

SPRINGER BRIEFS IN STATISTICS  
JSS RESEARCH SERIES IN STATISTICS

Hisayuki Tsukuma  
Tatsuya Kubokawa

# Shrinkage Estimation for Mean and Covariance Matrices



Springer

# **SpringerBriefs in Statistics**

## **JSS Research Series in Statistics**

### **Editors-in-Chief**

Naoto Kunitomo, Economics, Meiji University, Chiyoda-ku, Tokyo, Tokyo, Japan

Akimichi Takemura, The Center for Data Science Education and Research, Shiga University, Bunkyo-ku, Tokyo, Japan

### **Series Editors**

Genshiro Kitagawa, Meiji Institute for Advanced Study of Mathematical Sciences, Nakano-ku, Tokyo, Japan

Tomoyuki Higuchi, Faculty of Science and Engineering, Chuo University, Tokyo, Japan

Toshimitsu Hamasaki, Office of Biostatistics and Data Management, National Cerebral and Cardiovascular Center, Suita, Osaka, Japan

Shigeyuki Matsui, Graduate School of Medicine, Nagoya University, Nagoya, Aichi, Japan

Manabu Iwasaki, School of Data Science, Yokohama City University, Yokohama, Tokyo, Japan

Yasuhiro Omori, Graduate School of Economics, The University of Tokyo, Bunkyo-ku, Tokyo, Japan

Masafumi Akahira, Institute of Mathematics, University of Tsukuba, Tsukuba, Ibaraki, Japan

Takahiro Hoshino, Department of Economics, Keio University, Tokyo, Japan

Masanobu Taniguchi, Department of Mathematical Sciences/School, Waseda University/Science & Engineering, Shinjuku-ku, Japan

The current research of statistics in Japan has expanded in several directions in line with recent trends in academic activities in the area of statistics and statistical sciences over the globe. The core of these research activities in statistics in Japan has been the Japan Statistical Society (JSS). This society, the oldest and largest academic organization for statistics in Japan, was founded in 1931 by a handful of pioneer statisticians and economists and now has a history of about 80 years. Many distinguished scholars have been members, including the influential statistician Hirotugu Akaike, who was a past president of JSS, and the notable mathematician Kiyosi Itô, who was an earlier member of the Institute of Statistical Mathematics (ISM), which has been a closely related organization since the establishment of ISM. The society has two academic journals: the Journal of the Japan Statistical Society (English Series) and the Journal of the Japan Statistical Society (Japanese Series). The membership of JSS consists of researchers, teachers, and professional statisticians in many different fields including mathematics, statistics, engineering, medical sciences, government statistics, economics, business, psychology, education, and many other natural, biological, and social sciences. The JSS Series of Statistics aims to publish recent results of current research activities in the areas of statistics and statistical sciences in Japan that otherwise would not be available in English; they are complementary to the two JSS academic journals, both English and Japanese. Because the scope of a research paper in academic journals inevitably has become narrowly focused and condensed in recent years, this series is intended to fill the gap between academic research activities and the form of a single academic paper. The series will be of great interest to a wide audience of researchers, teachers, professional statisticians, and graduate students in many countries who are interested in statistics and statistical sciences, in statistical theory, and in various areas of statistical applications.

More information about this subseries at <http://www.springer.com/series/13497>

Hisayuki Tsukuma · Tatsuya Kubokawa

# Shrinkage Estimation for Mean and Covariance Matrices



Springer

Hisayuki Tsukuma  
Faculty of Medicine  
Toho University  
Tokyo, Japan

Tatsuya Kubokawa  
Faculty of Economics  
University of Tokyo  
Tokyo, Japan

ISSN 2191-544X  
SpringerBriefs in Statistics

ISSN 2364-0057  
JSS Research Series in Statistics

ISBN 978-981-15-1595-8

ISSN 2191-5458 (electronic)

ISSN 2364-0065 (electronic)

ISBN 978-981-15-1596-5 (eBook)

<https://doi.org/10.1007/978-981-15-1596-5>

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2020

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd.  
The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

# Preface

The rapid development of computer technology has started to yield many types of high-dimensional data and to enable us to deal with them well. Indeed, high-dimensional data appear in numerous fields such as web data science, genomics, telecommunication, atmospheric science, financial engineering, and others. With such a background, theory of statistical inference with high dimension has received much attention in recent years.

High-dimensional data in general are hard to handle, and ordinary or traditional methods in statistics are frequently inapplicable for them. This has inspired statisticians to develop new methodology in high dimension from both theoretical and practical aspects. Most statisticians' interests seem to be in development of efficient algorithms for statistical inference and in investigation of their asymptotic properties with the dimension going to infinity. On the other hand, there does not exist much literature in high-dimensional problems from a decision-theoretic point of view.

Statistical decision theory is the study of how to make decisions in the presence of statistical knowledge under uncertainty. It has been studied from around the 1940s and the researchers have already been produced many important and interesting results. Probably the most surprising result in decision-theoretic estimation is the inadmissibility of the sample mean vector to estimate a multivariate normal population mean. In the multivariate normal mean estimation, the sample mean vector is the maximum likelihood estimator and the uniformly minimum variance unbiased estimator, and thus it has been recognized to be optimal for a long time. However, in 1956, Charles Stein showed that the sample mean vector is admissible for the one- and two-dimensional cases but inadmissible for three or more dimensional cases. A little after that, a specific estimator, called a shrinkage estimator, was provided for exactly dominating the sample mean vector. To this day, various extensions of shrinkage estimation have been achieving in other statistical models.

The purpose of this book is to give a brief overview of shrinkage estimation in matrix-variate normal distribution model. More specifically, it includes recent techniques and results in estimation of mean and covariance matrices with a

high-dimensional setting that implies singularity of the sample covariance matrix. Such a high-dimensional model can really be analyzed by using the same arguments as for a low-dimensional model. Thus this book takes a unified approach to both high- and low-dimensional shrinkage estimation.

Theory of shrinkage estimation for matrix parameters needs many mathematical tools. In Chap. 1, we begin by briefly introducing basic terminology of decision-theoretic estimation and a mathematical technique in shrinkage estimation. Chapter 2 defines the notation with respect to matrix algebra and collects useful results in terms of the Moore-Penrose inverse, the Kronecker product and matrix decompositions. Chapter 3 provides the definition and some properties of matrix-variate normal distribution and related distributions, including the Wishart distribution and joint distributions corresponding to the Cholesky and the eigenvalue decompositions of the Wishart matrix. With a unified treatment for high- and low-dimensional cases, some related distributions are discussed. Chapter 4 introduces a multivariate linear model and derives its canonical form. To find decision-theoretically optimal estimators, we usually direct our attention to several classes of invariant estimators. Therefore Chap. 4 briefly explains group invariance in the canonical form as well. A key tool in shrinkage estimation is an integration by parts formula, called the Stein identity. Chapter 5 gives a generalized Stein identity on matrix-variate normal distribution. Moreover we list some results on matrix differential operators and in particular show useful differentiation formulae concerning the Moore-Penrose inverse. Chapter 6 addresses the problem of estimating the mean matrix in matrix-variate normal distribution model. A unified result on matricial shrinkage estimation is presented, and extensions and applications are given for more general models. Chapter 7 deals with the problem of estimating the covariance matrix relative to an extended Stein loss and provides various unified estimation procedures for high- and low-dimensional cases. Some related topics to covariance estimation are also touched.

The authors would like to thank Prof. M. Akahira for giving us the opportunity of publishing this book. The work of the first author was supported in part by Grant-in-Aid for Scientific Research (18K11201) from the Japan Society for the Promotion of Science (JSPS). The work of the second author was supported in part by Grant-in-Aid for Scientific Research (18K11188) from the JSPS.

Tokyo, Japan  
March 2020

Hisayuki Tsukuma  
Tatsuya Kubokawa

# Contents

<b>1</b>	<b>Decision-Theoretic Approach to Estimation</b>	<b>1</b>
1.1	Decision-Theoretic Framework for Estimation	1
1.2	James-Stein's Shrinkage Estimator	2
1.3	Unbiased Risk Estimate and Stein's Identity	3
	References	4
<b>2</b>	<b>Matrix Algebra</b>	<b>7</b>
2.1	Notation	7
2.2	Nonsingular Matrix and the Moore-Penrose Inverse	9
2.3	Kronecker Product and Vec Operator	10
2.4	Matrix Decompositions	11
	References	12
<b>3</b>	<b>Matrix-Variate Distributions</b>	<b>13</b>
3.1	Preliminaries	13
3.1.1	The Multivariate Normal Distribution	13
3.1.2	Jacobians of Matrix Transformations	14
3.1.3	The Multivariate Gamma Function	16
3.2	The Matrix-Variate Normal Distribution	17
3.3	The Wishart Distribution	21
3.4	The Cholesky Decomposition of the Wishart Matrix	23
	References	26
<b>4</b>	<b>Multivariate Linear Model and Group Invariance</b>	<b>27</b>
4.1	Multivariate Linear Model	27
4.2	A Canonical Form	30
4.3	Group Invariance	31
	References	33



<b>5</b>	<b>A Generalized Stein Identity and Matrix Differential Operators . . .</b>	<b>35</b>
5.1	Stein's Identity in Matrix-Variate Normal Distribution . . . . .	35
5.2	Some Useful Results on Matrix Differential Operators . . . . .	37
	Appendix . . . . .	40
	References . . . . .	42
<b>6</b>	<b>Estimation of the Mean Matrix . . . . .</b>	<b>45</b>
6.1	Introduction . . . . .	45
6.2	The Unified Efron-Morris Type Estimators Including Singular Cases . . . . .	48
6.2.1	Empirical Bayes Methods . . . . .	48
6.2.2	The Unified Efron-Morris Type Estimator . . . . .	49
6.3	A Unified Class of Matricial Shrinkage Estimators . . . . .	50
6.4	Unbiased Risk Estimate . . . . .	53
6.5	Examples for Specific Estimators . . . . .	55
6.5.1	The Unified Efron-Morris Type Estimator . . . . .	55
6.5.2	A Modified Stein-Type Estimator . . . . .	56
6.5.3	Modified Efron-Morris Type Estimator . . . . .	58
6.6	Related Topics . . . . .	59
6.6.1	Positive-Part Rule Estimators . . . . .	59
6.6.2	Shrinkage Estimation with a Loss Matrix . . . . .	62
6.6.3	Application to a GMANOVA Model . . . . .	63
6.6.4	Generalization in an Elliptically Contoured Model . . . . .	67
	Appendix . . . . .	68
	References . . . . .	74
<b>7</b>	<b>Estimation of the Covariance Matrix . . . . .</b>	<b>75</b>
7.1	Introduction . . . . .	75
7.2	Scale Invariant Estimators . . . . .	77
7.3	Triangular Invariant Estimators and the James-Stein Estimator . . .	79
7.3.1	The James-Stein Estimator . . . . .	79
7.3.2	Improvement Using a Subgroup Invariance . . . . .	81
7.4	Orthogonally Invariant Estimators . . . . .	84
7.4.1	Class of Orthogonally Invariant Estimators . . . . .	84
7.4.2	Unbiased Risk Estimate . . . . .	84
7.4.3	Examples . . . . .	86
7.5	Improvement Using Information on Mean Statistic . . . . .	96
7.5.1	A Class of Estimators and Its Risk Function . . . . .	97
7.5.2	Examples of Improved Estimators . . . . .	98
7.5.3	Further Improvements with a Truncation Rule . . . . .	100
7.6	Related Topics . . . . .	102
7.6.1	Decomposition of the Estimation Problem . . . . .	102
7.6.2	Decision-Theoretic Studies Under Quadratic Losses . . . . .	104

7.6.3	Estimation of the Generalized Variance . . . . .	105
7.6.4	Estimation of the Precision Matrix . . . . .	106
	References . . . . .	108
<b>Index</b>	. . . . .	<b>111</b>

# Chapter 1

## Decision-Theoretic Approach to Estimation



Statistical decision theory has been studied from around the 1940s and the researchers have already been producing many remarkable results. In the field of decision-theoretic estimation, the most surprising result is the inadmissibility of the sample mean vector in estimation of a mean vector of multivariate normal distribution. The inadmissibility result is closely relevant to the discovery of shrinkage estimator. This chapter summarizes basic terminology of decision-theoretic estimation and shrinkage estimators in the multivariate normal mean estimation. Also, Stein's unbiased estimate of risk is briefly explained as a general method of how to find better estimators. The unbiased risk estimate method is applied to estimation of mean and covariance matrices discussed in this book.

### 1.1 Decision-Theoretic Framework for Estimation

Let  $x$  be a random vector or matrix having a probability distribution characterized by an unknown parameter  $\theta$  (possibly,  $\theta$  can be a vector or matrix). Assume that we want to estimate  $\theta$  based on  $x$ . An estimator of  $\theta$  is denoted by  $\hat{\theta} = \hat{\theta}(x)$  which is a function of  $x$ . In the literature  $\hat{\theta}$  is also called the decision rule.

Let  $\mathcal{P}$  be the parameter space. There are usually many estimation procedures for the unknown parameter  $\theta \in \mathcal{P}$  and thus we need to decide how to select an optimal procedure. From an intuitive point of view, it seems reasonable to select an estimator minimizing a mathematical distance between  $\hat{\theta}$  and  $\theta$  or making it smaller as soon as possible. In statistical decision theory, such a distance is regarded as the loss induced from  $\hat{\theta}$  by estimating  $\theta$ . For this reason, the distance is called a loss function of  $\hat{\theta}$  and  $\theta$ . The loss function of  $\hat{\theta}$  and  $\theta$  is denoted by  $L(\hat{\theta}, \theta)$ , where it is nonnegative for any  $\hat{\theta}$  and  $\theta$ . We usually employ a loss function with properties that it takes zero when  $\hat{\theta}$  is equal to  $\theta$  and increases when  $\hat{\theta}$  goes away from  $\theta$ .

However, the loss function of  $\hat{\theta}$  and  $\theta$  is random, and practically a distance between  $\hat{\theta}$  and  $\theta$  is measured by an expected loss

$$R(\hat{\theta}, \theta) = E[L(\hat{\theta}, \theta)],$$

where  $E$  denotes the expectation taken with respect to the distribution of  $x$ . Here  $R(\hat{\theta}, \theta)$  is viewed as a quantified risk induced from  $\hat{\theta}$  by estimating  $\theta$  and it is called the risk function relative to  $L(\hat{\theta}, \theta)$ .

The risk function evaluates performance of estimators and can be employed to compare two estimators. If

$$R(\hat{\theta}_0, \theta) \leq R(\hat{\theta}, \theta)$$

for any  $\theta \in \mathcal{P}$ , with strict inequality for some  $\theta$ , then an estimator  $\hat{\theta}_0$  is said to be better than  $\hat{\theta}$ , or dominate  $\hat{\theta}$ , or improve on  $\hat{\theta}$ . An estimator  $\hat{\theta}$  is said to be inadmissible if there exists another estimator  $\hat{\theta}_0$  which dominates  $\hat{\theta}$  and to be admissible if no such estimator  $\hat{\theta}_0$  exists. For example, if we want to prove inadmissibility of an estimator then it simply suffices to find another dominating estimator.

Admissibility is a fundamental criterion in decision-theoretic estimation, but in general there exist many admissible estimators. Therefore we consider another criterion called minimaxity. The minimaxity of an estimator  $\hat{\theta}_0$  implies that  $\hat{\theta}_0$  minimizes a supremum of the risk function among any estimator. Namely,  $\hat{\theta}_0$  is said to be minimax if

$$\sup_{\theta \in \mathcal{P}} R(\hat{\theta}_0, \theta) \leq \sup_{\theta \in \mathcal{P}} R(\hat{\theta}, \theta)$$

for any estimator  $\hat{\theta}$ . If an estimator  $\hat{\theta}$  is better than a minimax estimator, then  $\hat{\theta}$  is also minimax.

For a general explanation on decision-theoretic estimation, see, for example, Ferguson (1967) and Lehmann and Casella (1998). Clearly, admissibility and minimaxity strongly depend on loss functions, and there exists some criticism concerning criteria based on loss functions. Berger (1985) and Robert (2007) discussed some criticism on decision-theoretic estimation and the use of loss functions.

## 1.2 James-Stein's Shrinkage Estimator

Now, we look at the problem of estimating the mean vector  $\theta$  of  $p$ -variate normal distribution with the identity covariance matrix relative to the quadratic loss  $L(\hat{\theta}, \theta) = \|\hat{\theta} - \theta\|^2$ , where  $\|\cdot\|$  denotes the Euclidean norm and  $\hat{\theta}$  is an estimator of  $\theta$ . A random vector drawn from the  $p$ -variate normal distribution is denoted by  $X$ .

In estimation of the normal mean vector  $\theta$ , the maximum likelihood estimator is  $\hat{\theta}^{ML} = X$ , which is the uniformly minimum variance unbiased estimator. Also,  $\hat{\theta}^{ML}$  is the best invariant estimator under the group of the affine transformations

$$X \rightarrow AX + b, \quad \theta \rightarrow A\theta + b, \quad (1.1)$$

where  $A$  is a  $p \times p$  orthogonal matrix and  $b$  is a  $p$ -dimensional vector, and it is minimax with the constant risk  $p$ . From the above facts,  $\hat{\theta}^{ML}$  has been recognized to be optimal for a long time. However, Stein (1956) showed that  $\hat{\theta}^{ML}$  is inadmissible for three or more dimensional cases. Further James and Stein (1961) succeeded in providing an explicit estimator which dominates  $\hat{\theta}^{ML}$ . The dominating estimator is of the form

$$\hat{\theta}^{JS} = \left(1 - \frac{p-2}{\|X\|^2}\right) X,$$

which is invariant under a subgroup of (1.1),  $X \rightarrow AX$  and  $\theta \rightarrow A\theta$  for a  $p \times p$  orthogonal matrix  $A$ . James and Stein's estimator  $\hat{\theta}^{JS}$  is called a shrinkage estimator since it is shrinking  $\hat{\theta}^{ML}$  toward the origin. James and Stein (1961) showed that the risk function of  $\hat{\theta}^{JS}$  can be expressed as

$$R(\hat{\theta}^{JS}, \theta) = p - (p-2)^2 E[(p-2+2K)^{-1}], \quad (1.2)$$

where  $K$  is a Poisson random variable with mean  $\|\theta\|^2/2$ . The expectation in the r.h.s. of (1.2) is finite when  $p \geq 3$ , and then  $R(\hat{\theta}^{JS}, \theta) \leq R(\hat{\theta}^{ML}, \theta)$  for any  $\theta$ , namely,  $\hat{\theta}^{ML}$  is inadmissible relative to the quadratic loss  $L$ .

James and Stein's shrinkage estimator can be characterized by an empirical Bayes method as shown by Efron and Morris (1973). See Gruber (1998), who compared the shrinkage and the ML estimators for some linear models from the Bayesian and Frequentist points of view. A broad survey on shrinkage estimation is presented by Kubokawa (1998) and a modern Bayesian approach is explained extensively by Fourdrinier et al. (2018). For geometrical interpretation of the shrinkage estimator, see Brown and Zhao (2012).

### 1.3 Unbiased Risk Estimate and Stein's Identity

When  $X$  follows the  $p$ -variate normal distribution with mean  $\theta$  and the identity covariance matrix  $\|X\|^2$  is distributed as the noncentral chi-square distribution with  $p$  degrees of freedom and the noncentrality parameter  $\|\theta\|^2$ . The p.d.f. of the noncentral chi-square distribution is written as a Poisson mixture and  $E[\|X\|^{-2}] = E[(p-2+2K)^{-1}]$ , where  $K$  is the Poisson random variable defined in the previous section. The risk function of  $\hat{\theta}^{JS}$  can be rewritten as  $R(\hat{\theta}^{JS}, \theta) = E[\hat{R}^{JS}(X)]$ , where

$$\widehat{R}^{JS}(\mathbf{x}) = p - \frac{(p-2)^2}{\|\mathbf{x}\|^2}.$$

Here  $\widehat{R}^{JS}(\mathbf{x})$  is a function of  $p$ -dimensional vector  $\mathbf{x}$  but independent of  $\boldsymbol{\theta}$ . Therefore it is called the unbiased risk estimate for  $R(\widehat{\boldsymbol{\theta}}^{JS}, \boldsymbol{\theta})$ . The unbiased risk estimate for  $R(\widehat{\boldsymbol{\theta}}^{ML}, \boldsymbol{\theta})$  is given by  $\widehat{R}^{ML}(\mathbf{x}) = p$ , so that  $\widehat{R}^{JS}(\mathbf{x}) \leq \widehat{R}^{ML}(\mathbf{x})$  for any  $\mathbf{x}$  except when  $\|\mathbf{x}\| = 0$ .

The unbiased risk estimate provides simple methods of proving the inadmissibility of estimators and of finding better estimators. In fact, if there exist unbiased risk estimates  $\widehat{R}_1(\mathbf{x})$  and  $\widehat{R}_2(\mathbf{x})$  with respect to two estimators  $\widehat{\boldsymbol{\theta}}_1$  and  $\widehat{\boldsymbol{\theta}}_2$ , respectively, such that  $\widehat{R}_1(\mathbf{x}) \leq \widehat{R}_2(\mathbf{x})$  for any  $\mathbf{x}$ , then  $R(\widehat{\boldsymbol{\theta}}_1, \boldsymbol{\theta}) \leq R(\widehat{\boldsymbol{\theta}}_2, \boldsymbol{\theta})$  uniformly for any  $\boldsymbol{\theta}$ .

In general, unbiased risk estimates are hard to derive. When we consider the class of estimators of the form  $\widehat{\boldsymbol{\theta}}^G = \widehat{\boldsymbol{\theta}}^{ML} + \mathbf{G}$  with a vector-valued function  $\mathbf{G}$  of  $\mathbf{X}$ , the risk function of  $\widehat{\boldsymbol{\theta}}^G$  is expressed by

$$R(\widehat{\boldsymbol{\theta}}^G, \boldsymbol{\theta}) = R(\widehat{\boldsymbol{\theta}}^{ML}, \boldsymbol{\theta}) + E[\|\mathbf{G}\|^2 + 2(\mathbf{X} - \boldsymbol{\theta})^\top \mathbf{G}].$$

Therefore we need to evaluate  $E[(\mathbf{X} - \boldsymbol{\theta})^\top \mathbf{G}]$  for deriving an unbiased risk estimate of  $R(\widehat{\boldsymbol{\theta}}^G, \boldsymbol{\theta})$ . To this end, Stein (1973, 1981) proved a useful integration by parts formula in terms of a normal distribution: Let  $x$  be a normal random variable with mean  $\theta$  and variance one. Let  $g$  be an absolutely continuous function such that  $E[(x - \theta)g(x)]$  and  $E[g'(x)]$  are finite. Then the integration by parts formula is given by

$$E[(x - \theta)g(x)] = E[g'(x)], \quad (1.3)$$

which enables us to evaluate  $E[(\mathbf{X} - \boldsymbol{\theta})^\top \mathbf{G}]$  given above.

The integration by parts formula (1.3) is named the Stein identity or the Stein lemma after Stein (1973, 1981). The Stein identity is nowadays a key tool for evaluating risk functions and for deriving unbiased risk estimates, and it has exerted an immeasurable influence on the development of shrinkage estimation.

## References

- J.O. Berger, *Statistical Decision Theory and Bayesian Analysis*, 2nd edn. (Springer, New York, 1985)
- L.D. Brown, L.H. Zhao, A geometrical explanation of Stein shrinkage. *Stat. Sci.* **27**, 24–30 (2012)
- B. Efron, C. Morris, Stein's estimation rule and its competitors—An empirical Bayes approach. *J. Am. Stat. Assoc.* **68**, 117–130 (1973)
- T.S. Ferguson, *Mathematical Statistics: A Decision Theoretic Approach* (Academic Press, New York, 1967)
- D. Fourdrinier, W.E. Strawderman, M.T. Wells, *Shrinkage Estimation* (Springer, New York, 2018)
- M.H.J. Gruber, *Improving Efficiency by Shrinkage* (Marcel Dekker, New York, 1998)
- W. James, C. Stein, Estimation with quadratic loss, in *proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, ed. by J. Neyman, vol. 1 (University of California Press, Berkeley, 1961), pp. 361–379

- T. Kubokawa, The Stein phenomenon in simultaneous estimation: A review, in *Applied Statistical Science III*, ed. by S.E. Ahmed, M. Ahsanullah, B.K. Sinha (Nova Science Publishers, New York, 1998), pp. 143–173
- E.L. Lehmann, G. Casella, *Theory of Point Estimation*, 2nd edn. (Springer, New York, 1998)
- C. Robert, *The Bayesian Choice*, 2nd edn. (Springer, New York, 2007)
- C. Stein, Inadmissibility of the usual estimator for the mean of a multivariate normal distribution, *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, ed. by J. Neyman, vol. 1 (University of California Press, Berkeley, 1956), pp. 197–206
- C. Stein, Estimation of the mean of a multivariate normal distribution, Technical Reports No.48 (Department of Statistics, Stanford University, Stanford, 1973)
- C. Stein, Estimation of the mean of a multivariate normal distribution. *Ann. Stat.* **9**, 1135–1151 (1981)

# Chapter 2

## Matrix Algebra



Matrix algebra is an important step in mathematical treatment of shrinkage estimation for matrix parameters, and in particular the Moore-Penrose inverse and some matrix decompositions are required for defining matricial shrinkage estimators. This chapter first explains the notation used in this book and subsequently lists helpful results in matrix algebra.

### 2.1 Notation

Let  $\mathbb{R}^n$  be the  $n$ -dimensional real vector space and in particular denote  $\mathbb{R} = \mathbb{R}^1$ . Let  $\mathbb{R}^{m \times n}$  be the set of all  $m \times n$  matrices with real elements. Note that  $\mathbb{R}^{m \times n} \neq \mathbb{R}^{mn}$  and  $\mathbb{R}^{mn}$  is the  $(mn)$ -dimensional real vector space. If  $\mathbf{A} \in \mathbb{R}^{m \times n}$  then  $\mathbf{A}$  is of the form

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix},$$

where all the  $a_{ij}$ 's belong to  $\mathbb{R}$  for  $i = 1, \dots, m$  and  $j = 1, \dots, n$ . In some cases,  $\mathbf{A}$  given above is written as  $\mathbf{A} = (a_{ij})$ , or  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$  where for  $j = 1, \dots, n$

$$\mathbf{a}_j = \begin{pmatrix} a_{1j} \\ \vdots \\ a_{mj} \end{pmatrix} \in \mathbb{R}^m.$$

Also, the  $(i, j)$ -th element of  $\mathbf{A}$  is sometimes expressed as  $a_{i,j}$  or  $\{\mathbf{A}\}_{ij}$ .



Let  $\mathbf{0}_{m \times n}$  be the zero matrix of size  $m \times n$ . The diagonal matrix of size  $n \times n$  is denoted by  $\text{diag}(d_1, \dots, d_n)$ , where  $d_1, \dots, d_n$  are diagonal elements from the upper left corner to the lower right one. Define the identity matrix of size  $n \times n$  as  $\mathbf{I}_n$ , namely,  $\mathbf{I}_n = \text{diag}(1, \dots, 1)$  consisting of  $n$  ones on the diagonal.

Let  $\mathbf{A}^\top$  be the transpose of a matrix  $\mathbf{A}$  and let  $\text{tr } \mathbf{A}$  and  $|\mathbf{A}|$  be, respectively, the trace and the determinant of a square matrix  $\mathbf{A}$ . Also, let  $\mathbf{A}^{-1}$  be the inverse of a nonsingular matrix  $\mathbf{A}$ . A matrix square root of a symmetric positive semi-definite matrix  $\mathbf{A}$  is written as  $\mathbf{A}^{1/2}$ , where  $\mathbf{A}^{1/2}$  is symmetric such that  $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$ . The inverse of  $\mathbf{A}^{1/2}$  is expressed as  $\mathbf{A}^{-1/2}$  if it exists.

We list the notation for special subsets in  $\mathbb{R}^{m \times n}$  as follows:

- $\mathbb{U}_n$ : Set of all  $n \times n$  nonsingular matrices;
- $\mathbb{O}_n$ : Set of all  $n \times n$  orthogonal matrices;
- $\mathbb{V}_{m,n}$ : Set of all  $m \times n$  matrices  $\mathbf{A}$  such that  $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_n$  for  $m \geq n$ ;
- $\mathbb{D}_n$ : Set of all  $n \times n$  diagonal matrices;
- $\mathbb{D}_n^{(\geq)}$ : Set of all  $n \times n$  diagonal matrices  $\text{diag}(d_1, \dots, d_n)$  such that  $d_1 \geq \dots \geq d_n$ ;
- $\mathbb{D}_n^{(\geq 0)}$ : Set of all  $n \times n$  diagonal matrices  $\text{diag}(d_1, \dots, d_n)$  such that  $d_1 \geq \dots \geq d_n \geq 0$ ;
- $\mathbb{L}_n^{(+)}$ : Set of all  $n \times n$  lower triangular matrices with positive diagonal elements;
- $\mathbb{L}_n^{(1)}$ : Set of all  $n \times n$  lower triangular matrices with ones on the diagonal;
- $\mathbb{L}_{m,n}^{(+)}$ : Set of all  $m \times n$  matrices of the form

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{pmatrix},$$

where  $m \geq n$ ,  $\mathbf{A}_1 \in \mathbb{L}_n^{(+)}$  and  $\mathbf{A}_2 \in \mathbb{R}^{(m-n) \times n}$ ;

- $\mathbb{L}_{m,n}^{(1)}$ : Set of all  $m \times n$  matrices  $\mathbf{A} = (a_{ij}) \in \mathbb{L}_{m,n}^{(+)}$  such that  $a_{ii} = 1$  for  $i = 1, \dots, n$ ;
- $\mathbb{S}_n$ : Set of all  $n \times n$  symmetric matrices;
- $\mathbb{S}_n^{(+)}$ : Set of all  $n \times n$  symmetric positive definite matrices;
- $\mathbb{S}_{n,r}^{(+)}$ : Set of all  $n \times n$  symmetric positive semi-definite matrices with rank  $r$ .

Note that  $\mathbb{O}_n = \mathbb{V}_{n,n}$ ,  $\mathbb{L}_n^{(+)} = \mathbb{L}_{n,n}^{(+)}$ ,  $\mathbb{L}_n^{(1)} = \mathbb{L}_{n,n}^{(1)}$  and  $\mathbb{S}_n^{(+)} = \mathbb{S}_{n,n}^{(+)}$ . Set  $\mathbb{V}_{m,n}$  is referred to as the Stiefel manifold. Also, it is important to note that  $\mathbb{O}_n$  and  $\mathbb{L}_n^{(+)}$  are groups, called respectively the orthogonal and the lower triangular groups, with respect to the group action by usual matrix multiplication.

Matrix inequality is defined in the Löwner sense. Namely, for  $\mathbf{A}_1, \mathbf{A}_2 \in \mathbb{S}_n$ , the matrix inequality  $\mathbf{A}_1 \preceq \mathbf{A}_2$  (or  $\mathbf{A}_1 \succeq \mathbf{A}_2$ ) means that  $\mathbf{A}_2 - \mathbf{A}_1$  (or  $\mathbf{A}_1 - \mathbf{A}_2$ ) is symmetric positive semi-definite.

For definition, concepts and applications in terms of matrix algebra, see, for example, Rao (1973), Golub and Van Loan (1996) and Harville (1997).

## 2.2 Nonsingular Matrix and the Moore-Penrose Inverse

A matrix  $A \in \mathbb{R}^{n \times n}$  is said to be invertible if there exists a matrix  $B \in \mathbb{R}^{n \times n}$  such that  $AB = BA = I_n$ . Such a matrix  $B$  is uniquely defined from  $A$  and is denoted by  $A^{-1}$ , called the inverse of  $A$ . A matrix  $A$  is invertible if and only if  $A$  belongs to  $\mathbb{U}_n$ , namely,  $|A| \neq 0$ , and thus an invertible matrix is also called a nonsingular matrix. Here we give an important result on the inverse of a partitioned matrix.

**Lemma 2.1** *Partition  $A \in \mathbb{R}^{n \times n}$  as*

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix},$$

where  $A_{11}$ ,  $A_{12}$ ,  $A_{21}$  and  $A_{22}$  are matrix subblocks of any size. If  $A_{11}$  and  $A_{22.1} = A_{22} - A_{21}A_{11}^{-1}A_{12}$  are squared and nonsingular then  $A$  is nonsingular and

$$A^{-1} = \begin{pmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}A_{22.1}^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}A_{22.1}^{-1} \\ -A_{22.1}^{-1}A_{21}A_{11}^{-1} & A_{22.1}^{-1} \end{pmatrix}.$$

In addition, if  $A_{22}$  and  $A_{11.2} = A_{11} - A_{12}A_{22}^{-1}A_{21}$  are also nonsingular then

$$A_{11}^{-1} + A_{11}^{-1}A_{12}A_{22.1}^{-1}A_{21}A_{11}^{-1} = A_{11.2}^{-1}.$$

In particular, if  $A \in \mathbb{S}_n^{(+)}$  then  $A_{11}$ ,  $A_{22}$ ,  $A_{11.2}$  and  $A_{22.1}$  are nonsingular.

If a matrix  $A$  is partitioned as in Lemma 2.1 and  $A_{11} \in \mathbb{U}_m$  with  $m < n$ , then  $A$  can be expressed as

$$A = \begin{pmatrix} I_m & \mathbf{0}_{m \times (n-m)} \\ A_{21}A_{11}^{-1} & I_{n-m} \end{pmatrix} \begin{pmatrix} A_{11} & \mathbf{0}_{m \times (n-m)} \\ \mathbf{0}_{(n-m) \times m} & A_{22.1} \end{pmatrix} \begin{pmatrix} I_m & A_{11}^{-1}A_{12} \\ \mathbf{0}_{(n-m) \times m} & I_{n-m} \end{pmatrix}. \quad (2.1)$$

Therefore Lemma 2.1 can easily be verified by (2.1).

Next, we describe some basic and useful properties of the Moore-Penrose inverse, which will be needed for a unified treatment of high- and low-dimensional shrinkage estimators. The Moore-Penrose inverse is an extension of invertibility to a singular square matrix and to a rectangular matrix, and further it is a special case of generalized inverse. Here, for a matrix  $A \in \mathbb{R}^{m \times n}$ , a matrix  $A^- (\in \mathbb{R}^{n \times m})$  is said to be a generalized inverse of  $A$  if  $A^-$  satisfies  $AA^-A = A$ . The generalized inverse is not unique, while the Moore-Penrose inverse is unique.

**Definition 2.1** For a matrix  $A (\in \mathbb{R}^{m \times n})$ , a matrix  $A^+ (\in \mathbb{R}^{n \times m})$  is called the Moore-Penrose inverse of  $A$  if  $A^+$  satisfies

- (i)  $AA^+A = A$ ,
- (ii)  $A^+AA^+ = A^+$ ,

- (iii)  $(\mathbf{A}\mathbf{A}^+)^\top = \mathbf{A}\mathbf{A}^+$ ,
- (iv)  $(\mathbf{A}^+\mathbf{A})^\top = \mathbf{A}^+\mathbf{A}$ .

**Lemma 2.2** *For any matrix  $\mathbf{A}$ ,  $\mathbf{A}^+$  always exists and it is unique.*

**Lemma 2.3** *For  $m \geq n$ , let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be of full column rank. Then we have*

- (i)  $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \in \mathbb{R}^{n \times m}$ , and in particular  $\mathbf{A}^+ = \mathbf{A}^\top$  if  $\mathbf{A} \in \mathbb{V}_{m,n}$ ,
- (ii)  $(\mathbf{A}^\top)^+ = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1}$ , and thus  $(\mathbf{A}^\top)^+ = (\mathbf{A}^+)^\top$ ,
- (iii)  $\mathbf{A}^+ \mathbf{A} = \mathbf{I}_n$ ,
- (iv)  $(\mathbf{A}\mathbf{B}\mathbf{A}^\top)^+ = (\mathbf{A}^\top)^+ \mathbf{B}^{-1} \mathbf{A}^+$  for any  $\mathbf{B} \in \mathbb{U}_n$ .

Parts (i) and (ii) of Lemma 2.3 can easily be verified by Definition 2.1 and Lemma 2.2. Part (iii) can be obtained from (i). Part (iv) follows from Definition 2.1 and Lemma 2.2 with (ii) and (iii).

Using (i) of Lemma 2.3, we can see that if  $\mathbf{A} \in \mathbb{U}_n$  then  $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top = \mathbf{A}^{-1}(\mathbf{A}^\top)^{-1} \mathbf{A}^\top = \mathbf{A}^{-1}$ . This implies that if  $\mathbf{A} \in \mathbb{U}_n$  and it is symmetric then  $\mathbf{A}^+$  is symmetric as well. For more general results on the Moore-Penrose inverse, see Harville (1997) and Magnus and Neudecker (1999).

## 2.3 Kronecker Product and Vec Operator

The notion of the Kronecker product and the vec operator is very important to discuss a clear theorization in terms of matrix-variate distributions. Here, we provide their definition and list some useful results without proofs. For the proofs and more details, see Harville (1997) and Muirhead (1982).

**Definition 2.2** The Kronecker product of two matrices  $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{p \times q}$  is denoted by  $\mathbf{A} \otimes \mathbf{B}$ , which is a block matrix of the form

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{mp \times nq}.$$

**Lemma 2.4** *Some results on the Kronecker product are given as follows:*

- (i)  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$  for  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{p \times q}$ ,  $\mathbf{C} \in \mathbb{R}^{n \times r}$  and  $\mathbf{D} \in \mathbb{R}^{q \times s}$ ,
- (ii)  $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$ ,
- (iii)  $|\mathbf{A} \otimes \mathbf{B}| = |\mathbf{A}|^n |\mathbf{B}|^m$  for  $\mathbf{A} \in \mathbb{R}^{m \times m}$  and  $\mathbf{B} \in \mathbb{R}^{n \times n}$ ,
- (iv)  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$  for nonsingular matrices  $\mathbf{A}$  and  $\mathbf{B}$ .

