

# Simultaneous testing of the mean vector and the covariance matrix with two-step monotone missing data

Miki Hosoya and Takashi Seo

(Received August 18, 2014; Revised November 15, 2014)

**Abstract.** In this paper, we consider the problem of simultaneous testing of the mean vector and the covariance matrix when the data have a two-step monotone pattern that is missing observations. We give the likelihood ratio test (LRT) statistic and propose an approximate upper percentile of the null distribution using linear interpolation based on an asymptotic expansion of the modified LRT statistic in the case of a complete data set. As another approach, we give the modified LRT statistics with a two-step monotone missing data pattern using the coefficient of the modified LRT statistic with complete data. Finally, we investigate the asymptotic behavior of the upper percentiles of these test statistics by Monte Carlo simulation.

*AMS 2010 Mathematics Subject Classification.* 62E20, 62H10.

*Key words and phrases.* Asymptotic expansion, linear interpolation, modified likelihood ratio test statistic, two-step monotone missing data.

## §1. Introduction

Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_1}$  be distributed as the  $p$ -dimensional normal distribution  $N_p(\boldsymbol{\mu}, \Sigma)$  and  $\mathbf{x}_{1, N_1+1}, \mathbf{x}_{1, N_1+2}, \dots, \mathbf{x}_{1, N}$  be distributed as the  $p_1$ -dimensional normal distribution  $N_{p_1}(\boldsymbol{\mu}_1, \Sigma_{11})$ , where

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

We partition  $\mathbf{x}_j$  into a  $p_1 \times 1$  random vector and a  $p_2 \times 1$  random vector as  $\mathbf{x}_j = (\mathbf{x}'_{1j}, \mathbf{x}'_{2j})'$ , where  $\mathbf{x}_{ij} : p_i \times 1, i = 1, 2, j = 1, 2, \dots, N_1$ .

Such a data set has two-step monotone missing data:

$$\left( \begin{array}{cc} \mathbf{x}'_{11} & \mathbf{x}'_{21} \\ \vdots & \vdots \\ \mathbf{x}'_{1N_1} & \mathbf{x}'_{2N_1} \\ \mathbf{x}'_{1,N_1+1} & * \cdots * \\ \vdots & \vdots \\ \mathbf{x}'_{1N} & * \cdots * \end{array} \right) \left. \vphantom{\begin{array}{c} \mathbf{x}'_{11} \\ \vdots \\ \mathbf{x}'_{1N_1} \\ \mathbf{x}'_{1,N_1+1} \\ \vdots \\ \mathbf{x}'_{1N} \end{array}} \right\}^{N_1}, \left. \vphantom{\begin{array}{c} \mathbf{x}'_{21} \\ \vdots \\ \mathbf{x}'_{2N_1} \\ * \cdots * \\ \vdots \\ * \cdots * \end{array}} \right\}^{N_2},$$

$$\underbrace{\hspace{1.5cm}}_{p_1} \quad \underbrace{\hspace{1.5cm}}_{p_2}$$

where  $N = N_1 + N_2$ ,  $p = p_1 + p_2$ ,  $N_1 > p$ , and “\*” indicates a missing observation.

Missing data is an important problem in statistical data analyses. A variety of statistical procedures to deal with missing data have been developed by many authors, including Anderson (1957), Bhargava (1962), McLachlan and Krishnan (1997), and Little and Rubin (2002). For a general missing pattern, Srivastava (1985) discussed the LRT for mean vectors in one-sample and two-sample problems. Seo and Srivastava (2000) derived a test of equality of means and the simultaneous confidence intervals for the monotone missing data in a one-sample problem. Anderson (1957) developed an approach to derive the MLEs of the mean vector and the covariance matrix by solving the likelihood equations for monotone missing data with several missing patterns. Anderson and Olkin (1985) derived the MLEs for two-step monotone missing data in a one-sample problem. For the related discussion of the MLEs in cases of general  $k$ -step monotone missing data, see Jinadasa and Tracy (1992) and Kanda and Fujikoshi (1998).

Further, by the use of the MLEs of the mean vector and the covariance matrix, the LRT statistic and Hotelling’s  $T^2$ -type statistic for tests of mean vectors with two or three-step monotone missing data has been discussed by Krishnamoorthy and Pannala (1999), Chang and Richards (2009), Seko et al. (2012), and Yagi and Seo (2014), among others. The problem of simultaneous testing of the mean and the variance under univariate and non-missing normality has been discussed by Choudhari et al. (2001) and Zhang et al. (2012). For non-missing and multivariate normality, Davis (1971) gave the modified LRT statistic (see Muirhead (1982) and Srivastava (2002)). In this paper, the LRT and modified LRT statistics are given under multivariate normality with a two-step monotone missing data pattern. Further, we assume that the data are missing completely at random (MCAR), see Hao and Krishnamoorthy (2001), and Little and Rubin (2002).

The remainder of this paper is organized as follows. In Section 2, we consider the case in which the missing observations are of the two-step monotone

type and provide an LRT statistic for the simultaneous testing of the mean vector and the covariance matrix. In Section 3, an approximation to the upper percentile of the LRT statistic and the modified LRT statistics are given. Finally, in Section 4, the accuracy of the approximation and the asymptotic behavior of modified statistics are investigated by Monte Carlo simulation.

