

# MINIMAX ESTIMATION OF THE MEAN MATRIX OF THE MATRIX VARIATE NORMAL DISTRIBUTION UNDER THE DIVERGENCE LOSS FUNCTION

Shokofeh Zinodiny

*Amirkabir University of Technology, Tehran, Iran*

Sadegh Rezaei

*Amirkabir University of Technology, Tehran, Iran*

Saralees Nadarajah <sup>1</sup>

*University of Manchester, Manchester M13 9PL, UK*

## 1. INTRODUCTION

Let  $\mathbf{X} = (x_{i,j})$  be a  $p \times m$  random matrix having a matrix variate normal distribution with mean matrix  $\Theta = (\theta_{i,j})$  and covariance matrix  $\mathbf{V} \otimes \mathbf{I}_m$ , where  $\mathbf{V}$  is known or unknown,  $\mathbf{I}_m$  is the  $m \times m$  identity matrix and  $\otimes$  denotes the Kronecker product. This note considers estimation of  $\Theta$  relative to the divergence loss function

$$L(\mathbf{a}; \Theta) = \frac{1}{\beta(1-\beta)} \left[ 1 - \int_{\mathbf{R}^{p \times m}} f^\beta(\mathbf{X}|\mathbf{a}) f^{1-\beta}(\mathbf{X}|\Theta) d\mathbf{x} \right], \quad (1)$$

where  $0 < \beta < 1$ ,  $f(\mathbf{X}|\Theta)$  is the conditional density function of  $\mathbf{X}$  given  $\Theta$  and  $\mathbf{a}$  is an estimator of  $\Theta$ .

Estimators for the mean matrix have been proposed under different loss functions. Efron and Morris (1972) proposed an empirical Bayes estimator outperforming the maximum likelihood estimator,  $X$ , for the case  $m > p + 1$ . Stein (1973), Zhang (1986a), Zhang (1986b), Bilodeau and Kariya (1989), Ghosh and Shieh (1991), Konno (1991), Tsukuma (2008), Tsukuma (2010) and others found minimax estimators better than the maximum likelihood estimator for quadratic loss and general quadratic loss functions.

The most recent minimax estimator for  $\Theta$  is that due to Zinodiny *et al.* (2017). They estimated  $\Theta$  under the balanced loss function. But Zinodiny *et al.* (2017) considered only the case  $\mathbf{V} = \sigma^2 \mathbf{I}_p$ , where  $\sigma^2$  is unknown. In this note, we consider three cases

---

<sup>1</sup> Corresponding author. E-mail: mbbssn2@manchester.ac.uk

as outlined in the abstract. Furthermore, the divergence loss function in (1) has attractive properties not shared by other loss functions, including the balanced loss function. Firstly, it is invariant to one-to-one coordinate transformation of  $\mathbf{a}$ . Secondly, the Fisher information function, used commonly to measure the sensitivity of estimates, is proportional to (1). Thirdly, many loss functions including the balanced loss function are symmetric, but (1) is not symmetric. If  $\beta$  is closer to 1 more weight is given to  $\mathbf{a}$ . If  $\beta$  is closer to 0 more weight is given to  $\Theta$ . Other attractive properties of (1) can be found in Kashyap (1974).

The divergence loss function is directly comparable to the densities  $f(\mathbf{X}|\Theta)$  and  $f(\mathbf{X}|\mathbf{a})$ . Robert (2001) refers to it as an intrinsic loss function, see also Amari (1982) and Cressie and Read (1984). The divergence loss function has been applied in many areas. Some examples include: discrimination between stationary Gaussian processes (Corduas, 1985); the Vocal Joystick engine, a real-time software library which can be used to map non-linguistic vocalizations into realizable continuous control signals (Malkin et al., 2011); risk-aware modeling framework for speech summarization (Chen and Lin, 2012); models for web image retrieval (Yang et al., 2012); models for coclustering the mouse brain atlas (Ji et al., 2013); properties of record values (Paul and Thomas, 2016).

The aim of this note is to estimate the mean matrix of the matrix variate normal distribution under the divergence loss function. The contents are organized as follows: Section 2 shows that the empirical Bayes estimators dominate the maximum likelihood estimator under (1) for  $m > p + 1$  and hence the latter is inadmissible for the case  $\mathbf{V}$  is known. We find a general class of minimax estimators using a technique due to Stein (1973) when  $\mathbf{V} = \mathbf{I}_p$ . Section 3 shows that  $X$  is inadmissible for  $m > p + 1$  when  $\mathbf{V} = \sigma^2 \mathbf{I}_p$ , where  $\sigma^2$  is unknown. Section 4 shows that  $X$  is inadmissible for the case  $\mathbf{V}$  is unknown. These sections in fact extend the results of Ghosh et al. (2008) and Ghosh and Mergel (2009) to the multivariate case. Section 5 performs a simulation study to compare one of the derived estimators with that in Zinodiny et al. (2017). Section 6 concludes the note.

## 2. ESTIMATION OF THE MEAN MATRIX WHEN $\mathbf{V}$ IS KNOWN

Here, we suppose

$$\mathbf{X} \sim N_{p \times m}(\Theta, \mathbf{V} \otimes \mathbf{I}_m),$$

where  $N_{p \times m}(\Theta, \mathbf{V} \otimes \mathbf{I}_m)$  denotes the matrix variate normal distribution with mean matrix  $\Theta$  and covariance matrix  $\mathbf{V} \otimes \mathbf{I}_m$ , where  $\mathbf{V}$  is assumed known. We consider the problem of estimating the mean matrix  $\Theta$  under the loss function (1), namely

$$\begin{aligned} L(\mathbf{a}; \Theta) &= \frac{1}{\beta(1-\beta)} \left[ 1 - \int_{\mathbf{R}^{p \times m}} f^\beta(\mathbf{X}|\mathbf{a}) f^{1-\beta}(\mathbf{X}|\Theta) d\mathbf{x} \right] \\ &= \frac{1}{\beta(1-\beta)} \left[ 1 - e^{-\frac{\beta(1-\beta)\text{tr}(\mathbf{a}-\Theta)'\mathbf{V}^{-1}(\mathbf{a}-\Theta)}{2}} \right], \end{aligned} \quad (2)$$

where  $\text{tr}(\mathbf{A})$  and  $\mathbf{A}'$  denote, respectively, the trace and the transpose of a matrix  $\mathbf{A}$ , and  $f(\mathbf{X}|\Theta)$  denotes the matrix variate normal density function. The usual estimator for  $\Theta$  is the maximum likelihood estimator,  $X$ .