**Definition 2.3** Let  $A = (a_1, \dots, a_n) \in \mathbb{R}^{m \times n}$ , where the  $a_i$ 's lie in  $\mathbb{R}^m$ . The vec operation on  $A$  is expressed by  $\text{vec}(A)$ , which is of the form

$$\text{vec}(A) = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \in \mathbb{R}^{mn}.$$

**Lemma 2.5** *Some results on the vec operation are given as follows:*

(i) *If  $A \in \mathbb{R}^{m \times n}$ ,  $X \in \mathbb{R}^{n \times p}$  and  $B \in \mathbb{R}^{p \times q}$  then*

$$\text{vec}(AXB) = (B^\top \otimes A)\text{vec}(X).$$

(ii) *If  $X \in \mathbb{R}^{m \times n}$ ,  $A \in \mathbb{R}^{m \times m}$  and  $B \in \mathbb{R}^{n \times n}$  then*

$$\text{tr} AXBX^\top = \text{vec}(X)^\top (B^\top \otimes A)\text{vec}(X) = \text{vec}(X)^\top (B \otimes A^\top)\text{vec}(X).$$

## 2.4 Matrix Decompositions

A matrix decomposition, or a matrix factorization, is to express a matrix as a product of some matrices. The matrix decomposition appears in various scenes of statistical analysis. In shrinkage estimation, it is often used for finding better estimators. This section provides some important matrix decompositions without proofs. For the proofs and properties of matrix decompositions, see Golub and Van Loan (1996) and Harville (1997).

**Lemma 2.6** (QR decomposition) *Let  $A \in \mathbb{R}^{m \times n}$  be of full column rank. Then there exist unique  $R \in \mathbb{L}_n^{(+)}$  and  $Q \in \mathbb{V}_{m,n}$  such that  $A = QR^\top$ .*

The QR decomposition  $A = QR^\top$  yields  $A^\top = RQ^\top$ , implying that for  $A \in \mathbb{R}^{m \times n}$  of full row rank there exist unique  $L \in \mathbb{L}_m^{(+)}$  and  $Q \in \mathbb{V}_{n,m}$  such that  $A = LQ^\top$ . The decomposition  $A = LQ^\top$  is called the LQ decomposition.

**Lemma 2.7** (Cholesky decomposition) *For any  $A \in \mathbb{S}_n^{(+)}$ , there exists a unique  $L \in \mathbb{L}_n^{(+)}$  such that  $A = LL^\top$ .*

If  $A \in \mathbb{S}_n^{(+)}$  then  $A$  can also be decomposed as  $A = LDL^\top$ , where  $L \in \mathbb{L}_n^{(1)}$ ,  $D \in \mathbb{D}_n$  with positive diagonals, and  $L$  and  $D$  are unique. This is known as the  $LDL^\top$  decomposition.

When  $A \in \mathbb{R}^{n \times n}$  can be decomposed as  $A = LU$ , where  $L$  and  $U$  are, respectively, lower and upper triangular matrices, the decomposition is called the LU decomposition of  $A$ .

**Lemma 2.8** (Eigenvalue decomposition) *For any  $A \in \mathbb{R}^{n \times n}$ , there exists  $P \in \mathbb{U}_n$  such that  $A = PDP^{-1}$ , where  $D \in \mathbb{D}_n^{(\geq)}$ . In particular, if  $A \in \mathbb{S}_{n,r}^{(+)}$  then there exist  $D \in \mathbb{D}_r^{(\geq 0)}$  and  $P \in \mathbb{V}_{n,r}$  such that  $A = PDP^\top$ .*

The eigenvalue decomposition is also named spectral decomposition. The diagonal elements of  $\mathbf{D}$  in Lemma 2.8 are called the eigenvalues of  $\mathbf{A}$  and the  $i$ -th column of  $\mathbf{P}$  is called the eigenvector corresponding to the  $i$ -th diagonal element of  $\mathbf{D}$ . The eigenvalue and eigenvector are also referred to as the characteristic root and characteristic vector, respectively. The eigenvalues of  $\mathbf{A}$  are arranged in descending order on the diagonal of  $\mathbf{D}$ , implying that  $\mathbf{D}$  is uniquely determined. When  $\mathbf{A}$  is a member of  $\mathbb{S}_{n,r}^{(+)}$ ,  $\mathbf{D}$  is unique and  $\mathbf{P}$  is also unique up to sign changes of columns of  $\mathbf{P}$ .

**Lemma 2.9** (Singular value decomposition) *For any  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with rank  $r$ , there exist  $\mathbf{D} \in \mathbb{D}_r^{(\geq 0)}$ ,  $\mathbf{U} \in \mathbb{V}_{m,r}$  and  $\mathbf{V} \in \mathbb{V}_{n,r}$  such that  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ .*

The diagonal elements of  $\mathbf{D}$  in Lemma 2.9 are called the singular values of  $\mathbf{A}$ . The singular value decomposition  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$  is unique up to sign changes of columns of  $\mathbf{V}$ .

Here we introduce the wedge and descending wedge symbols,  $\wedge$  and  $\vee$ , implying that, for numbers  $a$  and  $b$ ,

$$a \wedge b = \min(a, b), \quad a \vee b = \max(a, b).$$

Also for numbers  $a, b$  and  $c$ ,  $a \wedge b \wedge c = \min(a, b, c)$ . The following identities hold:

$$a \vee b - (a \wedge b) = |a - b| = a + b - 2(a \wedge b) = 2(a \vee b) - a - b.$$

If  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is of full rank, then its rank is  $m \wedge n$  and the singular value decomposition of  $\mathbf{A}$  can be expressed by  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ , where  $\mathbf{D} \in \mathbb{D}_{m \wedge n}^{(\geq 0)}$ ,  $\mathbf{U} \in \mathbb{V}_{m, m \wedge n}$  and  $\mathbf{V} \in \mathbb{V}_{n, m \wedge n}$ .

## References

- G.H. Golub, C.F. Van Loan, *Matrix Computations*, 3rd edn. (The Johns Hopkins University Press, Baltimore, 1996)
- D.A. Harville, *Matrix Algebra From a Statistician's Perspective* (Springer, New York, 1997)
- J.R. Magnus, H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, 2nd edn. (Wiley, New York, 1999)
- R.J. Muirhead, *Aspects of Multivariate Statistical Theory* (Wiley, New York, 1982)
- C.R. Rao, *Linear Statistical Inference and its Applications*, 2nd edn. (Wiley, New York, 1973)

# Chapter 3

## Matrix-Variate Distributions



This chapter provides the definition and some useful properties of a matrix-variate normal distribution, nonsingular and singular Wishart distributions, and other related distributions. The matrix-variate distributions considered here are based on the multivariate (vector-valued) normal distribution, and so we begin by briefly introducing the multivariate normal distribution and some Jacobians for matrix transformations used to obtain probability density functions of the matrix-variate distributions.

### 3.1 Preliminaries

#### 3.1.1 The Multivariate Normal Distribution

A  $p$ -dimensional random vector follows a multivariate (vector-valued) normal distribution if the probability density function (p.d.f.) is given by

$$\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right), \quad \mathbf{x} \in \mathbb{R}^p,$$

where  $\boldsymbol{\mu} \in \mathbb{R}^p$  and  $\Sigma \in \mathbb{S}_p^{(+)}$  are parameters. Such a multivariate ( $p$ -variate) normal distribution is denoted by  $\mathcal{N}_p(\boldsymbol{\mu}, \Sigma)$ . The multivariate normal distribution has the following properties.

**Lemma 3.1** *Let  $\mathbf{x} = (x_1, \dots, x_p)^\top \sim \mathcal{N}_p(\boldsymbol{\mu}, \Sigma)$  with  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top \in \mathbb{R}^p$  and  $\Sigma = (\sigma_{ij}) \in \mathbb{S}_p^{(+)}$ . Then*

- (i)  $E[x_i] = \mu_i$  and  $E[x_i x_j] = \mu_i \mu_j + \sigma_{ij}$  for all  $i, j \in \{1, \dots, p\}$ . In other words,

$$\begin{aligned} E[\mathbf{x}] &= (E[x_1], \dots, E[x_p])^\top = \boldsymbol{\mu}, \\ E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] &= (E[(x_i - \mu_i)(x_j - \mu_j)]) = \Sigma. \end{aligned}$$

- (ii) For a full row rank constant matrix  $A \in \mathbb{R}^{q \times p}$  and a constant vector  $\mathbf{b} \in \mathbb{R}^q$ ,  $A\mathbf{x} + \mathbf{b} \sim \mathcal{N}_q(A\boldsymbol{\mu} + \mathbf{b}, A\boldsymbol{\Sigma}A^\top)$ .
- (iii) If  $\mathbf{x}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are partitioned, respectively, as

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{21}^\top \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix},$$

where  $\mathbf{x}_1 \in \mathbb{R}^q$ ,  $\boldsymbol{\mu}_1 \in \mathbb{R}^q$  and  $\boldsymbol{\Sigma}_{11} \in \mathbb{S}_q^{(+)}$ , then

$$\begin{aligned} \mathbf{x}_1 &\sim \mathcal{N}_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}), \\ \mathbf{x}_2 | \mathbf{x}_1 &\sim \mathcal{N}_{p-q}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{21}^\top). \end{aligned}$$

Further,  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are independent if and only if  $\boldsymbol{\Sigma}_{21} = \mathbf{0}_{(p-q) \times q}$ .

From (i) of Lemma 3.1,  $\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is usually called the  $p$ -variate normal distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ . For theoretical properties other than Lemma 3.1, see Muirhead (1982) and Srivastava and Khatri (1979).

### 3.1.2 Jacobians of Matrix Transformations

Let  $X$  be a matrix such that it has  $k$  functionally independent variables  $x_i$ 's and let  $Y = F(X)$  be a matrix transformation, where  $Y$  is a matrix having  $k$  functionally independent variables  $y_i$ 's. The Jacobian of the transformation is given by  $J(X \rightarrow Y) = |J|$ , where  $J = (\partial x_i / \partial y_j)$  is the  $k \times k$  Jacobian matrix. Let  $(dX) = \bigwedge_{i=1}^k dx_i$  and  $(dY) = \bigwedge_{i=1}^k dy_i$ , where  $\bigwedge$  denotes the exterior product. Then  $(dX) = J(X \rightarrow Y)(dY)$ .

This section provides some specific Jacobians of matrix transformations. When considering a matrix transformation, we need to take attention to the number of functionally independent variables. For example,  $S = (s_{ij}) \in \mathbb{S}_n$  has  $\{n(n+1)/2\}$  distinct elements, and the exterior product of distinct elements of the differential  $dS = (ds_{ij}) \in \mathbb{S}_n$  is  $(dS) = \bigwedge_{i=1}^n \bigwedge_{j=1}^i ds_{ij}$ . See Muirhead (1982) and Mathai (1997) for more details of Jacobians and exterior products. Hereafter we ignore the signs of exterior differential forms and define only positive integrals.

**Lemma 3.2** Let  $X = (x_{ij}) \in \mathbb{R}^{m \times n}$  and  $Y = (y_{ij}) \in \mathbb{R}^{m \times n}$ . If  $X = AYB + C$  with  $A \in \mathbb{U}_m$ ,  $B \in \mathbb{U}_n$  and  $C \in \mathbb{R}^{m \times n}$ , then  $(dX) = |A|^n |B|^m (dY)$ , where  $(dX) = \bigwedge_{i=1}^m \bigwedge_{j=1}^n dx_{ij}$  and  $(dY) = \bigwedge_{i=1}^m \bigwedge_{j=1}^n dy_{ij}$ .

**Proof** See Muirhead (1982) and Mathai (1997). □

**Lemma 3.3** Let  $X \in \mathbb{R}^{p \times n}$  with full row rank. Denote by  $X = TQ^\top$  the LQ decomposition of  $X$ , where  $T = (t_{ij}) \in \mathbb{L}_p^{(+)}$  and  $Q \in \mathbb{V}_{n,p}$ . Then

$$(\mathrm{d}X) = \left( \prod_{i=1}^p t_{ii}^{n-i} \right) (\mathrm{d}T)(\mathbf{Q}^\top \mathrm{d}\mathbf{Q}),$$

where  $(\mathrm{d}T) = \bigwedge_{i=1}^p \bigwedge_{j=1}^i \mathrm{d}t_{ij}$  and  $(\mathbf{Q}^\top \mathrm{d}\mathbf{Q})$  is an unnormalized probability measure on  $\mathbb{V}_{n,p}$ .

**Proof** See Muirhead (1982, Theorem 2.1.13).  $\square$

Muirhead (1982, p. 69) pointed out that the measure  $(\mathbf{Q}^\top \mathrm{d}\mathbf{Q})$  on  $\mathbb{V}_{n,p}$  is invariant under the orthogonal transformations  $\mathbf{Q} \rightarrow \mathbf{A}\mathbf{Q}$  for  $\mathbf{A} \in \mathbb{O}_n$  and  $\mathbf{Q} \rightarrow \mathbf{Q}\mathbf{B}$  for  $\mathbf{B} \in \mathbb{O}_p$ . Also,

$$\int_{\mathbb{V}_{n,p}} (\mathbf{Q}^\top \mathrm{d}\mathbf{Q}) = \frac{2^p \pi^{np/2}}{\Gamma_p(n/2)}, \quad (3.1)$$

which is given in Muirhead (1982, Theorem 2.1.15), where  $\Gamma_p(n/2)$  is the multi-variate gamma function (see Definition 3.1 given below).

**Lemma 3.4** Let  $S \in \mathbb{S}_p^{(+)}$  and let  $S = \mathbf{T}\mathbf{T}^\top$  be the Cholesky decomposition of  $S$ , where  $\mathbf{T} = (t_{ij}) \in \mathbb{L}_p^{(+)}$ . Then  $(\mathrm{d}S) = 2^p (\prod_{i=1}^p t_{ii}^{p+1-i}) (\mathrm{d}T)$ .

**Proof** See Muirhead (1982, Theorem 2.1.9).  $\square$

**Lemma 3.5** Let  $S \in \mathbb{S}_{p,r}^{(+)}$ . Denote by  $S = \mathbf{H}\mathbf{L}\mathbf{H}^\top$  the eigenvalue decomposition of  $S$ , where  $\mathbf{L} = \text{diag}(\ell_1, \dots, \ell_r) \in \mathbb{D}_r^{(\geq 0)}$  and  $\mathbf{H} \in \mathbb{V}_{p,r}$ . Then

$$(\mathrm{d}S) = \frac{1}{2^r} \left( \prod_{i=1}^r \ell_i^{p-r} \right) \left( \prod_{1 \leq i < j \leq r} (\ell_i - \ell_j) \right) (\mathrm{d}L)(\mathbf{H}^\top \mathrm{d}\mathbf{H}),$$

where  $(\mathrm{d}L) = \bigwedge_{i=1}^r \mathrm{d}\ell_i$ .

**Proof** See Uhlig (1994, Theorem 2).  $\square$

**Lemma 3.6** Let  $X \in \mathbb{R}^{p \times n}$  with rank  $r$ . Denote by  $X = \mathbf{U}\mathbf{D}\mathbf{V}^\top$  the singular value decomposition of  $X$ , where  $\mathbf{D} = \text{diag}(d_1, \dots, d_r) \in \mathbb{D}_r^{(\geq 0)}$ ,  $\mathbf{U} \in \mathbb{V}_{p,r}$  and  $\mathbf{V} \in \mathbb{V}_{n,r}$ . Then

$$(\mathrm{d}X) = \frac{1}{2^r} \left( \prod_{i=1}^r d_i^{p+n-2r} \right) \left( \prod_{1 \leq i < j \leq r} (d_i^2 - d_j^2) \right) (\mathbf{U}^\top \mathrm{d}\mathbf{U})(\mathrm{d}\mathbf{D})(\mathbf{V}^\top \mathrm{d}\mathbf{V}).$$

**Proof** See Díaz-García et al. (1997, Theorem 3.1).  $\square$



### 3.1.3 The Multivariate Gamma Function

The multivariate gamma function is defined as a multivariate extension of the gamma function. It is convenient for clearly expressing normalizing constants of matrix-variate distributions.

**Definition 3.1** For  $a > (p - 1)/2$ , the multivariate gamma function  $\Gamma_p(a)$  is defined by

$$\Gamma_p(a) = \int_{\mathbb{S}_p^{(+)}} |\mathbf{W}|^{a-(p+1)/2} \exp(-\operatorname{tr} \mathbf{W}) (d\mathbf{W}).$$

When  $p = 1$ ,  $\Gamma(a) = \Gamma_1(a) = \int_0^\infty w^{a-1} e^{-w} dw$ , which is the usual gamma function. The multivariate gamma function can be rewritten as a product of the gamma functions.

**Proposition 3.1** For  $a > (p - 1)/2$ ,

$$\Gamma_p(a) = \pi^{p(p-1)/4} \prod_{i=1}^p \Gamma\left(a - \frac{i-1}{2}\right).$$

**Proof** For  $\mathbf{W} \in \mathbb{S}_p^{(+)}$ , let  $\mathbf{W} = \mathbf{T}\mathbf{T}^\top$  be the Cholesky decomposition of  $\mathbf{W}$ , where  $\mathbf{T} = (t_{ij}) \in \mathbb{L}_p^{(+)}$ . From Lemma 3.4, the Jacobian of transformation  $\mathbf{W} \rightarrow \mathbf{T}$  is given by  $J(\mathbf{W} \rightarrow \mathbf{T}) = 2^p \prod_{i=1}^p t_{ii}^{p-i+1}$ , so that

$$\Gamma_p(a) = 2^p \int_{\mathbb{L}_p^{(+)}} |\mathbf{T}\mathbf{T}^\top|^{a-(p+1)/2} \exp(-\operatorname{tr} \mathbf{T}\mathbf{T}^\top) \left( \prod_{i=1}^p t_{ii}^{p-i+1} \right) (d\mathbf{T}).$$

Since  $|\mathbf{T}\mathbf{T}^\top| = \prod_{i=1}^p t_{ii}^2$  and  $\operatorname{tr} \mathbf{T}\mathbf{T}^\top = \sum_{i=1}^p \sum_{j=1}^i t_{ij}^2$ , we have

$$\begin{aligned} \Gamma_p(a) &= 2^p \left( \prod_{i=2}^p \prod_{j=1}^{i-1} \int_{-\infty}^{\infty} e^{-t_{ij}^2} dt_{ij} \right) \left( \prod_{i=1}^p \int_0^{\infty} t_{ii}^{2a-i} e^{-t_{ii}^2} dt_{ii} \right) \\ &= 2^p \times (\sqrt{\pi})^{p(p-1)/2} \times \prod_{i=1}^p \frac{1}{2} \Gamma\left(a - \frac{i-1}{2}\right) \\ &= \pi^{p(p-1)/4} \prod_{i=1}^p \Gamma\left(a - \frac{i-1}{2}\right), \end{aligned}$$

which completes the proof. □

## 3.2 The Matrix-Variate Normal Distribution

This section provides the definition of matrix-variate normal distribution and its useful properties to analyze a multivariate linear model. Here, the matrix-variate normal distribution is defined as an extension of  $\mathcal{N}_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ .

**Definition 3.2** For a random matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , assume that

$$\text{vec}(\mathbf{X}^\top) \sim \mathcal{N}_{np}(\text{vec}(\mathbf{M}^\top), \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}).$$

That is,  $\text{vec}(\mathbf{X}^\top)$  follows the  $(np)$ -variate normal distribution with mean  $\text{vec}(\mathbf{M}^\top)$  and covariance  $\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}$ , where  $\mathbf{M} \in \mathbb{R}^{n \times p}$ , and  $\boldsymbol{\Omega} \in \mathbb{S}_n^{(+)}$  and  $\boldsymbol{\Sigma} \in \mathbb{S}_p^{(+)}$ . Then  $\mathbf{X}$  is said to follow the matrix-variate normal distribution with mean matrix  $\mathbf{M}$  and covariance matrix  $\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}$ , which is denoted by  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ .

Note that  $\mathcal{N}_{np}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$  means an  $(np)$ -variate (vector-valued) normal distribution and is distinguished from  $\mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ .

The p.d.f. of the matrix-variate normal distribution can be written as follows.

**Proposition 3.2** The p.d.f. of  $\mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$  is given by

$$\frac{1}{(2\pi)^{np/2} |\boldsymbol{\Omega}|^{p/2} |\boldsymbol{\Sigma}|^{n/2}} \exp \left( -\frac{1}{2} \text{tr} \boldsymbol{\Omega}^{-1} (\mathbf{X} - \mathbf{M}) \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{M})^\top \right).$$

**Proof** Using (iii) and (iv) of Lemma 2.4 and (ii) of Lemma 2.5, we observe  $|\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}| = |\boldsymbol{\Omega}|^p |\boldsymbol{\Sigma}|^n$  and

$$\begin{aligned} & \text{vec}((\mathbf{X} - \mathbf{M})^\top)^\top (\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})^{-1} \text{vec}((\mathbf{X} - \mathbf{M})^\top) \\ &= \text{vec}((\mathbf{X} - \mathbf{M})^\top)^\top (\boldsymbol{\Omega}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}((\mathbf{X} - \mathbf{M})^\top) \\ &= \text{tr} \boldsymbol{\Omega}^{-1} (\mathbf{X} - \mathbf{M}) \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{M})^\top. \end{aligned}$$

Since  $\text{vec}(\mathbf{X}^\top) \sim \mathcal{N}_{np}(\text{vec}(\mathbf{M}^\top), \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ , the p.d.f. is written as

$$\begin{aligned} & \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma}|^{1/2}} \exp \left( -\frac{1}{2} \text{vec}((\mathbf{X} - \mathbf{M})^\top)^\top (\boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})^{-1} \text{vec}((\mathbf{X} - \mathbf{M})^\top) \right) \\ &= \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Omega}|^{p/2} |\boldsymbol{\Sigma}|^{n/2}} \exp \left( -\frac{1}{2} \text{tr} \boldsymbol{\Omega}^{-1} (\mathbf{X} - \mathbf{M}) \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{M})^\top \right). \end{aligned}$$

Hence the proof is complete.  $\square$

Using properties of the Kronecker product and the vec operator, we have the following proposition.

**Proposition 3.3** If  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ , then  $\mathbf{X}^\top \sim \mathcal{N}_{p \times n}(\mathbf{M}^\top, \boldsymbol{\Sigma} \otimes \boldsymbol{\Omega})$ .

**Proof** Using (iv) of Lemma 2.4 and (ii) of Lemma 2.5 gives

$$\begin{aligned}\text{tr } \boldsymbol{\Omega}^{-1}(\mathbf{X} - \mathbf{M})\boldsymbol{\Sigma}^{-1}(\mathbf{X} - \mathbf{M})^\top &= \text{tr } \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \mathbf{M})^\top \boldsymbol{\Omega}^{-1}(\mathbf{X} - \mathbf{M}) \\ &= \text{vec}(\mathbf{X} - \mathbf{M})^\top (\boldsymbol{\Sigma} \otimes \boldsymbol{\Omega})^{-1} \text{vec}(\mathbf{X} - \mathbf{M}).\end{aligned}$$

This implies that  $\text{vec}(\mathbf{X}) \sim \mathcal{N}_{pn}(\text{vec}(\mathbf{M}), \boldsymbol{\Sigma} \otimes \boldsymbol{\Omega})$ .  $\square$

The first and second moments of the matrix-variate normal distribution are given in the following proposition.

**Proposition 3.4** *Let  $\mathbf{X} = (x_{ij}) \sim \mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$  with  $\mathbf{M} = (\mu_{ij})$ ,  $\boldsymbol{\Omega} = (\omega_{ij})$  and  $\boldsymbol{\Sigma} = (\sigma_{ij})$ . Then, for any  $i, k \in \{1, \dots, n\}$  and  $j, l \in \{1, \dots, p\}$ ,*

$$E[x_{ij}] = \mu_{ij}, \quad E[x_{ij}x_{kl}] = \omega_{ik}\sigma_{jl} + \mu_{ij}\mu_{kl}.$$

**Proof** Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$  and  $\mathbf{M} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n)^\top$ , where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top \in \mathbb{R}^p$  and  $\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{ip})^\top \in \mathbb{R}^p$  for  $i = 1, \dots, n$ . Then

$$\text{vec}(\mathbf{X}^\top) = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} \sim \mathcal{N}_{np} \left( \begin{pmatrix} \boldsymbol{\mu}_1 \\ \vdots \\ \boldsymbol{\mu}_n \end{pmatrix}, \begin{pmatrix} \omega_{11}\boldsymbol{\Sigma} & \dots & \omega_{1n}\boldsymbol{\Sigma} \\ \vdots & & \vdots \\ \omega_{n1}\boldsymbol{\Sigma} & \dots & \omega_{nn}\boldsymbol{\Sigma} \end{pmatrix} \right),$$

implying that, according to (i) of Lemma 3.1,

$$E[\mathbf{x}_i] = \boldsymbol{\mu}_i, \quad E[(\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_k - \boldsymbol{\mu}_k)^\top] = \omega_{ik}\boldsymbol{\Sigma},$$

further implying that

$$E[x_{ij}] = \mu_{ij}, \quad E[(x_{ij} - \mu_{ij})(x_{kl} - \mu_{kl})] = \omega_{ik}\sigma_{jl}.$$

Hence the proof is complete.  $\square$

Proposition 3.4 suggests that, if  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ , covariance of any two rows of  $\mathbf{X}$  is proportional to  $\boldsymbol{\Sigma}$  and also covariance of any two columns of  $\mathbf{X}$  is proportional to  $\boldsymbol{\Omega}$ . Further Proposition 3.4 yields the following corollary.

**Corollary 3.1** *If  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ , then*

$$\begin{aligned}E[\mathbf{X}] &= \mathbf{M}, \\ E[\mathbf{X}^\top \mathbf{A} \mathbf{X}] &= (\text{tr } \mathbf{A} \boldsymbol{\Omega}) \boldsymbol{\Sigma} + \mathbf{M}^\top \mathbf{A} \mathbf{M}, \\ E[\mathbf{X} \mathbf{B} \mathbf{X}^\top] &= (\text{tr } \mathbf{B} \boldsymbol{\Sigma}) \boldsymbol{\Omega} + \mathbf{M} \mathbf{B} \mathbf{M}^\top\end{aligned}$$

for constant matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{B} \in \mathbb{R}^{p \times p}$ .

**Proof** Using Proposition 3.4 immediately gives  $E[X] = (E[x_{ij}]) = (\mu_{ij}) = \mathbf{M}$ . Let  $\mathbf{A} = (a_{ij})$  and  $\mathbf{B} = (b_{ij})$ . For  $i, j \in \{1, \dots, p\}$ , the  $(i, j)$ -th element of  $\mathbf{X}^\top \mathbf{A} \mathbf{X}$  is  $\{\mathbf{X}^\top \mathbf{A} \mathbf{X}\}_{ij} = \sum_{k=1}^n \sum_{l=1}^n x_{ki} a_{kl} x_{lj}$ , so that

$$\begin{aligned} E[\{\mathbf{X}^\top \mathbf{A} \mathbf{X}\}_{ij}] &= \sum_{k=1}^n \sum_{l=1}^n a_{kl} E[x_{ki} x_{lj}] \\ &= \sum_{k=1}^n \sum_{l=1}^n a_{kl} (\omega_{kl} \sigma_{ij} + \mu_{ki} \mu_{lj}) \\ &= (\text{tr } \mathbf{A} \mathbf{\Omega}^\top) \sigma_{ij} + \{\mathbf{M}^\top \mathbf{A} \mathbf{M}\}_{ij}. \end{aligned}$$

It holds that  $\text{tr } \mathbf{A} \mathbf{\Omega}^\top = \text{tr } \mathbf{A} \mathbf{\Omega}$ , so that  $E[\mathbf{X}^\top \mathbf{A} \mathbf{X}] = (\text{tr } \mathbf{A} \mathbf{\Omega}) \mathbf{\Sigma} + \mathbf{M}^\top \mathbf{A} \mathbf{M}$ . Similarly, for  $i, j \in \{1, \dots, n\}$ ,

$$\begin{aligned} E[\{\mathbf{X} \mathbf{B} \mathbf{X}^\top\}_{ij}] &= \sum_{k=1}^p \sum_{l=1}^p b_{kl} E[x_{ik} x_{jl}] \\ &= \sum_{k=1}^p \sum_{l=1}^p b_{kl} (\omega_{ij} \sigma_{kl} + \mu_{ik} \mu_{jl}) \\ &= (\text{tr } \mathbf{B} \mathbf{\Sigma}) \omega_{ij} + \{\mathbf{M} \mathbf{B} \mathbf{M}^\top\}_{ij}. \end{aligned}$$

Hence the proof is complete.  $\square$

Next, we provide a result on linear transformation of the matrix-variate normal distribution. The following proposition is an extension of (ii) in Lemma 3.1.

**Proposition 3.5** *Let  $\mathbf{A}$  be a full row rank constant matrix in  $\mathbb{R}^{m \times n}$  and let  $\mathbf{B}$  be a full column rank constant matrix in  $\mathbb{R}^{p \times q}$ . Also, let  $\mathbf{C}$  be a constant matrix in  $\mathbb{R}^{m \times q}$ . If  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \mathbf{\Omega} \otimes \mathbf{\Sigma})$ , then*

$$\mathbf{A} \mathbf{X} \mathbf{B} + \mathbf{C} \sim \mathcal{N}_{m \times q}(\mathbf{A} \mathbf{M} \mathbf{B} + \mathbf{C}, \mathbf{A} \mathbf{\Omega} \mathbf{A}^\top \otimes \mathbf{B}^\top \mathbf{\Sigma} \mathbf{B}).$$

**Proof** From (i) of Lemma 2.5, it is seen that  $\text{vec}((\mathbf{A} \mathbf{X} \mathbf{B} + \mathbf{C})^\top) = \text{vec}(\mathbf{B}^\top \mathbf{X}^\top \mathbf{A}^\top) + \text{vec}(\mathbf{C}^\top) = (\mathbf{A} \otimes \mathbf{B}^\top) \text{vec}(\mathbf{X}^\top) + \text{vec}(\mathbf{C}^\top)$ . Using (ii) of Lemma 3.1 gives

$$\begin{aligned} &(\mathbf{A} \otimes \mathbf{B}^\top) \text{vec}(\mathbf{X}^\top) + \text{vec}(\mathbf{C}^\top) \\ &\sim \mathcal{N}_{mq}((\mathbf{A} \otimes \mathbf{B}^\top) \text{vec}(\mathbf{M}^\top) + \text{vec}(\mathbf{C}^\top), (\mathbf{A} \otimes \mathbf{B}^\top)(\mathbf{\Omega} \otimes \mathbf{\Sigma})(\mathbf{A} \otimes \mathbf{B}^\top)^\top), \end{aligned}$$

which implies from (i) and (ii) of Lemma 2.4 that

$$(\mathbf{A} \otimes \mathbf{B}^\top) \text{vec}(\mathbf{X}^\top) + \text{vec}(\mathbf{C}^\top) \sim \mathcal{N}_{mq}(\text{vec}((\mathbf{A} \mathbf{M} \mathbf{B} + \mathbf{C})^\top), \mathbf{A} \mathbf{\Omega} \mathbf{A}^\top \otimes \mathbf{B}^\top \mathbf{\Sigma} \mathbf{B}).$$

Hence the proof is complete.  $\square$

A partitioning of a random matrix is needed in various fields of statistical inference. Here we give a distributional property for a partitioned random matrix with respect to the matrix-variate normal distribution. The following proposition is a generalization from (iii) of Lemma 3.1.

**Proposition 3.6** *Let  $X \sim \mathcal{N}_{n \times p}(\mathbf{M}, \mathbf{\Omega} \otimes \mathbf{\Sigma})$ . Partition  $X$ ,  $\mathbf{M}$  and  $\mathbf{\Omega}$  as, respectively,*

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, \quad \mathbf{M} = \begin{pmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \end{pmatrix} \quad \text{and} \quad \mathbf{\Omega} = \begin{pmatrix} \mathbf{\Omega}_{11} & \mathbf{\Omega}_{21}^\top \\ \mathbf{\Omega}_{21} & \mathbf{\Omega}_{22} \end{pmatrix},$$

where  $X_1 \in \mathbb{R}^{m \times p}$ ,  $\mathbf{M}_1 \in \mathbb{R}^{m \times p}$  and  $\mathbf{\Omega}_{11} \in \mathbb{S}_m^{(+)}$ . Let  $\mathbf{\Omega}_{22 \cdot 1} = \mathbf{\Omega}_{22} - \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} \mathbf{\Omega}_{21}^\top$ . Then

$$\begin{aligned} X_1 &\sim \mathcal{N}_{m \times p}(\mathbf{M}_1, \mathbf{\Omega}_{11} \otimes \mathbf{\Sigma}), \\ X_2 | X_1 &\sim \mathcal{N}_{(n-m) \times p}(\mathbf{M}_2 + \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1), \mathbf{\Omega}_{22 \cdot 1} \otimes \mathbf{\Sigma}). \end{aligned}$$

Further, if  $\mathbf{\Omega}_{21} = \mathbf{0}_{(n-m) \times m}$  then  $X_1$  and  $X_2$  are independently distributed as  $X_1 \sim \mathcal{N}_{m \times p}(\mathbf{M}_1, \mathbf{\Omega}_{11} \otimes \mathbf{\Sigma})$  and  $X_2 \sim \mathcal{N}_{(n-m) \times p}(\mathbf{M}_2, \mathbf{\Omega}_{22} \otimes \mathbf{\Sigma})$ .

**Proof** Lemma 2.1 guarantees  $\mathbf{\Omega}_{11}$  and  $\mathbf{\Omega}_{22 \cdot 1}$  to be nonsingular. Using (2.1) gives that  $|\mathbf{\Omega}| = |\mathbf{\Omega}_{11}| \times |\mathbf{\Omega}_{22 \cdot 1}|$  and

$$\mathbf{\Omega}^{-1} = \begin{pmatrix} \mathbf{I}_m & -\mathbf{\Omega}_{11}^{-1} \mathbf{\Omega}_{21}^\top \\ \mathbf{0}_{(n-m) \times m} & \mathbf{I}_{n-m} \end{pmatrix} \begin{pmatrix} \mathbf{\Omega}_{11}^{-1} & \mathbf{0}_{m \times (n-m)} \\ \mathbf{0}_{(n-m) \times m} & \mathbf{\Omega}_{22 \cdot 1}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I}_m & \mathbf{0}_{m \times (n-m)} \\ -\mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} & \mathbf{I}_{n-m} \end{pmatrix}.$$

Here,  $(X - \mathbf{M})^\top \mathbf{\Omega}^{-1} (X - \mathbf{M})$  becomes

$$\begin{aligned} &(X_1 - \mathbf{M}_1)^\top \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1) \\ &+ \{X_2 - \mathbf{M}_2 - \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1)\}^\top \mathbf{\Omega}_{22 \cdot 1}^{-1} \{X_2 - \mathbf{M}_2 - \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1)\}. \end{aligned}$$

From Proposition 3.2, the p.d.f. of  $X$  is rewritten as

$$\begin{aligned} &c_1 \exp \left( -\frac{1}{2} \text{tr} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1) \mathbf{\Sigma}^{-1} (X_1 - \mathbf{M}_1)^\top \right) \\ &\times c_2 \exp \left( -\frac{1}{2} \text{tr} \mathbf{\Omega}_{22 \cdot 1}^{-1} \{X_2 - \mathbf{M}_2 - \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1)\} \right. \\ &\quad \left. \times \mathbf{\Sigma}^{-1} \{X_2 - \mathbf{M}_2 - \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1)\}^\top \right), \end{aligned}$$

where

$$c_1 = (2\pi)^{-mp/2} |\mathbf{\Omega}_{11}|^{-p/2} |\mathbf{\Sigma}|^{-m/2}, \quad c_2 = (2\pi)^{-(n-m)p/2} |\mathbf{\Omega}_{22 \cdot 1}|^{-p/2} |\mathbf{\Sigma}|^{-(n-m)/2}.$$

Hence the above joint p.d.f. of  $X_1$  and  $X_2$  suggests that  $X_1 \sim \mathcal{N}_{m \times p}(\mathbf{M}_1, \mathbf{\Omega}_{11} \otimes \mathbf{\Sigma})$  and  $X_2 | X_1 \sim \mathcal{N}_{(n-m) \times p}(\mathbf{M}_2 + \mathbf{\Omega}_{21} \mathbf{\Omega}_{11}^{-1} (X_1 - \mathbf{M}_1), \mathbf{\Omega}_{22 \cdot 1} \otimes \mathbf{\Sigma})$ .

When  $\mathbf{\Omega}_{21} = \mathbf{0}_{(n-m) \times m}$ , it is seen that  $\mathbf{\Omega}_{22 \cdot 1} = \mathbf{\Omega}_{22}$  and  $\mathbf{X}_2 \sim \mathcal{N}_{(n-m) \times p}(\mathbf{M}_2, \mathbf{\Omega}_{22} \otimes \mathbf{\Sigma})$ . As a consequence,  $\mathbf{X}_2$  does not depend on  $\mathbf{X}_1$ , and thus  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are mutually independent.  $\square$

Proposition 3.6 suggests that if  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{M}, \mathbf{\Omega} \otimes \mathbf{\Sigma})$  and  $\mathbf{\Omega}$  is a diagonal matrix then all the rows of  $\mathbf{X}$  are mutually independent.

Various properties of the matrix-variate normal distribution have been discovered and studied in addition to useful properties mentioned above. For other properties of the matrix-variate normal distribution, see Gupta and Nagar (1999) and Muirhead (1982).

### 3.3 The Wishart Distribution

The Wishart distribution is known as a distribution of the sample covariance matrix in a multivariate normal model and plays an important role in multivariate analysis. It is named for Wishart (1928). First, we provide the definition of the Wishart distribution.

**Definition 3.3** Let  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{\Sigma})$  with  $\mathbf{\Sigma} \in \mathbb{S}_p^{(+)}$  and denote  $\nu = n \wedge p$ . Then  $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$  is called the Wishart matrix of rank  $\nu$  and is said to follow the Wishart distribution with  $n$  degrees of freedom and scale matrix  $\mathbf{\Sigma}$ , which is denoted by  $\mathbf{S} \sim \mathcal{W}_p^\nu(n, \mathbf{\Sigma})$ .