## §2. Likelihood ratio test statistic

In order to derive the LRT statistic of the simultaneous testing of the mean vector and the covariance matrix in the case of a two-step monotone missing data pattern, we present their MLEs, which are given by

$$(2.1) \quad \hat{\boldsymbol{\mu}} = \begin{pmatrix} \hat{\boldsymbol{\mu}}_1 \\ \hat{\boldsymbol{\mu}}_2 \end{pmatrix} = \begin{pmatrix} \frac{1}{N}(N_1\bar{\mathbf{x}}_{(1)1} + N_2\bar{\mathbf{x}}_{(2)}) \\ \bar{\mathbf{x}}_{(1)2} - \hat{\Sigma}_{21}\hat{\Sigma}_{11}^{-1}(\bar{\mathbf{x}}_{(1)1} - \hat{\boldsymbol{\mu}}_1) \end{pmatrix},$$

$$(2.2) \quad \hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{11} & \hat{\Sigma}_{12} \\ \hat{\Sigma}_{21} & \hat{\Sigma}_{22} \end{pmatrix} \\ = \begin{pmatrix} \frac{1}{N}(W_{(1)11} + W_{(2)}) & \hat{\Sigma}_{11}W_{(1)11}^{-1}W_{(1)12} \\ W_{(1)21}W_{(1)11}^{-1}\hat{\Sigma}_{11} & \frac{1}{N_1}W_{(1)22\cdot1} + \hat{\Sigma}_{21}\hat{\Sigma}_{11}^{-1}\hat{\Sigma}_{12} \end{pmatrix},$$

where

$$\bar{\mathbf{x}}_{(1)} = \begin{pmatrix} \bar{\mathbf{x}}_{(1)1} \\ \bar{\mathbf{x}}_{(1)2} \end{pmatrix}, \quad \bar{\mathbf{x}}_{(1)1} = \frac{1}{N_1} \sum_{j=1}^{N_1} \mathbf{x}_{1j}, \quad \bar{\mathbf{x}}_{(1)2} = \frac{1}{N_1} \sum_{j=1}^{N_1} \mathbf{x}_{2j}, \\ \bar{\mathbf{x}}_{(2)} = \frac{1}{N_2} \sum_{j=N_1+1}^N \mathbf{x}_{1j},$$

and

$$W_{(1)} = \begin{pmatrix} W_{(1)11} & W_{(1)12} \\ W_{(1)21} & W_{(1)22} \end{pmatrix} = \sum_{j=1}^{N_1} (\mathbf{x}_j - \bar{\mathbf{x}}_{(1)})(\mathbf{x}_j - \bar{\mathbf{x}}_{(1)})', \\ W_{(2)} = \sum_{j=N_1+1}^N (\mathbf{x}_{1j} - \bar{\mathbf{x}}_{(2)})(\mathbf{x}_{1j} - \bar{\mathbf{x}}_{(2)})' + \frac{N_1N_2}{N}(\bar{\mathbf{x}}_{(1)1} - \bar{\mathbf{x}}_{(2)})(\bar{\mathbf{x}}_{(1)1} - \bar{\mathbf{x}}_{(2)})', \\ W_{(1)22\cdot1} = W_{(1)22} - W_{(1)21}W_{(1)11}^{-1}W_{(1)12}.$$

These results follow from the results in Anderson and Olkin (1985) and Kanda and Fujikoshi (1998).

In the derivation, we use the following transformed parameters  $(\boldsymbol{\eta}, \Delta)$  :

$$\boldsymbol{\eta} = \begin{pmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 - \Delta_{21}\boldsymbol{\mu}_1 \end{pmatrix},$$

$$\Delta = \begin{pmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{21} & \Delta_{22} \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{11}^{-1}\Sigma_{12} \\ \Sigma_{21}\Sigma_{11}^{-1} & \Sigma_{22\cdot 1} \end{pmatrix},$$

where  $\Sigma_{22\cdot 1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$ . We note that  $(\boldsymbol{\eta}, \Delta)$  are in one-to-one correspondence to  $(\boldsymbol{\mu}, \Sigma)$ . After multiplying the observation vector  $\boldsymbol{x}_j$  by the transformation matrix

$$A = \begin{pmatrix} I_{p_1} & O \\ -\Delta_{21} & I_{p_2} \end{pmatrix}$$

on the left side, the log likelihood function is derived, and the results can then be obtained by differentiation.

We consider the following hypothesis test when the data set is of a two-step monotone pattern.

$$(2.3) \quad H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0, \Sigma = \Sigma_0 \text{ vs. } H_1 : \text{not } H_0.$$

Without loss of generality, we can assume that  $\boldsymbol{\mu} = \mathbf{0}$  and  $\Sigma = I_p$ . Then, from the MLEs in (2.1) and (2.2), we obtain the following theorem.

**Theorem 2.1.** *Suppose that the data have a two-step monotone pattern that is missing observations and that  $\lambda_1$  is the likelihood ratio (LR) in the case of the two-step monotone missing data. Then, the LR of the hypothesis test (2.3) is given by*

$$\lambda_1 = |\widehat{\Sigma}_{11}|^{\frac{N}{2}} |\widehat{\Sigma}_{22\cdot 1}|^{\frac{N_1}{2}} \frac{\text{etr} \left( -\frac{1}{2} \sum_{j=1}^N \boldsymbol{x}_{1j} \boldsymbol{x}'_{1j} \right) \text{etr} \left( -\frac{1}{2} \sum_{j=1}^{N_1} \boldsymbol{x}_{2j} \boldsymbol{x}'_{2j} \right)}{\exp \left( -\frac{1}{2} N p_1 \right) \exp \left( -\frac{1}{2} N_1 p_2 \right)}.$$

Further, the LR can be expressed as

$$\begin{aligned} \lambda_1 &= \left( \frac{e}{N} \right)^{\frac{1}{2} N p_1} |W_{(1)11} + W_{(2)}|^{\frac{1}{2} N} \\ &\times \text{etr} \left[ -\frac{1}{2} \left\{ W_{(1)11} + W_{(2)} + \frac{1}{N} (N_1 \bar{\boldsymbol{x}}_{(1)1} + N_2 \bar{\boldsymbol{x}}_{(2)}) (N_1 \bar{\boldsymbol{x}}_{(1)1} + N_2 \bar{\boldsymbol{x}}_{(2)})' \right\} \right] \\ &\times \left( \frac{e}{N_1} \right)^{\frac{1}{2} N_1 p_2} |W_{(1)22\cdot 1}|^{\frac{1}{2} N_1} \text{etr} \left\{ -\frac{1}{2} (W_{(1)22} + N_1 \bar{\boldsymbol{x}}_{(1)2} \bar{\boldsymbol{x}}'_{(1)2}) \right\}. \end{aligned}$$