Lemma 1 shows that the maximum likelihood estimator is minimax under the loss function (2).

LEMMA 1. *The estimator  $X$  is minimax under the loss function (2).*

PROOF. Suppose the proper prior sequence  $\pi_n(\Theta)$  is distributed according to a matrix variate normal distribution with mean matrix zero and covariance matrix  $n\mathbf{C} \otimes \mathbf{I}_m$ , i.e.  $\Theta \sim N_{p \times m}(\mathbf{0}_{p \times m}, n\mathbf{C} \otimes \mathbf{I}_m)$ , where  $\mathbf{C}$  is a  $p \times p$  known positive definite matrix. Then, the posterior distribution is

$$\Theta|\mathbf{X} \sim N_{p \times m} \left( (\mathbf{I}_p - \mathbf{V}(n\mathbf{C} + \mathbf{V})^{-1})\mathbf{X}, (\mathbf{V}^{-1} + (n\mathbf{C})^{-1})^{-1} \otimes \mathbf{I}_m \right).$$

Thus, the Bayes estimator is

$$\delta_{\pi_n}(\mathbf{X}) = E[\Theta|\mathbf{X}] = [\mathbf{I}_p - \mathbf{V}(n\mathbf{C} + \mathbf{V})^{-1}]\mathbf{X}.$$

Moreover, the risk function of  $\mathbf{X}$  and Bayes risk function of  $\delta_{\pi_n}(\mathbf{X})$  are

$$R(\mathbf{X}; \Theta) = \frac{1 - [1 + \beta(1 - \beta)]^{-\frac{pm}{2}}}{\beta(1 - \beta)}$$

and

$$r_n = r(\pi_n, \delta_{\pi_n}(\mathbf{X})) = \frac{1 - |\mathbf{V}^{-1} + (n\mathbf{C})^{-1}|^{\frac{m}{2}} [1 + \beta(1 - \beta)] |\mathbf{V}^{-1} + (n\mathbf{C})^{-1}|^{-\frac{m}{2}}}{\beta(1 - \beta)},$$

respectively, where  $|\mathbf{A}|$  denotes the determinant of  $\mathbf{A}$ . Since the determinant of a matrix is a continuous function, we have

$$\lim_{n \rightarrow \infty} r_n = \frac{1 - [1 + \beta(1 - \beta)]^{-\frac{pm}{2}}}{\beta(1 - \beta)} = \sup_{\Theta} R(\mathbf{X}; \Theta).$$

Hence, using Theorem 1.12 on page 613 of Lehmann and Casella (1998) and results of Blyth (1951),  $X$  is minimax under the loss function (2). □

Now we construct a class of empirical Bayes estimators better than  $X$ . Assume we have some additional information about  $\Theta$  that can be written as

$$\Theta \sim N_{p \times m}(\mathbf{0}_{p \times m}, \mathbf{A} \otimes \mathbf{I}_m).$$

The conditional distribution of  $\Theta$  given  $\mathbf{X}$  is

$$N_{p \times m} \left( \left( \mathbf{I}_p - \mathbf{V}(\mathbf{V} + \mathbf{A})^{-1} \right) \mathbf{X}, (\mathbf{V}^{-1} + \mathbf{A}^{-1})^{-1} \otimes \mathbf{I}_m \right).$$

So, the Bayes estimator of  $\Theta$  under (2) is

$$\widehat{\Theta}_B = E[\Theta | \mathbf{X}] = \left( \mathbf{I}_p - \mathbf{V}\Sigma^{-1} \right) \mathbf{X},$$

where  $\Sigma = \mathbf{A} + \mathbf{V}$ . In an empirical Bayes scenario,  $\Sigma$  is unknown, and is estimated from the marginal distribution of  $\mathbf{X}$ . The marginal distribution of  $\mathbf{X}$  is  $N_{p \times m} \left( \mathbf{0}_{p \times m}, \Sigma \otimes \mathbf{I}_m \right)$ . So,  $\mathbf{S} = \mathbf{X}\mathbf{X}'$  is completely sufficient for  $\Sigma$  and an empirical Bayes estimator is

$$\widehat{\Theta}_{EB} = E[\Theta | \mathbf{X}] = \left[ \mathbf{I}_p - \mathbf{V}\widehat{\Sigma}^{-1}(\mathbf{S}) \right] \mathbf{X},$$

where  $\widehat{\Sigma}^{-1}(\mathbf{S})$  is an estimator for  $\Sigma^{-1}$  and  $\widehat{\Sigma}^{-1}(\mathbf{S})$  depends on  $\mathbf{X}$  only through  $\mathbf{S}$ . According to Ghosh and Shieh (1991), a natural candidate for  $\widehat{\Sigma}^{-1}(\mathbf{S})$  is  $(m - p - 1)\mathbf{S}^{-1}$ . Ghosh and Shieh (1991) obtained this empirical Bayes estimator when the loss function was

$$L_1(\delta; \Theta) = \text{tr} \left( (\delta - \Theta)' \mathbf{Q} (\delta - \Theta) \right), \tag{3}$$

where  $\mathbf{Q}$  is a  $p \times p$  known positive definite matrix.

To continue, we use the following notations borrowed from Ghosh and Shieh (1991): let  $E_{\Theta}$  denote the expectation conditional on  $\Theta$ ;  $\widetilde{E}$  denotes the expectation over the marginal distribution of  $\Theta$ ;  $E$  denotes the expectation over the joint distribution of  $\mathbf{X}$  and  $\Theta$ . Let  $R_i(\delta; \Theta) = E_{\Theta} [L_i(\delta; \Theta)]$ ,  $i = 1, 2, 3$  and  $R(\delta; \Theta) = E_{\Theta} [L(\delta; \Theta)]$ . We use the notation  $\mathbf{D} = (d_{i,j})$ , where  $d_{i,j} = \left( \frac{1 + \delta_{i,j}}{2} \right) \frac{\partial}{\partial s_{i,j}}$  (see the last line of page 308 in Ghosh and Shieh (1991)),  $\delta_{i,j}$  being the Kronecker deltas and  $s_{i,j}$  the  $(i, j)$ th element of  $\mathbf{S}$ . Note that  $\mathbf{D}$  is a differential operator in the form of a matrix. If we assume  $f(\mathbf{A})$  is a real-valued function of a  $p \times p$  matrix  $\mathbf{A}$ , then  $\mathbf{D}f(\mathbf{A})$  is a  $p \times p$  matrix with  $(i, j)$ th element  $\frac{1 + \delta_{i,j}}{2} \frac{\partial f(\mathbf{A})}{\partial s_{i,j}}$ . For example, if  $p = 2$  and  $f(\mathbf{A}) = s_{1,1}^2 + s_{2,2}^2$  then  $\mathbf{D}f(\mathbf{A}) = \begin{pmatrix} 2s_{1,1} & 0 \\ 0 & 2s_{2,2} \end{pmatrix}$ .