If  $n \geq p$ ,  $\mathcal{W}_p^\nu(n, \mathbf{\Sigma})$  is often abbreviated to  $\mathcal{W}_p(n, \mathbf{\Sigma})$  and the Wishart matrix  $\mathbf{S}$  lies in  $\mathbb{S}_p^{(+)}$  with probability one.

When  $p > n$ , the Wishart matrix is singular with probability one and belongs to  $\mathbb{S}_{p,n}^{(+)}$ . Then in the literature, the distribution of the Wishart matrix is called a pseudo-Wishart distribution. See Srivastava and Khatri (1979) and Díaz-García et al. (1997).

The following proposition provides a unified p.d.f. of  $\mathcal{W}_p^\nu(n, \mathbf{\Sigma})$  in the nonsingular and the singular cases of the Wishart matrix.

**Proposition 3.7** Let  $(d\mathbf{S})$  be defined as in Lemma 3.5 with  $r = \nu$ . The p.d.f. of  $\mathcal{W}_p^\nu(n, \mathbf{\Sigma})$  with respect to  $(d\mathbf{S})$  is expressed by

$$\frac{\pi^{n(\nu-p)/2}}{2^{np/2} |\mathbf{\Sigma}|^{n/2} \Gamma_\nu(n/2)} |\mathbf{L}|^{(n-p-1)/2} \exp\left(-\frac{1}{2} \text{tr } \mathbf{\Sigma}^{-1} \mathbf{S}\right),$$

where  $\mathbf{L} \in \mathbb{D}_\nu^{(\geq 0)}$  and the diagonal of  $\mathbf{L}$  consists of  $\nu$  positive eigenvalues of  $\mathbf{S}$ .

**Proof** Let  $\mathbf{X} = (x_{ij}) \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{\Sigma})$ . Note from Proposition 3.3 that  $\mathbf{X}^\top \sim \mathcal{N}_{p \times n}(\mathbf{0}_{p \times n}, \mathbf{\Sigma} \otimes \mathbf{I}_n)$ . Let  $\mathbf{X}^\top = \mathbf{V} \mathbf{D} \mathbf{U}^\top$  be the singular value decomposition of  $\mathbf{X}^\top$ , where  $\mathbf{V} \in \mathbb{V}_{p,\nu}$ ,  $\mathbf{D} = \text{diag}(d_1, \dots, d_\nu) \in \mathbb{D}_\nu^{(\geq 0)}$  and  $\mathbf{U} \in \mathbb{V}_{n,\nu}$ . Since the p.d.f. of  $\mathcal{N}_{p \times n}(\mathbf{0}_{p \times n}, \mathbf{\Sigma} \otimes \mathbf{I}_n)$  with respect to  $(d\mathbf{X}^\top) = \bigwedge_{j=1}^p \bigwedge_{i=1}^n dx_{ij}$  on  $\mathbb{R}^{p \times n}$  is given by

$$\frac{1}{(2\pi)^{np/2} |\mathbf{\Sigma}|^{n/2}} \exp\left(-\frac{1}{2} \text{tr } \mathbf{\Sigma}^{-1} \mathbf{X}^\top \mathbf{X}\right) (d\mathbf{X}^\top). \quad (3.2)$$

Lemma 3.6 is used to express the joint (unnormalized) p.d.f. of  $\mathbf{V}$ ,  $\mathbf{D}$  and  $\mathbf{U}$  as

$$\frac{1}{(2\pi)^{np/2}|\mathbf{\Sigma}|^{n/2}} \exp\left(-\frac{1}{2} \text{tr } \mathbf{\Sigma}^{-1} \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top\right) \\ \times 2^{-\nu} |\mathbf{D}|^{n+p-2\nu} \left( \prod_{1 \leq i < j \leq \nu} (d_i^2 - d_j^2) \right) (\mathbf{V}^\top d\mathbf{V})(d\mathbf{D})(\mathbf{U}^\top d\mathbf{U}).$$

Note from (3.1) that  $\int_{\mathbb{V}_{n,\nu}} (\mathbf{U}^\top d\mathbf{U}) = 2^\nu \pi^{n\nu/2} / \Gamma_\nu(n/2)$ . Making the transformation  $\mathbf{L} = \text{diag}(\ell_1, \dots, \ell_\nu) = \mathbf{D}^2$  and integrating out with respect to  $\mathbf{U}$ , we obtain the joint (unnormalized) p.d.f. of  $\mathbf{V}$  and  $\mathbf{L}$ :

$$\frac{\pi^{n\nu/2}}{(2\pi)^{np/2}|\mathbf{\Sigma}|^{n/2}\Gamma_\nu(n/2)} \exp\left(-\frac{1}{2} \text{tr } \mathbf{\Sigma}^{-1} \mathbf{V} \mathbf{L} \mathbf{V}^\top\right) \\ \times 2^{-\nu} |\mathbf{L}|^{(n+p-2\nu-1)/2} \left( \prod_{1 \leq i < j \leq \nu} (\ell_i - \ell_j) \right) (\mathbf{V}^\top d\mathbf{V})(d\mathbf{L}). \quad (3.3)$$

Let  $\mathbf{S} = \mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{L} \mathbf{V}^\top$ . From Lemma 3.5,

$$(d\mathbf{S}) = \frac{1}{2^\nu} |\mathbf{L}|^{p-\nu} \left( \prod_{1 \leq i < j \leq \nu} (\ell_i - \ell_j) \right) (d\mathbf{L})(\mathbf{V}^\top d\mathbf{V}),$$

which yields the p.d.f. of  $\mathbf{S}$ . □

Equation (3.3) is the joint (unnormalized) p.d.f. of nonzero eigenvalues and the corresponding eigenvectors for the Wishart matrix  $\mathbf{S} \in \mathbb{S}_{p,\nu}^{(+)}$ . Note that, for  $\nu = p$ ,  $\mathbf{S} \in \mathbb{S}_p^{(+)}$  and then the p.d.f. of  $\mathcal{W}_p(n, \mathbf{\Sigma})$  on  $\mathbb{S}_p^{(+)}$  can be expressed as

$$\frac{1}{2^{np/2}|\mathbf{\Sigma}|^{n/2}\Gamma_p(n/2)} |\mathbf{S}|^{(n-p-1)/2} \exp\left(-\frac{1}{2} \text{tr } \mathbf{\Sigma}^{-1} \mathbf{S}\right),$$

because  $|\mathbf{L}| = |\mathbf{S}|$  in Proposition 3.7. When  $p = 1$  and  $\mathbf{\Sigma} = 1$ , the p.d.f. above is the same as that of the chi-square distribution with  $n$  degrees of freedom. Thus the Wishart distribution is a multivariate generalization of the chi-square distribution.

The following proposition is an important result on expectation of the Wishart matrix, which is the basis for unbiased estimation of  $\mathbf{\Sigma}$ .

**Proposition 3.8** *Let  $\mathbf{S} \sim \mathcal{W}_p^\nu(n, \mathbf{\Sigma})$ . Then  $E[\mathbf{S}] = n\mathbf{\Sigma}$ .*

**Proof** Recall that  $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$ , where  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{\Sigma})$ . Thus, this proposition can immediately be verified by Corollary 3.1. □

From Proposition 3.8,  $\mathcal{W}_p^\nu(n, \mathbf{\Sigma})$  is also called the Wishart distribution with  $n$  degrees of freedom and mean  $n\mathbf{\Sigma}$ .

Many interesting results on the Wishart distribution have already been obtained in the literature. For other results on the Wishart distribution, see Gupta and Nagar (1999) and Muirhead (1982).

### 3.4 The Cholesky Decomposition of the Wishart Matrix

Here, we provide distribution theories related to the Cholesky decomposition of the Wishart matrix. This decomposition is also named the Bartlett decomposition in statistics.

Let  $\mathbf{X} = (x_{ij}) \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ . The Wishart matrix is defined by  $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$ . In the case of  $p > n$ , partition  $\mathbf{X}$  as  $\mathbf{X} = (\mathbf{X}_1^\top, \mathbf{X}_2^\top)$ , where  $\mathbf{X}_1 \in \mathbb{R}^{n \times n}$ , and denote by  $\mathbf{X}_1 = \mathbf{T}_1 \mathbf{Q}^\top$  uniquely the LQ decomposition of  $\mathbf{X}_1$ , where  $\mathbf{T}_1 \in \mathbb{L}_n^{(+)}$  and  $\mathbf{Q} \in \mathbb{O}_n$ . Here,

$$\mathbf{X}^\top = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{T}_1 \mathbf{Q}^\top \\ \mathbf{X}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{T}_1 \\ \mathbf{X}_2 \mathbf{Q} \end{pmatrix} \mathbf{Q}^\top \equiv \mathbf{T} \mathbf{Q}^\top, \quad (3.4)$$

where  $\mathbf{T} = (\mathbf{T}_1^\top, \mathbf{T}_2^\top)^\top \in \mathbb{L}_{p,n}^{(+)}$  and  $\mathbf{T}_2 = \mathbf{X}_2 \mathbf{Q} \in \mathbb{R}^{(p-n) \times n}$ , which yields  $\mathbf{S} = \mathbf{T} \mathbf{T}^\top$ . When  $n \geq p$ , the usual Cholesky decomposition of  $\mathbf{S}$  is given by  $\mathbf{S} = \mathbf{T} \mathbf{T}^\top$ , where  $\mathbf{T} \in \mathbb{L}_p^{(+)} = \mathbb{L}_{p,p}^{(+)}$ . Hence the above decompositions for the cases of  $n \geq p$  and  $p > n$  can be integrated into  $\mathbf{S} = \mathbf{T} \mathbf{T}^\top$ , where a unique  $\mathbf{T} \in \mathbb{L}_{p,v}^{(+)}$  with  $v = n \wedge p$ . Then we have the following proposition.

**Proposition 3.9** *The p.d.f. of  $\mathbf{T} = (t_{ij}) \in \mathbb{L}_{p,v}^{(+)}$  with  $v = n \wedge p$  is given by*

$$\frac{2^v \pi^{nv/2}}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2} \Gamma_v(n/2)} \exp\left(-\frac{1}{2} \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{T} \mathbf{T}^\top\right) \prod_{i=1}^v t_{ii}^{n-i}. \quad (3.5)$$

**Proof** For the  $n \geq p$  case, let  $\mathbf{X}^\top = \mathbf{T} \mathbf{Q}^\top$  be the LQ decomposition of  $\mathbf{X}^\top$ , where  $\mathbf{T} \in \mathbb{L}_p^{(+)}$  and  $\mathbf{Q} \in \mathbb{V}_{n,p}$ . Making the transformation  $\mathbf{X}^\top \rightarrow (\mathbf{T}, \mathbf{Q})$  in (3.2) and using Lemma 3.3, we can write the joint (unnormalized) p.d.f. of  $\mathbf{T}$  and  $\mathbf{Q}$  as

$$\frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} \exp\left(-\frac{1}{2} \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{T} \mathbf{T}^\top\right) \left(\prod_{i=1}^p t_{ii}^{n-i}\right) (d\mathbf{T})(\mathbf{Q}^\top d\mathbf{Q}).$$

From (3.1), integrating out with respect to  $\mathbf{Q}$  gives

$$\frac{2^p \pi^{np/2}}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2} \Gamma_p(n/2)} \exp\left(-\frac{1}{2} \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{T} \mathbf{T}^\top\right) \left(\prod_{i=1}^p t_{ii}^{n-i}\right) (d\mathbf{T}).$$

This is the p.d.f. of  $\mathbf{T}$  for the  $n \geq p$  case.



When  $p > n$ , noting from (3.4) that  $(dX_2) = (dT_2)$ , we see

$$(dX^\top) = (dX_1)(dX_2) = \left( \prod_{i=1}^n t_{ii}^{n-i} \right) (dT_1)(Q^\top dQ)(dX_2) = \left( \prod_{i=1}^n t_{ii}^{n-i} \right) (dT)(Q^\top dQ).$$

Hence the joint (unnormalized) p.d.f. of  $T \in \mathbb{L}_n^{(+)}$  and  $Q \in \mathbb{O}_n = \mathbb{V}_{n,n}$  can be expressed as

$$\frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} \exp \left( -\frac{1}{2} \text{tr} \Sigma^{-1} T T^\top \right) \left( \prod_{i=1}^n t_{ii}^{n-i} \right) (dT)(Q^\top dQ).$$

Using (3.1) gives  $\int_{\mathbb{V}_{n,n}} (Q^\top dQ) = 2^n \pi^{n^2/2} / \Gamma_n(n/2)$ , yielding the p.d.f. of  $T$  given in (3.5) with  $v = n$ .  $\square$

For deriving moments of  $T$ , we provide the distributions of nonzero elements in each column of  $T$ . Define the Cholesky decomposition of  $\Sigma$  as  $\Sigma = \Xi \Xi^\top$ , where  $\Xi = (\xi_{i,j}) \in \mathbb{L}_p^{(+)}$ . Denote  $\Xi_{(1)} = \Xi$ ,  $\Xi_{(p)} = \xi_{p,p}$  and, for  $i = 1, \dots, p-1$ ,

$$\Xi_{(i)} = \begin{pmatrix} \xi_{i,i} & \mathbf{0}_{p-i}^\top \\ \xi_{(i)} & \Xi_{(i+1)} \end{pmatrix} \in \mathbb{L}_{p-i+1}^{(+)},$$

where

$$\xi_{(i)} = \begin{pmatrix} \xi_{i+1,i} \\ \vdots \\ \xi_{p,i} \end{pmatrix} \in \mathbb{R}^{p-i}, \quad \Xi_{(i+1)} = \begin{pmatrix} \xi_{i+1,i+1} & 0 \\ \vdots & \ddots \\ \xi_{p,i+1} & \cdots \xi_{p,p} \end{pmatrix} \in \mathbb{L}_{p-i}^{(+)}.$$

Let  $\gamma_{(i)} = (\gamma_{i+1,i}, \dots, \gamma_{p,i})^\top = \xi_{i,i}^{-1} \xi_{(i)}$  for  $i = 1, \dots, p-1$  and let  $\sigma_i^2 = \xi_{i,i}^2$  for  $i = 1, \dots, p$ . For  $i = 1, \dots, p$ , let  $\Sigma_{(i)} = \Xi_{(i)} \Xi_{(i)}^\top$ , where  $\Sigma_{(1)} = \Sigma$  and  $\Sigma_{(p)} = \sigma_p^2$ . Note that for  $i = 1, \dots, p-1$

$$\Sigma_{(i)} = \begin{pmatrix} 1 & \mathbf{0}_{p-i}^\top \\ \gamma_{(i)} & \mathbf{I}_{p-i} \end{pmatrix} \begin{pmatrix} \sigma_i^2 & \mathbf{0}_{p-i}^\top \\ \mathbf{0}_{p-i} & \Sigma_{(i+1)} \end{pmatrix} \begin{pmatrix} 1 & \gamma_{(i)}^\top \\ \mathbf{0}_{p-i} & \mathbf{I}_{p-i} \end{pmatrix}. \quad (3.6)$$

Similarly, for  $i = 1, \dots, v$ , let  $T_{(i)}$  be submatrices obtained by removing the first  $(i-1)$  rows and columns of  $T = (t_{i,j}) \in \mathbb{L}_{p,v}^{(+)}$ . For  $i = 1, \dots, v-1$ , partition  $T_{(i)}$  into four blocks as

$$T_{(i)} = \begin{pmatrix} t_{i,i} & \mathbf{0}_{p-i}^\top \\ \mathbf{t}_{(i)} & T_{(i+1)} \end{pmatrix} \in \mathbb{L}_{p-i+1,v-i+1}^{(+)},$$

where  $\mathbf{t}_{(i)} = (t_{i+1,i}, \dots, t_{p,i})^\top$ . Note that for  $i = 1, \dots, v-1$

$$T_{(i)} T_{(i)}^\top = \begin{pmatrix} 1 & \mathbf{0}_{p-i}^\top \\ t_{i,i}^{-1} \mathbf{t}_{(i)} & \mathbf{I}_{p-i} \end{pmatrix} \begin{pmatrix} t_{i,i}^2 & \mathbf{0}_{p-i}^\top \\ \mathbf{0}_{p-i} & T_{(i+1)} T_{(i+1)}^\top \end{pmatrix} \begin{pmatrix} 1 & t_{i,i}^{-1} \mathbf{t}_{(i)}^\top \\ \mathbf{0}_{p-i} & \mathbf{I}_{p-i} \end{pmatrix} \quad (3.7)$$

and

$$\mathbf{T}_{(v)} = (t_{v,v}, t_{v+1,v}, \dots, t_{p,v})^\top = \begin{cases} t_{p,p} & \text{for } n \geq p \text{ } (v = p), \\ (t_{n,n}, t_{n+1,n}, \dots, t_{p,n})^\top & \text{for } p > n \text{ } (v = n). \end{cases}$$

In the  $p > n$  case, define additionally  $\mathbf{T}_{(n)} = (t_{n,n}, \mathbf{t}_{(n)}^\top)^\top$  with  $\mathbf{t}_{(n)} = (t_{n+1,n}, \dots, t_{p,n})^\top$ . Then we have the following proposition.

**Proposition 3.10** *The columns of  $\mathbf{T}$  are mutually independent and*

$$\begin{cases} t_{i,i}^2 \sim \sigma_i^2 \chi_{n-i+1}^2 & \text{for } i = 1, \dots, v, \\ \mathbf{t}_{(i)} | t_{i,i} \sim \mathcal{N}_{p-i}(t_{i,i} \boldsymbol{\gamma}_{(i)}, \boldsymbol{\Sigma}_{(i+1)}) & \text{for } i = 1, \dots, n \wedge (p-1). \end{cases} \quad (3.8)$$

**Proof** From Proposition 3.9, it is immediately seen that the columns of  $\mathbf{T}$  are mutually independent. Using (3.6) and (3.7) gives

$$\text{tr } \boldsymbol{\Sigma}_{(i)}^{-1} \mathbf{T}_{(i)} \mathbf{T}_{(i)}^\top = t_{i,i}^2 / \sigma_i^2 + (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)})^\top \boldsymbol{\Sigma}_{(i+1)}^{-1} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)}) + \text{tr } \boldsymbol{\Sigma}_{(i+1)}^{-1} \mathbf{T}_{(i+1)} \mathbf{T}_{(i+1)}^\top,$$

so that

$$\begin{aligned} \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{T} \mathbf{T}^\top &= t_{1,1}^2 / \sigma_1^2 + (\mathbf{t}_{(1)} - t_{1,1} \boldsymbol{\gamma}_{(1)})^\top \boldsymbol{\Sigma}_{(2)}^{-1} (\mathbf{t}_{(1)} - t_{1,1} \boldsymbol{\gamma}_{(1)}) + \text{tr } \boldsymbol{\Sigma}_{(2)}^{-1} \mathbf{T}_{(2)} \mathbf{T}_{(2)}^\top \\ &= \sum_{i=1}^2 \frac{t_{i,i}^2}{\sigma_i^2} + \sum_{i=1}^2 (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)})^\top \boldsymbol{\Sigma}_{(i+1)}^{-1} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)}) + \text{tr } \boldsymbol{\Sigma}_{(3)}^{-1} \mathbf{T}_{(3)} \mathbf{T}_{(3)}^\top \\ &= \dots \\ &= \sum_{i=1}^v \frac{t_{i,i}^2}{\sigma_i^2} + \sum_{i=1}^{n \wedge (p-1)} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)})^\top \boldsymbol{\Sigma}_{(i+1)}^{-1} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)}). \end{aligned}$$

Also, for  $i = 1, \dots, p-1$ ,

$$|\boldsymbol{\Sigma}| = |\boldsymbol{\Sigma}_{(2)}| \sigma_1^2 = |\boldsymbol{\Sigma}_{(3)}| \prod_{j=1}^2 \sigma_j^2 = \dots = |\boldsymbol{\Sigma}_{(i+1)}| \prod_{j=1}^i \sigma_j^2,$$

yielding

$$|\boldsymbol{\Sigma}|^{n/2} = \left( \prod_{i=1}^{n \wedge (p-1)} |\boldsymbol{\Sigma}_{(i+1)}|^{1/2} \right) \left( \prod_{i=1}^v (\sigma_i^2)^{(n-i+1)/2} \right).$$

It turns out from Proposition 3.1 that

$$\frac{2^v \pi^{nv/2}}{(2\pi)^{np/2} \Gamma_v(n/2)} = \left( \prod_{i=1}^{n \wedge (p-1)} \frac{1}{(2\pi)^{(p-i)/2}} \right) \left( \prod_{i=1}^v \frac{2}{2^{(n-i+1)/2} \Gamma((n-i+1)/2)} \right),$$

so that the p.d.f. of  $\mathbf{T}$ , given in Proposition 3.9, can be rewritten as

$$\prod_{i=1}^{n \wedge (p-1)} \frac{1}{(2\pi)^{(p-i)/2} |\boldsymbol{\Sigma}_{(i+1)}|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)})^\top \boldsymbol{\Sigma}_{(i+1)}^{-1} (\mathbf{t}_{(i)} - t_{i,i} \boldsymbol{\gamma}_{(i)}) \right) \\ \times \prod_{i=1}^v \frac{2}{2^{(n-i+1)/2} \Gamma((n-i+1)/2)} (\sigma_i^2)^{-(n-i+1)/2} t_{i,i}^{n-i} \exp \left( -\frac{t_{i,i}^2}{2\sigma_i^2} \right).$$

Thus, for  $i = 1, \dots, n \wedge (p-1)$ ,  $\mathbf{t}_{(i)} | t_{i,i} \sim \mathcal{N}_{p-i}(t_{i,i} \boldsymbol{\gamma}_{(i)}, \boldsymbol{\Sigma}_{(i+1)})$ . Finally, making the change of variables  $y_i = t_{i,i}^2$  for  $i = 1, \dots, v$  gives  $y_i \sim \sigma_i^2 \chi_{n-i+1}^2$ .  $\square$

The distributional decomposition (3.8) will be used in Sect. 7.6 for estimation of a covariance matrix.

Proposition 3.10 immediately provides the following corollary.

**Corollary 3.2** *Let  $\mathbf{X} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{I}_p)$ . Denote by  $\mathbf{X}^\top \mathbf{X} = \mathbf{T} \mathbf{T}^\top$  the Cholesky decomposition of  $\mathbf{X}^\top \mathbf{X}$ , where  $\mathbf{T} = (t_{i,j}) \in \mathbb{L}_{p,v}^{(+)}$  with  $v = n \wedge p$ . Then all the nonzero elements of  $\mathbf{T}$  are mutually independent and distributed as*

$$\begin{cases} t_{i,i}^2 \sim \chi_{n-i+1}^2 & \text{for } i = 1, \dots, v, \\ t_{i,j} \sim \mathcal{N}(0, 1) & \text{for } 2 \leq i \leq p \text{ and } 1 \leq j \leq n \wedge (i-1). \end{cases}$$

## References

- J.A. Díaz-García, R. Gutierrez-Jáimez, K.V. Mardia, Wishart and pseudo-Wishart distributions and some applications to shape theory. *J. Multivar. Anal.* **63**, 73–87 (1997)
- A.K. Gupta, D.K. Nagar, *Matrix Variate Distributions* (Chapman & Hall/CRC, New York, 1999)
- A.M. Mathai, *Jacobian of Matrix Transformations and Functions of Matrix Argument* (World Scientific, Singapore, 1997)
- R.J. Muirhead, *Aspects of Multivariate Statistical Theory* (Wiley, New York, 1982)
- M.S. Srivastava, C.G. Khatri, *An Introduction to Multivariate Statistics* (North Holland, New York, 1979)
- H. Uhlig, On singular Wishart and singular multivariate Beta distributions. *Ann. Stat.* **22**, 395–405 (1994)
- J. Wishart, The generalised product moment distribution in samples from a normal multivariate population. *Biometrika* **20**, 32–52 (1928)

# Chapter 4

## Multivariate Linear Model and Group Invariance



Multivariate linear model is a multivariate generalization for the dimension of response variable in traditional multiple linear regression model. This chapter provides some fundamental properties in terms of the multivariate linear model and the corresponding canonical form. The group invariance is also explained for shrinkage estimation in the multivariate linear model.

### 4.1 Multivariate Linear Model

The sample size is denoted by  $N$ . For  $i = 1, \dots, N$ , let  $\mathbf{y}_i$  be a  $p$ -dimensional column vector of response variables and let

$$\mathbf{y}_i = \mathbf{B}^\top \mathbf{x}_i + \boldsymbol{\varepsilon}_i,$$

where  $\mathbf{x}_i \in \mathbb{R}^m$  is a column vector of known explanatory variables ( $m \leq N$ ),  $\mathbf{B} \in \mathbb{R}^{m \times p}$  is a matrix of unknown parameters and  $\boldsymbol{\varepsilon}_i \in \mathbb{R}^p$  is a random vector. Define  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)^\top \in \mathbb{R}^{N \times p}$ ,  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top \in \mathbb{R}^{N \times m}$  and  $\mathbf{E} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N)^\top \in \mathbb{R}^{N \times p}$ . Then we obtain

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}. \quad (4.1)$$

Assume that for  $i = 1, \dots, N$  the  $\boldsymbol{\varepsilon}_i$ 's are independently distributed as  $\mathcal{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma})$ , namely, the error matrix  $\mathbf{E}$  follows  $\mathcal{N}_{N \times p}(\mathbf{0}_{N \times p}, \mathbf{I}_N \otimes \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma} \in \mathbb{S}_p^{(+)}$  is unknown. In addition,  $\mathbf{X}$  is assumed to be of full column rank.

In the literature, model (4.1) is called the multivariate linear model or multivariate linear regression model, and the parameter matrix  $\mathbf{B}$  is called the regression coefficient matrix or simply the regression matrix. The multivariate linear model (4.1) is closely relevant to various models from simple to complex, such as a simple

mean-variance model, the multivariate analysis of variance (MANOVA) model and the growth curve model.

We now consider estimation of  $\mathbf{B}$ . There are many procedures for estimation of  $\mathbf{B}$ , and we here introduce the least squares method that is one of the most well-known procedures. Let

$$g(\mathbf{B}) = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{B}^\top \mathbf{x}_i\|^2,$$

where  $\|\cdot\|$  denotes the usual Euclidean norm. Then the least squares method is the minimization of  $g(\mathbf{B})$  subject to  $\mathbf{B}$ . Noting that

$$\mathbf{Y}^\top \mathbf{Y} = \sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i^\top, \quad \mathbf{X}^\top \mathbf{Y} = \sum_{i=1}^N \mathbf{x}_i \mathbf{y}_i^\top, \quad \mathbf{X}^\top \mathbf{X} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top,$$

we observe that

$$\begin{aligned} g(\mathbf{B}) &= \sum_{i=1}^N \text{tr}(\mathbf{y}_i - \mathbf{B}^\top \mathbf{x}_i)(\mathbf{y}_i - \mathbf{B}^\top \mathbf{x}_i)^\top \\ &= \text{tr}(\mathbf{Y}^\top \mathbf{Y} - \mathbf{B}^\top \mathbf{X}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X} \mathbf{B} + \mathbf{B}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}). \end{aligned}$$

Recall that  $\mathbf{X}$  is of full column rank, and thus  $\mathbf{X}^\top \mathbf{X}$  is nonsingular. Completing the square with respect to  $\mathbf{B}$  gives

$$g(\mathbf{B}) = \text{tr}(\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) + \text{tr}(\mathbf{Y}^\top \mathbf{Y} - \widehat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \widehat{\mathbf{B}}),$$

where

$$\widehat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Clearly  $g(\mathbf{B})$  is minimized at  $\mathbf{B} = \widehat{\mathbf{B}}$ . The resulting  $\widehat{\mathbf{B}}$  is called the least squares (LS) estimator of  $\mathbf{B}$ . Since  $\mathbf{E} \sim \mathcal{N}_{N \times p}(\mathbf{0}_{N \times p}, \mathbf{I}_N \otimes \boldsymbol{\Sigma})$  in (4.1), namely,  $\mathbf{Y} \sim \mathcal{N}_{N \times p}(\mathbf{X} \mathbf{B}, \mathbf{I}_N \otimes \boldsymbol{\Sigma})$ , the likelihood without a normalizing constant can be written as

$$|\boldsymbol{\Sigma}|^{-N/2} \exp \left[ -\frac{1}{2} \text{tr} \boldsymbol{\Sigma}^{-1} \{(\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) + (\mathbf{Y}^\top \mathbf{Y} - \widehat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \widehat{\mathbf{B}})\} \right].$$

Thus the LS estimator  $\widehat{\mathbf{B}}$  is the maximum likelihood estimator as well. Also, using Proposition 3.5 leads to

$$\widehat{\mathbf{B}} \sim \mathcal{N}_{m \times p}(\mathbf{B}, (\mathbf{X}^\top \mathbf{X})^{-1} \otimes \boldsymbol{\Sigma}),$$

implying that  $\widehat{\mathbf{B}}$  is an unbiased estimator of  $\mathbf{B}$ .

Next we find a reasonable estimator for the covariance matrix  $\Sigma$  of random error terms  $\epsilon_i$ 's in the multivariate linear model (4.1). In traditional multiple linear regression model, a residual sum of squares is often used to estimate an error variance, which can be extended to estimating the error covariance  $\Sigma$ .

Let

$$S = \sum_{i=1}^N (y_i - \widehat{B}^\top x_i)(y_i - \widehat{B}^\top x_i)^\top,$$

which is called the residual sum of squares matrix. It turns out that

$$\begin{aligned} \sum_{i=1}^N (y_i - \widehat{B}^\top x_i)(y_i - \widehat{B}^\top x_i)^\top &= Y^\top Y - \widehat{B}^\top X^\top Y - Y^\top X \widehat{B} + \widehat{B}^\top X^\top X \widehat{B} \\ &= Y^\top \{I_N - X(X^\top X)^{-1} X^\top\} Y, \end{aligned}$$

so that

$$S = Y^\top (I_N - P_X) Y,$$

where  $P_X = X(X^\top X)^{-1} X^\top$ . Here,  $P_X$  is the orthogonal projection matrix onto the subspace spanned by columns of  $X$ . Now, we write the QR decomposition of  $X$  as  $X = Q_1 L^\top$ , where  $Q_1 \in \mathbb{V}_{N,m}$  and  $L \in \mathbb{L}_m^{(+)}$ . There exists  $Q_2 \in \mathbb{V}_{N,n}$  for  $n = N - m$  such that  $(Q_1, Q_2) \in \mathbb{O}_N$ , namely, the set of  $N$  columns of an  $N \times N$  matrix  $(Q_1, Q_2)$  forms an orthonormal basis of  $\mathbb{R}^N$ . Then

$$P_X = Q_1 L^\top (L Q_1^\top Q_1 L^\top)^{-1} L Q_1^\top = Q_1 Q_1^\top$$

and  $I_N - P_X = I_N - Q_1 Q_1^\top = Q_2 Q_2^\top$ . Since  $Q_2^\top X B = Q_2^\top Q_1 L^\top B = \mathbf{0}_{n \times p}$ , using Proposition 3.5 gives  $Q_2^\top Y \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, I_n \otimes \Sigma)$ . From Definition 3.3,

$$S = Y^\top Q_2 Q_2^\top Y \sim \mathcal{W}_p^v(n, \Sigma)$$

with  $v = n \wedge p$ . Hence, by Proposition 3.8,

$$\widehat{\Sigma}^{UB} = \frac{1}{n} S, \quad n = N - m,$$

is unbiased for  $\Sigma$ .

The unbiased estimator  $\widehat{\Sigma}^{UB}$  is reasonable, but not the maximum likelihood estimator of  $\Sigma$ . Note that  $(\widehat{B}, \widehat{\Sigma}^{UB})$ , or  $(\widehat{B}, S)$ , is a sufficient statistic for  $(B, \Sigma)$ . See Muirhead (1982, Theorem 10.1.1) or Anderson (2003, Corollary 8.2.1).

## 4.2 A Canonical Form

Let  $\mathbf{Q} = (\mathbf{Q}_1, \mathbf{Q}_2) \in \mathbb{O}_N$ , where  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$  are defined as in the previous section. Then

$$\mathbf{Q}^\top \mathbf{X} = \begin{pmatrix} \mathbf{Q}_1^\top \mathbf{X} \\ \mathbf{Q}_2^\top \mathbf{X} \end{pmatrix} = \begin{pmatrix} \mathbf{L}^\top \\ \mathbf{0}_{n \times m} \end{pmatrix}$$

for  $n = N - m$ . Let  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m)^\top = \mathbf{L}^\top \mathbf{B} \in \mathbb{R}^{m \times p}$ . Define

$$\begin{pmatrix} \mathbf{Z} \\ \mathbf{U} \end{pmatrix} = \mathbf{Q}^\top \mathbf{Y} = \begin{pmatrix} \mathbf{Q}_1^\top \mathbf{Y} \\ \mathbf{Q}_2^\top \mathbf{Y} \end{pmatrix},$$

where  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_m)^\top \in \mathbb{R}^{m \times p}$  and  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)^\top \in \mathbb{R}^{n \times p}$ .

Recall that  $\mathbf{Y} \sim \mathcal{N}_{N \times p}(\mathbf{X}\mathbf{B}, \mathbf{I}_N \otimes \boldsymbol{\Sigma})$ , so that by Proposition 3.5

$$\begin{pmatrix} \mathbf{Z} \\ \mathbf{U} \end{pmatrix} \sim \mathcal{N}_{N \times p} \left( \begin{pmatrix} \boldsymbol{\Theta} \\ \mathbf{0}_{n \times p} \end{pmatrix}, \mathbf{I}_N \otimes \boldsymbol{\Sigma} \right).$$

Hence from Proposition 3.6,  $\mathbf{Z}$  and  $\mathbf{U}$  are independently distributed as

$$\mathbf{Z} \sim \mathcal{N}_{m \times p}(\boldsymbol{\Theta}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}), \quad \mathbf{U} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \boldsymbol{\Sigma}), \quad (4.2)$$

which is a canonical form of the multivariate linear model (4.1). The canonical form (4.2) is equivalent to

$$\begin{aligned} \mathbf{z}_i &\sim \mathcal{N}_p(\boldsymbol{\theta}_i, \boldsymbol{\Sigma}), \quad i = 1, \dots, m, \\ \mathbf{u}_j &\sim \mathcal{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma}), \quad j = 1, \dots, n, \end{aligned}$$

where all of the  $\mathbf{z}_i$ 's and  $\mathbf{u}_j$ 's are mutually independent.

The QR decomposition of  $\mathbf{X}$  is denoted by  $\mathbf{X} = \mathbf{Q}_1 \mathbf{L}^\top$ , yielding

$$\widehat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = (\mathbf{L} \mathbf{Q}_1^\top \mathbf{Q}_1 \mathbf{L}^\top)^{-1} \mathbf{L} \mathbf{Q}_1^\top \mathbf{Y} = (\mathbf{L}^\top)^{-1} \mathbf{Z},$$

namely,  $\mathbf{Z} = \mathbf{L}^\top \widehat{\mathbf{B}}$ . Here,

$$\mathbf{S} = \mathbf{U}^\top \mathbf{U} = \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top$$

is also called the Wishart matrix, which follows  $\mathcal{W}_p^v(n, \boldsymbol{\Sigma})$  with  $v = n \wedge p$ . Thus  $\widehat{\mathbf{B}}$  and  $\mathbf{S}$  are, respectively, made from  $\mathbf{Z}$  and  $\mathbf{U}$ , and are mutually independent. From (4.2),  $(\mathbf{Z}, \mathbf{U})$  is a sufficient statistic of  $(\boldsymbol{\Theta}, \boldsymbol{\Sigma})$ .

In this book,  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Sigma}$  are hereinafter called the mean and the covariance matrices, respectively, and the estimation problem for  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Sigma}$  will be treated in the canonical form (4.2).

### 4.3 Group Invariance

If an estimation problem is invariant under a transformation such as translation, scaling and rotation, it seems reasonable to require all decision rules, namely, all possible estimators are invariant under the transformation. In this section, we briefly provide the definition of group invariance, or simply called invariance, and simple examples in the canonical form (4.2) of the multivariate linear model (4.1).

Let  $P_\theta$  be a probability distribution on a sample space  $\mathcal{X}$  parameterized by  $\theta$ . A statistical model  $\mathcal{M}$  is defined by  $\mathcal{M} = \{P_\theta : \theta \in \mathcal{P}\}$ , namely, it is a family of probability distributions  $P_\theta$ , where  $\mathcal{P}$  is a parameter space. Let  $x$  be an observed random variable having  $P_\theta$  and suppose  $\theta$  is estimated by using  $x$ . Denote by  $\hat{\theta} = \hat{\theta}(x)$  an estimator of  $\theta$  based on  $x$ , and by  $\mathcal{D}$  a decision space consisting of all possible estimators  $\hat{\theta}$ . A distance between  $\theta$  and its estimator  $\hat{\theta}$  is measured by a loss function  $L(\hat{\theta}, \theta)$ .

Let  $\mathbb{G}$  be a transformation group which acts on  $\mathcal{X}$ . For any  $g \in \mathbb{G}$ , the group action on  $x \in \mathcal{X}$  is written as  $gx$ . A statistical model  $\mathcal{M}$  is said to be invariant under  $\mathbb{G}$  if, for any  $g \in \mathbb{G}$  and  $\theta \in \mathcal{P}$ , there exists a unique  $\tilde{g}\theta \in \mathcal{P}$  such that a distribution of  $gx$  is  $P_{\tilde{g}\theta} \in \mathcal{M}$ . For an invariant statistical model  $\mathcal{M}$  under  $\mathbb{G}$ , all of the  $\tilde{g}$ 's form a group of transformations from  $\mathcal{P}$  into itself and it is called the group induced by  $\mathbb{G}$ . When a statistical model  $\mathcal{M}$  is invariant under  $\mathbb{G}$ , a loss function  $L(\hat{\theta}, \theta)$  is said to be invariant under  $\mathbb{G}$  if, for any  $g \in \mathbb{G}$  and  $\hat{\theta} \in \mathcal{D}$ , there exists an estimator  $\tilde{g}\hat{\theta}$  such that  $L(\tilde{g}\hat{\theta}, \tilde{g}\theta) = L(\hat{\theta}, \theta)$  for any  $\theta \in \mathcal{P}$ . Note that all of the  $\tilde{g}$ 's also form a group of transformations from  $\mathcal{D}$  into itself. An estimator  $\hat{\theta}(x)$  is said to be invariant under  $\mathbb{G}$  if  $\hat{\theta}(gx) = \tilde{g}\hat{\theta}(x)$  for any  $g \in \mathbb{G}$  and  $x \in \mathcal{X}$ . An estimation problem is said to be invariant under  $\mathbb{G}$  if the model  $\mathcal{M}$  and the loss function  $L$  are invariant under  $\mathbb{G}$ .