The result in Theorem 2.1 coincides with the result in Hao and Krishnamoorthy (2001). We note that under  $H_0$ ,  $-2 \log \lambda_1$  is asymptotically distributed as a  $\chi^2$  distribution with  $f = p(p+3)/2$  degrees of freedom when  $N_1, N \rightarrow \infty$  with  $N_1/N \rightarrow \delta \in (0, 1]$ . However, when the sample size is not large, the  $\chi^2$  distribution is not a good approximation to the upper percentile of  $-2 \log \lambda_1$ . Further, it is not easy to find the exact distribution of the LRT statistic  $-2 \log \lambda_1$ . In the next section, we give an approximate upper percentile of  $-2 \log \lambda_1$  and propose modified LRT statistics whose upper percentile is close to that of the  $\chi^2$  distribution even for small samples.

### §3. The modified LRT statistics and an approximate upper percentile of the LRT statistic

In this section, we propose an approximate upper percentile of the null distribution of  $-2 \log \lambda_1$  using linear interpolation based on an asymptotic expansion of the modified LRT statistic in the case of a complete data set. Further, as another approach, we give the modified LRT statistics using the coefficient of the modified LRT statistic for the complete data.

#### 3.1. Modified coefficient approximation procedure

We first consider the LR in the case of a complete data set. Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \sim N_p(\boldsymbol{\mu}, \Sigma)$ , and let  $\lambda_{c,N}$  be the LR for the complete data set. Then, the LR is given by

$$\lambda_{c,N} = \left( \frac{e}{N} \right)^{\frac{Np}{2}} |V|^{\frac{N}{2}} \text{etr} \left\{ -\frac{1}{2}(V + N\bar{\mathbf{x}}\bar{\mathbf{x}}') \right\},$$

where

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad V = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$$

Further, the modified LRT statistic is given by  $-2\rho_{c,N} \log \lambda_{c,N}$ , where  $\rho_{c,N} = 1 - (2p^2 + 9p + 11)/\{6N(p+3)\}$ , and its cumulative distribution function can be expanded as

$$(3.1) \quad \Pr(-2\rho_{c,N} \log \lambda_{c,N} \leq x) = G_f(x) + \frac{\gamma}{M^2} \{G_{f+4}(x) - G_f(x)\} + O(M^{-3}),$$

where

$$M = \rho_{c,N}N, \quad \gamma = \frac{p}{288(p+3)}(2p^4 + 18p^3 + 49p^2 + 36p - 13),$$

and  $G_f(x)$  and  $G_{f+4}(x)$  are the cumulative distribution functions of the  $\chi^2$  distribution with  $f(=p(p+3)/2)$  and  $f+4$  degrees of freedoms, respectively.

This result was derived by Davis (1971) (see Muirhead ((1982), p. 370) and Srivastava ((2002), p. 494)). This means that if the  $\chi^2$  distribution is used as an approximation to the distribution of  $-2\rho_{c,N} \log \lambda_{c,N}$ , the error involved is not of order  $M^{-1}$  but of order  $M^{-2}$ .

If we denote the coefficients of the modified LRT statistics in the case of complete data sets  $N$  and  $N_1$  by  $\rho_{c,N}$  and  $\rho_{c,N_1}$ , respectively, then it may be noted that  $\rho_{\text{miss}}$  is between  $\rho_{c,N}$  and  $\rho_{c,N_1}$ , where  $\rho_{\text{miss}}$  is the coefficient of the modified LRT statistic  $-2\rho_{\text{miss}} \log \lambda_1$ . From the linear interpolation, we propose an approximation to the modified LRT statistic in the case of two-step monotone missing data. Calculating the approximate coefficient  $\rho_L = (p_1\rho_{c,N} + p_2\rho_{c,N_1})/p$ , we can obtain an approximate modified LRT statistic  $-2\rho_L \log \lambda_1$ , where

$$\rho_L = 1 - \frac{1}{N} \left( 1 + \frac{N_2 p_2}{N_1 p} \right) \frac{2p^2 + 9p + 11}{6(p+3)}.$$

### 3.2. Asymptotic expansion approximation procedure

In this subsection, we give an approximate upper percentile of  $-2 \log \lambda_1$  when the data have a two-step monotone pattern that is missing observations. First, in the case of a complete data set, we obtain the following lemma.

**Lemma 3.1.** *Suppose that  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  are distributed as  $N_p(\boldsymbol{\mu}, \Sigma)$ . Then, under the null hypothesis  $H_0$  in (2.3), the upper percentile of the modified LRT statistic,  $-2\rho_{c,N} \log \lambda_{c,N}$ , can be expanded as*

$$q_{\text{MLR-c}}(\alpha) = \chi_f^2(\alpha) + \frac{1}{M^2} \frac{2\gamma}{f(f+2)} \chi_f^2(\alpha) \{ \chi_f^2(\alpha) + f + 2 \} + o(M^{-2}),$$

where

$$M = \rho_{c,N} N, \quad \rho_{c,N} = 1 - \frac{2p^2 + 9p + 11}{6N(p+3)}, \quad f = \frac{1}{2}p(p+3),$$

and  $\chi_f^2(\alpha)$  is the upper percentile of the  $\chi^2$  distribution with  $f$  degrees of freedom.