Also for any  $p \times p$  matrix  $\mathbf{T}$ ,  $\mathbf{D}(\mathbf{T})$  is a  $p \times p$  matrix with  $(i, j)$ th element  $\sum_{l=1}^p d_{i,l} t_{l,j}$ .

For example, if  $p = 2$  then  $\mathbf{D}(\mathbf{S}) = \begin{pmatrix} \frac{3}{2} & 0 \\ 0 & \frac{3}{2} \end{pmatrix}$ .

**THEOREM 2.** *The empirical Bayes estimator in (3) is minimax under (2) if it is minimax under the quadratic loss function (3) when  $\mathbf{Q} = \mathbf{V}^{-1}$ .*

**PROOF.** The difference of the risks of  $\widehat{\Theta}_{EB}$  and  $\mathbf{X}$  is

$$R(\widehat{\Theta}_{EB}; \Theta) - R(\mathbf{X}; \Theta) = \frac{1}{\beta(1 - \beta)} \left[ g(\mathbf{X}; \Theta) - g(\widehat{\Theta}_{EB}; \Theta) \right],$$

where

$$g(\boldsymbol{\delta}; \boldsymbol{\Theta}) = E \left[ e^{-\beta(1-\beta)\text{tr}((\boldsymbol{\delta}-\boldsymbol{\Theta})'\mathbf{V}^{-1}(\boldsymbol{\delta}-\boldsymbol{\Theta}))/2} \right].$$

Using the fact that  $e^{-x} \geq 1 - x$ , we see that

$$g(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) \geq g(\mathbf{X}; \boldsymbol{\Theta}) - \frac{\beta(1-\beta)}{2[1+\beta(1-\beta)]^{\frac{pm}{2}}} E_{\boldsymbol{\Theta}}[\text{tr}(\mathbf{X}'\widehat{\boldsymbol{\Sigma}}^{-1}\mathbf{V}'\mathbf{V}^{-1}\mathbf{V}\widehat{\boldsymbol{\Sigma}}^{-1}\mathbf{X} - 2\text{tr}(\mathbf{X}'\widehat{\boldsymbol{\Sigma}}^{-1}\mathbf{V}'\mathbf{V}^{-1}(\mathbf{X}-\boldsymbol{\Theta}))],$$

where  $E_{\boldsymbol{\Theta}}$  is taken with respect to  $N_{p \times m}(\boldsymbol{\Theta}, \frac{\mathbf{V}}{1+\beta(1-\beta)} \otimes \mathbf{I}_m)$ . Using Theorem 2 in Ghosh and Shieh (1991) and the notation  $R_1(\boldsymbol{\delta}; \boldsymbol{\Theta}) = E_{\boldsymbol{\Theta}}[L_1(\boldsymbol{\delta}; \boldsymbol{\Theta})]$ , we can write

$$g(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) - g(\mathbf{X}; \boldsymbol{\Theta}) \geq -\frac{\beta(1-\beta)}{2} [1+\beta(1-\beta)]^{-\frac{pm}{2}} [R_1(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) - R_1(\mathbf{X}; \boldsymbol{\Theta})]. \tag{4}$$

Hence, if  $R_1(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) \leq R_1(\mathbf{X}; \boldsymbol{\Theta})$ , then  $R(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) \leq R(\mathbf{X}; \boldsymbol{\Theta})$ . □

Similar to Ghosh and Shieh (1991), we can write (4) as

$$g(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) - g(\mathbf{X}; \boldsymbol{\Theta}) \geq -\frac{\beta(1-\beta)}{2} [1+\beta(1-\beta)]^{-\frac{pm}{2}} E \left[ \text{tr}(\mathbf{V}^{\frac{1}{2}}\mathbf{U}_1\mathbf{V}^{\frac{1}{2}}) \right]$$

and

$$\mathbf{U}_1 = \widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S})\mathbf{S}\widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S}) - 4\mathbf{D}(\mathbf{S}\widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S})) - 2(m-p-1)\widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S}),$$

and  $R_1(\widehat{\boldsymbol{\Theta}}_{EB}; \boldsymbol{\Theta}) \leq R_1(\mathbf{X}; \boldsymbol{\Theta})$  if  $\mathbf{U}_1 \leq 0$  has a positive probability for some  $\boldsymbol{\Theta}$ . If  $\widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S})$  is chosen in such a way that  $\widehat{\boldsymbol{\Theta}}_{EB}$  improves on  $X$  under the quadratic loss function (3), then it improves on  $X$  also under (2).

Let  $\mathbf{O}_p$  be the set of orthogonal matrices of order  $p$  and let  $\mathbf{V}_{m,p}$  be the Stiefel manifold, namely,  $\mathbf{V}_{m,p} = \{ \mathbf{V} \in \mathbf{R}^{m \times p}, \mathbf{V}'\mathbf{V} = \mathbf{I}_p \}$ . Write the singular value decomposition of  $\mathbf{X}$  as  $\mathbf{U}\mathbf{L}\mathbf{V}'$ , where  $\mathbf{U} \in \mathbf{O}_p$ ,  $\mathbf{V} \in \mathbf{V}_{m,p}$  and  $\mathbf{L} = \text{diag}(l_1, l_2, \dots, l_p)$  with  $l_1 > l_2 > \dots > l_p > 0$ . Ghosh and Shieh (1991) showed that if  $\widehat{\boldsymbol{\Sigma}}^{-1}(\mathbf{S}) = \mathbf{U}\mathbf{L}^{-1}\boldsymbol{\Psi}^*(\mathbf{L})\mathbf{U}'$ , where  $\boldsymbol{\Psi}^*(\mathbf{L})$  is a diagonal matrix with diagonal elements equal to  $\psi_i^*(\mathbf{L}), i = 1, \dots, p$ , then  $\widehat{\boldsymbol{\Theta}}_{EB}$  is minimax under the following conditions when the loss function is (3) with  $\mathbf{Q} = \mathbf{V}^{-1}$ .