Now, simple examples of invariance are given in terms of the canonical form (4.2) of the multivariate linear model (4.1), which is rewritten as

$$X \sim \mathcal{N}_{m \times p}(\Theta, I_m \otimes \Sigma), \quad Y \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, I_n \otimes \Sigma), \quad (4.3)$$

where  $\Theta \in \mathbb{R}^{m \times p}$  and  $\Sigma \in \mathbb{S}_p^{(+)}$ . Consider the problem of estimating the mean matrix  $\Theta$  relative to the quadratic loss  $L_Q(\hat{\Theta}, \theta) = \text{tr}(\hat{\Theta} - \theta)\Sigma^{-1}(\hat{\Theta} - \theta)^\top$ , where  $\hat{\Theta} = \hat{\Theta}(x)$  with  $x = (X, Y)$  and  $\theta = (\Theta, \Sigma)$ .

In the model (4.3), the sample and the parameter spaces are, respectively,  $\mathcal{X} = \mathbb{R}^{m \times p} \times \mathbb{R}^{n \times p}$  and  $\mathcal{P} = \mathbb{R}^{m \times p} \times \mathbb{S}_p^{(+)}$ . The decision space is denoted by  $\mathcal{D} = \{\hat{\Theta}(x) \in \mathbb{R}^{m \times p} : x \in \mathcal{X}\}$ . We define the transformation group  $\mathbb{G}$  as

$$\mathbb{G} = \{(\mathbf{O}, U) : \mathbf{O} \in \mathbb{O}_m \text{ and } U \in \mathbb{U}_p\}$$

with operation  $g_1 g_2 = (\mathbf{O}_1 \mathbf{O}_2, U_2 U_1)$  for any  $g_1 = (\mathbf{O}_1, U_1)$ ,  $g_2 = (\mathbf{O}_2, U_2) \in \mathbb{G}$ . The group action of  $g = (\mathbf{O}, U) \in \mathbb{G}$  on  $x \in \mathcal{X}$  and the induced group actions  $\tilde{g}$  on  $\theta \in \mathcal{P}$  and  $\tilde{g}$  on  $\hat{\Theta} \in \mathcal{D}$  are given, respectively, by scale transformations



$$\begin{aligned}
x &\rightarrow gx = (\mathbf{O}XU, YU), \\
\theta &\rightarrow \bar{g}\theta = (\mathbf{O}\Theta U, U^\top \Sigma U), \\
\hat{\Theta} &\rightarrow \tilde{g}\hat{\Theta} = \mathbf{O}\hat{\Theta}U.
\end{aligned} \tag{4.4}$$

From Proposition 3.5, it is easy to see that the model (4.3) is invariant under  $\mathbb{G}$ . Also, the invariance of the  $L_Q$ -loss can be verified because

$$\begin{aligned}
L_Q(\tilde{g}\hat{\Theta}, \bar{g}\theta) &= \text{tr}(\mathbf{O}\hat{\Theta}U - \mathbf{O}\Theta U)(U^\top \Sigma U)^{-1}(\mathbf{O}\hat{\Theta}U - \mathbf{O}\Theta U)^\top \\
&= \text{tr}(\hat{\Theta} - \Theta)\Sigma^{-1}(\hat{\Theta} - \Theta)^\top = L_Q(\hat{\Theta}, \theta).
\end{aligned}$$

Thus the estimation problem of  $\Theta$  in (4.3) relative to the  $L_Q$ -loss is invariant under  $\mathbb{G}$ .

Let  $S = Y^\top Y$ . For  $n \geq p$ , let an estimator of  $\Theta$  be

$$\hat{\Theta} = \hat{\Theta}(x) = \begin{cases} X(I_p - c_1(X^\top X)^{-1}S) & \text{for } m > p, \\ (I_m - c_2(XS^{-1}X^\top)^{-1})X & \text{for } p \geq m, \end{cases}$$

where  $c_1$  and  $c_2$  are positive constants. When  $n \geq p$ ,  $S$  belongs to  $\mathbb{S}_p^{(+)}$  with probability one. Here  $\hat{\Theta}$  is invariant under  $\mathbb{G}$ . Indeed, for  $m > p$ ,

$$\begin{aligned}
\hat{\Theta}(gx) &= \mathbf{O}XU[I_p - c_1\{(\mathbf{O}XU)^\top \mathbf{O}XU\}^{-1}U^\top SU] \\
&= \mathbf{O}X\{I_p - c_1(X^\top X)^{-1}S\}U \\
&= \mathbf{O}\hat{\Theta}(x)U = \tilde{g}\hat{\Theta}(x),
\end{aligned}$$

and for  $p \geq m$

$$\begin{aligned}
\hat{\Theta}(gx) &= [I_m - c_2\{\mathbf{O}XU(U^\top SU)^{-1}(\mathbf{O}XU)^\top\}^{-1}]\mathbf{O}XU \\
&= \mathbf{O}\{I_m - c_2(XS^{-1}X^\top)^{-1}\}XU \\
&= \mathbf{O}\hat{\Theta}(x)U = \tilde{g}\hat{\Theta}(x).
\end{aligned}$$

On the other hand, for example, an estimator  $\hat{\Theta}_0 = X - X/\text{tr} XX^\top$  is not invariant under  $\mathbb{G}$ . Further, a quadratic-type loss  $L_F(\hat{\Theta}, \theta) = \text{tr}(\hat{\Theta} - \Theta)(\hat{\Theta} - \Theta)^\top$  is not invariant under  $\mathbb{G}$  since

$$\begin{aligned}
L_F(\tilde{g}\hat{\Theta}, \bar{g}\theta) &= \text{tr}(\mathbf{O}\hat{\Theta}U - \mathbf{O}\Theta U)(\mathbf{O}\hat{\Theta}U - \mathbf{O}\Theta U)^\top \\
&= \text{tr}(\hat{\Theta} - \Theta)UU^\top(\hat{\Theta} - \Theta)^\top \neq L_F(\hat{\Theta}, \theta).
\end{aligned}$$

We may need to discuss invariance in terms of the original multivariate linear model (4.1), but this is omitted. For invariance in estimation of the covariance matrix, see Chap. 7.

A general theory of invariance and its applications in statistics are discussed by Eaton (1983, 1989). For finding decision-theoretically optimal estimators, it is standard tactics to focus on a restricted class of invariant estimators. See also Lehmann and Casella (1998) for the role of invariance in decision-theoretic estimation.

## References

- T.W. Anderson, *An Introduction to Multivariate Statistical Analysis*, 3rd edn. (Wiley, New York, 2003)
- M.L. Eaton, *Multivariate Statistics: A Vector Space Approach* (Wiley, New York, 1983)
- M.L. Eaton, *Group Invariance Application in Statistics*. Regional Conference Series in Probability and Statistics, vol. 1 (Institute of Mathematical Statistics, Hayward, 1989)
- E.L. Lehmann, G. Casella, *Theory of Point Estimation*, 2nd edn. (Springer, New York, 1998)
- R.J. Muirhead, *Aspects of Multivariate Statistical Theory* (Wiley, New York, 1982)

## Chapter 5

# A Generalized Stein Identity and Matrix Differential Operators



In shrinkage estimation, the Stein (1973, 1981) identity is known as an integration by parts formula for deriving unbiased risk estimates. It is a simple but very powerful mathematical tool and has contributed significantly to the development of shrinkage estimation. This chapter provides a generalized Stein identity in matrix-variate normal distribution model and also some useful results on matrix differential operators for a unified application of the identity to high- and low-dimensional normal models.

## 5.1 Stein's Identity in Matrix-Variate Normal Distribution

For  $\mathbf{X} = (x_{ij}) \in \mathbb{R}^{m \times p}$ , denote by  $d_{ij}^{\mathbf{X}} = \partial / \partial x_{ij}$  the differential operator with respect to the  $(i, j)$ -th element of  $\mathbf{X}$ . The matrix differential operator with respect to  $\mathbf{X}$  is defined by  $\nabla_{\mathbf{X}} = (d_{ij}^{\mathbf{X}})$ , which is an  $m \times p$  matrix. Let  $\mathbf{G} = (g_{ij})$  be a  $p \times a$  matrix-valued function such that all the elements of  $\mathbf{G}$ ,  $g_{ij}$ 's, are absolutely continuous functions of  $\mathbf{X}$ . Define  $\nabla_{\mathbf{X}} \mathbf{G}$  as a usual matrix product: For  $i = 1, \dots, m$  and  $j = 1, \dots, a$ , the  $(i, j)$ -th element of  $\nabla_{\mathbf{X}} \mathbf{G}$  is  $\{\nabla_{\mathbf{X}} \mathbf{G}\}_{ij} = \sum_{k=1}^p d_{ik}^{\mathbf{X}} g_{kj}$ . As for a scalar-valued function  $f$ , an  $m \times p$  matrix  $\nabla_{\mathbf{X}} f(\mathbf{X})$  is defined by  $\nabla_{\mathbf{X}} f(\mathbf{X}) = \nabla_{\mathbf{X}} \{f(\mathbf{X}) \mathbf{I}_p\}$ . Then a generalized Stein identity is given in the following theorem.

**Theorem 5.1** *Let  $\mathbf{X} = (x_{ij}) \sim \mathcal{N}_{m \times p}(\boldsymbol{\Theta}, \boldsymbol{\Omega} \otimes \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Theta} = (\theta_{ij}) \in \mathbb{R}^{m \times p}$ ,  $\boldsymbol{\Omega} \in \mathbb{S}_m^{(+)}$  and  $\boldsymbol{\Sigma} \in \mathbb{S}_p^{(+)}$ . Let  $\mathbf{G}_1 \in \mathbb{R}^{a \times m}$  and  $\mathbf{G}_2 \in \mathbb{R}^{p \times b}$  such that*

- (i) *all the elements of  $\mathbf{G}_1$  and  $\mathbf{G}_2$  are absolutely continuous functions of  $\mathbf{X}$ ,*
- (ii)  *$E[\|\{\mathbf{G}_1(\mathbf{X} - \boldsymbol{\Theta})\mathbf{G}_2\}_{ij}\}] < \infty$  for any  $i \in \{1, \dots, a\}$  and  $j \in \{1, \dots, b\}$ .*

*Then*

$$E[\mathbf{G}_1(\mathbf{X} - \boldsymbol{\Theta})\mathbf{G}_2] = E[\mathbf{G}_1\boldsymbol{\Omega}\nabla_{\mathbf{X}}\boldsymbol{\Sigma}\mathbf{G}_2 + (\mathbf{G}_2^{\top}\boldsymbol{\Sigma}\nabla_{\mathbf{X}}^{\top}\boldsymbol{\Omega}\mathbf{G}_1^{\top})^{\top}]. \quad (5.1)$$

**Proof** Let  $\mathbf{\Gamma} = (\gamma_{ij}) \in \mathbb{U}_m$  and  $\mathbf{\Lambda} = (\lambda_{ij}) \in \mathbb{U}_p$  such that  $\mathbf{\Omega} = \mathbf{\Gamma}\mathbf{\Gamma}^\top$  and  $\mathbf{\Sigma} = \mathbf{\Lambda}\mathbf{\Lambda}^\top$ . Let  $\mathbf{Z} = (z_{ij}) = \mathbf{\Gamma}^{-1}\mathbf{X}(\mathbf{\Lambda}^\top)^{-1}$  and  $\mathbf{\Xi} = (\xi_{ij}) = \mathbf{\Gamma}^{-1}\mathbf{\Theta}(\mathbf{\Lambda}^\top)^{-1}$ . Note here from Proposition 3.5 that  $\mathbf{Z} \sim \mathcal{N}_{m \times p}(\mathbf{\Xi}, \mathbf{I}_m \otimes \mathbf{I}_p)$ , namely, the  $z_{ij}$ 's are independently distributed as  $z_{ij} \sim \mathcal{N}(\xi_{ij}, 1)$ .

Denote  $\mathbf{G}_1 = \mathbf{G}_1(\mathbf{X})$  and  $\mathbf{G}_2 = \mathbf{G}_2(\mathbf{X})$ . For  $i = 1, \dots, a$  and  $j = 1, \dots, b$ , let  $h_{ij}$  be the  $(i, j)$ -th element of  $E[\mathbf{G}_1(\mathbf{X} - \mathbf{\Theta})\mathbf{G}_2]$ , which is given by

$$\begin{aligned} h_{ij} &= E[\{\mathbf{G}_1(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\mathbf{\Gamma}(\mathbf{Z} - \mathbf{\Xi})\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\}_{ij}] \\ &= \sum_{k=1}^m \sum_{l=1}^p E[\{\mathbf{G}_1(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\mathbf{\Gamma}\}_{ik}(z_{kl} - \xi_{kl})\{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\}_{lj}]. \end{aligned}$$

By the integration by parts formula, or by the Stein identity (1.3),

$$\begin{aligned} h_{ij} &= \sum_{k=1}^m \sum_{l=1}^p E[\mathbf{d}_{kl}^Z [\{\mathbf{G}_1(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\mathbf{\Gamma}\}_{ik}\{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\}_{lj}]] \\ &= \sum_{k=1}^m \sum_{l=1}^p E[\{\mathbf{G}_1(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\mathbf{\Gamma}\}_{ik}\mathbf{d}_{kl}^Z\{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\}_{lj} \\ &\quad + \{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\}_{lj}\mathbf{d}_{kl}^Z\{\mathbf{G}_1(\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top)\mathbf{\Gamma}\}_{ik}], \end{aligned}$$

where  $\mathbf{d}_{kl}^Z = \partial/\partial z_{kl}$ . Since  $x_{ij} = \{\mathbf{\Gamma}\mathbf{Z}\mathbf{\Lambda}^\top\}_{ij} = \sum_{q=1}^m \sum_{r=1}^p \gamma_{iq} z_{qr} \lambda_{jr}$ , using the chain rule gives

$$\mathbf{d}_{kl}^Z = \sum_{i=1}^m \sum_{j=1}^p [\mathbf{d}_{kl}^Z x_{ij}] \cdot \mathbf{d}_{ij}^X = \sum_{i=1}^m \sum_{j=1}^p \gamma_{ik} \lambda_{jl} \cdot \mathbf{d}_{ij}^X = \{\mathbf{\Gamma}^\top \nabla_X \mathbf{\Lambda}\}_{kl}, \quad (5.2)$$

so that

$$\begin{aligned} h_{ij} &= \sum_{k=1}^m \sum_{l=1}^p E[\{\mathbf{G}_1(\mathbf{X})\mathbf{\Gamma}\}_{ik}\{\mathbf{\Gamma}^\top \nabla_X \mathbf{\Lambda}\}_{kl}\{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{X})\}_{lj} \\ &\quad + \{\mathbf{\Lambda}^\top\mathbf{G}_2(\mathbf{X})\}_{lj}\{\mathbf{\Gamma}^\top \nabla_X \mathbf{\Lambda}\}_{kl}\{\mathbf{G}_1(\mathbf{X})\mathbf{\Gamma}\}_{ik}] \\ &= E[\{\mathbf{G}_1\mathbf{\Gamma}\mathbf{\Gamma}^\top \nabla_X \mathbf{\Lambda}\mathbf{\Lambda}^\top\mathbf{G}_2\}_{ij} + \{\mathbf{G}_2^\top \mathbf{\Lambda}\mathbf{\Lambda}^\top \nabla_X^\top \mathbf{\Gamma}\mathbf{\Gamma}^\top \mathbf{G}_1^\top\}_{ji}] \\ &= E[\{\mathbf{G}_1\mathbf{\Omega} \nabla_X \mathbf{\Sigma} \mathbf{G}_2\}_{ij} + \{(\mathbf{G}_2^\top \mathbf{\Sigma} \nabla_X^\top \mathbf{\Omega} \mathbf{G}_1^\top)^\top\}_{ij}]. \end{aligned}$$

Thus the proof is complete.  $\square$

In the proof of Theorem 5.1, the differentiability of elements of  $\mathbf{G}_1$  and  $\mathbf{G}_2$  and the interchangeability of integrals are guaranteed, respectively, by conditions (i) and (ii) of Theorem 5.1.

From Theorem 5.1, if  $\mathbf{X} \sim \mathcal{N}_{m \times p}(\mathbf{\Theta}, \mathbf{I}_m \otimes \mathbf{\Sigma})$  then, under some suitable conditions,

$$E[\operatorname{tr}(\mathbf{X} - \mathbf{\Theta})\mathbf{\Sigma}^{-1}\mathbf{G}^\top] = E[\operatorname{tr} \nabla_{\mathbf{X}} \mathbf{G}^\top], \quad (5.3)$$

where all the elements of  $\mathbf{G} (\in \mathbb{R}^{m \times p})$  are absolutely continuous functions of  $\mathbf{X}$ . The r.h.s. of (5.3) depends only on expectations of the diagonals of  $\nabla_{\mathbf{X}} \mathbf{G}^\top$ , but not on those of the off-diagonals. See also Bilodeau and Kariya (1989) and Konno (1992) for Stein-type identities on matrix-variate normal distribution. The appendix of this chapter will discuss a simple derivation of (5.3) via the Gauss divergence theorem.

The Stein identity (5.1), or (5.3), yields a useful identity on the chi-square distribution. Let  $\mathbf{x} = (x_i) \sim \mathcal{N}_n(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$  and let  $\nabla_{\mathbf{x}} = (\partial/\partial x_i)$  be the  $n$ -dimensional differential operator vector. For a differentiable function  $g(s)$  with  $s = \|\mathbf{x}\|^2$ , using the Stein identity leads to

$$E\left[\frac{g(s)}{\sigma^2}\right] = E\left[\mathbf{x}^\top \mathbf{x} \frac{g(s)}{\sigma^2 s}\right] = E\left[\nabla_{\mathbf{x}}^\top \mathbf{x} \frac{g(s)}{s}\right] = E\left[(n-2) \frac{g(s)}{s} + 2g'(s)\right], \quad (5.4)$$

where  $g'(s) = dg(s)/ds$ . Since  $s \sim \sigma^2 \chi_n^2$ , the identity (5.4) is named the chi-square identity, which was derived by Efron and Morris (1976).

The chi-square identity (5.4) can be extended to an identity on the nonsingular Wishart distribution. Let  $\mathbf{S} = (s_{ij}) \sim \mathcal{W}_p(n, \mathbf{\Sigma})$  and let  $D_{\mathbf{S}}$  be the  $p \times p$  matrix differential operator whose  $(i, j)$ -th element is  $(1/2)(1 + \delta_{ij})(\partial/\partial s_{ij})$ , where  $\delta_{ij}$  represents the Kronecker delta, namely,  $\delta_{ij} = 1$  if  $i = j$  and  $\delta_{ij} = 0$  otherwise. Under suitable conditions, we can obtain an identity

$$E[\operatorname{tr} \mathbf{\Sigma}^{-1} \mathbf{G}] = E[(n - p - 1) \operatorname{tr} \mathbf{S}^{-1} \mathbf{G} + 2 \operatorname{tr} D_{\mathbf{S}} \mathbf{G}], \quad (5.5)$$

where all the elements of  $\mathbf{G} (\in \mathbb{S}_p)$  are absolutely continuous functions of  $\mathbf{S}$ . In the literature, the identity (5.5) is called the Haff (1977, 1979) identity, and it is a useful tool for estimation of  $\mathbf{\Sigma}$ . Clearly, when  $p = 1$ , the Haff identity (5.5) is equivalent to the chi-square identity (5.4). See the appendix of this chapter for a brief derivation of the Haff identity (5.5).

## 5.2 Some Useful Results on Matrix Differential Operators

Let  $\mathbf{Y} = (y_{ab}) \in \mathbb{R}^{n \times p}$ . Denote the  $n \times p$  matrix differential operator with respect to  $\mathbf{Y}$  by  $\nabla_{\mathbf{Y}} = (d_{ab}^{\mathbf{Y}})$  with  $d_{ab}^{\mathbf{Y}} = \partial/\partial y_{ab}$ . Here, we provide calculus formulas for  $\mathbf{S} = (s_{ij}) = \mathbf{Y}^\top \mathbf{Y} \in \mathbb{S}_{p,v}^{(+)}$  with  $v = n \wedge p$  and its Moore-Penrose inverse  $\mathbf{S}^+ = (s_{ij}^+)$ .

**Lemma 5.1** *Abbreviate  $d_{ab}^{\mathbf{Y}}$  to  $d$ . Denote  $d\mathbf{S} = (ds_{ij})$  and  $d\mathbf{S}^+ = (ds_{ij}^+)$ . Then*

$$d\mathbf{S}^+ = -\mathbf{S}^+ [d\mathbf{S}] \mathbf{S}^+ + (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+) [d\mathbf{S}] \mathbf{S}^+ \mathbf{S}^+ + \mathbf{S}^+ \mathbf{S}^+ [d\mathbf{S}] (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+).$$

**Proof** Note that  $S^+ = S^+ S S^+$ ,  $S S^+ = S^+ S$  and  $S = S S S^+$ . Differentiating both sides of  $S^+ = S^+ \times S \times S^+$ , we have  $dS^+ = [dS^+] S S^+ + S^+ [dS] S^+ + S S^+ dS^+$ , so that

$$[dS^+] S S^+ = -S^+ [dS] S^+ + (I_p - S S^+) dS^+.$$

Thus

$$\begin{aligned} dS^+ &= [dS^+] \{S S^+ + (I_p - S S^+)\} \\ &= [dS^+] S S^+ + [dS^+] (I_p - S S^+) \\ &= -S^+ [dS] S^+ + (I_p - S S^+) dS^+ + \{(I_p - S S^+) dS^+\}^\top. \end{aligned} \quad (5.6)$$

Differentiating both sides of  $S = S S^+ \times S$  yields  $dS = [d(S S^+)] S + S S^+ dS$ , which implies that  $[d(S S^+)] S = (I_p - S S^+) dS$ , which further implies that

$$[d(S S^+)] S^+ = (I_p - S S^+) [dS] S^+ S^+. \quad (5.7)$$

Differentiating both sides of  $S^+ = S S^+ \times S^+$ , we obtain  $dS^+ = [d(S S^+)] S^+ + S S^+ dS^+$ , namely,

$$(I_p - S S^+) dS^+ = [d(S S^+)] S^+ = (I_p - S S^+) [dS] S^+ S^+, \quad (5.8)$$

where the second equality follows from (5.7). Substituting (5.8) into (5.6) completes the proof.  $\square$

**Lemma 5.2** Denote the Kronecker delta by  $\delta_{ij}$ , namely,  $\delta_{ij} = 1$  for  $i = j$  and  $\delta_{ij} = 0$  for  $i \neq j$ . For  $a, i \in \{1, \dots, n\}$  and  $b, j, k \in \{1, \dots, p\}$ , we have

- (i)  $d_{ab}^Y y_{ij} = \delta_{ai} \delta_{bj}$ ,
- (ii)  $d_{ab}^Y s_{jk} = \delta_{bj} y_{ak} + \delta_{bk} y_{aj}$ ,
- (iii)  $d_{ab}^Y s_{jk}^+ = -\{Y S^+\}_{ak} s_{bj}^+ - \{Y S^+\}_{aj} s_{bk}^+ + \{Y S^+ S^+\}_{ak} \{I_p - S S^+\}_{bj} + \{Y S^+ S^+\}_{aj} \{I_p - S S^+\}_{bk}$ ,
- (iv)  $d_{ab}^Y \{Y S^+\}_{ik} = \{I_n - Y S^+ Y^\top\}_{ai} s_{bk}^+ + \{Y S^+ S^+ Y^\top\}_{ai} \{I_p - S S^+\}_{bk} - \{Y S^+\}_{ak} \{Y S^+\}_{ib}$ .

**Proof** Obviously, (i) holds. Differentiating  $s_{jk} = \sum_{c=1}^n y_{cj} y_{ck}$  with respect to  $y_{ab}$  yields

$$\begin{aligned} d_{ab}^Y s_{jk} &= \sum_{c=1}^n (y_{ck} d_{ab}^Y y_{cj} + y_{cj} d_{ab}^Y y_{ck}) \\ &= \sum_{c=1}^n (y_{ck} \delta_{ac} \delta_{bj} + y_{cj} \delta_{ac} \delta_{bk}) = \delta_{bj} y_{ak} + y_{aj} \delta_{bk}, \end{aligned}$$

which shows (ii).

From Lemma 5.1 and (ii), it is observed that

$$\begin{aligned}
 \mathbf{d}_{ab}^Y \mathbf{s}_{jk}^+ &= \{\mathbf{d}_{ab}^Y \mathbf{S}^+\}_{jk} \\
 &= \sum_{c=1}^p \sum_{d=1}^p \left[ -s_{jc}^+ [\mathbf{d}_{ab}^Y s_{cd}] s_{dk}^+ + \{\mathbf{I}_p - \mathbf{S} \mathbf{S}^+\}_{jc} [\mathbf{d}_{ab}^Y s_{cd}] \{\mathbf{S}^+ \mathbf{S}^+\}_{dk} \right. \\
 &\quad \left. + \{\mathbf{S}^+ \mathbf{S}^+\}_{jc} [\mathbf{d}_{ab}^Y s_{cd}] \{\mathbf{I}_p - \mathbf{S} \mathbf{S}^+\}_{dk} \right] \\
 &= -s_{bj}^+ \{\mathbf{Y} \mathbf{S}^+\}_{ak} - \{\mathbf{Y} \mathbf{S}^+\}_{aj} s_{bk}^+ \\
 &\quad + \{\mathbf{I}_p - \mathbf{S} \mathbf{S}^+\}_{bj} \{\mathbf{Y} \mathbf{S}^+ \mathbf{S}^+\}_{ak} + \{\mathbf{Y} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{aj} \{\mathbf{S}^+ \mathbf{S}^+\}_{bk} \\
 &\quad + \{\mathbf{S}^+ \mathbf{S}^+\}_{bj} \{\mathbf{Y} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{ak} + \{\mathbf{Y} \mathbf{S}^+ \mathbf{S}^+\}_{aj} \{\mathbf{I}_p - \mathbf{S} \mathbf{S}^+\}_{bk}.
 \end{aligned}$$

Noting that  $\mathbf{Y} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+) = \mathbf{0}_{n \times p}$  gives (iii).

In view of the product rule,

$$\mathbf{d}_{ab}^Y \{\mathbf{Y} \mathbf{S}^+\}_{ik} = \sum_{c=1}^p \left[ [\mathbf{d}_{ab}^Y y_{ic}] s_{ck}^+ + y_{ic} \mathbf{d}_{ab}^Y s_{ck}^+ \right].$$

Using (i) and (iii) and subsequently summing over all  $c$ , we obtain (iv).  $\square$

The following lemma is useful in estimation of a covariance matrix in multivariate normal distribution model. The proof of the lemma is similar to that of the  $p > n$  case given in Konno (2009). See also Stein (1977), Sheena (1995) and Kubokawa and Srivastava (2008).

**Lemma 5.3** *Denote by  $\mathbf{S} = \mathbf{H} \mathbf{L} \mathbf{H}^\top$  the eigenvalue decomposition of  $\mathbf{S} = \mathbf{Y}^\top \mathbf{Y}$ , where  $\mathbf{L} = \text{diag}(\ell_1, \dots, \ell_\nu) \in \mathbb{D}_\nu^{(\geq 0)}$  and  $\mathbf{H} \in \mathbb{V}_{p,\nu}$  with  $\nu = n \wedge p$ . Let  $\Phi = \text{diag}(\phi_1, \dots, \phi_\nu) \in \mathbb{D}_\nu$  such that all the diagonal elements  $\phi_i$ 's are differentiable functions of  $\mathbf{L}$ . Then*

$$\nabla_Y^\top \mathbf{Y} \mathbf{S}^+ \mathbf{H} \Phi \mathbf{H}^\top = \mathbf{H} \Phi^* \mathbf{H}^\top + (\text{tr } \mathbf{L}^{-1} \Phi) (\mathbf{I}_p - \mathbf{H} \mathbf{H}^\top),$$

where  $\Phi^* = \text{diag}(\phi_1^*, \dots, \phi_\nu^*)$  and for  $i = 1, \dots, \nu$

$$\phi_i^* = (n - \nu - 1) \frac{\phi_i}{\ell_i} + 2 \frac{\partial \phi_i}{\partial \ell_i} + \sum_{j \neq i}^{\nu} \frac{\phi_i - \phi_j}{\ell_i - \ell_j}.$$

In particular,

$$\text{tr } \nabla_Y^\top \mathbf{Y} \mathbf{S}^+ \mathbf{H} \Phi \mathbf{H}^\top = \sum_{i=1}^{\nu} \left\{ (|n - p| - 1) \frac{\phi_i}{\ell_i} + 2 \frac{\partial \phi_i}{\partial \ell_i} + 2 \sum_{j>i}^{\nu} \frac{\phi_i - \phi_j}{\ell_i - \ell_j} \right\}.$$

## Appendix

In this appendix, we first give another derivation of the Stein identity (5.3). Denote by  $f(\mathbf{X})$  the p.d.f. of  $\mathcal{N}_{m \times p}(\boldsymbol{\Theta}, \mathbf{I}_m \otimes \boldsymbol{\Sigma})$ . Let

$$I_{ST}(\mathbf{G}) = \int_{\mathbb{R}^{m \times p}} \text{tr} \nabla_{\mathbf{X}} \{\mathbf{G}^{\top} f(\mathbf{X})\} (d\mathbf{X}).$$

It follows that  $\nabla_{\mathbf{X}} f(\mathbf{X}) = -(\mathbf{X} - \boldsymbol{\Theta})\boldsymbol{\Sigma}^{-1} f(\mathbf{X})$ , so that

$$I_{ST}(\mathbf{G}) = E[\text{tr} \nabla_{\mathbf{X}} \mathbf{G}^{\top}] - E[\text{tr} (\mathbf{X} - \boldsymbol{\Theta})\boldsymbol{\Sigma}^{-1} \mathbf{G}^{\top}]$$

provided the expectations exist. Hence the Stein identity (5.3) can be verified if  $I_{ST}(\mathbf{G}) = 0$ .

For  $r > 0$ , let  $\mathbb{B}_r = \{\mathbf{X} \in \mathbb{R}^{m \times p} : \|\text{vec}(\mathbf{X})\| \leq r\}$ , where  $\|\cdot\|$  is the usual Euclidean norm and  $\text{vec}(\cdot)$  is defined in Definition 2.3. Then  $\mathbb{B}_r \rightarrow \mathbb{R}^{m \times p}$  as  $r \rightarrow \infty$  and

$$I_{ST}(\mathbf{G}) = \lim_{r \rightarrow \infty} \int_{\mathbb{B}_r} \text{vec}(\nabla_{\mathbf{X}})^{\top} \text{vec}(\mathbf{G} f(\mathbf{X})) (d\mathbf{X}).$$

The boundary of  $\mathbb{B}_r$  is expressed by  $\partial \mathbb{B}_r = \{\text{vec}(\mathbf{X}) \in \mathbb{R}^{mp} : \|\text{vec}(\mathbf{X})\| = r\}$ . Denote by  $\mathbf{u}$  an outward unit normal vector at a point  $\text{vec}(\mathbf{X}) \in \partial \mathbb{B}_r$ . Let  $\lambda_{\partial \mathbb{B}_r}$  be Lebesgue measure on  $\partial \mathbb{B}_r$ . By the Gauss divergence theorem,

$$I_{ST}(\mathbf{G}) = \lim_{r \rightarrow \infty} \int_{\partial \mathbb{B}_r} \mathbf{u}^{\top} \text{vec}(\mathbf{G}) f(\mathbf{X}) (d\lambda_{\partial \mathbb{B}_r}).$$

For details of the Gauss divergence theorem, see Fleming (1977).

Let  $o(\cdot)$  be the Landau symbol, namely, for real-valued functions  $f(x)$  and  $g(x)$  with  $g(x) \neq 0$ , we write  $f(x) = o(g(x))$  when  $\lim_{x \rightarrow c} |f(x)/g(x)| = 0$  for an extended real number  $c$ . If

$$\sup_{\text{vec}(\mathbf{X}) \in \partial \mathbb{B}_r} \|\text{vec}(\mathbf{G})\| f(\mathbf{X}) = o(r^{1-mp}) \quad \text{as } r \rightarrow \infty,$$

then  $I_{ST}(\mathbf{G}) = 0$ . In fact,

$$\begin{aligned} \int_{\partial \mathbb{B}_r} |\mathbf{u}^{\top} \text{vec}(\mathbf{G})| f(\mathbf{X}) (d\lambda_{\partial \mathbb{B}_r}) &\leq \int_{\partial \mathbb{B}_r} \|\text{vec}(\mathbf{G})\| f(\mathbf{X}) (d\lambda_{\partial \mathbb{B}_r}) \\ &\leq o(r^{1-mp}) \int_{\partial \mathbb{B}_r} (d\lambda_{\partial \mathbb{B}_r}) = o(1), \end{aligned}$$

because  $\int_{\partial \mathbb{B}_r} (d\lambda_{\partial \mathbb{B}_r})$  is the surface area of the  $(mp - 1)$ -sphere of radius  $r$  in  $\mathbb{R}^{mp}$ , namely,  $\int_{\partial \mathbb{B}_r} (d\lambda_{\partial \mathbb{B}_r}) \approx r^{mp-1}$ .



Next, a simple derivation of the Haff identity (5.5) is provided by using the Gauss divergence theorem. The derivation is essentially the same as Haff (1977, 1979). Let  $f(\mathbf{S})$  be the p.d.f. of  $\mathcal{W}_p(n, \mathbf{\Sigma})$ . For a differentiable matrix-valued function  $\mathbf{G} \in \mathbb{S}_p$ , let

$$I_{HF}(\mathbf{G}) = \int_{\mathbb{S}_p^{(+)}} \text{tr } D_S \{ \mathbf{G} f(\mathbf{S}) \} (d\mathbf{S})$$

Since  $D_S |\mathbf{S}| = \mathbf{S}^{-1} |\mathbf{S}|$  and  $D_S \text{tr } \mathbf{\Sigma}^{-1} \mathbf{S} = \mathbf{\Sigma}^{-1}$ , we get  $D_S f(\mathbf{S}) = \{(n - p - 1) \mathbf{S}^{-1} - \mathbf{\Sigma}^{-1}\} f(\mathbf{S})/2$ , implying that

$$I_{HF}(\mathbf{G}) = E[\text{tr } D_S \mathbf{G}] + \frac{n - p - 1}{2} E[\text{tr } \mathbf{S}^{-1} \mathbf{G}] - \frac{1}{2} E[\text{tr } \mathbf{\Sigma}^{-1} \mathbf{G}]$$

provided the expectations exist. Hence the Haff identity (5.5) follows if  $I_{HF}(\mathbf{G}) = 0$ .