*Proof.* Putting the upper percentile of  $-2\rho_{c,N} \log \lambda_{c,N}$  with

$$q_{\text{MLR-c}}(\alpha) = \chi_f^2(\alpha) + \frac{1}{M^2} h + o(M^{-2}),$$

where  $h$  is a constant, we have

$$(3.2) \quad 1 - \alpha = G_f(q_{\text{MLR}\cdot\text{c}}(\alpha)) - g_f(\chi_f^2(\alpha)) \frac{1}{M^2} h + o(M^{-2}),$$

where  $G_f(x)$  and  $g_f(x)$  are, respectively, the cumulative distribution function and the density function of the  $\chi^2$  distribution with  $f$  degrees of freedom. On the other hand, from (3.1), we can write

$$(3.3) \quad \begin{aligned} 1 - \alpha &= \Pr \{-2\rho_{\text{c},N} \log \lambda_{\text{c},N} \leq q_{\text{MLR}\cdot\text{c}}(\alpha)\} \\ &= G_f(q_{\text{MLR}\cdot\text{c}}(\alpha)) + \frac{\gamma}{M^2} \{G_{f+4}(q_{\text{MLR}\cdot\text{c}}(\alpha)) - G_f(q_{\text{MLR}\cdot\text{c}}(\alpha))\} \\ &\quad + o(M^{-2}). \end{aligned}$$

Therefore, using  $G_{f+2j}(x) = -2g_{f+2j}(x) + G_{f+2(j-1)}(x)$ ,  $j = 0, 1, 2$  and comparing (3.2) with (3.3), we obtain

$$h = \frac{2\gamma}{f(f+2)} \chi_f^2(\alpha) \{\chi_f^2(\alpha) + f + 2\} + o(M^{-2}).$$

□

From Lemma 3.1 and  $M^{-2} = N^{-2} + O(N^{-3})$ , we can expand the upper percentile of  $-2 \log \lambda_{\text{c},N}$  as

$$q_{\text{LR}\cdot\text{c}}(\alpha) = \chi_f^2(\alpha) + \frac{\nu}{N} \chi_f^2(\alpha) + \frac{1}{N^2} \chi_f^2(\alpha) \left\{ \nu^2 + \frac{2\gamma}{f} + \frac{2\gamma}{f(f+2)} \chi_f^2(\alpha) \right\} + o(N^{-2}),$$

where

$$\nu = \frac{2p^2 + 9p + 11}{6(p+3)}.$$

From the linear interpolation, letting  $q_{\text{LR}\cdot\text{m}}(\alpha)$  be the upper percentile of  $-2 \log \lambda_1$ , an approximate upper percentile of  $-2 \log \lambda_1$  can be obtained as

$$\begin{aligned} q_{\text{LR}\cdot\text{m}}^*(\alpha) &= \chi_f^2(\alpha) + \frac{1}{N} \left( p_1 + \frac{1}{c_1} p_2 \right) \frac{\nu}{p} \chi_f^2(\alpha) \\ &\quad + \frac{1}{N^2} \left( p_1 + \frac{1}{c_1^2} p_2 \right) \frac{\chi_f^2(\alpha)}{p} \left\{ \nu^2 + \frac{2\gamma}{f} + \frac{2\gamma}{f(f+2)} \chi_f^2(\alpha) \right\} + o(N^{-2}), \end{aligned}$$

where  $c_1 = N_1/N$ .

### 3.3. The LRT statistic's decomposition procedure

In this section, we give other modified LRT statistics by the decomposition of  $\lambda_1$ . We first consider the following test problem for  $\Sigma$ .

$$H_{01} : \Sigma = \Sigma_0 = I \text{ vs. } H_{11} : \Sigma \neq I.$$

Hao and Krishnamoorthy (2001) derived the modified LRT statistic  $\lambda_\Sigma^*$  in the case of two-step monotone missing data, which is given by

$$\begin{aligned} \lambda_\Sigma^* &= \left(\frac{e}{n}\right)^{\frac{1}{2}np_1} |W_{(1)11} + W_{(2)}|^{\frac{1}{2}n} \exp\left\{-\frac{1}{2}\text{tr}(W_{(1)11} + W_{(2)})\right\} \\ &\quad \times \left(\frac{e}{n_1}\right)^{\frac{1}{2}n_1p_2} |W_{(1)22\cdot 1}|^{\frac{1}{2}n_1} \exp\left\{-\frac{1}{2}\text{tr}W_{(1)22\cdot 1}\right\} \\ &\quad \times \exp\left\{-\frac{1}{2}\text{tr}(W_{(1)21}W_{(1)11}^{-1}W_{(1)12})\right\}, \end{aligned}$$

where  $n = N - 1$ ,  $n_1 = N_1 - p_1 - 1$ . We note that the modified LRT statistic  $-2\log\lambda_\Sigma^*$  is an unbiased test statistic (see Hao and Krishnamoorthy (2001) and Chang and Richards (2010)). Further, after modifying and rearranging some terms, Hao and Krishnamoorthy (2001) expressed the modified LR for  $H_0$  in (2.3) as  $\lambda_\Sigma^*\omega_1\omega_2$ , where

$$\begin{aligned} \omega_1 &= \exp\left\{-\frac{1}{2N}(N_1\bar{\mathbf{x}}_{(1)1} + N_2\bar{\mathbf{x}}_{(2)})'(N_1\bar{\mathbf{x}}_{(1)1} + N_2\bar{\mathbf{x}}_{(2)})\right\}, \\ \omega_2 &= \exp\left\{-\frac{1}{2}N_1\bar{\mathbf{x}}'_{(1)2}\bar{\mathbf{x}}_{(1)2}\right\}. \end{aligned}$$