**THEOREM 3.** *Suppose that  $\psi_i^*(\mathbf{L}), i = 1, \dots, p$  satisfy*

$$I. \ 0 < \psi_i^*(\mathbf{L}) < 2(m-p-1), \ i = 1, \dots, p,$$

II. For every  $i = 1, \dots, p$ ,  $\frac{\partial \psi_i^*(\mathbf{L})}{\partial l_i} \geq 0$ ,

III.  $\psi_i^*(\mathbf{L})$  are similarly ordered to  $l_i$ , that is  $[\psi_i^*(\mathbf{L}) - \psi_t^*(\mathbf{L})](l_i - l_t) \geq 0$  for every  $1 \leq i, t \leq p$ .

Then,  $\hat{\Theta}_{EB} = [\mathbf{I}_p - \mathbf{VUL}^{-1}\Psi^*(\mathbf{L})\mathbf{U}']\mathbf{X}$  improves on  $X$  for  $m > p + 1$  under the loss function (3).

PROOF. See Ghosh and Shieh (1991). □

Using Theorem 2, we can see that  $\hat{\Theta}_{EB} = [\mathbf{I}_p - \mathbf{VUL}^{-1}\Psi^*(\mathbf{L})\mathbf{U}']\mathbf{X}$  is minimax for  $m > p + 1$  under the conditions of Theorem 3 when the loss function is (2). This means,  $X$  is inadmissible for  $m > p + 1$  under the loss function (2). Examples (1) and (2) in Ghosh and Shieh (1991) illustrated the above result. Of course, one may obtain estimators of the form (3) which dominate  $X$  when the conditions of Theorem 3 are not met. Some examples are given in Ghosh and Shieh (1991).

In the rest of this section, we consider the problem of estimating  $\Theta$  under the loss function (2) for the case  $\mathbf{V} = \mathbf{I}_p$ . Let  $\delta = \mathbf{X} + \mathbf{G}$ , where  $\mathbf{G} = \mathbf{G}(\mathbf{X})$  is a  $p \times m$  matrix valued function of  $\mathbf{X}$ . Also let  $\nabla$  be a  $p \times m$  matrix with  $(i, j)$  element equal to the differential operator  $\frac{\partial}{\partial x_{i,j}}$ . We obtain a general condition for minimaxity.

THEOREM 4. Suppose

$$\delta = \mathbf{X} + \mathbf{G}$$

is minimax under (3). Then it is also minimax under (2).

PROOF. The difference of the risks of  $\delta$  and  $\mathbf{X}$  is

$$R(\delta; \Theta) - R(\mathbf{X}; \Theta) = g(\mathbf{X}; \Theta) - g(\delta; \Theta),$$

where

$$g(\delta; \Theta) = E_{\Theta} \left[ e^{-\frac{\beta(1-\beta)\text{tr}(\delta-\Theta)'(\delta-\Theta)}{2}} \right].$$

Using  $e^{-x} \geq 1 - x$ , we see that

$$g(\delta; \Theta) \geq g(\mathbf{X}; \Theta) - \frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} E [\text{tr}(\mathbf{G}'\mathbf{G}) + 2\text{tr}(\mathbf{G}'(\mathbf{X} - \Theta))],$$

where  $E_{\Theta}$  is taken with respect to  $N_{p \times m} \left( \Theta, \frac{\mathbf{I}_p}{1 + \beta(1-\beta)} \otimes \mathbf{I}_m \right)$ . Using Theorem 2 in Ghosh and Shieh (1991), we can write

$$g(\delta; \Theta) - g(\mathbf{X}; \Theta) \geq -\frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} [R_1(\delta; \Theta) - R_1(\mathbf{X}; \Theta)].$$

Hence, if  $R_1(\delta; \Theta) \leq R_1(\mathbf{X}; \Theta)$ , then  $R(\delta; \Theta) \leq R(\mathbf{X}; \Theta)$ . □

By Theorem 4, if  $E_{\Theta}[\text{tr}(\mathbf{G}'\mathbf{G}) + 2\text{tr}(\mathbf{G}'(\mathbf{X} - \Theta))] \leq 0$  then  $R(\boldsymbol{\delta}; \Theta) \leq R(\mathbf{X}; \Theta)$ . Using Stein (1973)'s identities,  $E_{\Theta} \left[ \frac{\text{tr}(\mathbf{G}'(\mathbf{X} - \Theta))}{1 + \beta(1 - \beta)} \right] = E_{\Theta}[\text{tr}(\nabla \mathbf{G}')]$ , so we can write

$$E_{\Theta}[\text{tr}(\mathbf{G}'\mathbf{G}) + 2\text{tr}(\mathbf{G}'(\mathbf{X} - \Theta))] = E \{ \text{tr}(\mathbf{G}'\mathbf{G}) + 2[1 + \beta(1 - \beta)]\text{tr}(\nabla \mathbf{G}') \}.$$

So, if  $\text{tr}(\mathbf{G}'\mathbf{G}) + 2[1 + \beta(1 - \beta)]\text{tr}(\nabla \mathbf{G}') \leq 0$  with a positive probability for some  $\Theta$ , then  $\boldsymbol{\delta}$  is minimax under the loss function (2).

Now, assume the class of shrinkage estimators

$$\boldsymbol{\delta} = [\mathbf{I}_p - \mathbf{U}\mathbf{F}^{-1}\boldsymbol{\Psi}(\mathbf{F})\mathbf{U}']\mathbf{X}, \tag{5}$$

where  $\boldsymbol{\Psi}(\mathbf{F}) = \text{diag}(\psi_1, \dots, \psi_p)$  is a diagonal matrix with  $\psi_i$  being functions of  $\mathbf{F} = \mathbf{L}^2$ . We obtain the following by applying Theorem 4.

COROLLARY 5. Suppose  $\psi_i$  satisfy

I. For fixed  $f_1, \dots, f_{i-1}, f_{i+1}, \dots, f_p$ ,  $\frac{\partial \psi_i}{\partial f_i} \geq 0$ , where  $i = 1, \dots, p$ ,

II.  $0 \leq \psi_p \leq \dots \leq \psi_1 \leq 2(m - p - 1)$ .

Then,  $\boldsymbol{\delta}$  in (5) is minimax under the loss function (2).

