276 & 2.1751 & -.03129 & -.03129 & -.03129 \\ .06790 \dots .06790 & -3.6812 & 2.9369 & .03123 & .03123 \end{pmatrix},$$

$$A_2 = \begin{pmatrix} -.18556 \dots -.18556 & -0.8535 & -.08535 & 2.1194 & -.08535 \\ .14143 \dots .14143 & .07098 & .06192 & -3.5814 & 3.0313 \end{pmatrix}.$$

Then

$$z = (7.314, 5.244)'.$$

From the above numerical results, we can see that the choice of a lower triangular or symmetric one as a solution of G will affect the basic random vectors, but little affect.

Acknowledgements

The author wishes to thank Professor M. Siotani and Professor Y. Fujikoshi, Hiroshima University, for their valuable comments and encouragement.

HIROSHIMA UNIVERSITY

REFERENCES

- [1] Chatterjee, S. K. (1959). On an extension of Stein's two-sample procedure to the multi-normal problem, *Calcutta Statist. Assoc. Bull.*, 8, 121-148.
- [2] Dudewicz, E. J. and Bishop, T. A. (1979). The heteroscedastic method, *Optimizing Methods in Statistics*, (ed. J. S. Rustagi), Academic Press, New York, 183-203.
- [3] Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems, *Ann.*

- Eugen.*, 7, 179-188.
- [4] Hyakutake, H. and Siotani, M. (1984). Distributions of some statistics in heteroscedastic method, *Technical Report*, No. 108, Statistical Research Group, Hiroshima University, Japan.
 - [5] Hyakutake, H., Siotani, M., Li, C. and Mustafid (1984). Distributions of some statistics in heteroscedastic inference method II: Tables of percentage points and power functions, *Technical Report*, No. 142, Statistical Research Group, Hiroshima University, Japan.
 - [6] Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*, John Wiley, New York.
 - [7] Stein, C. (1945). A two-sample test for a linear hypothesis whose power is independent of the variance, *Ann. Math. Statist.*, 16, 243-258.