If we denote

$$\begin{aligned} \omega_3 &= \left(\frac{e}{N}\right)^{\frac{1}{2}Np_1} |W_{(1)11} + W_{(2)}|^{\frac{1}{2}N} \exp\left\{-\frac{1}{2}\text{tr}(W_{(1)11} + W_{(2)})\right\}, \\ \omega_4 &= \left(\frac{e}{N_1}\right)^{\frac{1}{2}N_1p_2} |W_{(1)22\cdot 1}|^{\frac{1}{2}N_1} \exp\left\{-\frac{1}{2}\text{tr}W_{(1)22\cdot 1}\right\}, \\ \omega_5 &= \exp\left\{-\frac{1}{2}\text{tr}(W_{(1)21}W_{(1)11}^{-1}W_{(1)12})\right\}, \end{aligned}$$

we can express  $\lambda_1 = \prod_{i=1}^5 \omega_i$ . Since  $\omega_1\omega_3$  and  $\omega_2\omega_4$  are of the form of LR for  $H_0$  under non-missing normality, we can give the modified LRT statistics,  $-2\rho_{13}\log\omega_1\omega_3$  and  $-2\rho_{24}\log\omega_2\omega_4$ , respectively, where

$$\rho_{13} = 1 - \frac{2p_1^2 + 9p_1 + 11}{6N(p_1 + 3)}, \quad \rho_{24} = 1 - \frac{2p_2^2 + 9p_2 + 11}{6N_1(p_2 + 3)}.$$



Thus, we propose a new modified LRT statistic given by  $-2 \log \tau$ , where

$$\tau = (\omega_1 \omega_3)^{\rho_{13}} (\omega_2 \omega_4)^{\rho_{24}} \omega_5 .$$

In addition, we denote

$$\omega_3^* = \left( \frac{e}{n} \right)^{\frac{1}{2} n p_1} |W_{(1)11} + W_{(2)}|^{\frac{1}{2} n} \exp \left\{ -\frac{1}{2} \text{tr}(W_{(1)11} + W_{(2)}) \right\},$$

$$\omega_4^* = \left( \frac{e}{n_1} \right)^{\frac{1}{2} n_1 p_2} |W_{(1)22 \cdot 1}|^{\frac{1}{2} n_1} \exp \left\{ -\frac{1}{2} \text{tr} W_{(1)22 \cdot 1} \right\}.$$

Then, since  $\omega_3^*$  and  $\omega_4^*$  are of the form of LR for  $H_{01}$  under non-missing normality, we can propose the modified LRT statistic  $-2 \log \varphi^*$ , where

$$\varphi^* = \omega_1 \omega_2 (\omega_3^*)^{\rho_3^*} (\omega_4^*)^{\rho_4^*} \omega_5$$

and

$$\rho_3^* = 1 - \frac{2p_1^2 + 3p_1 - 1}{6n(p_1 + 1)}, \quad \rho_4^* = 1 - \frac{2p_2^2 + 3p_2 - 1}{6n_1(p_2 + 1)}.$$

#### §4. Simulation studies

We evaluate the accuracy and the asymptotic behaviors of the  $\chi^2$  approximations by Monte Carlo simulation ( $10^6$  runs).

In Table 1, we provide the simulated upper  $100\alpha$  percentiles of  $-2 \log \lambda_1$  and  $-2\rho_L \log \lambda_1$  and the approximate upper percentiles of  $-2 \log \lambda_1$ , that is,  $q_{\text{LR},m}^*(\alpha)$  for  $(p_1, p_2) = (8, 4)$ ;  $\alpha = 0.05, 0.01$ ; and for the following three cases of  $(N_1, N_2)$ ,

$$(N_1, N_2) = \begin{cases} (m, m), & m = 20, 40, 80, 160, 320, \\ (2m, m), & m = 10, 20, 40, 80, 160, \\ (m, 2m), & m = 20, 40, 80, 160. \end{cases}$$

In Table 2, we provide the same upper percentiles as those given in Table 1 for  $(p_1, p_2) = (8, 4)$ ;  $\alpha = 0.05, 0.01$ ;  $(N_1, N_2) = (m_1, m_2)$ ,  $m_1 = 40, 80, 160, 320$ ,  $m_2 = 10, 30, 60, 120$ , where the sets of  $(N_1, N_2)$  are combinations of  $m_1$  and  $m_2$ .

It may be noted from Tables 1 and 2 that the simulated values are closer to the upper percentile of the  $\chi^2$  distribution when the sample size becomes

large. In addition, it can be seen from both tables that the upper percentile of  $-2\rho_L \log \lambda_1$  is considerably better than that of  $-2 \log \lambda_1$  even for small sample sizes. Further, Tables 1 and 2 list the actual type I error rates for the upper percentiles of  $-2 \log \lambda_1$  and  $-2\rho_L \log \lambda_1$  as well as  $q_{\text{LR}\cdot\text{m}}^*(\alpha)$ , which are given by

$$\alpha_1 = \Pr \{ -2 \log \lambda_1 > \chi_f^2(\alpha) \},$$

$$\alpha_{\rho_L} = \Pr \{ -2\rho_L \log \lambda_1 > \chi_f^2(\alpha) \},$$

and

$$\alpha_{q_{\text{LR}\cdot\text{m}}^*} = \Pr \{ -2 \log \lambda_1 > q_{\text{LR}\cdot\text{m}}^*(\alpha) \},$$

respectively. It appears from the simulated results that the approximate value  $q_{\text{LR}\cdot\text{m}}^*(\alpha)$  based on the asymptotic expansion is good for all cases, even when  $N_1 < N_2$ . Therefore, it can be concluded that our approximation procedures are very accurate for most of the cases.

In Tables 3 and 4, we provide the simulated upper percentiles of  $-2 \log \tau$  and  $-2 \log \varphi^*$  for the same cases as those in Tables 1 and 2. It may also be noted that the upper percentiles of  $-2 \log \varphi^*$  are considerably good even for small sample sizes. Tables 3 and 4 list the actual type I error rates for the upper percentiles of  $-2 \log \tau$  and  $-2 \log \varphi^*$ , which are given by

$$\alpha_\tau = \Pr \{ -2 \log \tau > \chi_f^2(\alpha) \}$$

and

$$\alpha_{\varphi^*} = \Pr \{ -2 \log \varphi^* > \chi_f^2(\alpha) \},$$

respectively. The results for actual type I error rates also show that our modified LRT statistic  $-2 \log \varphi^*$  yields considerably good  $\chi^2$  approximations for cases in which the sample size is small.