PROOF. With the notation  $\boldsymbol{\Phi}(\mathbf{F}) = (\mathbf{F}')^{-1}\boldsymbol{\Psi}(\mathbf{F})$ , Stein (1973) showed that  $\boldsymbol{\delta} = [\mathbf{I}_p - \mathbf{U}\boldsymbol{\Phi}(\mathbf{F})\mathbf{U}']\mathbf{X}$  is minimax under conditions (I) and (II) when the loss function is (3). Using Theorem 4, we can see that  $\boldsymbol{\delta} = \mathbf{X} + \mathbf{G}$  is minimax under (2), where  $\mathbf{G} = -\mathbf{U}\boldsymbol{\Phi}(\mathbf{F})\mathbf{U}'\mathbf{X}$ . □

Tsukuma (2008) showed that the proper Bayes estimators with respect to the following prior

$$\Theta \sim N_{p \times m}(\mathbf{0}_{p \times m}, \Lambda^{-1}(\mathbf{I}_p - \Lambda) \otimes \mathbf{I}_m)$$

are of the form  $\boldsymbol{\delta} = [\mathbf{I}_p - \mathbf{U}\mathbf{F}^{-1}\boldsymbol{\Psi}(\mathbf{F})\mathbf{U}']\mathbf{X}$ , where  $\Lambda$  is a  $p \times m$  random matrix with the density function

$$\pi(\Lambda) \propto |\Lambda|^{\frac{a}{2}-1} \mathbf{I}(\mathbf{0}_{p \times p} < \Lambda < \mathbf{I}_p).$$

Under certain  $a$  and loss function (3), these estimators are minimax, thus we can see from Corollary 5 that  $\boldsymbol{\delta}$  is minimax under (2).

3. ESTIMATION OF THE MEAN MATRIX WHEN  $\mathbf{V} = \sigma^2 \mathbf{I}_p$  AND  $\sigma^2$  IS UNKNOWN

Here, we assume that

$$\mathbf{X} \sim N_{p \times m}(\boldsymbol{\Theta}, \sigma^2 \mathbf{I}_p \otimes \mathbf{I}_m),$$

where  $\sigma^2$  is unknown. So,  $S$  is independent of  $X$  and

$$S \sim \sigma^2 \chi_n^2,$$

where  $\chi_n^2$  denotes a chi-square random variable with  $n$  degrees of freedom. We consider estimation of  $\boldsymbol{\Theta}$  under

$$\begin{aligned} L(\mathbf{a}; \boldsymbol{\Theta}) &= \frac{1}{\beta(1-\beta)} \left[ 1 - \int_{\mathbf{R}^{p \times m}} f^\beta(\mathbf{X}|\mathbf{a}) f^{1-\beta}(\mathbf{X}|\boldsymbol{\Theta}) d\mathbf{x} \right] \\ &= \frac{1}{\beta(1-\beta)} \left[ 1 - e^{-\frac{\beta(1-\beta)\text{tr}((\mathbf{a}-\boldsymbol{\Theta})'(\mathbf{a}-\boldsymbol{\Theta}))}{2\sigma^2}} \right]. \end{aligned} \tag{6}$$

Similar to Lemma 1, one can show that the maximum likelihood estimator,  $X$ , is minimax under the loss function (6). Tsukuma (2010) obtained general conditions for minimaxity of estimators having the form  $\boldsymbol{\delta}_1 = \mathbf{X} + \mathbf{G}(\mathbf{X}, S)$ , where  $\mathbf{G}(\mathbf{X}, S)$  is a  $p \times m$  matrix valued function of  $\mathbf{X}$  and  $S$  under the quadratic loss function

$$L_2(\mathbf{a}; \boldsymbol{\Theta}) = \frac{\text{tr}((\mathbf{a}-\boldsymbol{\Theta})'(\mathbf{a}-\boldsymbol{\Theta}))}{\sigma^2}. \tag{7}$$

**THEOREM 6.** *The estimator  $\boldsymbol{\delta}_1 = \mathbf{X} + \mathbf{G}(\mathbf{X}, S)$  is minimax under the quadratic loss function (6) if it is minimax under (7) for  $m > p + 1$ .*

**PROOF.** We have

$$R(\boldsymbol{\delta}_1; \boldsymbol{\Theta}) - R(\mathbf{X}; \boldsymbol{\Theta}) = \frac{1}{\beta(1-\beta)} [g(\mathbf{X}; \boldsymbol{\Theta}) - g(\boldsymbol{\delta}_1; \boldsymbol{\Theta})].$$

Using  $e^{-x} \geq 1 - x$ ,

$$g(\boldsymbol{\delta}_1; \boldsymbol{\Theta}) \geq g(\mathbf{X}; \boldsymbol{\Theta}) - \frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} E_{\boldsymbol{\Theta}} \left[ \frac{\text{tr}(\mathbf{G}'\mathbf{G}) + 2\text{tr}(\mathbf{G}'(\mathbf{X}-\boldsymbol{\Theta}))}{\sigma^2} \right]$$

and

$$g(\boldsymbol{\delta}_1; \boldsymbol{\Theta}) - g(\mathbf{X}; \boldsymbol{\Theta}) \geq -\frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} [R_2(\boldsymbol{\delta}_1; \boldsymbol{\Theta}) - R_2(\mathbf{X}; \boldsymbol{\Theta})],$$

where

$$R_2(\boldsymbol{\delta}; \boldsymbol{\Theta}) = E_{\boldsymbol{\Theta}} [L_2(\boldsymbol{\delta}; \boldsymbol{\Theta})].$$

So, if

$$R_2(\delta_1; \Theta) \leq R_2(\mathbf{X}; \Theta),$$

then

$$R(\delta_1; \Theta) \leq R(\mathbf{X}; \Theta).$$

The proof is complete. □

Tsukuma (2010) showed that if

$$\frac{(n-2)\text{tr}(\mathbf{G}'\mathbf{G})}{s} + 2\text{tr}(\nabla\mathbf{G}') + \frac{\partial \text{tr}(\mathbf{G}'\mathbf{G})}{\partial s} \leq 0$$

then  $\delta_1 = \mathbf{X} + \mathbf{G}(\mathbf{X}, S)$  is minimax under the loss function (7). So, using Theorem 6, we see that  $\delta$  is minimax under the loss function (6).