Denote by  $\partial/\partial \mathbf{S} = (\partial/\partial s_{ij})$  the  $p \times p$  matrix differential operator with respect to  $\mathbf{S} \in \mathbb{S}_p$ . For  $\mathbf{A} = (a_{ij}) \in \mathbb{S}_p$ , define

$$\text{Vec}(\mathbf{A}) = (a_{11}, a_{21}, \dots, a_{p1}, a_{22}, a_{32}, \dots, a_{p2}, \dots, a_{p-1,p-1}, a_{p,p-1}, a_{pp})^\top \in \mathbb{R}^q,$$

where  $q = p(p + 1)/2$ . From symmetry of  $\mathbf{G}$ , it holds that

$$\begin{aligned} \text{tr } D_S \{ \mathbf{G} f(\mathbf{S}) \} &= \sum_{i=1}^p \sum_{j=1}^p \frac{1 + \delta_{ij}}{2} \frac{\partial}{\partial s_{ij}} \{ g_{ji} f(\mathbf{S}) \} = \sum_{i=1}^p \sum_{j=1}^i \frac{\partial}{\partial s_{ij}} \{ g_{ij} f(\mathbf{S}) \} \\ &= \text{Vec}(\partial/\partial \mathbf{S})^\top \text{Vec}(\mathbf{G} f(\mathbf{S})), \end{aligned}$$

so that

$$I_{HF}(\mathbf{G}) = \int_{\mathbb{S}_p^{(+)}} \text{Vec}(\partial/\partial \mathbf{S})^\top \text{Vec}(\mathbf{G} f(\mathbf{S})) (d\mathbf{S}).$$

For  $r > 0$ , let  $\partial \mathbb{B}_r^q = \{\text{Vec}(\mathbf{S}) \in \mathbb{R}^q : \|\text{Vec}(\mathbf{S})\| = r\}$  and, for  $0 < r_1 \leq r_2 < \infty$ , let  $\mathbb{C}_{r_1, r_2} = \{\text{Vec}(\mathbf{S}) \in \mathbb{R}^q : r_1 \leq \|\text{Vec}(\mathbf{S})\| \leq r_2\}$ . Then  $\mathbb{C}_{r_1, r_2} \cap \mathbb{S}_p^{(+)} \rightarrow \mathbb{S}_p^{(+)}$  as  $r_1 \rightarrow 0$  and  $r_2 \rightarrow \infty$ . The boundary of  $\mathbb{C}_{r_1, r_2} \cap \mathbb{S}_p^{(+)}$  can be expressed as  $\bigcup_{i=1}^3 \partial \mathbb{B}_i$ , where  $\partial \mathbb{B}_1$ ,  $\partial \mathbb{B}_2$  and  $\partial \mathbb{B}_3$  are certain sets satisfying  $\partial \mathbb{B}_1 \subset \partial \mathbb{B}_{r_1}^q$ ,  $\partial \mathbb{B}_2 \subset \partial \mathbb{B}_{r_2}^q$  and  $\partial \mathbb{B}_3 \subset \partial \mathbb{S}_p^{(+)}$ . Note that, for any point  $\mathbf{S} \in \partial \mathbb{S}_p^{(+)}$ ,  $|\mathbf{S}| = 0$ , namely,  $f(\mathbf{S}) = 0$  when  $n - p - 1 > 0$ . Let  $\mathbf{u}_1 = -\text{Vec}(\mathbf{S})/\|\text{Vec}(\mathbf{S})\|$  for  $\text{Vec}(\mathbf{S}) \in \partial \mathbb{B}_{r_1}^q$  and  $\mathbf{u}_2 = \text{Vec}(\mathbf{S})/\|\text{Vec}(\mathbf{S})\|$  for  $\text{Vec}(\mathbf{S}) \in \partial \mathbb{B}_{r_2}^q$ . Denote by  $\lambda_{\partial \mathbb{B}_i^q}$  Lebesgue measure on  $\partial \mathbb{B}_i^q$ . Using the Gauss divergence theorem gives

$$I_{HF}(\mathbf{G}) = \lim_{r_1 \rightarrow 0} \int_{\partial \mathbb{B}_1} \mathbf{u}_1^\top \text{Vec}(\mathbf{G}) f(\mathbf{S}) (d\lambda_{\partial \mathbb{B}_{r_1}^q}) + \lim_{r_2 \rightarrow \infty} \int_{\partial \mathbb{B}_2} \mathbf{u}_2^\top \text{Vec}(\mathbf{G}) f(\mathbf{S}) (d\lambda_{\partial \mathbb{B}_{r_2}^q}).$$

Using the Landau symbol  $o(\cdot)$ , we assume that

$$\sup_{\text{Vec}(\mathbf{S}) \in \partial \mathbb{B}_1} \|\text{Vec}(\mathbf{G})\| f(\mathbf{S}) = o(r_1^{1-q}) \quad \text{as } r_1 \rightarrow 0$$

and

$$\sup_{\text{Vec}(\mathbf{S}) \in \partial \mathbb{B}_2} \|\text{Vec}(\mathbf{G})\| f(\mathbf{S}) = o(r_2^{1-q}) \quad \text{as } r_2 \rightarrow \infty.$$

Under these assumptions, we can see that

$$\begin{aligned} \int_{\partial \mathbb{B}_1} |\mathbf{u}_1^\top \text{Vec}(\mathbf{G})| f(\mathbf{S}) (d\lambda_{\partial \mathbb{B}_1^q}) &\leq \int_{\partial \mathbb{B}_1} \|\text{Vec}(\mathbf{G})\| f(\mathbf{S}) (d\lambda_{\partial \mathbb{B}_1^q}) \\ &\leq o(r_1^{1-q}) \int_{\partial \mathbb{B}_1^q} (d\lambda_{\partial \mathbb{B}_1^q}) = o(1) \quad \text{as } r_1 \rightarrow 0 \end{aligned}$$

and also

$$\int_{\partial \mathbb{B}_2} |\mathbf{u}_2^\top \text{Vec}(\mathbf{G})| f(\mathbf{S}) (d\lambda_{\partial \mathbb{B}_2^q}) \leq o(r_2^{1-q}) \int_{\partial \mathbb{B}_2^q} (d\lambda_{\partial \mathbb{B}_2^q}) = o(1) \quad \text{as } r_2 \rightarrow \infty,$$

so that  $I_{HF}(\mathbf{G}) = 0$ .

## References

- M. Bilodeau, T. Kariya, Minimax estimators in the normal MANOVA model. *J. Multivar. Anal.* **28**, 260–270 (1989)
- B. Efron, C. Morris, Families of minimax estimators of the mean of a multivariate normal distribution. *Ann. Stat.* **4**, 11–21 (1976)
- W. Fleming, *Functions of Several Variables*, 2nd edn. (Springer, New York, 1977)
- L.R. Haff, Minimax estimators for a multinormal precision matrix. *J. Multivar. Anal.* **7**, 374–385 (1977)
- L.R. Haff, An identity for the Wishart distribution with applications. *J. Multivar. Anal.* **9**, 531–544 (1979)
- Y. Konno, Improved estimation of matrix of normal mean and eigenvalues in the multivariate  $F$ -distribution. Doctoral dissertation, Institute of Mathematics, University of Tsukuba, 1992. <http://mcm-www.jwu.ac.jp/~konno/>
- Y. Konno, Shrinkage estimators for large covariance matrices in multivariate real and complex normal distributions under an invariant quadratic loss. *J. Multivar. Anal.* **100**, 2237–2253 (2009)
- T. Kubokawa, M.S. Srivastava, Estimation of the precision matrix of a singular Wishart distribution and its application in high-dimensional data. *J. Multivar. Anal.* **99**, 1906–1928 (2008)
- Y. Sheena, Unbiased estimator of risk for an orthogonally invariant estimator of a covariance matrix. *J. Jpn. Stat. Soc.* **25**, 35–48 (1995)
- C. Stein, Estimation of the mean of a multivariate normal distribution. Technical Reports No.48 (Department of Statistics, Stanford University, Stanford, 1973)

- C. Stein, Lectures on the theory of estimation of many parameters, in *Proceedings of Scientific Seminars of the Steklov Institute Studies in the Statistical Theory of Estimation, Part I*, vol. 74, ed. by I.A. Ibragimov, M.S. Nikulin (Leningrad Division, 1977), pp. 4–65
- C. Stein, Estimation of the mean of a multivariate normal distribution. *Ann. Stat.* **9**, 1135–1151 (1981)

# Chapter 6

## Estimation of the Mean Matrix



This chapter introduces a unified approach to high- and low-dimensional cases for matricial shrinkage estimation of a normal mean matrix with unknown covariance matrix. A historical background is briefly explained and matricial shrinkage estimators are motivated from an empirical Bayes method. An unbiased risk estimate is unifiedly developed for a class of estimators corresponding to all possible orderings of sample size and dimensions. Specific examples of matricial shrinkage estimators are provided and also some related topics are discussed.

### 6.1 Introduction

Matricial shrinkage estimation of a mean matrix of a matrix-variate normal distribution has been studied since Efron and Morris (1972, 1976) and Stein (1973). Assume now that an  $m \times p$  observed data matrix  $X$  follows  $\mathcal{N}_{m \times p}(\Xi, I_m \otimes I_p)$ , where  $p \geq m$  and  $\Xi$  is unknown, and consider estimation of the mean matrix  $\Xi$ . The performance of its estimator  $\hat{\Xi}$  is evaluated by a risk function relative to the squared Frobenius norm loss

$$L_F(\hat{\Xi}, \Xi) = \|\hat{\Xi} - \Xi\|_F^2 = \text{tr}(\hat{\Xi} - \Xi)(\hat{\Xi} - \Xi)^\top.$$

The maximum likelihood (ML) estimator of  $\Xi$  is  $\hat{\Xi}^{ML} = X$ . It is unbiased and minimax. Efron and Morris (1972) considered empirical Bayes estimation for  $\Xi$  and showed that  $\hat{\Xi}^{ML}$  is dominated by the resulting empirical Bayes estimator of the form

$$\hat{\Xi}^{EM} = \{I_m - c_0(XX^\top)^{-1}\}X, \quad c_0 = p - m - 1. \quad (6.1)$$

This is equivalent to  $\{\mathbf{I}_m - c_0(\mathbf{X}\mathbf{X}^\top)^{-1}\}\widehat{\mathbf{\Xi}}^{ML}$ , which is a matrix multiple of  $\widehat{\mathbf{\Xi}}^{ML}$ , while the James-Stein (1961) type shrinkage estimator can be defined by

$$\widehat{\mathbf{\Xi}}^{JS} = \left(1 - \frac{mp - 2}{\|\mathbf{X}\|_F^2}\right)\widehat{\mathbf{\Xi}}^{ML} = \left(1 - \frac{mp - 2}{\text{tr } \mathbf{X}\mathbf{X}^\top}\right)\widehat{\mathbf{\Xi}}^{ML},$$

which is a scalar multiple of  $\widehat{\mathbf{\Xi}}^{ML}$ . Therefore  $\widehat{\mathbf{\Xi}}^{EM}$  and  $\widehat{\mathbf{\Xi}}^{JS}$  are called, respectively, matricial and scalar shrinkage estimators.

The Efron-Morris estimator  $\widehat{\mathbf{\Xi}}^{EM}$  can be written as  $\widehat{\mathbf{\Xi}}^{EM} = (\mathbf{I}_m - c_0\mathbf{B}\mathbf{L}^{-1}\mathbf{B}^\top)\mathbf{X}$ , where  $\mathbf{X}\mathbf{X}^\top = \mathbf{B}\mathbf{L}\mathbf{B}^\top$  such that  $\mathbf{L} = \text{diag}(l_1, \dots, l_m) \in \mathbb{D}_m^{(\geq 0)}$  and  $\mathbf{B} \in \mathbb{O}_m$ . Efron and Morris (1972) also pointed out an interesting relationship between the mean matrix estimation and estimating a restricted precision matrix. The relationship suggests a positive-part rule for  $\widehat{\mathbf{\Xi}}^{EM}$  and in fact they showed that

$$\widehat{\mathbf{\Xi}}^{PEM} = (\mathbf{I}_m - \mathbf{B}\mathbf{C}\mathbf{L}^{-1}\mathbf{B}^\top)\mathbf{X}, \quad \mathbf{C} = \text{diag}(c_1, \dots, c_m), \quad c_i = \min(l_i, c_0), \quad (6.2)$$

dominates  $\widehat{\mathbf{\Xi}}^{EM}$  under the  $L_F$ -loss. Efron and Morris (1976) presented another improved estimator of the form

$$\widehat{\mathbf{\Xi}}^{MEM} = \widehat{\mathbf{\Xi}}^{EM} - \frac{(m-1)(m+2)}{\|\mathbf{X}\|_F^2}\mathbf{X}, \quad (6.3)$$

which is uniformly better than  $\widehat{\mathbf{\Xi}}^{EM}$  under the  $L_F$ -loss. Meanwhile, Stein (1973) considered a multivariate generalization of Baranchik's (1970) estimator for a mean vector of multivariate normal distribution. Stein's class of estimators is given by

$$\widehat{\mathbf{\Xi}}_\Phi = \{\mathbf{I}_m - \mathbf{B}\Phi(\mathbf{L})\mathbf{B}^\top\}\mathbf{X}, \quad \Phi(\mathbf{L}) = \text{diag}(\phi_1(\mathbf{L}), \dots, \phi_m(\mathbf{L})), \quad (6.4)$$

where the diagonals of  $\Phi(\mathbf{L})$  are functions of  $\mathbf{L}$ . Stein (1973) derived an unbiased risk estimate of  $\widehat{\mathbf{\Xi}}_\Phi$  to provide alternative estimators. For example, he proposed

$$\widehat{\mathbf{\Xi}}^{ST} = (\mathbf{I}_m - \mathbf{B}\mathbf{D}\mathbf{L}^{-1}\mathbf{B}^\top)\mathbf{X}, \quad \mathbf{D} = \text{diag}(d_1, \dots, d_m), \quad d_i = m + p - 2i - 1, \quad (6.5)$$

which dominates  $\widehat{\mathbf{\Xi}}^{EM}$  relative to the  $L_F$ -loss.

The purpose of this chapter is to extend these results above to the unknown covariance case: Assume that  $\mathbf{X}$  and  $\mathbf{Y}$  are mutually independent random matrices distributed as, respectively,

$$\mathbf{X} \sim \mathcal{N}_{m \times p}(\mathbf{\Theta}, \mathbf{I}_m \otimes \mathbf{\Sigma}), \quad \mathbf{Y} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{\Sigma}), \quad (6.6)$$

where  $\mathbf{\Theta} \in \mathbb{R}^{m \times p}$  and  $\mathbf{\Sigma} \in \mathbb{S}_p^{(+)}$  are unknown. The model (6.6) is a canonical form of the multivariate linear regression model (4.1). Denote by  $\widehat{\mathbf{\Theta}}$  an estimator based on

$X$  and  $S = Y^\top Y$ . The problem we consider in this chapter is estimation of the mean matrix  $\Theta$  relative to invariant quadratic loss

$$L(\hat{\Theta}, \Theta | \Sigma) = \text{tr}(\hat{\Theta} - \Theta)\Sigma^{-1}(\hat{\Theta} - \Theta)^\top. \quad (6.7)$$

The invariance of (6.7) follows under (4.4) and, more generally, under the group of transformations  $\hat{\Theta} \rightarrow O\hat{\Theta}U + A$ ,  $\Theta \rightarrow O\Theta U + A$  and  $\Sigma \rightarrow U^\top \Sigma U$  for any  $O \in \mathbb{O}_m$ ,  $U \in \mathbb{U}_p$  and  $A \in \mathbb{R}^{m \times p}$ . The performance of  $\hat{\Theta}$  is measured by the risk function  $R(\hat{\Theta}, \Theta) = E[L(\hat{\Theta}, \Theta | \Sigma)]$ , where  $E$  is expectation taken with respect to (6.6).

The ML estimator of  $\Theta$  in (6.6) is  $\hat{\Theta}^{ML} = X$ , which is a minimax estimator with the constant risk  $mp$ . Thus all estimators dominating  $\hat{\Theta}^{ML}$  are minimax relative to the quadratic loss (6.7). When  $n \geq p$  in (6.6) and then the Wishart matrix  $S$  is nonsingular with probability one, some studies of improving  $\hat{\Theta}^{ML}$  via matricial shrinkage estimation can be found in Bilodeau and Kariya (1989), Honda (1991), Konno (1990, 1991, 1992) and Tsukuma and Kubokawa (2007). Bilodeau and Kariya (1989) and Konno (1992) studied general classes of matricial shrinkage estimators which can be regarded as an extension of (6.4). In particular, Konno's (1992) class has the form

$$\hat{\Theta}^K = \begin{cases} X\{I_p - Q\Psi(F)Q^{-1}\} & \text{for } m > p, \\ \{I_m - R\Psi(F)R^\top\}X & \text{for } p \geq m, \end{cases} \quad (6.8)$$

where  $F \in \mathbb{D}_{m \wedge p}^{(\geq 0)}$ ,  $Q \in \mathbb{U}_p$  and  $R \in \mathbb{O}_m$  satisfy

$$\begin{cases} Q^\top S Q = I_p \text{ and } Q^\top X^\top X Q = F & \text{for } m > p, \\ X S^{-1} X^\top = R F R^\top & \text{for } p \geq m, \end{cases}$$

and  $\Psi(F) \in \mathbb{D}_{m \wedge p}$  whose diagonal elements are functions of  $F$ . Konno (1992) showed that an unbiased risk estimate of the class (6.8) is only a function of  $F$  for both cases of  $m > p$  and  $p \geq m$ . Numerical comparison of shrinkage estimators has been carried out by Tsukuma and Kubokawa (2007). The numerical results suggest that when true eigenvalues of  $\Sigma^{-1}\Theta^\top\Theta$  are scattered, matricial shrinkage estimators outperform a James-Stein (1961) type scalar one, which may motivate us to study matricial shrinkage estimation.

This chapter is based in part on Tsukuma and Kubokawa (2015). The structure of this chapter is as follows. In Sect. 6.2, via an empirical Bayes method, we begin by deriving a unified Efron-Morris estimator for any possible ordering on  $m$ ,  $p$  and  $n$ . Section 6.3 yields a unified class of matricial shrinkage estimators from the class (6.8) and studies some properties of the unified class. In Sect. 6.4 we derive an unbiased risk estimate of the unified class, and Sect. 6.5 gives specific examples of matricial shrinkage estimators. Section 6.6 discusses some related topics including a method of positive-part rule and an extension to the GMANOVA model. In the Appendix, we supplement some results on matrix differential operators.

## 6.2 The Unified Efron-Morris Type Estimators Including Singular Cases

### 6.2.1 Empirical Bayes Methods

First, when  $n \geq p$ , namely, when  $S = Y^\top Y$  is nonsingular, we will derive empirical Bayes estimators of  $\Theta$  in (6.6) individually in the cases of  $m > p$  and  $p \geq m$ .

In the case of  $m > p$ , the prior distribution of  $\Theta$  is assumed to be  $\mathcal{N}_{m \times p}(\mathbf{0}_{m \times p}, I_m \otimes A)$ , where  $A \in \mathbb{S}_p^{(+)}$  is unknown. Then the posterior distribution of  $\Theta$  and the marginal distribution of  $X$  are, respectively,

$$\begin{aligned}\Theta|X &\sim \mathcal{N}_{m \times p}(X(I_p - \Omega), I_m \otimes (\Sigma^{-1} + A^{-1})^{-1}), \\ X &\sim \mathcal{N}_{m \times p}(\mathbf{0}_{m \times p}, I_m \otimes (\Sigma + A)),\end{aligned}$$

where  $\Omega = (\Sigma + A)^{-1}\Sigma$ . Thus the posterior mean of  $\Theta$  is  $\hat{\Theta}^B = X(I_p - \Omega)$ . Since  $\Omega$  is unknown, it may be estimated from the marginal distributions of  $S$  and  $W = X^\top X$ , which are given by  $S \sim \mathcal{W}_p(n, \Sigma)$  and  $W \sim \mathcal{W}_p(m, \Sigma + A)$ , respectively. It is reasonable that, as an estimator of  $\Omega$ , we take  $\hat{\Omega} = cW^{-1}S$  for a suitable constant  $c$ . Substituting  $\hat{\Omega}$  for  $\Omega$  in  $\hat{\Theta}^B$  yields an Efron-Morris (1972) type empirical Bayes estimator of the form

$$\hat{\Theta}_1^{EMK} = X(I_p - \hat{\Omega}) = X\{I_p - c(X^\top X)^{-1}S\}.$$

Next, we treat the case of  $p \geq m$ . Assume that the prior distribution of  $\Theta$  is  $\mathcal{N}_{m \times p}(\mathbf{0}_{m \times p}, B \otimes \Sigma)$ , where  $B \in \mathbb{S}_m^{(+)}$  is unknown. Then the posterior distribution of  $\Theta$  and the marginal distribution of  $X$  are, respectively,

$$\begin{aligned}\Theta|X &\sim \mathcal{N}_{m \times p}((I_m - \Omega)X, (I_m - \Omega) \otimes \Sigma), \\ X &\sim \mathcal{N}_{m \times p}(\mathbf{0}_{m \times p}, \Omega^{-1} \otimes \Sigma),\end{aligned}$$

where  $\Omega = (I_m + B)^{-1}$ . The resulting posterior mean of  $\Theta$  becomes  $\hat{\Theta}^B = (I_m - \Omega)X$ . Since we need to estimate  $\Omega$ , it will be estimated from the marginal distributions of  $X$  and  $S$ . Then from Corollary 3.1,  $E[X\Sigma^{-1}X^\top] = p\Omega^{-1}$ , and we think of  $\hat{\Omega} = c(XS^{-1}X^\top)^{-1}$  as an estimator of  $\Omega$ , where  $c$  is a positive constant. Thus the resulting Efron-Morris (1972) type empirical Bayes estimator of  $\Theta$  for  $p \geq m$  can be expressed as

$$\hat{\Theta}_2^{EMK} = (I_m - \hat{\Omega})X = \{I_m - c(XS^{-1}X^\top)^{-1}\}X.$$

The Efron-Morris type estimators given here have been studied by Konno (1991, 1992). For other empirical Bayes approaches, see Tsukuma and Kubokawa (2007).

### 6.2.2 The Unified Efron-Morris Type Estimator

Let us here give a unified form of empirical Bayes estimators  $\widehat{\Theta}_1^{EMK}$  and  $\widehat{\Theta}_2^{EMK}$  with properties of the Moore-Penrose inverse. When  $m > p$  with  $n \geq p$ , using (i), (ii) and (iv) of Lemma 2.3 yields

$$(XS^{-1}X^\top)^+ = (X^\top)^+SX^+ = X(X^\top X)^{-1}S(X^\top X)^{-1}X^\top,$$

so that  $(XS^{-1}X^\top)^+X = X(X^\top X)^{-1}S$ . When  $p \geq m$  with  $n \geq p$ ,  $XS^{-1}X^\top$  is of full rank and its Moore-Penrose inverse becomes  $(XS^{-1}X^\top)^+ = (XS^{-1}X^\top)^{-1}$ . Hence for both cases of  $m > p$  and  $p \geq m$  with  $n \geq p$ , the Efron-Morris type empirical Bayes estimators  $\widehat{\Theta}_1^{EMK}$  and  $\widehat{\Theta}_2^{EMK}$  can be unified into  $\widehat{\Theta}^{EMK} = X - c(XS^{-1}X^\top)^+X$ , where  $c$  is a constant.

In the case of  $p > n$ , the rank of  $S$  is deficient and its inverse does not exist. Therefore, we replace  $S^{-1}$  with  $S^+$ . This leads to  $\widehat{\Theta}^{EMK} = X - c(XS^+X^\top)^+X$ .

On the other hand, in the case where  $m = 1$ , Ch  telat and Wells (2012) suggested the shrinkage estimator

$$\widehat{\Theta}^{CW} = X - \frac{c}{XS^+X^\top}XSS^+.$$

An important problem is how to extend  $\widehat{\Theta}^{CW}$  to the framework of estimation of the mean matrix  $\Theta$ . In particular, the Efron-Morris type estimator seems to take various variants which depend on possible orderings among  $m$ ,  $n$  and  $p$ . One of interesting results provided here is that we can develop a unified form for the Efron-Morris type estimators, given by

$$\widehat{\Theta}^{EMK} = X - c(XS^+X^\top)^+XSS^+, \quad (6.9)$$

for any set of  $(m, n, p)$ . Of course, the expression (6.9) includes  $\widehat{\Theta}_1^{EMK}$ ,  $\widehat{\Theta}_2^{EMK}$  and  $\widehat{\Theta}^{CW}$  as special cases.

The matrix  $XS^+X^\top$  is nonsingular for  $n \wedge p \geq m$ , while it is singular for  $m > n \wedge p$ . In fact,  $(XS^+X^\top)^+$  for  $n \wedge p \geq m$  can be rewritten as

$$(XS^+X^\top)^+ = \begin{cases} (XS^{-1}X^\top)^{-1} & \text{for } n \geq p \geq m, \\ (XS^+X^\top)^{-1} & \text{for } p > n \geq m, \end{cases}$$

and the corresponding Efron-Morris type estimators are provided. In the case of  $m > n \wedge p$ ,  $\widehat{\Theta}^{EMK}$  in (6.9) can be expressed as in the following lemma.

**Lemma 6.1** *In the case of  $m > n \wedge p$ , the Efron-Morris type estimators in (6.9) can be expressed as*

$$\widehat{\Theta}^{EMK} = \begin{cases} X - cX(X^\top X)^{-1}S & \text{for } n \geq m > p \text{ or for } m > n \geq p, \\ X - cX(SS^+X^\top XSS^+)^+S & \text{for } m > p > n \text{ or for } p \geq m > n. \end{cases} \quad (6.10)$$



**Proof** When  $n \geq m > p$  or when  $m > n \geq p$ , using (i) of Lemma 2.3 gives  $S^+ = S^{-1}$ . Further from Lemma 2.3,

$$(XS^+X^\top)^+ \cdot XSS^+ = (X^\top)^+ SX^+ \cdot X = X(X^\top X)^{-1}S,$$

which verifies the expression (6.10) when  $n \geq m > p$  or when  $m > n \geq p$ .

When  $m > p > n$  or when  $p \geq m > n$ , we denote the eigenvalue decomposition of  $S$  by  $S = H L H^\top$ , where  $H \in \mathbb{V}_{p,n}$  and  $L \in \mathbb{D}_n^{(\geq 0)}$ . From (i) and (iv) of Lemma 2.3,  $S^+ = H L^{-1} H^\top$ , so that  $SS^+ = H H^\top = S^+ S$ . Since  $XH \in \mathbb{R}^{m \times n}$  is of rank  $n$ , it is observed that

$$\begin{aligned} (XS^+X^\top)^+ XSS^+ &= (XHL^{-1}H^\top X^\top)^+ XHH^\top \\ &= (H^\top X^\top)^+ L(XH)^+ XHH^\top & (\because \text{(iv) of Lemma 2.3}) \\ &= (H^\top X^\top)^+ L H^\top & (\because \text{(iii) of Lemma 2.3}) \\ &= XH(H^\top X^\top XH)^{-1} L H^\top & (\because \text{(ii) of Lemma 2.3}) \\ &= XH(H^\top X^\top XH)^{-1} H^\top H L H^\top & (\because H^\top H = I_n) \\ &= XH(H^\top X^\top XH)^{-1} H^\top S & (\because S = H L H^\top). \end{aligned}$$

Again from (i), (ii) and (iv) of Lemma 2.3,

$$H(H^\top X^\top XH)^{-1} H^\top = (H H^\top X^\top X H H^\top)^+ = (SS^+ X^\top X SS^+)^+.$$

Hence the expression (6.10) is obtained for the case where  $m > p > n$  or  $p \geq m > n$ .  $\square$

### 6.3 A Unified Class of Matricial Shrinkage Estimators

To define the class (6.8), Konno (1992) separately considered two cases, where  $m > p$  and where  $p \geq m$ , under  $n \geq p$ . The arguments stated in the previous section suggest that we can construct a well-defined class of matricial shrinkage estimators for all possible orders on  $m, n$  and  $p$ .

Hereafter in this chapter, we denote

$$\tau = m \wedge n \wedge p.$$

Define the eigenvalue decomposition of  $S$  as  $S = H L H^\top$ , where  $H \in \mathbb{V}_{p,n \wedge p}$  and  $L \in \mathbb{D}_{n \wedge p}^{(\geq 0)}$ . Let  $L^{1/2} = \text{diag}(\sqrt{\ell_1}, \dots, \sqrt{\ell_{n \wedge p}})$  and  $L^{-1/2} = (L^{1/2})^{-1}$ . Denote the singular value decomposition of  $XHL^{-1/2}$  by

$$XHL^{-1/2} = R F^{1/2} V^\top,$$

where  $\mathbf{R} \in \mathbb{V}_{m,\tau}$ ,  $\mathbf{V} \in \mathbb{V}_{n \wedge p, \tau}$  and  $\mathbf{F}^{1/2} = \text{diag}(\sqrt{f_1}, \dots, \sqrt{f_\tau}) \in \mathbb{D}_\tau^{(\geq 0)}$ . It is clear that  $\mathbf{X}\mathbf{S}^+\mathbf{X}^\top = \mathbf{X}\mathbf{H}\mathbf{L}^{-1}\mathbf{H}^\top\mathbf{X}^\top = \mathbf{R}\mathbf{F}\mathbf{R}^\top$ , which is the eigenvalue decomposition of  $\mathbf{X}\mathbf{S}^+\mathbf{X}^\top$ . For any possible triplet  $(m, n, p)$ , a unified class of matricial shrinkage estimators is defined by

$$\widehat{\Theta}^{SH} = \widehat{\Theta}^{SH}(\mathbf{X}, \mathbf{S}) = \mathbf{X} - \mathbf{R}\Phi(\mathbf{F})\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+, \quad (6.11)$$

where  $\Phi(\mathbf{F}) = \text{diag}(\phi_1(\mathbf{F}), \dots, \phi_\tau(\mathbf{F})) \in \mathbb{D}_\tau$  and the  $\phi_i(\mathbf{F})$ 's are absolutely continuous functions of  $\mathbf{F}$ .

Let us here discuss invariance of the unified class (6.11) under the orthogonal transformations  $\mathbf{X} \rightarrow \mathbf{O}\mathbf{X}\mathbf{P}$ ,  $\mathbf{S} \rightarrow \mathbf{P}^\top\mathbf{S}\mathbf{P}$ ,  $\Theta \rightarrow \mathbf{O}\Theta\mathbf{P}$  and  $\Sigma \rightarrow \mathbf{P}^\top\Sigma\mathbf{P}$  for any  $\mathbf{O} \in \mathbb{O}_m$  and  $\mathbf{P} \in \mathbb{O}_p$ . Then for an estimator  $\widehat{\Theta} = \widehat{\Theta}(\mathbf{X}, \mathbf{S})$ , it seems natural to require  $\widehat{\Theta}(\mathbf{O}\mathbf{X}\mathbf{P}, \mathbf{P}^\top\mathbf{S}\mathbf{P}) = \mathbf{O}\widehat{\Theta}(\mathbf{X}, \mathbf{S})\mathbf{P}$ . Since the eigenvalue decomposition of  $\mathbf{P}^\top\mathbf{S}\mathbf{P}$  is  $\mathbf{P}^\top\mathbf{H}\mathbf{L}\mathbf{H}^\top\mathbf{P}$ , it turns out that, due to (i), (ii) and (iv) of Lemma 2.3,

$$\begin{aligned} (\mathbf{P}^\top\mathbf{S}\mathbf{P})^+ &= (\mathbf{H}^\top\mathbf{P})^+\mathbf{L}^{-1}(\mathbf{P}^\top\mathbf{H})^+ \\ &= \mathbf{P}^\top\mathbf{H}(\mathbf{H}^\top\mathbf{P}\mathbf{P}^\top\mathbf{H})^{-1}\mathbf{L}^{-1}(\mathbf{H}^\top\mathbf{P}\mathbf{P}^\top\mathbf{H})^{-1}\mathbf{H}^\top\mathbf{P} \\ &= \mathbf{P}^\top\mathbf{H}\mathbf{L}^{-1}\mathbf{H}^\top\mathbf{P} = \mathbf{P}^\top\mathbf{S}^+\mathbf{P}. \end{aligned}$$

Thus,  $\mathbf{O}\mathbf{X}\mathbf{P}(\mathbf{P}^\top\mathbf{S}\mathbf{P})^+(\mathbf{O}\mathbf{X}\mathbf{P})^\top = \mathbf{O}\mathbf{X}\mathbf{S}^+\mathbf{X}^\top\mathbf{O}^\top$ , whose eigenvalue decomposition is  $\mathbf{O}\mathbf{R}\mathbf{F}\mathbf{R}^\top\mathbf{O}^\top$ . This yields, for any  $\mathbf{O} \in \mathbb{O}_m$  and  $\mathbf{P} \in \mathbb{O}_p$ ,

$$\begin{aligned} \widehat{\Theta}^{SH}(\mathbf{O}\mathbf{X}\mathbf{P}, \mathbf{P}^\top\mathbf{S}\mathbf{P}) &= \mathbf{O}\mathbf{X}\mathbf{P} - \mathbf{O}\mathbf{R} \cdot \Phi(\mathbf{F}) \cdot \mathbf{R}^\top\mathbf{O}^\top \cdot \mathbf{O}\mathbf{X}\mathbf{P} \cdot \mathbf{P}^\top\mathbf{S}\mathbf{P} \cdot (\mathbf{P}^\top\mathbf{S}\mathbf{P})^+ \\ &= \mathbf{O}\widehat{\Theta}^{SH}(\mathbf{X}, \mathbf{S})\mathbf{P}, \end{aligned}$$

which shows invariance of  $\widehat{\Theta}^{SH}$ . Note that if  $\mathbf{P} \in \mathbb{U}_p$  and  $\mathbf{P} \notin \mathbb{O}_p$  with  $p > n$  then  $(\mathbf{P}^\top\mathbf{H})^+ \neq \mathbf{H}^\top\mathbf{P}$  and  $\widehat{\Theta}^{SH}$  does not retain invariance, namely, it is not invariant under the scale transformations (4.4).

The Efron-Morris type estimator (6.9) lies in the unified class (6.11). It indeed holds that, according to (i), (ii) and (iv) of Lemma 2.3,  $(\mathbf{X}\mathbf{S}^+\mathbf{X}^\top)^+ = (\mathbf{R}\mathbf{F}\mathbf{R}^\top)^+ = \mathbf{R}\mathbf{F}^{-1}\mathbf{R}^\top$ , yielding

$$\widehat{\Theta}^{EMK} = \mathbf{X} - \mathbf{R}\Phi^{EMK}(\mathbf{F})\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+, \quad \Phi^{EMK}(\mathbf{F}) = c\mathbf{F}^{-1}. \quad (6.12)$$

Next, we give some convenient representations for (6.11).

**Lemma 6.2** *The unified class in (6.11) can be rewritten by*

$$\widehat{\Theta}^{SH} = \mathbf{X}(\mathbf{I}_p - \mathbf{S}\mathbf{S}^+) + \mathbf{R}\{\mathbf{I}_\tau - \Phi(\mathbf{F})\}\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+.$$

*Proof* Since  $\mathbf{S} = \mathbf{H}\mathbf{L}\mathbf{H}^\top$  and  $\mathbf{X}\mathbf{H}\mathbf{L}^{-1/2} = \mathbf{R}\mathbf{F}^{1/2}\mathbf{V}^\top$  where  $\mathbf{R} \in \mathbb{V}_{m,\tau}$ , it is seen that

$$\begin{aligned}
(\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top)\mathbf{X}\mathbf{S}\mathbf{S}^+ &= (\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top)\mathbf{X}\mathbf{H}\mathbf{H}^\top \\
&= (\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top)\mathbf{R}\mathbf{F}^{1/2}\mathbf{V}^\top\mathbf{L}^{1/2}\mathbf{H}^\top = \mathbf{0}_{m \times p},
\end{aligned}$$

which is used to rewrite the class (6.11) as

$$\begin{aligned}
\hat{\Theta}^{SH} &= \mathbf{X} - \mathbf{X}\mathbf{S}\mathbf{S}^+ + \mathbf{X}\mathbf{S}\mathbf{S}^+ - \mathbf{R}\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+ + \mathbf{R}\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+ - \mathbf{R}\Phi(\mathbf{F})\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+ \\
&= \mathbf{X}(\mathbf{I}_p - \mathbf{S}\mathbf{S}^+) + \mathbf{R}\{\mathbf{I}_\tau - \Phi(\mathbf{F})\}\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+.
\end{aligned}$$

Hence the proof is complete.  $\square$

When  $p > n$ , because of (i), (ii) and (iv) in Lemma 2.3,

$$\mathbf{S}\mathbf{S}^+ = \mathbf{Y}^\top\mathbf{Y} \cdot (\mathbf{Y}^\top\mathbf{Y})^+ = \mathbf{Y}^\top\mathbf{Y} \cdot \mathbf{Y}^\top(\mathbf{Y}\mathbf{Y}^\top)^{-2}\mathbf{Y} = \mathbf{Y}^\top(\mathbf{Y}\mathbf{Y}^\top)^{-1}\mathbf{Y}$$

is the orthogonal projection matrix onto the subspace spanned by rows of  $\mathbf{Y}$ . The ML estimator is rewritten as  $\hat{\Theta}^{ML} = \mathbf{X}(\mathbf{I}_p - \mathbf{S}\mathbf{S}^+) + \mathbf{X}\mathbf{S}\mathbf{S}^+$ . Thus Lemma 6.2 implies that  $\hat{\Theta}^{SH}$  is shrinking with respect only to the projections of rows of  $\mathbf{X}$  onto the subspace spanned by rows of  $\mathbf{Y}$ .

Further, the unified class (6.11) can be rewritten as in the following lemma which is an extension for Konno's (1992) class (6.8) with  $m > p$ .

**Lemma 6.3** *Let  $\mathbf{Q} = \mathbf{H}\mathbf{L}^{-1/2}\mathbf{V} \in \mathbb{R}^{p \times \tau}$ . Then  $\mathbf{Q}^- = \mathbf{V}^\top\mathbf{L}^{1/2}\mathbf{H}^\top \in \mathbb{R}^{\tau \times p}$  is the generalized inverse of  $\mathbf{Q}$ . Further the unified class in (6.11) can be rewritten as*

$$\begin{aligned}
\hat{\Theta}^{SH} &= \mathbf{X}\{\mathbf{I}_p - \mathbf{Q}\Phi(\mathbf{F})\mathbf{Q}^-\} \\
&= \mathbf{X}(\mathbf{I}_p - \mathbf{S}\mathbf{S}^+) + \mathbf{X}\mathbf{S}\mathbf{S}^+\mathbf{Q}\{\mathbf{I}_\tau - \Phi(\mathbf{F})\}\mathbf{Q}^-.
\end{aligned} \tag{6.13}$$

**Proof** It is seen that

$$\mathbf{Q}\mathbf{Q}^-\mathbf{Q} = \mathbf{H}\mathbf{L}^{-1/2}\mathbf{V}\mathbf{V}^\top\mathbf{L}^{1/2}\mathbf{H}^\top\mathbf{H}\mathbf{L}^{-1/2}\mathbf{V} = \mathbf{H}\mathbf{L}^{-1/2}\mathbf{V}\mathbf{V}^\top\mathbf{V} = \mathbf{H}\mathbf{L}^{-1/2}\mathbf{V} = \mathbf{Q},$$

and consequently  $\mathbf{Q}^-$  is the generalized inverse of  $\mathbf{Q}$ . Since

$$\begin{aligned}
\mathbf{R} &= \mathbf{X}\mathbf{H}\mathbf{L}^{-1/2}\mathbf{V}\mathbf{F}^{-1/2} = \mathbf{X}\mathbf{Q}\mathbf{F}^{-1/2}, \\
\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+ &= \mathbf{R}^\top\mathbf{R}\mathbf{F}^{1/2}\mathbf{V}^\top\mathbf{L}^{1/2}\mathbf{H}^\top = \mathbf{F}^{1/2}\mathbf{Q}^-,
\end{aligned} \tag{6.14}$$

it follows that

$$\mathbf{R}\Phi(\mathbf{F})\mathbf{R}^\top\mathbf{X}\mathbf{S}\mathbf{S}^+ = \mathbf{X}\mathbf{Q}\mathbf{F}^{-1/2}\Phi(\mathbf{F})\mathbf{F}^{1/2}\mathbf{Q}^- = \mathbf{X}\mathbf{Q}\Phi(\mathbf{F})\mathbf{Q}^-.$$

Hence, we get the first equality in (6.13). The second equality in (6.13) can be obtained similarly by using Lemma 6.2 with the fact that  $\mathbf{Q} = \mathbf{S}\mathbf{S}^+\mathbf{Q}$ .  $\square$

As for  $\mathbf{Q}$  in Lemma 6.3, it is easy to check that  $\mathbf{Q} = \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \mathbf{F}^{-1/2}$  and  $\mathbf{Q}^- = \mathbf{F}^{-1/2} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+$ . Also,  $\mathbf{Q}^\top \mathbf{S} \mathbf{Q} = \mathbf{I}_\tau$  and  $\mathbf{Q}^\top \mathbf{X}^\top \mathbf{X} \mathbf{Q} = \mathbf{F}$ , while

$$(\mathbf{Q}^-)^\top \mathbf{Q}^- = \begin{cases} \mathbf{S} & \text{for } m > n \wedge p, \\ \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{S} \mathbf{S}^+ & \text{for } m \leq n \wedge p, \end{cases}$$

and

$$(\mathbf{Q}^-)^\top \mathbf{F} \mathbf{Q}^- = \begin{cases} \mathbf{X}^\top \mathbf{X} & \text{for } n \geq p \geq m, \\ \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top) (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)^+ \mathbf{X} \mathbf{S} \mathbf{S}^+ & \text{otherwise.} \end{cases}$$

Lemmas 6.2 and 6.3 suggest that  $\widehat{\Theta}^{SH}$  shrinks not only columns, but also rows of  $\mathbf{X} \mathbf{S} \mathbf{S}^+$  in terms of  $\widehat{\Theta}^{ML}$ .