In conclusion, we have developed the approximate upper percentiles of the LRT statistic  $-2 \log \lambda_1$  and some modified LRT statistics for simultaneous testing of the mean vector and the covariance matrix for the case of two-step monotone missing data. The null distribution of the modified LRT statistic  $-2 \log \varphi^*$  proposed in this paper has considerably good approximation to the  $\chi^2$  distribution even when the sample size is small.

Table 1: The simulated values for  $-2 \log \lambda_1$  and  $-2\rho_L \log \lambda_1$ , and the approximate value for  $-2 \log \lambda_1$ , and the type I error rates when  $(p_1, p_2) = (8, 4)$ 

Sample Size		Upper Percentile			Type I Error Rate		
$N_1$	$N_2$	$-2 \log \lambda_1$	$-2\rho_L \log \lambda_1$	$q_{LR-m}^*(\alpha)$	$\alpha_1$	$\alpha_{\rho_L}$	$\alpha_{q_{LR-m}^*}$
$\alpha = 0.05$							
20	20	148.24	125.89	134.65	0.562	0.180	<b>0.162</b>
40	40	126.08	116.58	122.79	0.190	0.076	<b>0.073</b>
80	80	119.01	114.52	117.69	0.099	<b>0.059</b>	<b>0.059</b>
160	160	115.92	113.74	115.35	0.071	<b>0.054</b>	<b>0.054</b>
320	320	114.48	113.40	114.23	0.059	<b>0.052</b>	<b>0.052</b>
20	10	150.21	123.79	138.65	0.596	0.152	<b>0.137</b>
40	20	127.02	115.85	124.50	0.203	0.070	<b>0.067</b>
80	40	119.38	114.14	118.47	0.104	0.057	<b>0.056</b>
160	80	116.11	113.55	115.72	0.073	<b>0.053</b>	<b>0.053</b>
320	160	114.61	113.35	114.41	0.060	<b>0.051</b>	<b>0.051</b>
20	40	146.54	128.13	130.99	0.531	0.212	<b>0.189</b>
40	80	125.28	117.41	121.16	0.177	0.084	<b>0.080</b>
80	160	118.60	114.87	116.93	0.095	<b>0.062</b>	<b>0.062</b>
160	320	115.67	113.86	114.98	0.069	<b>0.055</b>	<b>0.055</b>
$\alpha = 0.01$							
20	20	163.40	138.77	147.80	0.328	0.061	<b>0.052</b>
40	40	138.45	128.01	134.71	0.063	0.018	<b>0.017</b>
80	80	130.49	125.57	129.11	0.025	0.013	<b>0.012</b>
160	160	127.20	124.80	126.53	0.016	<b>0.011</b>	<b>0.011</b>
320	320	125.63	124.44	125.31	0.013	<b>0.011</b>	<b>0.011</b>
20	10	165.49	136.39	152.21	0.360	0.048	<b>0.041</b>
40	20	139.43	127.17	136.60	0.069	0.016	<b>0.015</b>
80	40	130.92	125.16	129.97	0.027	<b>0.012</b>	<b>0.012</b>
160	80	127.23	124.43	126.94	0.016	0.011	<b>0.010</b>
320	160	125.67	124.29	125.51	0.013	<b>0.010</b>	<b>0.010</b>
20	40	161.66	141.35	143.76	0.300	0.077	<b>0.065</b>
40	80	137.56	128.92	132.93	0.057	0.020	<b>0.019</b>
80	160	130.07	125.99	128.27	0.024	<b>0.013</b>	<b>0.013</b>
160	320	126.86	124.86	126.13	0.015	<b>0.011</b>	<b>0.011</b>

Note. The closest to  $\alpha$  in the values  $\alpha_1$ ,  $\alpha_{\rho_L}$ , and  $\alpha_{q_{LR-m}^*}$  of each row is in bold.

$$\chi_f^2(0.05) = 113.145, \chi_f^2(0.01) = 124.116.$$

Table 2: The simulated values for  $-2 \log \lambda_1$  and  $-2\rho_L \log \lambda_1$ , and the approximate value for  $-2 \log \lambda_1$ , and the type I error rates when  $(p_1, p_2) = (8, 4)$