**THEOREM 7.** *Consider the estimator*

$$\delta_2(\mathbf{X}, S) = [\mathbf{I}_p - \mathbf{U}\mathbf{F}^{-1}\mathbf{\Psi}(\mathbf{F}, S)\mathbf{U}']\mathbf{X},$$

where  $\mathbf{F} = \frac{\mathbf{I}}{s} = \text{diag}(f_1, \dots, f_p)$  and  $\mathbf{\Psi}(\mathbf{F}, s) = \text{diag}(\psi_1, \dots, \psi_p)$  are diagonal matrices with  $\psi_i$  being functions of  $\mathbf{F}$  and  $s$ . If the following conditions hold, then  $\delta_2$  is minimax under the loss function (7)

- I. For fixed  $\mathbf{F}$ ,  $\frac{\partial \psi_i}{\partial s} \leq 0$ ,
- II. For fixed  $f_1, \dots, f_{j-1}, f_{j+1}, \dots, f_p$ ,  $\frac{\partial \psi_i}{\partial f_j} \geq 0$ , where  $i, j = 1, \dots, p$ ,
- III.  $0 \leq \psi_p \leq \dots \leq \psi_1 \leq \frac{2(m-p-1)}{n+2}$ .

**PROOF.** Tsukuma (2010) showed that  $\delta_2$  is minimax under conditions I-III when the loss function is (7). So, using Theorem 6, we see that  $\delta_2$  is minimax under the loss function (6) if conditions I-III hold. □

#### 4. ESTIMATION OF THE MEAN MATRIX WHEN $\mathbf{V}$ IS UNKNOWN

Here, we assume

$$\mathbf{X} \sim N_{p \times m}(\Theta, \mathbf{V} \otimes \mathbf{I}_m),$$

where  $\mathbf{V}$  is unknown. So,  $\hat{\mathbf{V}}$  is independent of  $X$  and

$$\hat{\mathbf{V}} \sim W_p(\mathbf{V}, n),$$

where  $W_p(\mathbf{V}, n)$  denotes a Wishart random vector with  $n$  degrees of freedom and mean  $n\mathbf{V}$ . We consider estimating  $\Theta$  under the loss function

$$\begin{aligned} L(\mathbf{a}; \Theta) &= \frac{1}{\beta(1-\beta)} \left[ 1 - \int_{\mathbf{R}^{p \times m}} f^\beta(\mathbf{X}|\mathbf{a}) f^{1-\beta}(\mathbf{X}|\Theta) d\mathbf{x} \right] \\ &= \frac{1}{\beta(1-\beta)} \left[ 1 - e^{-\frac{\beta(1-\beta)\text{tr}((\mathbf{a}-\Theta)' \mathbf{V}^{-1}(\mathbf{a}-\Theta))}{2}} \right]. \end{aligned} \tag{8}$$

Similar to Lemma 1, one can show that the maximum likelihood estimator,  $X$ , is minimax under the loss function (8).

We now consider the relationship between this loss function and the following quadratic loss function

$$L_3(\mathbf{a}; \Theta, \mathbf{V}) = \text{tr}((\mathbf{a}-\Theta)' \mathbf{Q}^*(\mathbf{a}-\Theta)), \tag{9}$$

where  $\mathbf{Q}^* = \mathbf{V}^{-\frac{1}{2}} \mathbf{Q} \mathbf{V}^{-\frac{1}{2}}$ . Assume  $\delta_3 = \mathbf{X} + \mathbf{G}(\mathbf{X}, \widehat{\mathbf{V}})$ , where  $\mathbf{G}(\mathbf{X}, \widehat{\mathbf{V}})$  is a  $p \times m$  matrix valued function of  $\mathbf{X}$  and  $\widehat{\mathbf{V}}$ .

**THEOREM 8.** *The estimator  $\delta_3 = \mathbf{X} + \mathbf{G}(\mathbf{X}, \widehat{\mathbf{V}})$  is minimax under the loss function (8) if it is minimax under (9).*

**PROOF.** We have

$$R(\delta_3; \Theta) - R(\mathbf{X}; \Theta) = \frac{1}{\beta(1-\beta)} [g(\mathbf{X}; \Theta) - g(\delta_3; \Theta)],$$

where

$$g(\delta_3; \Theta) = E \left[ e^{-\frac{\beta(1-\beta)\text{tr}((\delta_3-\Theta)' \mathbf{V}^{-1}(\delta_3-\Theta))}{2}} \right].$$

So,

$$\begin{aligned} g(\delta_3; \Theta) \geq g(\mathbf{X}; \Theta) - \frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} \\ E_\Theta [\text{tr}(\mathbf{G}' \mathbf{Q}^* \mathbf{G}) + 2 \text{tr}(\mathbf{G}' \mathbf{Q}^*(\mathbf{X} - \Theta))], \end{aligned}$$

where  $E_\Theta$  is taken with respect to  $N_{p \times m}(\Theta, \frac{\mathbf{V}}{1+\beta(1-\beta)} \otimes \mathbf{I}_m)$ . Using

$$R_3(\delta; \Theta) = E_\Theta [L_3(\delta; \Theta)],$$

we can write

$$g(\delta_3; \Theta) - g(\mathbf{X}; \Theta) \geq -\frac{\beta(1-\beta)}{2} [1 + \beta(1-\beta)]^{-\frac{pm}{2}} [R_3(\delta_3; \Theta) - R_3(\mathbf{X}; \Theta)].$$

Hence, if  $R_3(\delta_3; \Theta) \leq R_3(\mathbf{X}; \Theta)$ , then  $R(\delta_3; \Theta) \leq R(\mathbf{X}; \Theta)$ . □

Shieh (1993) considered empirical Bayes estimators of the form

$$\widehat{\Theta}_{1EB} = [\mathbf{I}_p - \widehat{\mathbf{V}}\mathbf{S}^{-1}\tau(\widehat{\mathbf{V}}, \mathbf{S})]\mathbf{X},$$

where  $\mathbf{S} = \mathbf{X}\mathbf{X}'$  and  $\tau(\widehat{\mathbf{V}}, \mathbf{S})$  is a symmetric matrix. Shieh (1993) showed that  $\widehat{\Theta}_{1EB}$  is better than the maximum likelihood estimator,  $X$ , under (8). For example, he showed that if  $\tau(\widehat{\mathbf{V}}, \mathbf{S}) = (m - p - 1)\mathbf{I}_p$ , then  $\widehat{\Theta}_{EB}$  improves on  $X$  for  $m > p + 1$ .