If all the diagonals of  $\Phi(\mathbf{F})$  are positive and the matrix square root of  $\Phi(\mathbf{F})$  is denoted by  $\Phi^{1/2}$  then

$$\mathbf{R} \Phi \mathbf{F} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ = \mathbf{R} \Phi^{1/2} \mathbf{R}^\top \mathbf{R} \Phi^{1/2} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ = \mathbf{R} \Phi^{1/2} \mathbf{R}^\top \mathbf{X} \mathbf{Q} \Phi^{1/2} \mathbf{Q}^-.$$

Since  $\mathbf{Q} = \mathbf{S} \mathbf{S}^+ \mathbf{Q}$ , we get the following lemma as well.

**Lemma 6.4** *If all the diagonals of  $\Phi(\mathbf{F})$  are positive and the matrix square root of  $\Phi(\mathbf{F})$  is denoted by  $\Phi^{1/2}$  then the unified class (6.11) can be rewritten by*

$$\widehat{\Theta}^{SH} = \mathbf{X} - \mathbf{R} \Phi^{1/2} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ \mathbf{Q} \Phi^{1/2} \mathbf{Q}^-.$$

*If all the diagonals of  $\mathbf{I}_\tau - \Phi(\mathbf{F})$  are positive and the matrix square root of  $\mathbf{I}_\tau - \Phi(\mathbf{F})$  is denoted by  $(\mathbf{I}_\tau - \Phi)^{1/2}$  then the unified class (6.11) can be rewritten by*

$$\widehat{\Theta}^{SH} = \mathbf{X}(\mathbf{I}_p - \mathbf{S} \mathbf{S}^+) + \mathbf{R}(\mathbf{I}_\tau - \Phi)^{1/2} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ \mathbf{Q}(\mathbf{I}_\tau - \Phi)^{1/2} \mathbf{Q}^-.$$

As seen in the above lemmas,  $\widehat{\Theta}^{SH}$  takes several different forms. In the following, we will employ the different forms for different purposes.

## 6.4 Unbiased Risk Estimate

Abbreviate  $\Phi(\mathbf{F})$  to  $\Phi$ . The quadratic loss (6.7) of  $\widehat{\Theta}^{SH}$  in (6.11) is expanded to

$$\begin{aligned} L(\widehat{\Theta}^{SH}, \Theta | \Sigma) &= \text{tr}(\mathbf{X} - \Theta) \Sigma^{-1} (\mathbf{X} - \Theta)^\top - 2 \text{tr}(\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top \\ &\quad + \text{tr} \mathbf{R} \Phi \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ \Sigma^{-1} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top. \end{aligned}$$

Recalling that  $R(\mathbf{X}, \Theta) = mp$ , we obtain  $R(\widehat{\Theta}^{SH}, \Theta) = mp + E_2 - 2E_1$ , where

$$E_1 = E[\text{tr}(\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top], \quad E_2 = E[\text{tr} \Sigma^{-1} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi^2 \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+].$$

Here Theorem 5.1, or (5.3), is used to evaluate  $E_1$ . If the conditions in Theorem 5.1 are satisfied for  $\mathbf{G}_1 = \mathbf{I}_m$  and  $\mathbf{G}_2 = \boldsymbol{\Sigma}^{-1} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{R}^\top$ , then  $E_1$  can be expressed as

$$E_1 = E[\text{tr } \nabla_{\mathbf{X}} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{R}^\top].$$

Similarly, since  $\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ , Theorem 5.1 is used to rewrite  $E_2$  as

$$\begin{aligned} E_2 &= E[\text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{Y}^\top \mathbf{Y} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi}^2 \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+] \\ &= E[\text{tr } \nabla_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi}^2 \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+]. \end{aligned}$$

A sufficient condition for applying Theorem 5.1 to  $E_1$  and  $E_2$  is

$$E[\text{tr } \mathbf{S} \cdot \text{tr } \mathbf{F} \boldsymbol{\Phi}^2] < \infty. \quad (6.15)$$

For more details, see Tsukuma and Kubokawa (2015).

From the above observation, an unbiased risk estimate of  $\widehat{\boldsymbol{\Theta}}^{SH}$  is given by

$$\widehat{R}(\widehat{\boldsymbol{\Theta}}^{SH}) = mp + \text{tr } \nabla_{\mathbf{Y}}^\top \mathbf{Y} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi}^2 \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+ - 2 \text{tr } \nabla_{\mathbf{X}} \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{R}^\top.$$

Using Lemmas 6.6 and 6.7 in the Appendix yields

**Theorem 6.1** *Let  $\phi_i = \phi_i(\mathbf{F})$  for  $i = 1, \dots, \tau$ . Assume that (6.15) is satisfied. For any possible ordering on  $m, n$  and  $p$ , an unbiased risk estimate of  $\widehat{\boldsymbol{\Theta}}^{SH}$  is*

$$\begin{aligned} \widehat{R}(\widehat{\boldsymbol{\Theta}}^{SH}) &= mp + \sum_{i=1}^{\tau} \left\{ a f_i \phi_i^2 - 2b \phi_i - 4 f_i^2 \phi_i \frac{\partial \phi_i}{\partial f_i} - 4 f_i \frac{\partial \phi_i}{\partial f_i} \right. \\ &\quad \left. - 2 \sum_{j>i}^{\tau} \frac{f_i^2 \phi_i^2 - f_j^2 \phi_j^2}{f_i - f_j} - 4 \sum_{j>i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j} \right\}, \end{aligned}$$

where  $a = a_{m,p,n} = (|n - p| + 2m) \wedge (n + p) - 3$  and  $b = b_{m,p,n} = |n \wedge p - m| + 1$ .

The unbiased risk estimate  $\widehat{R}(\widehat{\boldsymbol{\Theta}}^{SH})$  in Theorem 6.1 depends on  $\mathbf{F} \in \mathbb{D}_{\tau}^{(\geq 0)}$ . Let  $\widehat{\boldsymbol{\Theta}}_0^{SH}$  and  $\widehat{\boldsymbol{\Theta}}_1^{SH}$  be estimators belonging to the unified class (6.11). If  $\widehat{R}(\widehat{\boldsymbol{\Theta}}_0^{SH}) \leq \widehat{R}(\widehat{\boldsymbol{\Theta}}_1^{SH})$  for any  $\mathbf{F} \in \mathbb{D}_{\tau}^{(\geq 0)}$  and fixed  $(m, n, p)$ , then  $\widehat{\boldsymbol{\Theta}}_0^{SH}$  dominates  $\widehat{\boldsymbol{\Theta}}_1^{SH}$  relative to the quadratic loss (6.7). For example, we have

**Corollary 6.1** *If  $\widehat{R}(\widehat{\boldsymbol{\Theta}}^{SH}) \leq \widehat{R}(\widehat{\boldsymbol{\Theta}}^{ML}) = mp$  for any  $\mathbf{F} \in \mathbb{D}_{\tau}^{(\geq 0)}$  and fixed  $(m, n, p)$  then  $\widehat{\boldsymbol{\Theta}}^{SH}$  is a minimax estimator dominating  $\widehat{\boldsymbol{\Theta}}^{ML}$  relative to the quadratic loss (6.7).*

## 6.5 Examples for Specific Estimators

### 6.5.1 The Unified Efron-Morris Type Estimator

The unified Efron-Morris type estimator  $\hat{\Theta}^{EMK}$  could be rewritten as in (6.12). To apply Theorem 6.1 to  $\hat{\Theta}^{EMK}$ , we put  $\phi_i = cf_i^{-1}$ . Then condition (6.15) is expressed by  $c^2 E[\text{tr } \mathbf{S} \text{tr } \mathbf{F}^{-1}] < \infty$ , which is satisfied if  $b \geq 3$  or, equivalently,  $|n \wedge p - m| \geq 2$  (see Lemma 6.8 of Tsukuma and Kubokawa 2015). The unbiased risk estimate of  $\hat{\Theta}^{EMK}$  is

$$\hat{R}(\hat{\Theta}^{EMK}) = mp + \{(a+4)c - 2(b-2)\}c \sum_{i=1}^{\tau} \frac{1}{f_i},$$

implying that if  $0 < c \leq 2(b-2)/(a+4)$  and  $b \geq 3$  then  $\hat{\Theta}^{EMK}$  dominates  $\hat{\Theta}^{ML}$  relative to the quadratic loss (6.7).

The unbiased risk estimate  $\hat{R}(\hat{\Theta}^{EMK})$  is a quadratic function of  $c$  and attains its minimum at

$$c^{EM} = \frac{b-2}{a+4} = \frac{|n \wedge p - m| - 1}{(|n - p| + 2m) \wedge (n + p) + 1}. \quad (6.16)$$

Define

$$\hat{\Theta}^{EM} = \mathbf{X} - c^{EM} \mathbf{R} \mathbf{F}^{-1} \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^+ = \mathbf{X} - c^{EM} (\mathbf{X} \mathbf{S}^+ \mathbf{X}^{\top})^+ \mathbf{X} \mathbf{S} \mathbf{S}^+. \quad (6.17)$$

This is an extension of Efron and Morris' (1972) original estimator in (6.1). The unbiased risk estimate of  $\hat{\Theta}^{EM}$  has the form

$$\hat{R}(\hat{\Theta}^{EM}) = mp - (b-2)c^{EM} \sum_{i=1}^{\tau} \frac{1}{f_i} = \hat{R}(\hat{\Theta}^{ML}) - (b-2)c^{EM} \text{tr } \mathbf{F}^{-1}. \quad (6.18)$$

When  $m, n$  and  $p$  are given, the corresponding specific values of  $a$  and  $b$  in  $c^{EM}$  are determined. Noting that  $(|n - p| + 2m) \wedge (n + p) = n + p$  for  $m > n \wedge p$ , we can obtain specific values of  $c^{EM}$ ,

$$c^{EM} = \begin{cases} (p - m - 1)/(n - p + 2m + 1) & \text{for } n \geq p \geq m, \\ (n - m - 1)/(p - n + 2m + 1) & \text{for } p > n \geq m, \\ (m - p - 1)/(n + p + 1) & \text{for } n \geq m > p \text{ and } m > n \geq p, \\ (m - n - 1)/(n + p + 1) & \text{for } m > p > n \text{ and } p \geq m > n. \end{cases}$$

The cases satisfying  $n \geq p$ , namely,  $n \geq p \geq m$ ,  $n \geq m > p$  and  $m > n \geq p$ , are provided by Konno (1992).

### 6.5.2 A Modified Stein-Type Estimator

A modified Stein-type estimator is defined by

$$\widehat{\Theta}^{mST} = \widehat{\Theta}^{ST} - \frac{d}{\text{tr } \mathbf{X} \mathbf{S}^+ \mathbf{X}^\top} \mathbf{R} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+,$$

where  $\widehat{\Theta}^{ST} = \mathbf{X} - \mathbf{R} \mathbf{C} \mathbf{F}^{-1} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+$  for  $\mathbf{C} = \text{diag}(c_1, \dots, c_\tau)$  with  $c_1 \geq \dots \geq c_\tau$ . This corresponds to the form

$$\phi_i = \frac{c_i}{f_i} + \frac{d}{\sum_{j=1}^\tau f_j} = \frac{c_i}{f_i} + \frac{d}{\text{tr } \mathbf{F}}.$$

Then, from Theorem 6.1, it follows that

$$\begin{aligned} \widehat{\Delta} &= \widehat{R}(\widehat{\Theta}^{mST}) - \widehat{R}(\widehat{\Theta}^{ML}) \\ &= \sum_{i=1}^\tau \frac{1}{f_i} (ac_i^2 - 2bc_i + 4c_i + 4c_i^2) \\ &\quad + \frac{1}{\text{tr } \mathbf{F}} \left\{ (a - 2\tau + 2)d^2 - 2\tau bd - 2\tau(\tau - 1)d + 4d + 2(a + 2)d \sum_{i=1}^\tau c_i \right\} \\ &\quad + 4d \frac{\text{tr } \mathbf{C} \mathbf{F}}{(\text{tr } \mathbf{F})^2} + 4d^2 \frac{\text{tr } \mathbf{F}^2}{(\text{tr } \mathbf{F})^3} - 2 \sum_{i=1}^\tau \sum_{j>i}^\tau \frac{(c_i - c_j)(c_i + c_j + 2)}{f_i - f_j} \\ &\quad - \frac{4d}{\text{tr } \mathbf{F}} \sum_{i=1}^\tau \sum_{j>i}^\tau \frac{c_i f_i - c_j f_j}{f_i - f_j}, \end{aligned} \tag{6.19}$$

because  $\sum_{i=1}^\tau \sum_{j>i}^\tau (f_i + f_j) = (\tau - 1) \text{tr } \mathbf{F}$ . The condition for obtaining (6.19), namely, for satisfying (6.15), is  $b \geq 3$ . It is noted that  $\text{tr } \mathbf{F}^2 \leq (\text{tr } \mathbf{F})^2$ ,  $\text{tr } \mathbf{C} \mathbf{F} / (\text{tr } \mathbf{F})^2 \leq c_1 / \text{tr } \mathbf{F}$ ,

$$\begin{aligned} \sum_{i=1}^\tau \sum_{j>i}^\tau \frac{(c_i - c_j)(c_i + c_j + 2)}{f_i - f_j} &\geq \sum_{i=1}^\tau \frac{1}{f_i} \sum_{j>i}^\tau (c_i - c_j)(c_i + c_j + 2), \\ \sum_{i=1}^\tau \sum_{j>i}^\tau \frac{c_i f_i - c_j f_j}{f_i - f_j} &\geq \sum_{i=1}^\tau (\tau - i) c_i. \end{aligned}$$

Thus,  $\widehat{\Delta} \leq \sum_{i=1}^\tau h_c(i)/f_i + h_d/\text{tr } \mathbf{F}$ , where

$$h_c(i) = (a + 4 - 2\tau + 2i)c_i^2 - 2(b - 2 + 2\tau - 2i)c_i + 2 \sum_{j>i}^{\tau} c_j(c_j + 2),$$

$$h_d = (a - 2\tau + 6)d^2 - 2 \left\{ b\tau + \tau(\tau - 1) - 2 - 2c_1 - (a + 2) \sum_{i=1}^{\tau} c_i + 2 \sum_{i=1}^{\tau} (\tau - i)c_i \right\} d.$$

For  $i = 1, \dots, \tau$ , put

$$c_i = \frac{b - 2 + 2\tau - 2i}{a + 4 - 2\tau + 2i}, \quad (6.20)$$

which satisfy  $c_1 \geq \dots \geq c_{\tau}$  and  $c_{\tau} = c^{EM}$  given in (6.16). Since  $h_c(i)$  is a quadratic function of  $c_i$ , and  $(a + 4 - 2\tau + 2i)c_i^2 - 2(b - 2 + 2\tau - 2i)c_i \leq (a + 4 - 2\tau + 2i)c_{i+1}^2 - 2(b - 2 + 2\tau - 2i)c_{i+1}$  for each  $i$ , it is observed that

$$\begin{aligned} h_c(i) &= (a + 4 - 2\tau + 2i)c_i^2 - 2(b - 2 + 2\tau - 2i)c_i + 2c_{i+1}(c_{i+1} + 2) + 2 \sum_{j>i+1}^{\tau} c_j(c_j + 2) \\ &\leq \{a + 4 - 2\tau + 2(i + 1)\}c_{i+1}^2 - 2\{b - 2 + 2\tau - 2(i + 1)\}c_{i+1} + 2 \sum_{j>i+1}^{\tau} c_j(c_j + 2) \\ &\leq \dots \leq (a + 2)c_{\tau-1}^2 - 2bc_{\tau-1} + 2c_{\tau}(c_{\tau} + 2) \leq (a + 4)c_{\tau}^2 - 2(b - 2)c_{\tau} \\ &= -(b - 2)c_{\tau} = -(b - 2)c^{EM}. \end{aligned}$$

It is also seen that  $2 \sum_{i=1}^{\tau} (\tau - i)c_i = -(b + \tau - 3)\tau + (a + 4) \sum_{i=1}^{\tau} c_i$ , which provides

$$h_d = (a - 2\tau + 6)d^2 - 4 \left( \tau - 1 + \sum_{i=2}^{\tau} c_i \right) d.$$

Hence, from (6.18),

$$\widehat{\Delta} \leq \widehat{R}(\widehat{\Theta}^{EM}) - \widehat{R}(\widehat{\Theta}^{ML}) + \left\{ (a - 2\tau + 6)d^2 - 4 \left( \tau - 1 + \sum_{i=2}^{\tau} c_i \right) d \right\} \frac{1}{\text{tr } \mathbf{F}}.$$

These observations imply that  $\widehat{\Theta}^{ML}$  and  $\widehat{\Theta}^{EM}$  are improved on by  $\widehat{\Theta}^{ST} = \mathbf{X} - \mathbf{RCF}^{-1}\mathbf{R}^{\top}\mathbf{XSS}^{+}$  with constants  $c_i$ 's given in (6.20). With these  $c_i$ 's,  $\widehat{\Theta}^{ST}$  is an extension of Stein's (1973) estimator in (6.5) and further improved on by the modified Stein-type estimator

$$\widehat{\Theta}^{mST} = \mathbf{X} - \mathbf{RCF}^{-1}\mathbf{R}^{\top}\mathbf{XSS}^{+} - \frac{d}{\text{tr } \mathbf{XS}^{+}\mathbf{X}^{\top}} \mathbf{RR}^{\top}\mathbf{XSS}^{+}$$

if  $d$  satisfies  $0 < d \leq 4 \left\{ \tau - 1 + \sum_{i=2}^{\tau} c_i \right\} / (a - 2\tau + 6)$ . This is a generalization of Tsukuma and Kubokawa (2007) for any possible ordering on  $m, n$  and  $p$ .



### 6.5.3 Modified Efron-Morris Type Estimator

Next, we extend the modified Efron-Morris (1976) estimator in (6.3) to the unknown covariance case. Let

$$\widehat{\Theta}^{mEM} = \widehat{\Theta}^{EM} - \frac{d}{\text{tr } XS^+X^\top} \mathbf{R} \mathbf{R}^\top X S S^+,$$

where  $\widehat{\Theta}^{EM}$  is given in (6.17). This corresponds to the form

$$\phi_i = \frac{c^{EM}}{f_i} + \frac{d}{\sum_{j=1}^{\tau} f_j} = \frac{c^{EM}}{f_i} + \frac{d}{\text{tr } \mathbf{F}}.$$

Letting  $c_i = c^{EM}$  in (6.19) for all  $i$  and using the fact that  $\text{tr } \mathbf{F}^2 \leq (\text{tr } \mathbf{F})^2$ , we get

$$\begin{aligned} & \widehat{R}(\widehat{\Theta}^{mEM}) - \widehat{R}(\widehat{\Theta}^{ML}) \\ & \leq -(b-2)c^{EM} \text{tr } \mathbf{F}^{-1} \\ & \quad + \left\{ (a-2\tau+2)d^2 - 2\tau bd - 2\tau(\tau-1)d + 4d + 2\tau(a+2)c^{EM}d \right\} \frac{1}{\text{tr } \mathbf{F}} \\ & \quad + \frac{4c^{EM}d}{\text{tr } \mathbf{F}} + \frac{4d^2}{\text{tr } \mathbf{F}} - \frac{4c^{EM}d}{\text{tr } \mathbf{F}} \sum_{i=1}^{\tau} (\tau-i). \end{aligned}$$

With some algebraic manipulation,

$$\widehat{R}(\widehat{\Theta}^{mEM}) - \widehat{R}(\widehat{\Theta}^{EM}) \leq \left[ (a-2\tau+6)d^2 - 2d \frac{(a+b+2)(\tau-1)(\tau+2)}{a+4} \right] \frac{1}{\text{tr } \mathbf{F}}. \quad (6.21)$$

Therefore  $\widehat{\Theta}^{mEM}$  improves on  $\widehat{\Theta}^{EM}$  if

$$0 < d \leq 2 \frac{(a+b+2)(\tau-1)(\tau+2)}{(a+4)(a-2\tau+6)}$$

for  $b \geq 3$  and  $\tau \geq 2$ . The r.h.s. of (6.21) is quadratic in  $d$  and attains its minimum at

$$d_0 = \frac{(a+b+2)(\tau-1)(\tau+2)}{(a+4)(a-2\tau+6)} = \frac{(m+n \vee p)(m \wedge n \wedge p - 1)(m \wedge n \wedge p + 2)}{\{(|n-p|+2m) \wedge (n+p)+1\}(|n-p|+3)}.$$

That is, the Efron-Morris type estimator  $\widehat{\Theta}^{EM}$  is dominated by the modified Efron-Morris type estimator

$$\widehat{\Theta}^{mEM} = \widehat{\Theta}^{EM} - \frac{d_0}{\text{tr } \mathbf{X} \mathbf{S}^+ \mathbf{X}^\top} \mathbf{R} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+$$

for  $b \geq 3$  and  $\tau \geq 2$ .

## 6.6 Related Topics

### 6.6.1 Positive-Part Rule Estimators

In the field of estimating a mean vector of the multivariate normal distribution, a positive-part rule for shrinkage estimators is well known as an improving method of risk. For an analytical proof of the improvement, see Baranchik (1970). Gruber (1998) provided in detail numerical examples to compare risk performance of the James-Stein (1961) shrinkage and the corresponding positive-part rule estimators. In this section, we will give a positive-part rule for matricial shrinkage estimator (6.11).

With the help of Lemma 6.2, denote

$$\widehat{\Theta}^{SH} = \mathbf{X}(\mathbf{I}_p - \mathbf{S} \mathbf{S}^+) + \mathbf{R} \Psi(\mathbf{F}) \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+,$$

where  $\Psi(\mathbf{F}) = \text{diag}(\psi_1(\mathbf{F}), \dots, \psi_\tau(\mathbf{F})) = \mathbf{I}_\tau - \Phi(\mathbf{F})$ . Instead of  $\Psi(\mathbf{F})$ , we here use  $\Psi_+(\mathbf{F}) = \text{diag}(\psi_1^+(\mathbf{F}), \dots, \psi_\tau^+(\mathbf{F}))$  for  $\psi_i^+(\mathbf{F}) = \max\{0, \psi_i(\mathbf{F})\}$ . The resulting estimator is denoted by

$$\widehat{\Theta}_+^{SH} = \mathbf{X}(\mathbf{I}_p - \mathbf{S} \mathbf{S}^+) + \mathbf{R} \Psi_+(\mathbf{F}) \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+.$$

If  $m = 1$  with  $n \geq p$ , then  $\widehat{\Theta}_+^{SH}$  is the same as Baranchik's (1970) positive-part rule estimator. An analytical dominance result for  $m > 1$  with  $n \geq p$  was given by Tsukuma (2010). When  $m = 1$  with  $p > n$ ,  $\widehat{\Theta}_+^{SH}$  was suggested by Chételat and Wells (2012), who showed by simulation that  $\widehat{\Theta}_+^{SH}$  outperforms  $\widehat{\Theta}^{SH}$ . These kind of dominance results can be proved unifiedly for any set  $(m, n, p)$ .

**Theorem 6.2** Assume that  $\Pr(\psi_i(\mathbf{F}) < 0) > 0$  for some  $i$ . Then  $\widehat{\Theta}_+^{SH}$  dominates  $\widehat{\Theta}^{SH}$  relative to the quadratic loss (6.7) regardless of an order relation among  $m, n$  and  $p$ .

**Proof** Abbreviate  $\Psi(\mathbf{F}) = \text{diag}(\psi_1(\mathbf{F}), \dots, \psi_\tau(\mathbf{F}))$  to  $\Psi = \text{diag}(\psi_1, \dots, \psi_\tau)$  and  $\Psi_+(\mathbf{F}) = \text{diag}(\psi_1^+(\mathbf{F}), \dots, \psi_\tau^+(\mathbf{F}))$  to  $\Psi_+ = \text{diag}(\psi_1^+, \dots, \psi_\tau^+)$ , respectively. Put  $v = n \wedge p$ . Let  $\mathbf{H}_0 \in \mathbb{R}^{p \times (p-v)}$  such that  $(\mathbf{H}, \mathbf{H}_0) \in \mathbb{O}_p$ . We can express  $\widehat{\Theta}^{SH}$  as  $\widehat{\Theta}^{SH} = \mathbf{X} \mathbf{H}_0 \mathbf{H}_0^\top + \mathbf{R} \Psi \mathbf{R}^\top \mathbf{X} \mathbf{H} \mathbf{H}^\top$ . Note that

$$\begin{aligned}
\text{tr}(\widehat{\Theta}^{SH} - \Theta)\Sigma^{-1}(\widehat{\Theta}^{SH} - \Theta)^\top &= \text{tr}(XH_0H_0^\top - \Theta)\Sigma^{-1}(XH_0H_0^\top - \Theta)^\top \\
&\quad + 2\text{tr}\Psi R^\top XHH^\top \Sigma^{-1}(XH_0H_0^\top - \Theta)^\top R \\
&\quad + \text{tr}\Psi^2 R^\top XHH^\top \Sigma^{-1}HH^\top X^\top R.
\end{aligned}$$

Thus the difference in risk of  $\widehat{\Theta}_+^{SH}$  and  $\widehat{\Theta}^{SH}$  is

$$\begin{aligned}
&R(\widehat{\Theta}_+^{SH}, \Theta) - R(\widehat{\Theta}^{SH}, \Theta) \\
&= E[\text{tr}(\Psi_+^2 - \Psi^2)R^\top XHH^\top \Sigma^{-1}HH^\top X^\top R] \\
&\quad + 2E[\text{tr}(\Psi_+ - \Psi)R^\top XHH^\top \Sigma^{-1}(XH_0H_0^\top - \Theta)^\top R]. \tag{6.22}
\end{aligned}$$

The first expectation in the r.h.s. of (6.22) is not positive because  $(\psi_i^+)^2 \leq \psi_i^2$  for all  $i$ .

Recall that  $S = H L H^\top$  is the eigenvalue decomposition, where  $H \in \mathbb{V}_{p,v}$  and  $L = \text{diag}(\ell_1, \dots, \ell_v) \in \mathbb{D}_v^{(\geq 0)}$ . From Proposition 3.2 and Equation (3.3), the joint (unnormalized) p.d.f. of  $(X, L, H)$  without a normalizing constant can be written as

$$\exp\left(-\frac{1}{2}\text{tr}(X - \Theta)\Sigma^{-1}(X - \Theta)^\top - \frac{1}{2}\text{tr}\Sigma^{-1}H L H^\top\right)g_{n,p}(L),$$

where

$$g_{n,p}(L) = |L|^{(|n-p|-1)/2} \prod_{1 \leq i < j \leq v} (\ell_i - \ell_j).$$

Noting that

$$\begin{aligned}
\text{tr}(X - \Theta)\Sigma^{-1}(X - \Theta)^\top &= \text{tr}(XH_0H_0^\top - \Theta)\Sigma^{-1}(XH_0H_0^\top - \Theta)^\top \\
&\quad + 2\text{tr}XHH^\top \Sigma^{-1}(XH_0H_0^\top - \Theta)^\top \\
&\quad + \text{tr}XHH^\top \Sigma^{-1}HH^\top X^\top,
\end{aligned}$$

we make the transformation  $(Z, Z_0) = (X H L^{-1/2}, X H_0)$ . Since, by Lemma 3.2, the Jacobian of the transformation is given by  $J[X \rightarrow (Z, Z_0)] = |L|^{m/2}$ , the joint (unnormalized) p.d.f. of  $(Z, Z_0, L, H)$  without a normalizing constant is proportional to

$$\begin{aligned}
&\exp\left(-\frac{1}{2}\text{tr}(Z_0H_0^\top - \Theta)\Sigma^{-1}(Z_0H_0^\top - \Theta)^\top - \text{tr}Z L^{1/2}H^\top \Sigma^{-1}(Z_0H_0^\top - \Theta)^\top\right. \\
&\quad \left.- \frac{1}{2}\text{tr}Z L^{1/2}H^\top \Sigma^{-1}H L^{1/2}Z^\top - \frac{1}{2}\text{tr}\Sigma^{-1}H L H^\top\right)|L|^{m/2}g_{n,p}(L).
\end{aligned}$$

Then the second expectation in the r.h.s. of (6.22) becomes

$$K_0 \iiint_{\mathbb{R}^{m \times (p-v)} \times \mathbb{D}_v^{(\geq 0)} \times \mathbb{V}_{p,v}} I \times f(\mathbf{Z}_0, \mathbf{L}, \mathbf{H}) |\mathbf{L}|^{m/2} g_{n,p}(\mathbf{L}) (d\mathbf{Z}_0)(d\mathbf{L})(\mathbf{H}^\top d\mathbf{H}),$$

where  $K_0$  is a normalizing constant,

$$I = \int_{\mathbb{R}^{m \times v}} \text{tr}(\Psi_+ - \Psi) \mathbf{R}^\top \mathbf{Z} \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top \mathbf{R} \\ \times \exp \left( -\text{tr} \mathbf{Z} \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top - \frac{1}{2} \text{tr} \mathbf{Z} \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2} \mathbf{Z}^\top \right) (d\mathbf{Z})$$

and

$$f(\mathbf{Z}_0, \mathbf{L}, \mathbf{H}) = \exp \left( -\frac{1}{2} \text{tr} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta) \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top - \frac{1}{2} \text{tr} \Sigma^{-1} \mathbf{H} \mathbf{L} \mathbf{H}^\top \right).$$

Hence if it is shown that  $I \leq 0$ , the proof of Theorem 6.2 will be complete.

We next consider the singular value decomposition  $\mathbf{Z} = \mathbf{R} \mathbf{D} \mathbf{V}^\top$ , where  $\mathbf{R} \in \mathbb{V}_{m,\tau}$ ,  $\mathbf{D} = \text{diag}(d_1, \dots, d_\tau) = \mathbf{F}^{1/2} \in \mathbb{D}_\tau^{(\geq 0)}$ ,  $\mathbf{V} \in \mathbb{V}_{v,\tau}$  and  $\tau = m \wedge (n \wedge p) = m \wedge v$ . From Lemma 3.6, we have

$$(d\mathbf{Z}) = \frac{1}{2^\tau} |\mathbf{D}|^{m-v} \prod_{1 \leq i < j \leq \tau} (d_i^2 - d_j^2) (\mathbf{R}^\top d\mathbf{R}) (d\mathbf{D}) (\mathbf{V}^\top d\mathbf{V}) \\ = \frac{1}{2^{2\tau}} g_{m,v}(\mathbf{F}) (\mathbf{R}^\top d\mathbf{R}) (d\mathbf{F}) (\mathbf{V}^\top d\mathbf{V}),$$

where the second equality is verified by the transformation  $\mathbf{F} = \mathbf{D}^2$ . Recall that  $(\mathbf{R}^\top d\mathbf{R})$  and  $(\mathbf{V}^\top d\mathbf{V})$  are invariant with respect to any orthogonal transformation. For  $i = 1, \dots, \tau$ , it is observed that

$$\{\mathbf{R}^\top \mathbf{Z} \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top \mathbf{R}\}_{ii} = f_i^{1/2} \mathbf{v}_i^\top \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top \mathbf{r}_i,$$

where  $\mathbf{v}_i$  and  $\mathbf{r}_i$  are the  $i$ -th column vectors of  $\mathbf{V}$  and  $\mathbf{R}$ , respectively. Letting  $\mathbf{a}_i^\top = f_i^{1/2} \mathbf{v}_i^\top \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} (\mathbf{Z}_0 \mathbf{H}_0^\top - \Theta)^\top$ , we obtain

$$I = \sum_{i=1}^{\tau} \iiint_{\mathbb{V}_{m,\tau} \times \mathbb{D}_\tau^{(\geq 0)} \times \mathbb{V}_{v,\tau}} (\psi_i^+ - \psi_i) \mathbf{a}_i^\top \mathbf{r}_i e^{-\mathbf{a}_i^\top \mathbf{r}_i} G_i(\mathbf{R}^\top d\mathbf{R}) (d\mathbf{F}) (\mathbf{V}^\top d\mathbf{V}), \quad (6.23)$$

where

$$G_i = \exp \left( -\sum_{j \neq i} \mathbf{a}_j^\top \mathbf{r}_j - \frac{1}{2} \text{tr} \mathbf{F} \mathbf{V}^\top \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2} \mathbf{V} \right) \times \frac{1}{2^{2\tau}} g_{m,v}(\mathbf{F}).$$

For each  $i \in \{1, \dots, \tau\}$ , we make the transformation  $\mathbf{r}_i \rightarrow -\mathbf{r}_i$ . This transformation is equivalent to the orthogonal transformation  $\mathbf{R} \rightarrow \mathbf{R} \mathbf{O}_i$ , where  $\mathbf{O}_i \in \mathbb{D}_\tau$  such that

the  $i$ -th diagonal is minus one and the other diagonals are ones. Since  $(\mathbf{R}^\top \mathbf{d}\mathbf{R})$  is invariant with respect to the orthogonal transformation, (6.23) is rewritten as

$$I = \sum_{i=1}^{\tau} \iiint_{\mathbb{V}_{m,\tau} \times \mathbb{D}_{\tau}^{(\geq 0)} \times \mathbb{V}_{v,\tau}} (\psi_i^+ - \psi_i) (-\mathbf{a}_i^\top \mathbf{r}_i e^{\mathbf{a}_i^\top \mathbf{r}_i}) G_i(\mathbf{R}^\top \mathbf{d}\mathbf{R}) (\mathbf{d}\mathbf{F}) (\mathbf{V}^\top \mathbf{d}\mathbf{V}). \quad (6.24)$$

Adding each side of (6.23) and (6.24) yields

$$2I = \sum_{i=1}^{\tau} \iiint_{\mathbb{V}_{m,\tau} \times \mathbb{D}_{\tau}^{(\geq 0)} \times \mathbb{V}_{v,\tau}} (\psi_i^+ - \psi_i) \mathbf{a}_i^\top \mathbf{r}_i (e^{-\mathbf{a}_i^\top \mathbf{r}_i} - e^{\mathbf{a}_i^\top \mathbf{r}_i}) G_i(\mathbf{R}^\top \mathbf{d}\mathbf{R}) (\mathbf{d}\mathbf{F}) (\mathbf{V}^\top \mathbf{d}\mathbf{V}).$$

Since  $\psi_i^+ \geq \psi_i$  and  $\mathbf{a}_i^\top \mathbf{r}_i (e^{-\mathbf{a}_i^\top \mathbf{r}_i} - e^{\mathbf{a}_i^\top \mathbf{r}_i}) \leq 0$  for any value of  $\mathbf{a}_i^\top \mathbf{r}_i$ , it always holds that  $I \leq 0$ . Thus the proof of Theorem 6.2 is complete.  $\square$

For example, the Efron-Morris estimator  $\hat{\Theta}^{EM}$  is dominated by

$$\hat{\Theta}_+^{EM} = \mathbf{X}(\mathbf{I}_p - \mathbf{S}\mathbf{S}^+) + \mathbf{R}\Psi_+^{EM}(\mathbf{F})\mathbf{R}^\top \mathbf{X}\mathbf{S}\mathbf{S}^+,$$

where the  $i$ -th diagonal element of  $\Psi_+^{EM}(\mathbf{F})$  is  $\max(0, 1 - c^{EM}/f_i)$ . This positive-part rule is extending (6.2) to the unknown covariance case. Also, Theorem 6.2 can be applied to  $\hat{\Theta}^{mEM}$ ,  $\hat{\Theta}^{ST}$  and  $\hat{\Theta}^{mST}$  given in Sect. 6.5, but the applications are omitted.

### 6.6.2 Shrinkage Estimation with a Loss Matrix

Next, we look at shrinkage estimation under a loss matrix of the form

$$L_M(\hat{\Theta}, \Theta | \Sigma) = (\hat{\Theta} - \Theta) \Sigma^{-1} (\hat{\Theta} - \Theta)^\top, \quad (6.25)$$

which is an  $m \times m$  symmetric positive semi-definite matrix. The corresponding risk matrix is defined by  $R_M(\hat{\Theta}, \Theta) = E[L_M(\hat{\Theta}, \Theta | \Sigma)]$ . The loss matrix (6.25) is used in Bilodeau and Kariya (1989). For a more general loss matrix, see Honda (1991).

Using Theorem 5.1 gives

$$R_M(\hat{\Theta}^{ML}, \Theta) = E[(\mathbf{X} - \Theta) \Sigma^{-1} (\mathbf{X} - \Theta)^\top] = E[\nabla_X (\mathbf{X} - \Theta)] = p \mathbf{I}_m.$$

Thus an estimator  $\hat{\Theta}$  is said to be better than  $\hat{\Theta}^{ML}$  relative to the loss matrix (6.25) if  $\hat{\Theta}$  has a smaller risk matrix than  $p \mathbf{I}_m$  in the Löwner sense, namely,  $R_M(\hat{\Theta}, \Theta) \preceq p \mathbf{I}_m$ .