Sample Size		Upper Percentile			Type I Error Rate		
$N_1$	$N_2$	$-2 \log \lambda_1$	$-2\rho_L \log \lambda_1$	$q_{LR-m}^*(\alpha)$	$\alpha_1$	$\alpha_{\rho_L}$	$\alpha_{q_{LR-m}^*}$
$\alpha = 0.05$							
40	10	127.69	115.18	125.92	0.214	0.064	<b>0.061</b>
80	10	119.81	113.54	119.55	0.109	0.053	<b>0.052</b>
160	10	116.45	113.29	116.35	0.075	<b>0.051</b>	<b>0.051</b>
320	10	114.78	113.19	114.75	0.061	<b>0.050</b>	<b>0.050</b>
40	30	126.47	116.26	123.51	0.196	0.074	<b>0.070</b>
80	30	119.49	113.97	118.76	0.105	0.056	<b>0.055</b>
160	30	116.28	113.34	116.12	0.074	<b>0.051</b>	<b>0.051</b>
320	30	114.70	113.17	114.68	0.061	<b>0.050</b>	<b>0.050</b>
40	60	125.61	117.09	121.80	0.182	0.081	<b>0.077</b>
80	60	119.14	114.33	118.02	0.101	0.058	<b>0.057</b>
160	60	116.16	113.48	115.86	0.073	<b>0.052</b>	<b>0.052</b>
320	60	114.70	113.25	114.60	0.061	<b>0.051</b>	<b>0.051</b>
40	120	124.90	117.84	120.38	0.172	0.088	<b>0.084</b>
80	120	118.72	114.70	117.23	0.097	0.061	<b>0.060</b>
160	120	115.95	113.61	115.51	0.071	<b>0.053</b>	<b>0.053</b>
320	120	114.61	113.29	114.48	0.061	<b>0.051</b>	<b>0.051</b>
$\alpha = 0.01$							
40	10	140.21	126.48	138.17	0.075	0.014	<b>0.013</b>
80	10	131.45	124.57	131.15	0.029	0.011	<b>0.010</b>
160	10	127.73	124.26	127.63	0.017	<b>0.010</b>	<b>0.010</b>
320	10	125.85	124.10	125.87	0.013	<b>0.010</b>	<b>0.010</b>
40	30	138.67	127.47	135.51	0.066	0.017	<b>0.016</b>
80	30	131.24	125.17	130.28	0.028	<b>0.012</b>	<b>0.012</b>
160	30	127.53	124.31	127.38	0.017	<b>0.010</b>	<b>0.010</b>
320	30	125.76	124.09	125.81	0.013	<b>0.010</b>	<b>0.010</b>
40	60	137.82	128.47	133.63	0.060	0.019	<b>0.018</b>
80	60	130.74	125.46	129.47	0.026	<b>0.012</b>	<b>0.012</b>
160	60	127.43	124.48	127.09	0.016	<b>0.011</b>	<b>0.011</b>
320	60	125.79	124.20	125.72	0.013	<b>0.010</b>	<b>0.010</b>
40	120	137.08	129.34	132.07	0.055	0.021	<b>0.020</b>
80	120	130.35	125.93	128.60	0.024	<b>0.013</b>	<b>0.013</b>
160	120	127.22	124.65	126.71	0.016	<b>0.011</b>	<b>0.011</b>
320	120	125.60	124.15	125.58	0.013	<b>0.010</b>	<b>0.010</b>

Note. The closest to  $\alpha$  in the values  $\alpha_1$ ,  $\alpha_{\rho_L}$ , and  $\alpha_{q_{LR-m}^*}$  of each low is in bold.

$$\chi_f^2(0.05) = 113.145, \chi_f^2(0.01) = 124.116.$$

Table 3: The simulated values for  $-2 \log \tau$  and  $-2 \log \varphi^*$ , and the type I error rates when  $(p_1, p_2) = (8, 4)$ 

Sample Size		Upper Percentile		Type I Error Rate	
$N_1$	$N_2$	$-2 \log \tau$	$-2 \log \varphi^*$	$\alpha_\tau$	$\alpha_{\varphi^*}$
$\alpha = 0.05$					
20	20	138.73	113.43	0.397	<b>0.052</b>
40	40	122.40	113.16	0.139	<b>0.050</b>
80	80	117.25	113.12	0.083	<b>0.050</b>
160	160	115.12	113.15	0.064	<b>0.050</b>
320	320	114.09	113.15	0.056	<b>0.050</b>
20	10	138.91	113.56	0.400	<b>0.053</b>
40	20	122.49	113.22	0.140	<b>0.050</b>
80	40	117.36	113.21	0.084	<b>0.050</b>
160	80	115.12	113.12	0.064	<b>0.050</b>
320	160	114.06	113.13	0.056	<b>0.050</b>
20	40	138.59	113.35	0.395	<b>0.051</b>
40	80	122.38	113.20	0.139	<b>0.050</b>
80	160	117.27	113.12	0.083	<b>0.050</b>
160	320	115.07	113.12	0.064	<b>0.050</b>
$\alpha = 0.01$					
20	20	152.69	124.43	0.189	<b>0.011</b>
40	40	134.29	124.05	0.040	<b>0.010</b>
80	80	128.61	124.07	0.020	<b>0.010</b>
160	160	126.33	124.18	0.014	<b>0.010</b>
320	320	125.24	124.18	0.012	<b>0.010</b>
20	10	153.09	124.62	0.192	<b>0.011</b>
40	20	134.44	124.22	0.041	<b>0.010</b>
80	40	128.74	124.18	0.020	<b>0.010</b>
160	80	126.30	124.11	0.014	<b>0.010</b>
320	160	125.20	124.14	0.012	<b>0.010</b>
20	40	152.57	124.45	0.188	<b>0.010</b>
40	80	134.26	124.16	0.040	<b>0.010</b>
80	160	128.62	124.07	0.020	<b>0.010</b>
160	320	126.28	124.18	0.014	<b>0.010</b>

Note. The closer to  $\alpha$  in the values  $\alpha_\tau$  and  $\alpha_{\varphi^*}$  of each row is in bold.

$$\chi_f^2(0.05) = 113.145, \chi_f^2(0.01) = 124.116.$$

Table 4: The simulated values for  $-2 \log \tau$  and  $-2 \log \varphi^*$ , and the type I error rates when  $(p_1, p_2) = (8, 4)$ 