For the case  $\mathbf{Q} = \mathbf{V} = \mathbf{I}_p$ , Konno (1990), Konno (1991), Konno (1992) considered the estimators

$$\delta_4 = \begin{cases} [\mathbf{I}_p - \mathbf{R}\mathbf{F}^{-1}\Phi(\mathbf{F})\mathbf{R}']\mathbf{X}, & \text{if } m < p, \\ [\mathbf{I}_p - \mathbf{Q}\mathbf{F}^{-1}\Phi(\mathbf{F})\mathbf{R}\mathbf{Q}^{-1}]\mathbf{X}, & \text{if } m \geq p, \end{cases}$$

where, for  $p < m$ ,  $\mathbf{X}\widehat{\mathbf{V}}^{-1}\mathbf{X}' = \mathbf{R}\mathbf{F}\mathbf{R}'$ ,  $\mathbf{R} \in \mathbf{O}_p$  and, for  $p \geq m$ ,  $\mathbf{Q}\widehat{\mathbf{V}}^{-1}\mathbf{Q}' = \mathbf{I}_p$ ,  $\mathbf{Q}'\mathbf{X}\mathbf{X}'\mathbf{Q} = \mathbf{F} = \text{diag}(f_1, \dots, f_{m \wedge p})$ ,  $f_1 \geq \dots \geq f_{m \wedge p} > 0$ . Under the following conditions, Konno (1990), Konno (1991), Konno (1992) showed that  $\delta_4$  improves on  $X$  when the loss function is (8):

- i) For  $i = 1, \dots, p$ ,  $\phi_i(\mathbf{F})$  is nondecreasing in  $f_i$ ;
- ii)  $0 \leq \phi_{m \wedge p}(\mathbf{F}) \leq \phi_{m \wedge p - 1}(\mathbf{F}) \leq \dots \leq \phi_1(\mathbf{F}) \leq \frac{2(m \vee p - m \wedge p - 1)}{n + (2m - p) \wedge p + 1}$ .

So, using Theorem 8, we see that these estimators are minimax under (8).

### 5. SIMULATION STUDY

As mentioned in Section 1, Zinodiny *et al.* (2017) estimated  $\Theta$  under the balanced loss function when  $\mathbf{V} = \sigma^2\mathbf{I}_p$ , where  $\sigma^2$  is unknown. Here, we perform a simulation study to compare this estimator with the estimator under the divergence loss function, see Theorem 8. The simulation study was performed as follows:

1. simulate 10000 random samples each of size  $n$  from a matrix variate normal distribution with zero means and  $\mathbf{V} = \sigma^2\mathbf{I}_p$ ;
2. compute Zinodiny *et al.* (2017)'s estimator and the estimator in Theorem 8 for each of the samples, say  $\{\widehat{\Theta}_1, \widehat{\Theta}_2, \dots, \widehat{\Theta}_{10000}\}$  and  $\{\widetilde{\Theta}_1, \widetilde{\Theta}_2, \dots, \widetilde{\Theta}_{10000}\}$ ;
3. compute the biases of Zinodiny *et al.* (2017)'s estimator and the estimator in Theorem 8 as

$$\frac{1}{10000mp} \sum_{i=1}^{10000} \sum_{j=1}^p \sum_{k=1}^m \widehat{\Theta}_{i,j,k}$$

and

$$\frac{1}{10000mp} \sum_{i=1}^{10000} \sum_{j=1}^p \sum_{k=1}^m \tilde{\Theta}_{i,j,k};$$

4. compute the mean squared errors of Zinodiny *et al.* (2017)'s estimator and the estimator in Theorem 8 as

$$\frac{1}{10000mp} \sum_{i=1}^{10000} \sum_{j=1}^p \sum_{k=1}^m \hat{\Theta}_{i,j,k}^2$$

and

$$\frac{1}{10000mp} \sum_{i=1}^{10000} \sum_{j=1}^p \sum_{k=1}^m \tilde{\Theta}_{i,j,k}^2.$$

We repeated this procedure for  $n = 10, 12, \dots, 109$ . We chose  $\sigma = 1$ ,  $m = 10$  and  $p = 5$ .

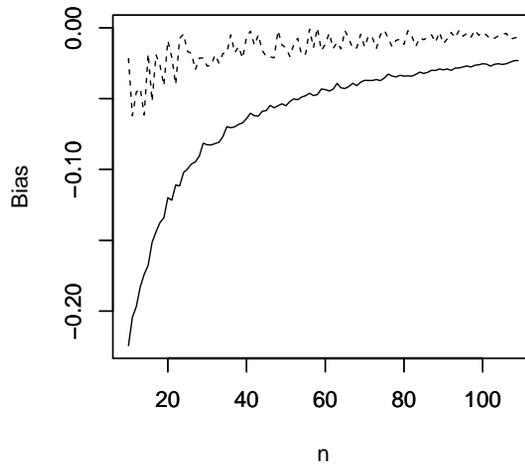


Figure 1 – Biases of Zinodiny *et al.* (2017)'s estimator (solid curve) and the estimator in Theorem 8 (broken curve) versus  $n$ .

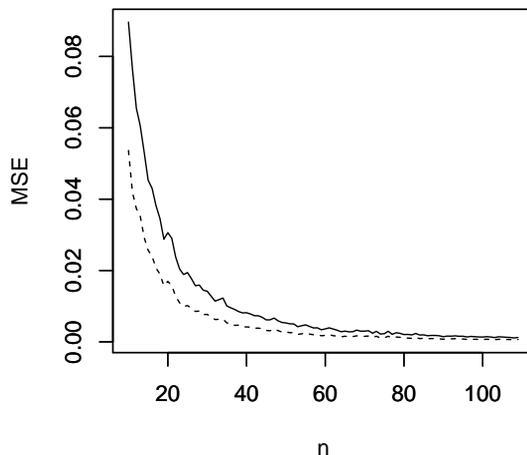


Figure 2 – Mean squared errors of Zinodiny *et al.* (2017)'s estimator (solid curve) and the estimator in Theorem 8 (broken curve) versus  $n$ .