Here, we focus our attention on  $\hat{\Theta}^{SH}$  in (6.11). Denote  $\mathbf{G} = \mathbf{R}\Phi\mathbf{R}^\top \mathbf{X}$ . The risk matrix of  $\hat{\Theta}^{SH}$  is written as

$$\begin{aligned} R_M(\hat{\Theta}^{SH}, \Theta) &= R_M(\hat{\Theta}^{ML}, \Theta) - E[(\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{S}\mathbf{S}^+ \mathbf{G}^\top] \\ &\quad - E[(\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{S}\mathbf{S}^+ \mathbf{G}^\top]^\top + E[\mathbf{G}\mathbf{S}\mathbf{S}^+ \Sigma^{-1} \mathbf{S}\mathbf{S}^+ \mathbf{G}^\top]. \end{aligned}$$

Using the Stein identity (5.1) gives

$$E[(X - \Theta)\Sigma^{-1}SS^+G^\top] = E[\nabla_X SS^+G^\top]$$

and

$$\begin{aligned} E[GSS^+\Sigma^{-1}SS^+G^\top] &= E[GS^+Y^\top Y\Sigma^{-1}SS^+G^\top] \\ &= E[GS^+Y^\top \nabla_Y SS^+G^\top] + E[\{GSS^+\nabla_Y^\top YS^+G^\top\}^\top]. \end{aligned}$$

Thus the unbiased risk estimate of  $\hat{\Theta}^{SH}$  relative to the loss matrix (6.25) becomes

$$\begin{aligned} \hat{R}_M(\hat{\Theta}^{SH}) &= pI_m - \nabla_X SS^+G^\top - \{\nabla_X SS^+G^\top\}^\top \\ &\quad + GS^+Y^\top \nabla_Y SS^+G^\top + \{GSS^+\nabla_Y^\top YS^+G^\top\}^\top, \end{aligned}$$

implying that, according to Lemmas 6.6 and 6.7 in the Appendix,

$$\hat{R}_M(\hat{\Theta}^{SH}) = pI_m - 2(\text{tr } \Phi)(I_m - \mathbf{R}\mathbf{R}^\top) + \mathbf{R}\Phi^*\mathbf{R}^\top,$$

where  $\Phi^* = \text{diag}(\phi_1^*, \dots, \phi_\tau^*)$  and for  $i = 1, \dots, \tau$

$$\begin{aligned} \phi_i^* &= af_i\phi_i^2 - 4f_i^2\phi_i \frac{\partial \phi_i}{\partial f_i} - 2 \sum_{j \neq i}^{\tau} \frac{f_i^2 \phi_i^2}{f_i - f_j} + 2 \sum_{j \neq i}^{\tau} \frac{f_i \phi_i f_j \phi_j}{f_i - f_j} \\ &\quad - 2(n \wedge p - \tau + 1)\phi_i - 4f_i \frac{\partial \phi_i}{\partial f_i} - 2 \sum_{j \neq i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j} \end{aligned}$$

with  $a = n + p - 2(n \wedge p) + 2\tau - 3 = (|n - p| + 2m) \wedge (n + p) - 3$ . The above discussion is summarized as follows.

**Proposition 6.1** *If  $\sum_{i=1}^{\tau} \phi_i \geq 0$  and  $\phi_i^* \leq 0$  for  $i = 1, \dots, \tau$ , then  $\hat{\Theta}^{SH}$  dominates  $\hat{\Theta}^{ML}$  relative to the loss matrix (6.25).*

As a specific example, we consider improvement on the Efron-Morris type estimator (6.9). Putting  $\phi_i = c/f_i$  gives  $\phi_i^* = \{(a+4)c^2 - 2(n \wedge p - \tau - 1)\}f_i^{-1}$  for each  $i$ . Hence if  $\tau = m$  and  $0 < c \leq 2(n \wedge p - m - 1)/(a+4)$ , then  $\hat{R}_M(\hat{\Theta}^{EMK}) \leq pI_m$ .

### 6.6.3 Application to a GMANOVA Model

Here, we treat shrinkage estimation in a generalized MANOVA model (Potthoff and Roy, 1964) of the form

$$\mathbf{Z} = \mathbf{ABC} + \mathbf{E}, \quad (6.26)$$

where  $\mathbf{Z} \in \mathbb{R}^{N \times q}$  is an observation matrix,  $\mathbf{A} \in \mathbb{R}^{N \times m_1}$  and  $\mathbf{C} \in \mathbb{R}^{p \times q}$  are constant matrices of full rank with  $N \geq m_1$  and  $q \geq p$ ,  $\mathbf{B} \in \mathbb{R}^{m_1 \times p}$  is an unknown regression coefficient matrix and  $\mathbf{E} \in \mathbb{R}^{N \times q}$  is a random error matrix. Assume that  $\mathbf{E} \sim \mathcal{N}_{N \times q}(\mathbf{0}_{N \times q}, \mathbf{I}_N \otimes \mathbf{\Sigma}_0)$  and  $\mathbf{\Sigma}_0 \in \mathbb{S}_q^{(+)}$  is unknown. The generalized MANOVA model is abbreviated by the GMANOVA model and it is also called the growth curve model. The purpose of this section is to present a shrinkage estimator of  $\mathbf{B}$  improving the ML estimator.

To simplify the estimation problem, we first derive a canonical form of (6.26). Let  $\mathbf{\Gamma}_A \in \mathbb{O}_N$  and  $\mathbf{\Gamma}_C \in \mathbb{O}_q$  such that

$$\mathbf{\Gamma}_A \mathbf{A} = \begin{pmatrix} (\mathbf{A}^\top \mathbf{A})^{1/2} \\ \mathbf{0}_{(N-m_1) \times m_1} \end{pmatrix}, \quad \mathbf{C} \mathbf{\Gamma}_C = ((\mathbf{C} \mathbf{C}^\top)^{1/2}, \mathbf{0}_{p \times m_2})$$

for  $m_2 = q - p$ . Denote  $\mathbf{\Theta} = (\mathbf{A}^\top \mathbf{A})^{1/2} \mathbf{B} (\mathbf{C} \mathbf{C}^\top)^{1/2}$ ,

$$\mathbf{\Gamma}_A \mathbf{Z} \mathbf{\Gamma}_C = \begin{pmatrix} \mathbf{X}_1 & \mathbf{U} \\ \mathbf{Z}_1 & \mathbf{Z}_2 \end{pmatrix}, \quad \mathbf{\Gamma}_C^\top \mathbf{\Sigma}_0 \mathbf{\Gamma}_C = \begin{pmatrix} \mathbf{I}_p & \mathbf{\Xi}^\top \\ \mathbf{0}_{m_2 \times p} & \mathbf{I}_{m_2} \end{pmatrix} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0}_{p \times m_2} \\ \mathbf{0}_{m_2 \times p} & \mathbf{\Omega} \end{pmatrix} \begin{pmatrix} \mathbf{I}_p & \mathbf{0}_{p \times m_2} \\ \mathbf{\Xi} & \mathbf{I}_{m_2} \end{pmatrix},$$

where  $\mathbf{X}_1 \in \mathbb{R}^{m_1 \times p}$ ,  $\mathbf{U} \in \mathbb{R}^{m_1 \times m_2}$ ,  $\mathbf{\Sigma} \in \mathbb{S}_p^{(+)}$ ,  $\mathbf{\Omega} \in \mathbb{S}_{m_2}^{(+)}$  and  $\mathbf{\Xi} \in \mathbb{R}^{m_2 \times p}$ . Further let  $\mathbf{\Gamma}_Z \in \mathbb{O}_{N-m_1}$  such that

$$\mathbf{\Gamma}_Z \mathbf{Z}_2 = \begin{pmatrix} \mathbf{W}^{1/2} \\ \mathbf{0}_{n \times m_2} \end{pmatrix}$$

with  $n = N - m_1 - m_2$  and  $\mathbf{W} = \mathbf{Z}_2^\top \mathbf{Z}_2$ . Denote  $\mathbf{\Gamma}_Z \mathbf{Z}_1 = (\mathbf{V}^\top \mathbf{W}^{1/2}, \mathbf{Y}^\top)^\top$ , where  $\mathbf{V} \in \mathbb{R}^{m_2 \times p}$ . Then a canonical form of (6.26) is given as follows:

$$\mathbf{X}_1 | \mathbf{U} \sim \mathcal{N}_{m_1 \times p}(\mathbf{\Theta} + \mathbf{U} \mathbf{\Xi}, \mathbf{I}_{m_1} \otimes \mathbf{\Sigma}), \quad (6.27)$$

$$\mathbf{U} \sim \mathcal{N}_{m_1 \times m_2}(\mathbf{0}_{m_1 \times m_2}, \mathbf{I}_{m_1} \otimes \mathbf{\Omega}), \quad (6.28)$$

$$\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{\Sigma}), \quad (6.29)$$

$$\mathbf{V} | \mathbf{W} \sim \mathcal{N}_{m_2 \times p}(\mathbf{\Xi}, \mathbf{W}^{-1} \otimes \mathbf{\Sigma}), \quad (6.30)$$

$$\mathbf{W} \sim \mathcal{W}_{m_2}(N - m_1, \mathbf{\Omega}), \quad (6.31)$$

where  $(\mathbf{X}_1, \mathbf{U})$ ,  $\mathbf{Y}$  and  $(\mathbf{V}, \mathbf{W})$  are independent. For derivation of the above canonical form, see Srivastava and Khatri (1979, pp.192–193).

Here we view the estimation problem of  $\mathbf{\Theta}$  relative to the quadratic loss

$$L(\widehat{\mathbf{\Theta}}, \mathbf{\Theta} | \mathbf{\Sigma}) = \text{tr}(\widehat{\mathbf{\Theta}} - \mathbf{\Theta}) \mathbf{\Sigma}^{-1} (\widehat{\mathbf{\Theta}} - \mathbf{\Theta})^\top. \quad (6.32)$$

The risk is defined by  $R(\widehat{\mathbf{\Theta}}, \mathbf{\Theta}) = E[L(\widehat{\mathbf{\Theta}}, \mathbf{\Theta} | \mathbf{\Sigma})]$ , where the expectation is taken with respect to (6.27)–(6.31).

The ML estimator of  $\mathbf{\Theta}$  is

$$\widehat{\mathbf{\Theta}}^{ML} = \mathbf{X}_1 - \mathbf{U} \mathbf{V}.$$

Denote  $S = Y^\top Y$ . To improve  $\hat{\Theta}^{ML}$ , we consider a double shrinkage estimator (Kariya et al. 1996, 1999) having the form

$$\hat{\Theta}^{DSH} = X_1 - G_1 - U(V - G_2),$$

where  $G_1 = G_1(X_1, S) \in \mathbb{R}^{m_1 \times p}$ , and  $G_2 = G_2(V, S|U, W) \in \mathbb{R}^{m_2 \times p}$  satisfies

$$G_2(V, S|U, W) = G_2(V, S| - U, W). \quad (6.33)$$

The risk of  $\hat{\Theta}^{DSH}$  can be written as

$$\begin{aligned} R(\hat{\Theta}^{DSH}, \Theta) &= E[\text{tr}(X_1 - G_1 - \Theta - U\Xi)\Sigma^{-1}(X_1 - G_1 - \Theta - U\Xi)] \\ &\quad - 2E[\text{tr}(X_1 - G_1 - \Theta - U\Xi)\Sigma^{-1}(V - G_2 - \Xi)^\top U^\top] \\ &\quad + E[\text{tr}U(V - G_2 - \Xi)\Sigma^{-1}(V - G_2 - \Xi)^\top U^\top]. \end{aligned} \quad (6.34)$$

The second term of the r.h.s. in (6.34) is zero because the distributions (6.27), (6.28) and (6.30) are symmetric and  $G_2$  has the symmetry assumption (6.33), so that

$$R(\hat{\Theta}^{DSH}, \Theta) = E^U[R_1(G_1)] + E^{U,W}[R_2(G_2)],$$

where

$$\begin{aligned} R_1(G_1) &= E[\text{tr}(X_1 - G_1 - \Theta - U\Xi)\Sigma^{-1}(X_1 - G_1 - \Theta - U\Xi)|U], \\ R_2(G_2) &= E[\text{tr}U^\top U(V - G_2 - \Xi)\Sigma^{-1}(V - G_2 - \Xi)^\top |U, W]. \end{aligned}$$

This suggests the possibility of double shrinkage estimation in both distributions of  $X_1$  and  $V$ . The risk of the ML estimator can be expressed as  $R(\hat{\Theta}^{ML}, \Theta) = E^U[R_1(\mathbf{0}_{m_1 \times p})] + E^{U,W}[R_2(\mathbf{0}_{m_2 \times p})]$ .

For example, we consider the case of  $m_1 \geq m_2$ . Let  $\tau_1 = m_1 \wedge n \wedge p$ ,

$$X_1 S^+ X_1^\top = R_1 F_1 R_1^\top, \quad G_1^{EM} = c_1 R_1 F_1^{-1} R_1^\top X_1 S S^+,$$

where  $F_1 \in \mathbb{D}_{\tau_1}^{(\geq 0)}$ ,  $R_1 \in \mathbb{V}_{m_1 \times \tau_1}$  and  $c_1$  is a positive constant. In addition, let  $\tau_2 = m_2 \wedge n \wedge p$ ,  $X_2 = (U^\top U)^{-1/2} W V$ ,

$$X_2 S^+ X_2^\top = R_2 F_2 R_2^\top, \quad G_2^{EM} = c_2 (U^\top U)^{-1/2} R_2 F_2^{-1} R_2^\top X_2 S S^+,$$

where  $F_2 \in \mathbb{D}_{\tau_2}^{(\geq 0)}$ ,  $R_2 \in \mathbb{V}_{m_2 \times \tau_2}$  and  $c_2$  is a positive constant. Then the Efron-Morris type estimator is defined by

$$\hat{\Theta}^{EM} = X_1 - G_1^{EM} - U(V - G_2^{EM}).$$



Using the same arguments as in Sect. 6.5.1 immediately gives  $R_1(\mathbf{G}_1^{EM}) \leq R_1(\mathbf{0}_{m_1 \times p})$  if

$$0 < c_1 \leq \frac{2(|n \wedge p - m_1| - 1)}{(|n - p| + 2m_1) \wedge (n + p) + 1}. \quad (6.35)$$

For evaluating  $R_2(\mathbf{G}_2^{EM})$ , let  $\nabla_V$  and  $\nabla_{X_2}$  be matrix differential operators with respect to  $V$  and  $X_2$ , respectively. Using the same arguments as in (5.2) gives  $\nabla_{X_2} = (\mathbf{U}^\top \mathbf{U})^{1/2} \mathbf{W}^{-1} \nabla_V$ . The Stein identity (5.1) in terms of (6.30) is used to obtain

$$\begin{aligned} & E[\text{tr } \mathbf{U}^\top \mathbf{U}(\mathbf{V} - \mathbf{\Xi}) \mathbf{\Sigma}^{-1} (\mathbf{G}_2^{EM})^\top | \mathbf{U}, \mathbf{W}] \\ &= E[\text{tr } \mathbf{U}^\top \mathbf{U} \mathbf{W}^{-1} \nabla_V (\mathbf{G}_2^{EM})^\top | \mathbf{U}, \mathbf{W}] \\ &= c_2 E[\text{tr } (\mathbf{U}^\top \mathbf{U})^{1/2} \mathbf{W}^{-1} \nabla_V \mathbf{S} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-1} \mathbf{R}_2^\top | \mathbf{U}, \mathbf{W}] \\ &= c_2 E[\text{tr } \nabla_{X_2} \mathbf{S} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-1} \mathbf{R}_2^\top | \mathbf{U}, \mathbf{W}]. \end{aligned}$$

Note also that

$$\begin{aligned} E[\text{tr } \mathbf{U}^\top \mathbf{U} \mathbf{G}_2^{EM} \mathbf{\Sigma}^{-1} (\mathbf{G}_2^{EM})^\top | \mathbf{U}, \mathbf{W}] &= E[\text{tr } \mathbf{\Sigma}^{-1} \mathbf{S} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-2} \mathbf{R}_2^\top \mathbf{X}_2 \mathbf{S} \mathbf{S}^\top | \mathbf{U}, \mathbf{W}] \\ &= E[\text{tr } \nabla_Y^\top \mathbf{Y} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-2} \mathbf{R}_2^\top \mathbf{X}_2 \mathbf{S} \mathbf{S}^\top | \mathbf{U}, \mathbf{W}]. \end{aligned}$$

Hence,

$$\begin{aligned} R_2(\mathbf{G}_2^{EM}) &= R_2(\mathbf{0}_{m_2 \times p}) - 2c_2 E[\text{tr } \nabla_{X_2} \mathbf{S} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-1} \mathbf{R}_2^\top | \mathbf{U}, \mathbf{W}] \\ &\quad + c_2^2 E[\text{tr } \nabla_Y^\top \mathbf{Y} \mathbf{S}^\top \mathbf{X}_2^\top \mathbf{R}_2 \mathbf{F}_2^{-2} \mathbf{R}_2^\top \mathbf{X}_2 \mathbf{S} \mathbf{S}^\top | \mathbf{U}, \mathbf{W}]. \end{aligned}$$

Using the same arguments as in Sects. 6.4 and 6.5.1 gives  $R_2(\mathbf{G}_2^{EM}) \leq R_2(\mathbf{0}_{m_2 \times p})$  if

$$0 < c_2 \leq \frac{2(|n \wedge p - m_2| - 1)}{(|n - p| + 2m_2) \wedge (n + p) + 1}. \quad (6.36)$$

From the abovementioned,  $\hat{\mathbf{\Theta}}^{EM}$  dominates  $\hat{\mathbf{\Theta}}^{ML}$  relative to the quadratic loss (6.32) if  $c_1$  and  $c_2$  satisfy, respectively, (6.35) and (6.36).

A unified dominance result in the case of  $m_2 > m_1$  can be established in a similar way to the case of  $m_1 \geq m_2$ , but it is omitted. Other approaches to decision-theoretic estimation in the GMANOVA model have been studied by Tan (1991), Kubokawa et al. (1992) and Kariya et al. (1996, 1999).

### 6.6.4 Generalization in an Elliptically Contoured Model

Consider the multivariate linear model (4.1), but the  $N \times p$  random error matrix  $\mathbf{E}$  is assumed to have a p.d.f. of the form

$$|\Sigma|^{-N/2} g(\text{tr } \Sigma^{-1} \mathbf{E}^\top \mathbf{E}), \quad (6.37)$$

where  $g$  is a nonnegative and nonincreasing function on the nonnegative real line. In general, a probability distribution having the p.d.f. (6.37) is commonly called the elliptically contoured distribution. This section introduces shrinkage estimation in the elliptically contoured distribution model.

Using the orthogonal transformation as in Sect. 4.2, we can rewrite (6.37) as

$$|\Sigma|^{-N/2} g(\text{tr } (\mathbf{X} - \Theta) \Sigma^{-1} (\mathbf{X} - \Theta)^\top + \text{tr } \Sigma^{-1} \mathbf{Y}^\top \mathbf{Y}), \quad (6.38)$$

where  $N = m + n$ ,  $\mathbf{X} \in \mathbb{R}^{m \times p}$ ,  $\mathbf{Y} \in \mathbb{R}^{n \times p}$ ,  $\Theta \in \mathbb{R}^{m \times p}$  and  $\Sigma \in \mathbb{S}_p^{(+)}$ . Suppose that  $\Theta$  and  $\Sigma$  are unknown and that, using  $\mathbf{X}$  and  $\mathbf{Y}$ , we want to decision-theoretically estimate  $\Theta$  relative to the quadratic loss (6.7).

Let

$$G(x) = \frac{1}{2} \int_x^\infty g(t) dt.$$

Denote

$$\begin{aligned} E^g[u(\mathbf{X}, \mathbf{Y})] &= \iint_{\mathbb{R}^{m \times p} \times \mathbb{R}^{n \times p}} u(\mathbf{X}, \mathbf{Y}) |\Sigma|^{-N/2} g(w) (d\mathbf{X})(d\mathbf{Y}), \\ E^G[u(\mathbf{X}, \mathbf{Y})] &= \iint_{\mathbb{R}^{m \times p} \times \mathbb{R}^{n \times p}} u(\mathbf{X}, \mathbf{Y}) |\Sigma|^{-N/2} G(w) (d\mathbf{X})(d\mathbf{Y}), \end{aligned}$$

where  $w = \text{tr } (\mathbf{X} - \Theta) \Sigma^{-1} (\mathbf{X} - \Theta)^\top + \text{tr } \Sigma^{-1} \mathbf{Y}^\top \mathbf{Y}$  and  $u$  is an integrable function. Let  $\mathbf{U} \in \mathbb{R}^{m \times p}$  such that all elements of  $\mathbf{U}$  are absolutely continuous functions of  $\mathbf{X}$ . Then, under some suitable conditions,

$$E^g[\text{tr } (\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{U}^\top] = E^G[\text{tr } \nabla_{\mathbf{X}} \mathbf{U}^\top]. \quad (6.39)$$

For details of the conditions, see Kubokawa and Srivastava (2001). The identity (6.39) is an extension of the Stein identity (5.3).

Applying the Stein identity (6.39) to the risk function of  $\hat{\Theta}^{ML} = \mathbf{X}$ , we obtain

$$\begin{aligned} R_g(\hat{\Theta}^{ML}, \Theta) &= E^g[\text{tr } (\mathbf{X} - \Theta) \Sigma^{-1} (\mathbf{X} - \Theta)^\top] = E^G[\text{tr } \nabla_{\mathbf{X}} (\mathbf{X} - \Theta)^\top] \\ &= E^G[mp]. \end{aligned}$$

The risk of  $\widehat{\Theta}^{SH}$  in (6.11) is expanded to

$$R_g(\widehat{\Theta}^{SH}, \Theta | \Sigma) = R_g(\widehat{\Theta}^{ML}, \Theta | \Sigma) + E^g[\text{tr } \mathbf{R} \Phi \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^\top \Sigma^{-1} \mathbf{S} \mathbf{S}^\top \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top - 2 \text{tr } (\mathbf{X} - \Theta) \Sigma^{-1} \mathbf{S} \mathbf{S}^\top \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top],$$

so that, by the Stein identity (6.39),  $R_g(\widehat{\Theta}^{SH}, \Theta | \Sigma) = E^G[\widehat{\Delta}_G^{SH}]$ , where

$$\widehat{\Delta}_G^{SH} = mp + \text{tr } \nabla_Y^\top \mathbf{Y} \mathbf{S}^\top \mathbf{X}^\top \mathbf{R} \Phi^2 \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^\top - 2 \text{tr } \nabla_X \mathbf{S} \mathbf{S}^\top \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top.$$

Although  $\widehat{\Delta}_G^{SH}$  is not an unbiased risk estimator for  $\widehat{\Theta}^{SH}$ , we can see that  $\widehat{\Theta}^{SH}$  dominates  $\widehat{\Theta}^{ML}$  if  $\widehat{\Delta}_G^{SH} \leq mp$ . Hence the improving procedures in Sect. 6.5 can be applied to estimation of  $\Theta$  in (6.38), and the corresponding dominance results hold without depending on the underlying function  $g$ .

## Appendix

This appendix provides some brief proofs of useful results on matrix differential operators that were previously applied to Theorem 6.1.

Let  $\mathbf{X} = (x_{ab}) \in \mathbb{R}^{m \times p}$  and  $\mathbf{Y} = (y_{ab}) \in \mathbb{R}^{n \times p}$ . Denote the matrix differential operators with respect to  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively, by  $\nabla_X = (d_{ab}^X)$  with  $d_{ab}^X = \partial/\partial x_{ab}$  and by  $\nabla_Y = (d_{ab}^Y)$  with  $d_{ab}^Y = \partial/\partial y_{ab}$ . Let

$$\tau = m \wedge n \wedge p.$$

Here, the eigenvalue decomposition of  $\mathbf{X} \mathbf{S}^\top \mathbf{X}^\top$  is

$$\mathbf{X} \mathbf{S}^\top \mathbf{X}^\top = \mathbf{R} \mathbf{F} \mathbf{R}^\top,$$

where  $\mathbf{F} = \text{diag}(f_1, \dots, f_\tau) \in \mathbb{D}_\tau^{(\geq 0)}$  and  $\mathbf{R} = (r_{ij}) \in \mathbb{V}_{m, \tau}$ .

**Lemma 6.5** For  $i = 1, \dots, \tau$ ,  $k = 1, \dots, m$ ,  $a = 1, \dots, m$  and  $b = 1, \dots, p$ , we have

$$\begin{aligned} \text{(i)} \quad d_{ab}^X f_i &= A_{ab}^{ii}, \\ \text{(ii)} \quad d_{ab}^X r_{ki} &= \sum_{j \neq i} \frac{r_{kj} A_{ab}^{ij}}{f_i - f_j} + f_i^{-1} \{ \mathbf{I}_m - \mathbf{R} \mathbf{R}^\top \}_{ka} \{ \mathbf{R}^\top \mathbf{X} \mathbf{S}^\top \}_{ib}, \end{aligned}$$

where  $A_{ab}^{ij} = r_{aj} \{ \mathbf{R}^\top \mathbf{X} \mathbf{S}^\top \}_{ib} + r_{ai} \{ \mathbf{R}^\top \mathbf{X} \mathbf{S}^\top \}_{jb}$ .

For  $i = 1, \dots, \tau$ ,  $k = 1, \dots, m$ ,  $a = 1, \dots, n$  and  $b = 1, \dots, p$ , we have

$$\begin{aligned} \text{(iii)} \quad d_{ab}^Y f_i &= B_{ab}^{ii}, \\ \text{(iv)} \quad d_{ab}^Y r_{ki} &= \sum_{j \neq i} \frac{r_{kj} B_{ab}^{ij}}{f_i - f_j} + f_i^{-1} \{ \mathbf{R}^\top \mathbf{X} \mathbf{S}^\top \mathbf{S}^\top \mathbf{Y}^\top \}_{ia} \{ (\mathbf{I}_m - \mathbf{R} \mathbf{R}^\top) \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^\top) \}_{kb}, \end{aligned}$$

where

$$\begin{aligned} B_{ab}^{ij} = & -\{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{Y}^\top\}_{ia} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{jb} - \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{Y}^\top\}_{ja} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{ib} \\ & + \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{S}^+ \mathbf{Y}^\top\}_{ia} \{\mathbf{R}^\top \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{jb} \\ & + \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{S}^+ \mathbf{Y}^\top\}_{ja} \{\mathbf{R}^\top \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{ib}. \end{aligned}$$

**Proof** Take  $\mathbf{R}_* \in \mathbb{V}_{m,m-\tau}$  such that  $\mathbf{R}_*^\top \mathbf{R} = \mathbf{0}_{(m-\tau) \times \tau}$ . Define  $\mathbf{R}_0 = (\mathbf{R}, \mathbf{R}_*) \in \mathbb{O}_m$ . Denote  $\mathbf{F}_0 = \text{diag}(f_1, \dots, f_\tau, 0, \dots, 0) (\in \mathbb{D}_m^{\geq 0})$ . Then  $\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top = \mathbf{R}_0 \mathbf{F}_0 \mathbf{R}_0^\top$ .

Differentiating both sides of  $\mathbf{R}_0^\top \mathbf{R}_0 = \mathbf{I}_m$  gives  $[\mathbf{d} \mathbf{R}_0^\top] \mathbf{R}_0 + \mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0 = \mathbf{0}_{m \times m}$ , implying that  $\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0$  is skew-symmetric in  $\mathbb{R}^{m \times m}$ . Thus, for  $j, i \in \{1, \dots, m\}$ ,  $\{\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0\}_{ji} = 0$  if  $j = i$  and  $\{\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0\}_{ji} = -\{\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0\}_{ij} = -\{[\mathbf{d} \mathbf{R}_0^\top] \mathbf{R}_0\}_{ji}$  otherwise. Differentiating both sides of  $\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top = \mathbf{R}_0 \mathbf{F}_0 \mathbf{R}_0^\top$  gives

$$\mathbf{d}(\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top) = [\mathbf{d} \mathbf{R}_0] \mathbf{F}_0 \mathbf{R}_0^\top + \mathbf{R}_0 [\mathbf{d} \mathbf{F}_0] \mathbf{R}_0^\top + \mathbf{R}_0 \mathbf{F}_0 \mathbf{d} \mathbf{R}_0^\top,$$

and then

$$\begin{aligned} \mathbf{R}_0^\top [\mathbf{d}(\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}_0 &= [\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0] \mathbf{F}_0 + \mathbf{d} \mathbf{F}_0 + \mathbf{F}_0 [\mathbf{d} \mathbf{R}_0^\top] \mathbf{R}_0 \\ &= [\mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0] \mathbf{F}_0 + \mathbf{d} \mathbf{F}_0 - \mathbf{F}_0 \mathbf{R}_0^\top \mathbf{d} \mathbf{R}_0. \end{aligned}$$

Comparing each element in both sides of the above identity, we have

$$\begin{aligned} \mathbf{d} f_i &= \{\mathbf{R}^\top [\mathbf{d}(\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ii} \quad \text{for } i \in \{1, \dots, \tau\}, \\ \{\mathbf{R}^\top \mathbf{d} \mathbf{R}\}_{ji} &= \frac{\{\mathbf{R}^\top [\mathbf{d}(\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ji}}{f_i - f_j} \quad \text{for } j, i \in \{1, \dots, \tau\} \text{ with } j \neq i, \\ \{\mathbf{R}_*^\top \mathbf{d} \mathbf{R}\}_{ji} &= \frac{\{\mathbf{R}_*^\top [\mathbf{d}(\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ji}}{f_i} \quad \text{for } j \in \{1, \dots, m - \tau\} \text{ and } i \in \{1, \dots, \tau\}. \end{aligned}$$

Note that  $\mathbf{d}_{ab}^X \mathbf{X} = \mathbf{E}_{ab}$ , where  $\mathbf{E}_{ab} \in \mathbb{R}^{m \times p}$  such that the  $(a, b)$ -th element is one and the other elements are zeros. Since  $\mathbf{d}_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top) = [\mathbf{d}_{ab}^X \mathbf{X}] \mathbf{S}^+ \mathbf{X}^\top + \mathbf{X} \mathbf{S}^+ [\mathbf{d}_{ab}^X \mathbf{X}^\top] = \mathbf{E}_{ab} \mathbf{S}^+ \mathbf{X}^\top + \mathbf{X} \mathbf{S}^+ \mathbf{E}_{ab}^\top$ , we observe that, for  $j, i \in \{1, \dots, \tau\}$ ,

$$\begin{aligned} \{\mathbf{R}^\top [\mathbf{d}_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ji} &= \{\mathbf{R}^\top \mathbf{E}_{ab} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{ji} + \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{E}_{ab}^\top \mathbf{R}\}_{ji} \\ &= r_{aj} \{\mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bi} + \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{jb} r_{ai} = A_{ab}^{ij}. \end{aligned} \quad (6.40)$$

Thus, for  $i = 1, \dots, \tau$ ,  $\mathbf{d}_{ab}^X f_i = \{\mathbf{R}^\top [\mathbf{d}_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ii} = A_{ab}^{ii}$ , which shows (i).

On the other hand, it is observed that for  $k = 1, \dots, m$  and  $i = 1, \dots, \tau$

$$\begin{aligned}
d_{ab}^X r_{ki} &= \{d_{ab}^X \mathbf{R}\}_{ki} = \{(\mathbf{R}\mathbf{R}^\top + \mathbf{R}_* \mathbf{R}_*^\top) d_{ab}^X \mathbf{R}\}_{ki} \\
&= \sum_{j \neq i}^{\tau} r_{kj} \{\mathbf{R}^\top d_{ab}^X \mathbf{R}\}_{ji} + \{\mathbf{R}_* \mathbf{R}_*^\top d_{ab}^X \mathbf{R}\}_{ki} \\
&= \sum_{j \neq i}^{\tau} \frac{r_{kj} A_{ab}^{ij}}{f_i - f_j} + \frac{\{\mathbf{R}_* \mathbf{R}_*^\top [d_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ki}}{f_i}. \tag{6.41}
\end{aligned}$$

In a similar way to (6.40), for  $k = 1, \dots, m$  and  $i = 1, \dots, \tau$ ,

$$\{\mathbf{R}_* \mathbf{R}_*^\top [d_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ki} = \{\mathbf{R}_* \mathbf{R}_*^\top\}_{ka} \{\mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bi} + \{\mathbf{R}_* \mathbf{R}_*^\top \mathbf{X} \mathbf{S}^+\}_{kb} r_{ai}.$$

Here,  $\mathbf{R}_*^\top \mathbf{X} \mathbf{S}^+ = \mathbf{0}_{(m-\tau) \times p}$ , so that

$$\{\mathbf{R}_* \mathbf{R}_*^\top [d_{ab}^X (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ki} = \{\mathbf{I}_m - \mathbf{R} \mathbf{R}^\top\}_{ka} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{ib}. \tag{6.42}$$

Substituting (6.42) into (6.41), we obtain (ii).

Since  $\{\mathbf{R}^\top [d_{ab}^Y (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ji} = \{\mathbf{R}^\top \mathbf{X} [d_{ab}^Y \mathbf{S}^+] \mathbf{X}^\top \mathbf{R}\}_{ji}$  for  $j, i \in \{1, \dots, \tau\}$ , it is observed from (iii) of Lemma 5.2 that

$$\begin{aligned}
&\{\mathbf{R}^\top [d_{ab}^Y (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ji} \\
&= -\{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{jb} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{Y}^\top\}_{ia} - \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{Y}^\top\}_{ja} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{ib} \\
&\quad + \{\mathbf{R}^\top \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{jb} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{S}^+ \mathbf{Y}^\top\}_{ia} \\
&\quad + \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{S}^+ \mathbf{Y}^\top\}_{ja} \{\mathbf{R}^\top \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{ib} = B_{ab}^{ij}.
\end{aligned}$$

Similarly,

$$\{\mathbf{R}_* \mathbf{R}_*^\top [d_{ab}^Y (\mathbf{X} \mathbf{S}^+ \mathbf{X}^\top)] \mathbf{R}\}_{ki} = \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{S}^+ \mathbf{Y}^\top\}_{ia} \{(\mathbf{I}_m - \mathbf{R} \mathbf{R}^\top) \mathbf{X} (\mathbf{I}_p - \mathbf{S} \mathbf{S}^+)\}_{kb}$$

for  $k = 1, \dots, m$  and  $i = 1, \dots, \tau$ . Hence using the same arguments as in the proofs of (i) and (ii) yields (iii) and (iv).  $\square$

**Lemma 6.6** Let  $c_0 = |n \wedge p - m| + 1$ . Define  $\Phi = \text{diag}(\phi_1, \dots, \phi_\tau) \in \mathbb{D}_\tau$  such that the diagonals of  $\Phi$  are absolutely continuous functions of  $\mathbf{F}$ . Then

$$\nabla_X \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top = \mathbf{R} \Phi^* \mathbf{R}^\top + (\text{tr } \Phi) (\mathbf{I}_m - \mathbf{R} \mathbf{R}^\top), \tag{6.43}$$

where  $\Phi^* = \text{diag}(\phi_1^*, \dots, \phi_\tau^*)$  and for  $i = 1, \dots, \tau$

$$\phi_i^* = (n \wedge p - \tau + 1) \phi_i + 2 f_i \frac{\partial \phi_i}{\partial f_i} + \sum_{j \neq i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j}.$$

In particular,

$$\text{tr } \nabla_X \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top = \sum_{i=1}^{\tau} \left\{ c_0 \phi_i + 2 f_i \frac{\partial \phi_i}{\partial f_i} + 2 \sum_{j>i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j} \right\}. \quad (6.44)$$

**Proof** For  $a, c \in \{1, \dots, \tau\}$ , the  $(a, c)$ -th element of  $\nabla_X \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top$  is

$$\{\nabla_X \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top\}_{ac} = \sum_{b=1}^p \sum_{d=1}^p \sum_{k=1}^m \sum_{i=1}^{\tau} d_{ab}^X [\{\mathbf{S} \mathbf{S}^+\}_{bd} x_{kd} r_{ki} \phi_i r_{ci}],$$

and then

$$\{\nabla_X \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \Phi \mathbf{R}^\top\}_{ac} = D_{ac}^{(1)} + D_{ac}^{(2)} + D_{ac}^{(3)}, \quad (6.45)$$

where

$$\begin{aligned} D_{ac}^{(1)} &= \sum_{b=1}^p \sum_{d=1}^p \sum_{k=1}^m \{\mathbf{S} \mathbf{S}^+\}_{bd} \{\mathbf{R} \Phi \mathbf{R}^\top\}_{kc} d_{ab}^X x_{kd}, \\ D_{ac}^{(2)} &= \sum_{b=1}^p \sum_{i=1}^{\tau} \{\mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bi} r_{ci} d_{ab}^X \phi_i, \\ D_{ac}^{(3)} &= \sum_{b=1}^p \sum_{k=1}^m \sum_{i=1}^{\tau} \{\mathbf{S} \mathbf{S}^+ \mathbf{X}^\top\}_{bk} \phi_i (r_{ci} d_{ab}^X r_{ki} + r_{ki} d_{ab}^X r_{ci}). \end{aligned}$$

Since  $d_{ab}^X x_{kd} = \delta_{ak} \delta_{bd}$  and  $\mathbf{S} \mathbf{S}^+$  is idempotent with rank  $n \wedge p$ , it follows that

$$D_{ac}^{(1)} = \sum_{b=1}^p \{\mathbf{S} \mathbf{S}^+\}_{bb} \{\mathbf{R} \Phi \mathbf{R}^\top\}_{ac} = (n \wedge p) \{\mathbf{R} \Phi \mathbf{R}^\top\}_{ac}. \quad (6.46)$$

To evaluate  $D_{ac}^{(2)}$ , we first use the chain rule and (i) of Lemma 6.5 to obtain

$$d_{ab}^X \phi_i = \sum_{j=1}^{\tau} [d_{ab}^X f_j] \frac{\partial \phi_i}{\partial f_j} = \sum_{j=1}^{\tau} A_{ab}^{ij} \frac{\partial \phi_i}{\partial f_j}.$$

Note that  $\mathbf{S}^+ \mathbf{S} \mathbf{S}^+ = \mathbf{S}^+$  and

$$\sum_{b=1}^p \{\mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bi} A_{ab}^{ij} = 2r_{aj} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{ji} = 2r_{aj} \{\mathbf{F}\}_{ji},$$

so that

$$D_{ac}^{(2)} = 2 \sum_{i=1}^{\tau} \sum_{j=1}^{\tau} r_{aj} r_{ci} \{\mathbf{F}\}_{ji} \frac{\partial \phi_i}{\partial f_j} = 2 \sum_{i=1}^{\tau} r_{ai} r_{ci} f_i \frac{\partial \phi_i}{\partial f_i}. \quad (6.47)$$

Finally, we consider  $D_{ac}^{(3)}$ . Using (ii) of Lemma 6.5 yields

$$\begin{aligned} \sum_{k=1}^m \{\mathbf{S}\mathbf{S}^+ \mathbf{X}^\top\}_{bk} d_{ab}^X r_{ki} &= \sum_{j \neq i}^{\tau} \frac{\{\mathbf{S}\mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bj} A_{ab}^{ij}}{f_i - f_j} \\ &\quad + f_i^{-1} \{(\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top) \mathbf{X} \mathbf{S} \mathbf{S}^+\}_{ab} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{ib}. \end{aligned}$$

Since

$$\begin{aligned} \sum_{b=1}^p \{\mathbf{S}\mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{bj} A_{ab}^{ij} &= r_{aj} \{\mathbf{F}\}_{ij} + r_{ai} f_j, \\ \sum_{b=1}^p \{(\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top) \mathbf{X} \mathbf{S} \mathbf{S}^+\}_{ab} \{\mathbf{R}^\top \mathbf{X} \mathbf{S}^+\}_{ib} &= \{(\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top) \mathbf{X} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R}\}_{ai} = 0, \end{aligned}$$

it is seen that

$$\begin{aligned} &\sum_{b=1}^p \sum_{k=1}^m \sum_{i=1}^{\tau} \{\mathbf{S}\mathbf{S}^+ \mathbf{X}^\top\}_{bk} \phi_i r_{ci} d_{ab}^X r_{ki} \\ &= \sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{r_{ai} r_{ci} f_j \phi_i}{f_i - f_j} = \sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{r_{ai} r_{ci} (f_j - f_i + f_i) \phi_i}{f_i - f_j} \\ &= -(\tau - 1) \sum_{i=1}^{\tau} r_{ai} r_{ci} \phi_i + \sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{r_{ai} r_{ci} f_i \phi_i}{f_i - f_j}. \end{aligned} \quad (6.48)$$

Similarly,

$$\begin{aligned} \sum_{b=1}^p \sum_{k=1}^m \sum_{i=1}^{\tau} \{\mathbf{S}\mathbf{S}^+ \mathbf{X}^\top\}_{bk} \phi_i r_{ki} d_{ab}^X r_{ci} &= \sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{r_{aj} r_{ci} f_i \phi_i}{f_i - f_j} + (\text{tr } \Phi) \{\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top\}_{ac} \\ &= -\sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{r_{ai} r_{ci} f_j \phi_j}{f_i - f_j} + (\text{tr } \Phi) \{\mathbf{I}_m - \mathbf{R}\mathbf{R}^\top\}_{ac}. \end{aligned} \quad (6.49)$$

Combining (6.48) and (6.49) gives

$$D_{ac}^{(3)} = \sum_{i=1}^{\tau} r_{ai} r_{ci} \left\{ -(\tau - 1)\phi_i + \sum_{j \neq i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j} \right\} + (\text{tr } \Phi) \{I_m - \mathbf{R} \mathbf{R}^{\top}\}_{ac}. \quad (6.50)$$

Substituting (6.46), (6.47) and (6.50) into (6.45) yields (6.43).