Sample Size		Upper Percentile		Type I Error Rate	
$N_1$	$N_2$	$-2 \log \tau$	$-2 \log \varphi^*$	$\alpha_\tau$	$\alpha_{\varphi^*}$
$\alpha = 0.05$					
40	10	122.57	113.28	0.141	<b>0.051</b>
80	10	117.35	113.16	0.083	<b>0.050</b>
160	10	115.12	113.14	0.064	<b>0.050</b>
320	10	114.11	113.14	0.057	<b>0.050</b>
40	30	122.45	113.22	0.140	<b>0.050</b>
80	30	117.32	113.17	0.083	<b>0.050</b>
160	30	115.13	113.15	0.064	<b>0.050</b>
320	30	114.14	113.17	0.057	<b>0.050</b>
40	60	122.37	113.20	0.139	<b>0.050</b>
80	60	117.25	113.11	0.083	<b>0.050</b>
160	60	115.11	113.11	0.064	<b>0.050</b>
320	60	114.13	113.17	0.057	<b>0.050</b>
40	120	122.39	113.19	0.139	<b>0.050</b>
80	120	117.24	113.11	0.082	<b>0.050</b>
160	120	115.18	113.18	0.065	<b>0.050</b>
320	120	114.19	113.24	0.057	<b>0.051</b>
$\alpha = 0.01$					
40	10	134.47	124.24	0.041	<b>0.010</b>
80	10	128.76	124.17	0.020	<b>0.010</b>
160	10	126.29	124.12	0.014	<b>0.010</b>
320	10	125.09	124.07	0.012	<b>0.010</b>
40	30	134.45	124.20	0.041	<b>0.010</b>
80	30	128.76	124.19	0.020	<b>0.010</b>
160	30	126.23	124.15	0.014	<b>0.010</b>
320	30	125.29	124.25	0.012	<b>0.010</b>
40	60	134.33	124.18	0.040	<b>0.010</b>
80	60	128.59	124.00	0.020	<b>0.010</b>
160	60	126.25	124.05	0.014	<b>0.010</b>
320	60	125.17	124.15	0.012	<b>0.010</b>
40	120	134.31	124.19	0.040	<b>0.010</b>
80	120	128.65	124.08	0.020	<b>0.010</b>
160	120	126.38	124.22	0.014	<b>0.010</b>
320	120	125.21	124.15	0.012	<b>0.010</b>

Note. The closer to  $\alpha$  in the values  $\alpha_\tau$  and  $\alpha_{\varphi^*}$  of each low is in bold.

$$\chi_f^2(0.05) = 113.145, \chi_f^2(0.01) = 124.116.$$

### Acknowledgments

The authors would like to thank the referee for helpful comments and suggestions. Second author's research was in part supported by Grant-in-Aid for Scientific Research (C) (26330050).

### References

- [1] Anderson, T. W. (1957). Maximum likelihood estimates for a multivariate normal distribution when some observations are missing, *J. Amer. Statist. Assoc.*, **52**, 200–203.
- [2] Anderson, T. W. and Olkin, I. (1985). Maximum-likelihood estimation of the parameters of a multivariate normal distribution, *Linear Algebra Appl.*, **70**, 147–171.
- [3] Bhargava, R. (1962). Multivariate tests of hypotheses with incomplete data, *Technical Report No.3, Applied Mathematics and Statistics Laboratories, Stanford University*, Stanford, California.
- [4] Chang, W. -Y. and Richards, D. St. P. (2009). Finite-sample inference with monotone incomplete multivariate normal data, I, *J. Multivariate Anal.*, **100**, 1883–1899.
- [5] Chang, W. -Y. and Richards, D. St. P. (2010). Finite-sample inference with monotone incomplete multivariate normal data, II, *J. Multivariate Anal.*, **101**, 603–620.
- [6] Choudhari, P., Kundu, D. and Misra, N. (2001). Likelihood ratio test for simultaneous testing of the mean and the variance of a normal distribution, *J. Statist. Comput. Simul.*, **71**, 313–333.
- [7] Davis, A. W. (1971). Percentile approximations for a class of likelihood ratio criteria, *Biometrika*, **58**, 349–356.
- [8] Hao, J. and Krishnamoorthy, K. (2001). Inferences on a normal covariance matrix and generalized variance with monotone missing data, *J. Multivariate Anal.*, **78**, 62–82.
- [9] Jinadasa, K. G. and Tracy, D. S. (1992). Maximum likelihood estimation for multivariate normal distribution with monotone sample, *Comm. Statist. Theory Methods*, **21**, 41–50.
- [10] Kanda, T. and Fujikoshi, Y. (1998). Some basic properties of the MLE's for a multivariate normal distribution with monotone missing data, *Amer. J. Math. Management Sci.*, **18**, 161–190.
- [11] Krishnamoorthy, K. and Pannala, M. K. (1999). Confidence estimation of a normal mean vector with incomplete data, *Canad. J. Statist.*, **27**, 395–407.

- [12] Little, R. J. and Rubin, D. B. (2002). *Statistical Analysis with Missing Data*, 2nd ed., Hoboken, NJ: Wiley.
- [13] McLachlan, J. G. and Krishnan, T. (1997). *The EM Algorithm and Extensions*, New York, NY: Wiley.
- [14] Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*, New York, NY: Wiley.
- [15] Seo, T. and Srivastava, M. S. (2000). Testing equality of means and simultaneous confidence intervals in repeated measures with missing data, *Biometrical J.*, **42**, 981–993.
- [16] Seko, N., Yamazaki, A. and Seo, T. (2012). Tests for mean vector with two-step monotone missing data, *SUT J. Math.*, **48**, 13–36.
- [17] Srivastava, M. S. (1985). Multivariate data with missing observations, *Comm. Statist. Theory Methods*, **14**, 775–792.
- [18] Srivastava, M. S. (2002). *Methods of Multivariate Statistics*, New York, NY: Wiley.
- [19] Yagi, A. and Seo, T. (2014). A test for mean vector and simultaneous confidence intervals with three-step monotone missing data, *Amer. J. Math. Management Sci.*, **33**, 161–175.
- [20] Zhang, L., Xu, X. and Chen, G. (2012). The exact likelihood ratio test for equality of two normal populations, *Amer. Statist.*, **66**, 180–184.

Miki Hosoya  
Department of Mathematical Information Science  
Tokyo University of Science  
1-3, Kagurazaka, Shinjuku-ku, Tokyo 162-8601, Japan  
*E-mail*: 1414622@ed.tus.ac.jp

Takashi Seo  
Department of Mathematical Information Science  
Tokyo University of Science  
1-3, Kagurazaka, Shinjuku-ku, Tokyo 162-8601, Japan  
*E-mail*: seo@rs.tus.ac.jp