The biases versus  $n$  are drawn in Figure 1. The mean squared errors versus  $n$  are drawn in Figure 2. We see that the estimator in Theorem 8 has consistently smaller bias and consistently smaller mean squared error for every  $n$ . We took  $\sigma = 1$ ,  $m = 10$  and  $p = 5$  in the simulations. But the same observation was noted for a wide range of other values of  $\sigma$ ,  $m$  and  $p$ . Hence, the estimator in Theorem 8 is a better estimator with respect to both bias and mean squared error.

## 6. CONCLUSIONS

We have considered the problem of estimating the mean of a matrix variate normal distribution. We have provided minimax estimators under the divergence loss function. The divergence loss function has several attractive features compared to other loss functions considered in the literature.

We have given minimax estimators for the mean considering different forms for the covariance matrix. The forms considered include the cases that the covariance matrix is known and the covariance matrix is completely unknown.

We have performed a simulation study to compare one of our estimators with the most recently proposed estimator. The simulation shows that our estimator has consistently smaller bias and consistently smaller mean squared error.

## ACKNOWLEDGEMENTS

The authors would like to thank the referee and the Editor for careful reading and comments which greatly improved the paper.

## REFERENCES

- S. AMARI (1982). *Differential geometry of curved exponential families-curvatures and information loss*. Annals of Statistics, 10, pp. 357–387.
- M. BILODEAU, T. KARIYA (1989). *Minimax estimators in the normal manova model*. Journal of Multivariate Analysis, 28, pp. 260–270.
- C. R. BLYTH (1951). *On minimax statistical decision procedures and their admissibility*. Annals of Mathematical Statistics, 22, pp. 22–42.
- B. CHEN, S. H. LIN (2012). *A risk-aware modeling framework for speech summarization*. IEEE Transactions on Audio, Speech, and Language Processing, 20, pp. 211–222.
- M. CORDUAS (1985). *On the divergence between linear processes*. Statistica, 45, pp. 393–401.
- N. CRESSIE, T. R. C. READ (1984). *Multinomial goodness-of-fit tests*. Journal of the Royal Statistical Society Series B, 46, pp. 440–464.
- B. EFRON, C. MORRIS (1972). *Empirical Bayes on vector observations: An extension of Stein's method*. Biometrika, 59, pp. 335–347.
- M. GHOSH, V. MERGEL (2009). *On the stein phenomenon under divergence loss and an unknown variance-covariance matrix*. Journal of Multivariate Analysis, 100, pp. 2331–2336.
- M. GHOSH, V. MERGEL, G. S. DATTA (2008). *Estimation, prediction and the stein phenomenon under divergence loss*. Journal of Multivariate Analysis, 99, pp. 1941–1961.
- M. GHOSH, G. SHIEH (1991). *Empirical bayes minimax estimators of matrix normal means*. Journal of Multivariate Analysis, 38, pp. 306–318.
- S. JI, W. ZHANG, R. LI (2013). *A probabilistic latent semantic analysis model for coclustering the mouse brain atlas*. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 10, pp. 1460–1468.
- R. L. KASHYAP (1974). *Minimax estimation with divergence loss function*. Information Sciences, 7, pp. 341–364.
- Y. KONNO (1990). *Families of minimax estimators of matrix of normal means with unknown covariance matrix*. Journal of the Japan Statistical Society, 20, pp. 191–201.

- Y. KONNO (1991). *On estimation of a matrix of normal means with unknown covariance matrix*. Journal of Multivariate Analysis, 36, pp. 44–55.
- Y. KONNO (1992). *Improved estimation of matrix of normal mean and eigenvalues in the multivariate F-distribution*. Ph.D. thesis, Institute of Mathematics, University of Tsukuba.
- E. L. LEHMANN, G. CASELLA (1998). *Theory of Point Estimation*. Springer Verlag, New York, 2 ed.
- J. MALKIN, X. LI, S. HARADA, J. LANDAY, J. BILMES (2011). *The vocal joystick engine v1.0*. Computer Speech and Language, 25, pp. 535–555.
- J. PAUL, P. Y. THOMAS (2016). *Sharma-mittal entropy properties on record values*. Statistica, 76, pp. 273–287.
- C. P. ROBERT (2001). *The Bayesian Choice, second edition*. Springer Verlag, New York.
- G. SHIEH (1993). *Empirical bayes minimax estimators of matrix normal means for arbitrary quadratic loss and unknown covariance matrix*. Statistics and Decisions, 11, pp. 317–341.
- C. STEIN (1973). *Estimation of the mean of a multivariate normal distribution*. In J. HÁJEK (ed.), *Proceedings of the Prague Symposium on Asymptotic statistics*. Universita Karlova, Prague, pp. 345–381.
- H. TSUKUMA (2008). *Admissibility and minimaxity of Bayes estimators for a normal mean matrix*. Journal of Multivariate Analysis, 99, pp. 2251–2264.
- H. TSUKUMA (2010). *Proper Bayes minimax estimators of the normal mean matrix with common unknown variances*. Journal of Statistical Planning and Inference, 140, pp. 2596–2606.
- L. YANG, B. GENG, A. HANJALIC, X. S. HUA (2012). *A unified context model for web image retrieval*. ACM Transactions on Multimedia Computing, Communications, and Applications, 8. Article No. 28.
- Z. ZHANG (1986a). *On estimation of matrix of normal mean*. Journal of Multivariate Analysis, 18, pp. 70–82.
- Z. ZHANG (1986b). *Selecting a minimax estimator doing well, at a point*. Journal of Multivariate Analysis, 19, pp. 14–23.
- S. ZINODINY, S. REZAEI, S. NADARAJAH (2017). *Bayes minimax estimation of the mean matrix of matrix-variate normal distribution under balanced loss function*. Statistics and Probability Letters, 125, pp. 110–120.

## SUMMARY

The problem of estimating the mean matrix of a matrix-variate normal distribution with a covariance matrix is considered under two loss functions. We construct a class of empirical Bayes estimators which are better than the maximum likelihood estimator under the first loss function and hence show that the maximum likelihood estimator is inadmissible. We find a general class of minimax estimators. Also we give a class of estimators that improve on the maximum likelihood estimator under the second loss function and hence show that the maximum likelihood estimator is inadmissible.

*Keywords:* Empirical Bayes estimation; Matrix variate normal distribution; Mean matrix; Minimax estimation.