Note that  $n \wedge p - \tau + 1 + \text{tr } (I_m - \mathbf{R} \mathbf{R}^{\top}) = n \wedge p + m - 2\tau + 1 = |n \wedge p - m| + 1 = c_0$  and also that

$$\sum_{i=1}^{\tau} \sum_{j \neq i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j} = 2 \sum_{i=1}^{\tau} \sum_{j > i}^{\tau} \frac{f_i \phi_i - f_j \phi_j}{f_i - f_j}.$$

Hence taking the trace of (6.43) yields (6.44), which completes the proof.  $\square$

**Lemma 6.7** Let  $c_1 = n - (n \wedge p) + \tau - 2$ ,  $c_2 = p - (n \wedge p) + \tau - 1$  and  $c_0 = c_1 + c_2$ . Let  $\Phi = \Phi(F) = \text{diag}(\phi_1, \dots, \phi_{\tau})$ , where the  $\phi_i$ 's are absolutely continuous functions of  $F$ . Then

$$\mathbf{R} \Phi \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^{\top} \nabla_Y^{\top} \mathbf{Y} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi \mathbf{R}^{\top} = \mathbf{R} \Phi^{*1} \mathbf{R}^{\top}, \quad (6.51)$$

$$\mathbf{R} \Phi \mathbf{R}^{\top} \mathbf{X} \mathbf{S}^{\top} \mathbf{Y}^{\top} \nabla_Y \mathbf{S} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi \mathbf{R}^{\top} = \mathbf{R} \Phi^{*2} \mathbf{R}^{\top}, \quad (6.52)$$

where  $\Phi^{*k} = \text{diag}(\phi_1^{*k}, \dots, \phi_{\tau}^{*k})$  for  $k = 1, 2$  and, for  $i = 1, \dots, \tau$ ,

$$\phi_i^{*k} = c_k f_i \phi_i^2 - 2 f_i^2 \phi_i \frac{\partial \phi_i}{\partial f_i} - \sum_{j \neq i}^{\tau} \frac{f_i^2 \phi_i^2}{f_i - f_j} + \sum_{j \neq i}^{\tau} \frac{f_i \phi_i f_j \phi_j}{f_i - f_j}.$$

In particular,

$$\text{tr } \nabla_Y^{\top} \mathbf{Y} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi^2 \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^{\top} = \sum_{i=1}^{\tau} \left\{ c_0 f_i \phi_i^2 - 2 f_i^2 \frac{\partial(\phi_i^2)}{\partial f_i} - 2 \sum_{j > i}^{\tau} \frac{f_i^2 \phi_i^2 - f_j^2 \phi_j^2}{f_i - f_j} \right\}. \quad (6.53)$$

**Proof** The proofs of (6.51) and (6.52) can be done by using the same arguments as in the proof of Lemma 6.6. Since

$$\begin{aligned} \text{tr } \nabla_Y^{\top} \mathbf{Y} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi^2 \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^{\top} &= \text{tr } \nabla_Y^{\top} \{ \mathbf{Y} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi \mathbf{R}^{\top} \cdot \mathbf{R} \Phi \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^{\top} \} \\ &= \text{tr } \mathbf{R} \Phi \mathbf{R}^{\top} \mathbf{X} \mathbf{S} \mathbf{S}^{\top} \nabla_Y^{\top} \mathbf{Y} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi \mathbf{R}^{\top} \\ &\quad + \text{tr } \mathbf{R} \Phi \mathbf{R}^{\top} \mathbf{X} \mathbf{S}^{\top} \mathbf{Y}^{\top} \nabla_Y \mathbf{S} \mathbf{S}^{\top} \mathbf{X}^{\top} \mathbf{R} \Phi \mathbf{R}^{\top}, \end{aligned}$$

the identity (6.53) can be verified by combining (6.51) and (6.52).  $\square$



## References

- A.J. Baranchik, A family of minimax estimators of the mean of a multivariate normal distribution. *Ann. Math. Stat.* **41**, 642–645 (1970)
- M. Bilodeau, T. Kariya, Minimax estimators in the normal MANOVA model. *J. Multivar. Anal.* **28**, 260–270 (1989)
- D. Chételat, M.T. Wells, Improved multivariate normal mean estimation with unknown covariance when  $p$  is greater than  $n$ . *Ann. Stat.* **40**, 3137–3160 (2012)
- B. Efron, C. Morris, Empirical Bayes on vector observations: an extension of Stein's method. *Biometrika* **59**, 335–347 (1972)
- B. Efron, C. Morris, Multivariate empirical Bayes and estimation of covariance matrices. *Ann. Stat.* **4**, 22–32 (1976)
- M.H.J. Gruber, *Improving Efficiency by Shrinkage* (Marcel Dekker, New York, 1998)
- T. Honda, Minimax estimators in the manova model for arbitrary quadratic loss and unknown covariance matrix. *J. Multivar. Anal.* **36**, 113–120 (1991)
- James, W. and Stein, C. (1961). Estimation with quadratic loss, in *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, ed. by J. Neyman (University of California Press, Berkeley), pp. 361–379
- T. Kariya, Y. Konno, W.E. Strawderman, Double shrinkage estimators in the GMANOVA model. *J. Multivar. Anal.* **56**, 245–258 (1996)
- T. Kariya, Y. Konno, W.E. Strawderman, Construction of shrinkage estimators for the regression coefficient matrix in the GMANOVA model. *Commun. Stat.—Theory Methods* **28**, 597–611 (1999)
- Y. Konno, Families of minimax estimators of matrix of normal means with unknown covariance matrix. *J. Japan Stat. Soc.* **20**, 191–201 (1990)
- Y. Konno, On estimation of a matrix of normal means with unknown covariance matrix. *J. Multivar. Anal.* **36**, 44–55 (1991)
- Y. Konno, Improved estimation of matrix of normal mean and eigenvalues in the multivariate  $F$ -distribution. Doctoral dissertation, Institute of Mathematics, University of Tsukuba, 1992. (<http://mcm-www.jwu.ac.jp/~konno/>)
- T. Kubokawa, AKMdE Saleh, K. Morita, Improving on MLE of coefficient matrix in a growth curve model. *J. Stat. Plann. Infer.* **31**, 169–177 (1992)
- T. Kubokawa, M.S. Srivastava, Robust improvement in estimation of a mean matrix in an elliptically contoured distribution. *J. Multivar. Anal.* **76**, 138–152 (2001)
- R.F. Potthoff, S.N. Roy, A generalized multivariate analysis of variance model useful especially for growth curve problems. *Biometrika* **51**, 313–326 (1964)
- M.S. Srivastava, C.G. Khatri, *An Introduction to Multivariate Statistics* (North Holland, New York, 1979)
- C. Stein, Estimation of the mean of a multivariate normal distribution. Technical Reports No. 48 (Department of Statistics, Stanford University, Stanford, 1973)
- M. Tan, Improved estimators for the GMANOVA problem with application to Monte Carlo simulation. *J. Multivar. Anal.* **38**, 262–274 (1991)
- H. Tsukuma, Shrinkage minimax estimation and positive-part rule for a mean matrix in an elliptically contoured distribution. *Stat. Probab. Lett.* **80**, 215–220 (2010)
- H. Tsukuma, T. Kubokawa, Methods for improvement in estimation of a normal mean matrix. *J. Multivar. Anal.* **98**, 1592–1610 (2007)
- H. Tsukuma, T. Kubokawa, A unified approach to estimating a normal mean matrix in high and low dimensions. *J. Multivar. Anal.* **139**, 312–328 (2015)

# Chapter 7

## Estimation of the Covariance Matrix



This chapter addresses decision-theoretic estimation of an error covariance matrix in a multivariate linear model relative to a Stein-type entropy loss. With a unified treatment for high- and low-dimensions, some important improving methods of the best scale and the best triangular invariant estimators are discussed by using the residual sum of squares matrix only. Also this chapter provides interesting dominance results by using the information on both the residual sum of squares matrix and the least squares estimator of the regression coefficient matrix.

### 7.1 Introduction

As seen in (4.2), a canonical form of multivariate linear model (4.1) is given by

$$Y \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \Sigma), \quad X \sim \mathcal{N}_{m \times p}(\Theta, \mathbf{I}_m \otimes \Sigma),$$

where  $Y$  and  $X$  are mutually independent, and  $\Theta \in \mathbb{R}^{m \times p}$  and  $\Sigma \in \mathbb{S}_p^{(+)}$  are matrices of unknown parameters. Throughout this chapter, we use the following notation:

$$v = n \wedge p, \quad \kappa = n \vee p.$$

Here the Wishart matrix  $S = Y^T Y$  is of rank  $v$  and the above canonical form is replaced by

$$S \sim \mathcal{W}_p^v(n, \Sigma), \quad X \sim \mathcal{N}_{m \times p}(\Theta, \mathbf{I}_m \otimes \Sigma). \quad (7.1)$$

The problem of estimating the covariance matrix  $\Sigma$  in (7.1) is looked at from a decision-theoretic point of view.

In the literature, several loss functions have been employed for decision-theoretic estimation of  $\Sigma$  with  $n \geq p$ . One of such loss functions is Stein's (1956) entropy loss

$$L_S(\widehat{\Sigma}, \Sigma) = \text{tr } \Sigma^{-1} \widehat{\Sigma} - \log |\Sigma^{-1} \widehat{\Sigma}| - p. \quad (7.2)$$

Stein (1956) focused on triangular invariant estimators and succeeded in deriving a minimax estimator improving the sample covariance matrix relative to the loss (7.2). Note that Stein's (1956) results were summarized in James and Stein (1961) and, in this chapter, the minimax estimator is called James and Stein's minimax estimator. Although James and Stein's minimax estimator dominates the sample covariance matrix, the minimax estimator depends on the coordinate system and the dependence causes inadmissibility of the minimax estimator. Typical improved estimators on James and Stein's minimax estimator are orthogonally invariant estimators, which are not influenced by the coordinate system. The orthogonally invariant estimators have been studied since Stein (1975, 1977). For other studies on orthogonally invariant estimators, see Takemura (1984), Dey and Srinivasan (1985), Sheena and Takemura (1992) and Perron (1992).

James and Stein's minimax estimator and its improved estimators mentioned above are based only on  $S = Y^T Y$ , while truncation rules have been proposed for improving the existing estimators by using the information contained in  $X$ . Such a truncation rule was first derived by Stein (1964) in estimation of variance of a normal distribution, and several extensions to multivariate models were studied by Sinha and Ghosh (1987), Perron (1990), Kubokawa et al. (1992) and Kubokawa and Srivastava (2003) in the  $n \geq p$  case. These articles applied conditional arguments to deriving dominance results, but Kubokawa and Tsai (2006) used the Stein identity (5.3) to suggest an alternative truncation rule with shrinkage.

When  $p > n$ , Konno (2009) studied decision-theoretic covariance estimation relative to a quadratic loss. In the  $p > n$  case, the Stein loss (7.2) is not available for singular estimators such as the unbiased estimator  $S/n$  since  $|\Sigma^{-1} S| = 0$ . An extended Stein-type entropy loss applicable to singular estimators was treated by Tsukuma (2016a) and Tsukuma and Kubokawa (2016).

This chapter will take a unified approach to both cases of  $n \geq p$  and  $p > n$ . We assume that any estimator lies in  $\mathbb{S}_{p,v}^{(+)}$  and it is of the same rank as  $S$ . More specifically, any estimator is of rank  $p$  when  $n \geq p$  and is of rank  $n$  when  $p > n$ .

Now, let  $\widehat{\Sigma}$  be an estimator of  $\Sigma$  based on  $S$  and  $X$  in (7.1), where  $\widehat{\Sigma} \in \mathbb{S}_{p,v}^{(+)}$ . Since  $\Sigma^{-1}$  is positive definite,  $\Sigma^{-1} \widehat{\Sigma}$  has  $v$  nonzero eigenvalues and they are all positive. Let  $\text{Ch}(\Sigma^{-1} \widehat{\Sigma}) \in \mathbb{D}_v$  such that its diagonal elements consist of  $v$  positive eigenvalues of  $\Sigma^{-1} \widehat{\Sigma}$ . The extended Stein loss is defined by

$$L_{ES}(\widehat{\Sigma}, \Sigma) = \text{tr} [\text{Ch}(\Sigma^{-1} \widehat{\Sigma})] - \log |\text{Ch}(\Sigma^{-1} \widehat{\Sigma})| - v. \quad (7.3)$$

If  $n \geq p$ ,  $L_{ES}$  is the same as the ordinary Stein loss (7.2). The accuracy of estimators is measured by the risk function  $R_{ES}(\widehat{\Sigma}, \Sigma) = E[L_{ES}(\widehat{\Sigma}, \Sigma)]$ , where the expectation

$E$  is taken with respect to the model (7.1). With the extended Stein loss (7.3) used, this chapter gives some dominance results unifying both cases of  $n \geq p$  and  $p > n$ .

First in Sect. 7.2, we deal with the best scale invariant estimator that forms a scalar multiple of  $S$ . Section 7.3 considers the class of triangular invariant estimators and gives a unified expression of the James-Stein (1961) type estimators for the  $n \geq p$  and  $p > n$  cases. Section 7.4 provides some orthogonally invariant estimators improving on the best scale invariant and the unified James-Stein-type estimators relative to the extended Stein loss (7.3). Section 7.5 gives alternative unified estimators using information on the mean statistic  $X$ . In Sect. 7.6, we point out some relevant topics on decision-theoretic covariance estimation.

## 7.2 Scale Invariant Estimators

Recall that  $\nu = n \wedge p$  and  $\kappa = n \vee p$ . Consider a simple class of estimators which forms a constant multiple of  $S$ . The simple class is denoted by

$$\widehat{\Sigma}_c = \widehat{\Sigma}_c(S) = cS, \quad (7.4)$$

where  $c$  is a positive constant. The class (7.4) satisfies  $P\widehat{\Sigma}_c(S)P^\top = \widehat{\Sigma}_c(PSP^\top)$  for any  $P \in \mathbb{U}_p$ , namely,  $\widehat{\Sigma}_c$  is invariant under the scale transformations  $S \rightarrow PSP^\top$  and  $\widehat{\Sigma} \rightarrow P\widehat{\Sigma}P^\top$  for any  $P \in \mathbb{U}_p$ . The unbiased estimator of  $\Sigma$ ,

$$\widehat{\Sigma}^{UB} = \frac{1}{n}S,$$

belongs to the class (7.4), but  $\widehat{\Sigma}^{UB}$  is not the best estimator among the class (7.4) relative to the extended Stein loss (7.3). The best estimator is given in the following proposition.

**Proposition 7.1** *Among the class (7.4), the best estimator relative to the extended Stein loss (7.3) is given by*

$$\widehat{\Sigma}^{BS} = \widehat{\Sigma}_{c_0}$$

with  $c_0 = 1/\kappa$ . Hence for  $p > n$ ,  $\widehat{\Sigma}^{BS} = S/p$  dominates  $\widehat{\Sigma}^{UB} = S/n$  relative to the extended Stein loss (7.3).

**Proof** Recall that  $S = Y^\top Y$  and  $Y \in \mathbb{R}^{n \times p}$ . The positive eigenvalues of  $\Sigma^{-1}S$  are identical to those of  $Y\Sigma^{-1}Y^\top$ , so that the positive eigenvalues of  $\Sigma^{-1}S$  are identical to those of the full-rank matrix

$$\begin{cases} \Sigma^{-1}S \ (\in \mathbb{R}^{p \times p}) & \text{for } n \geq p, \\ Y\Sigma^{-1}Y^\top \ (\in \mathbb{S}_n^{(+)} \subset \mathbb{R}^{n \times n}) & \text{for } p > n. \end{cases}$$

Note that  $\text{tr} [\text{Ch}(\Sigma^{-1} \widehat{\Sigma}_c)] = \text{tr} \Sigma^{-1} \widehat{\Sigma}_c$ . Since  $\Sigma^{-1} \mathbf{S}$  has  $\nu$  positive eigenvalues with probability one, we obtain  $|\text{Ch}(c \Sigma^{-1} \mathbf{S})| = c^\nu |\text{Ch}(\Sigma^{-1} \mathbf{S})|$ . The risk of  $\widehat{\Sigma}_c$  with respect to the extended Stein loss (7.3) is expressed as

$$\begin{aligned} R_{ES}(\widehat{\Sigma}_c, \Sigma) &= nc \text{tr} \Sigma^{-1} \Sigma - \nu \log c - E[\log |\text{Ch}(\Sigma^{-1} \mathbf{S})|] - \nu \\ &= npc - \nu \log c - E[\log |\text{Ch}(\Sigma^{-1} \mathbf{S})|] - \nu, \end{aligned}$$

which is minimized at  $c = c_0$  with

$$c_0 = \frac{\nu}{np} = \frac{1}{\kappa}.$$

Thus the proof is complete.  $\square$

Denote  $r_{\kappa, \nu} = E[\log |\text{Ch}(\Sigma^{-1} \mathbf{S})|]$ . The risk of  $\widehat{\Sigma}^{BS}$  is

$$R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) = \nu \log \kappa - r_{\kappa, \nu}. \quad (7.5)$$

Now, it follows that

$$|\text{Ch}(\Sigma^{-1} \mathbf{S})| = \begin{cases} |\Sigma^{-1} \mathbf{Y}^\top \mathbf{Y}| & \text{for } n \geq p, \\ |\mathbf{Y} \Sigma^{-1} \mathbf{Y}^\top| & \text{for } p > n, \end{cases}$$

and consequently

$$r_{\kappa, \nu} = \begin{cases} E[\log |\mathbf{Z}^\top \mathbf{Z}|] & \text{for } n \geq p, \\ E[\log |\mathbf{Z} \mathbf{Z}^\top|] & \text{for } p > n, \end{cases}$$

where  $\mathbf{Z} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \mathbf{I}_p)$ . Since  $\mathbf{Z}^\top \mathbf{Z} \sim \mathcal{W}_p(n, \mathbf{I}_p)$  for  $n \geq p$  and  $\mathbf{Z} \mathbf{Z}^\top \sim \mathcal{W}_n(p, \mathbf{I}_n)$  for  $p > n$ , using Corollary 3.2 gives  $r_{\kappa, \nu} = \sum_{i=1}^{\nu} E[\log s_i]$ , where  $s_i \sim \chi_{\kappa-i+1}^2$  for  $i = 1, \dots, \nu$ . Denoting the digamma function by

$$F(t) = \frac{d}{dt} \log \Gamma(t) = \frac{\Gamma'(t)}{\Gamma(t)},$$

we observe  $E[\log s_i] = F((\kappa - i + 1)/2) + \log 2$ , so that

$$r_{\kappa, \nu} = \sum_{i=1}^{\nu} F\left(\frac{\kappa - i + 1}{2}\right) + \nu \log 2.$$

Hence  $r_{\kappa, \nu}$  is a constant and  $\widehat{\Sigma}^{BS}$  has a constant risk.

## 7.3 Triangular Invariant Estimators and the James-Stein Estimator

### 7.3.1 The James-Stein Estimator

As seen in Sect. 3.4, the Cholesky decomposition of  $S$  is written as

$$S = T T^\top = \begin{pmatrix} T_1 \\ T_2 \end{pmatrix} (T_1^\top, T_2^\top),$$

where  $T = (T_1^\top, T_2^\top)^\top \in \mathbb{L}_{p,v}^{(+)}$ ,  $T_1 \in \mathbb{L}_p^{(+)}$  and  $T_2 \in \mathbb{R}^{(p-v) \times v}$ . Define a class of estimators as

$$\widehat{\Sigma}^T = \widehat{\Sigma}^T(S) = T D_v T^\top, \quad (7.6)$$

where  $D_v = \text{diag}(d_1, \dots, d_v)$  and the  $d_i$ 's are positive constants. The unbiased estimator  $\widehat{\Sigma}^{UB}$  and the best scale invariant estimator  $\widehat{\Sigma}^{BS}$  are members of the class (7.6).

The class (7.6) is invariant under the scale transformation with respect to the lower triangular group  $\mathbb{L}_p^{(+)}$ . Indeed, this can be verified as follows: Denote by  $T_* T_*^\top$  the Cholesky decomposition of  $LSL^\top$  for  $L \in \mathbb{L}_p^{(+)}$ . It is observed that  $LSL^\top = LTT^\top L^\top$  and  $LT \in \mathbb{L}_{p,v}^{(+)}$ . As discussed in the beginning of Sect. 3.4, the Cholesky decomposition of a symmetric positive semi-definite matrix  $S$  is unique. Thus, we obtain  $T_* = LT$ , so that

$$\widehat{\Sigma}^T(LSL^\top) = T_* D_v T_*^\top = LT D_v T^\top L^\top = L \widehat{\Sigma}^T(S) L^\top.$$

Here  $\widehat{\Sigma}^T$  is named triangular invariant estimator. We will now investigate the risk function of  $\widehat{\Sigma}^T$ . The Cholesky decomposition of  $\Sigma$  is expressed as  $\Sigma = \Xi \Xi^\top$ , where  $\Xi \in \mathbb{L}_p^{(+)}$ . Let  $U = \Xi^{-1}T$ . Then

$$\begin{aligned} E[\text{tr}[\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^T)]] &= E[\text{tr} D_v T^\top \Sigma^{-1} T] \\ &= E[\text{tr} D_v (\Xi^{-1} T)^\top \Xi^{-1} T] \\ &= E[\text{tr} D_v U^\top U]. \end{aligned} \quad (7.7)$$

The distributions of nonzero elements of  $U$  are given in Corollary 3.2. In the  $p > n$  case, we partition  $U$  as  $U = (u_{ij}) = (U_1^\top, U_2^\top)^\top$ , where  $U_1 \in \mathbb{L}_n^{(+)}$ . Since  $U_2 \sim \mathcal{N}_{(p-n) \times n}(\mathbf{0}_{(p-n) \times n}, I_{(p-n)} \otimes I_n)$ , Corollary 3.1 leads to  $E[U_2^\top U_2] = (p - n)I_n$ . For  $i = 1, \dots, n$ , the  $i$ -th diagonal element of  $E[U_1^\top U_1]$  is

$$\begin{aligned} \sum_{j=i}^n E[u_{ji}^2] &= E[u_{ii}^2] + \sum_{j>i}^n E[u_{ji}^2] \\ &= (n - i + 1) + (n - i) = 2n - 2i + 1, \end{aligned}$$

so that

$$\begin{aligned}
 E[\operatorname{tr} \mathbf{D}_n \mathbf{U}^\top \mathbf{U}] &= \operatorname{tr} \mathbf{D}_n E[\mathbf{U}_1^\top \mathbf{U}_1] + \operatorname{tr} \mathbf{D}_n E[\mathbf{U}_2^\top \mathbf{U}_2] \\
 &= \sum_{i=1}^n \{(2n - 2i + 1)d_i + (p - n)d_i\} \\
 &= \sum_{i=1}^n (n + p - 2i + 1)d_i.
 \end{aligned} \tag{7.8}$$

When  $n \geq p$ , it follows that

$$E[\operatorname{tr} \mathbf{D}_p \mathbf{U}^\top \mathbf{U}] = \sum_{i=1}^p \sum_{j=i}^p E[d_i u_{ji}^2] = \sum_{i=1}^p (n + p - 2i + 1)d_i. \tag{7.9}$$

Combining (7.7), (7.8) and (7.9) gives

$$E[\operatorname{tr} [\operatorname{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^T)]] = \sum_{i=1}^v (n + p - 2i + 1)d_i. \tag{7.10}$$

It is seen that  $\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^T$  has the same positive eigenvalues as  $\mathbf{D}_v \mathbf{T}^\top \boldsymbol{\Sigma}^{-1} \mathbf{T}$ , implying that

$$|\operatorname{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^T)| = |\mathbf{D}_v \mathbf{T}^\top \boldsymbol{\Sigma}^{-1} \mathbf{T}|.$$

Since  $\mathbf{T}^\top \boldsymbol{\Sigma}^{-1} \mathbf{T} \in \mathbb{S}_v^{(+)}$ , it follows that

$$E[\log |\operatorname{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^T)|] = \log |\mathbf{D}_v| + E[\log |\mathbf{T}^\top \boldsymbol{\Sigma}^{-1} \mathbf{T}|] = \sum_{i=1}^v \log d_i + r_{\kappa, v}, \tag{7.11}$$

where  $r_{\kappa, v}$  is the same as in (7.5). Using (7.10) and (7.11), we can write the risk of  $\widehat{\boldsymbol{\Sigma}}^T$  under the extended Stein loss (7.3) as

$$R_{ES}(\widehat{\boldsymbol{\Sigma}}^T, \boldsymbol{\Sigma}) = \sum_{i=1}^v \{(n + p - 2i + 1)d_i - \log d_i\} - r_{\kappa, v} - v.$$

Hence the triangular invariant estimator  $\widehat{\boldsymbol{\Sigma}}^T$  has a constant risk.

Clearly, the  $d_i$ 's minimizing the risk  $R_{ES}(\widehat{\boldsymbol{\Sigma}}^T, \boldsymbol{\Sigma})$  are given by

$$d_i^{JS} = \frac{1}{n + p - 2i + 1}$$

for  $i = 1, \dots, v$ . Thus the best triangular invariant estimator can be expressed as

$$\widehat{\Sigma}^{JS} = \mathbf{T} \mathbf{D}_v^{JS} \mathbf{T}^\top, \quad (7.12)$$

where  $\mathbf{D}_v^{JS} = \text{diag}(d_1^{JS}, \dots, d_v^{JS})$ , which is named the James-Stein (1961) estimator. Since  $\widehat{\Sigma}^{BS}$  belongs to the class (7.6),  $\widehat{\Sigma}^{JS}$  dominates  $\widehat{\Sigma}^{BS}$  relative to the extended Stein loss (7.3). In fact,  $\widehat{\Sigma}^{JS}$  has the constant risk

$$R_{ES}(\widehat{\Sigma}^{JS}, \Sigma) = \sum_{i=1}^v \log(n + p - 2i + 1) - r_{\kappa, v}, \quad (7.13)$$

which implies by (7.5) that

$$R_{ES}(\widehat{\Sigma}^{JS}, \Sigma) - R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) = \sum_{i=1}^v \log(n + p - 2i + 1) - v \log \kappa < 0,$$

where the inequality follows immediately from concavity of the logarithmic function.

The abovementioned can be summarized as follows.

**Proposition 7.2** *The James-Stein estimator  $\widehat{\Sigma}^{JS}$ , namely, the best triangular invariant estimator dominates the best scale invariant estimator  $\widehat{\Sigma}^{BS}$  relative to the extended Stein loss (7.3).*

As pointed out by Stein (1956) and James and Stein (1961),  $\widehat{\Sigma}^{JS}$  is minimax when  $n \geq p$ . The proof of minimaxity comes from the invariance approach. A general theory of the invariance approach was studied in Kiefer (1957). For proving minimaxity of a specific estimator, the least favorable prior approach are also well known (Strawderman, 2000). See Tsukuma and Kubokawa (2015) for the minimaxity proof of  $\widehat{\Sigma}^{JS}$  by using the least favorable prior approach.

### 7.3.2 Improvement Using a Subgroup Invariance

In the literature, various estimators have been proposed for improving the James-Stein estimator  $\widehat{\Sigma}^{JS}$  in (7.12). Here we introduce an invariant estimator under the commutator subgroup of  $\mathbb{L}_p^{(+)}$ .

For two elements  $\mathbf{A}$  and  $\mathbf{B}$  of the group  $\mathbb{L}_p^{(+)}$ , the commutator of  $\mathbf{A}$  and  $\mathbf{B}$  is defined by  $\mathbf{A}^{-1} \mathbf{B}^{-1} \mathbf{A} \mathbf{B}$ . The commutator subgroup of  $\mathbb{L}_p^{(+)}$  is generated by all the commutators of  $\mathbb{L}_p^{(+)}$  and coincides with  $\mathbb{L}_p^{(1)}$ , where  $\mathbb{L}_p^{(1)}$  consists of all  $p \times p$  lower triangular matrices with ones on the diagonal.

Let  $\mathbf{S} = \mathbf{T}_1 \mathbf{T}_0 \mathbf{T}_1^\top$ , where  $\mathbf{T}_0$  and  $\mathbf{T}_1$  are, respectively, unique elements of  $\mathbb{D}_v$  and  $\mathbb{L}_{p,v}^{(1)}$ . Note that, when  $n \geq p$ ,  $\mathbf{T}_1 \mathbf{T}_0 \mathbf{T}_1^\top$  is the LDL<sup>T</sup> decomposition of  $\mathbf{S}$ . Here we define a class of estimators as

$$\widehat{\Sigma}^I = \widehat{\Sigma}^I(\mathbf{T}_0, \mathbf{T}_1) = \mathbf{T}_1 \Phi(\mathbf{T}_0) \mathbf{T}_1^\top, \quad (7.14)$$



where  $\Phi(T_0) \in \mathbb{D}_v$  and each diagonal element of  $\Phi(T_0)$  is an absolutely continuous function of  $T_0$ . The class (7.14) is invariant under the transformations  $S \rightarrow ASA^\top$  and  $\hat{\Sigma} \rightarrow A\hat{\Sigma}A^\top$  for any  $A \in \mathbb{L}_p^{(1)}$ .

The risk of the class (7.14) can be expressed as follows.

**Theorem 7.1** Denote  $T_0 = \text{diag}(t_1, \dots, t_v)$  and  $\Phi(T_0) = \text{diag}(\phi_1, \dots, \phi_v)$ . Then the risk function of  $\hat{\Sigma}^I$  with respect to the extended Stein loss (7.3) is expressed by

$$R_{ES}(\hat{\Sigma}^I, \Sigma) = E \left[ \sum_{i=1}^v \left\{ (n + p - 2i - 1) \frac{\phi_i}{t_i} + 2 \frac{\partial \phi_i}{\partial t_i} - \log \frac{\phi_i}{t_i} \right\} \right] - r_{\kappa, v} - v,$$

where  $r_{\kappa, v}$  is given by (7.5).

**Proof** It is observed that

$$\begin{aligned} \log |\text{Ch}(\Sigma^{-1} \hat{\Sigma}^I)| &= \log |\text{Ch}(T_1^\top \Sigma^{-1} T_1 T_0 T_0^{-1} \Phi(T_0))| \\ &= \log |\text{Ch}(T_1^\top \Sigma^{-1} T_1 T_0)| + \log |\text{Ch}(T_0^{-1} \Phi(T_0))| \\ &= \log |\text{Ch}(\Sigma^{-1} S)| + \sum_{i=1}^v \log \frac{\phi_i}{t_i}, \end{aligned}$$

so that

$$E[\log |\text{Ch}(\Sigma^{-1} \hat{\Sigma}^I)|] = E \left[ \sum_{i=1}^v \log \frac{\phi_i}{t_i} \right] + r_{\kappa, v}.$$

Thus,

$$\begin{aligned} R_{ES}(\hat{\Sigma}^I, \Sigma) &= E[\text{tr} \Sigma^{-1} T_1 \Phi(T_0) T_1^\top - \log |\text{Ch}(\Sigma^{-1} T_1 \Phi(T_0) T_1^\top)| - v] \\ &= E \left[ \text{tr} \Sigma^{-1} T_1 \Phi(T_0) T_1^\top - \sum_{i=1}^v \log \frac{\phi_i}{t_i} \right] - r_{\kappa, v} - v. \end{aligned} \quad (7.15)$$

Denote by  $\Sigma^{-1} = \Gamma^\top \Sigma_0^{-1} \Gamma$  the LDL<sup>⊤</sup> decomposition of  $\Sigma^{-1}$ , where  $\Sigma_0 = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$  and  $\Gamma$  are, respectively, unique elements of  $\mathbb{D}_p$  and  $\mathbb{L}_p^{(1)}$ . Making the transformation  $U = (u_{ij}) = \Gamma T_1$  ( $\in \mathbb{L}_{p, v}^{(1)}$ ) yields

$$\begin{aligned} E[\text{tr} \Sigma^{-1} \hat{\Sigma}^I] &= E[\text{tr} \Sigma_0^{-1} U \Phi(T_0) U^\top] = E \left[ \sum_{i=1}^v \{U^\top \Sigma_0^{-1} U\}_{ii} \phi_i \right] \\ &= E \left[ \sum_{i=1}^v \phi_i \sum_{j=i}^p \frac{u_{ji}^2}{\sigma_j^2} \right]. \end{aligned}$$

Using Proposition 3.10 with some manipulation, we can see that  $t_i \sim \sigma_i^2 \chi_{n-i+1}^2$  for  $i = 1, \dots, v$  and  $u_{ji}|t_i \sim \mathcal{N}(0, (t_i/\sigma_j^2)^{-1})$  for  $j > i$ . Noting that  $\mathbf{U} \in \mathbb{L}_{p,v}^{(1)}$ , namely,  $u_{ii} = 1$ , we obtain

$$E[\text{tr } \Sigma^{-1} \widehat{\Sigma}^I] = E\left[\sum_{i=1}^v \phi_i (1/\sigma_i^2 + (p-i)/t_i)\right],$$

which implies by the chi-square identity (5.4) that

$$E[\text{tr } \Sigma^{-1} \widehat{\Sigma}^I] = E\left[\sum_{i=1}^v \left\{(n+p-2i-1)\frac{\phi_i}{t_i} + 2\frac{\partial \phi_i}{\partial t_i}\right\}\right]. \quad (7.16)$$

Hence combining (7.15) and (7.16) completes the proof.  $\square$

The James-Stein estimator  $\widehat{\Sigma}^{JS}$  belongs to the class (7.14). Using Theorem 7.1 with  $\phi_i^{JS} = d_i^{JS} t_i$ , we can obtain the same expression for risk of  $\widehat{\Sigma}^{JS}$  as in (7.13).

Next, we provide an improved estimator on  $\widehat{\Sigma}^{JS}$ . Let  $\Phi^M(\mathbf{T}_0) = \text{diag}(\phi_1^M, \dots, \phi_v^M)$  with

$$\phi_i^M = d_i^{JS} t_i - \frac{(t_i \log t_i) g(w)}{b + w}, \quad w = \sum_{i=1}^v (\log t_i)^2,$$

where  $b$  is a suitable constant and  $g(w)$  is a differentiable function of  $w$ .

**Proposition 7.3** Suppose  $v \geq 3$  and  $b \geq 144(v-2)^2/\{25(n+p-1)^2\}$ . If  $g(w)$  is nondecreasing in  $w$  and  $0 < g(w) \leq 12(v-2)/\{5(n+p-1)^2\}$ , then

$$\widehat{\Sigma}^M = \widehat{\Sigma}^M(\mathbf{T}_0, \mathbf{T}_1) = \mathbf{T}_1 \Phi^M(\mathbf{T}_0) \mathbf{T}_1^\top$$

dominates  $\widehat{\Sigma}^{JS}$  relative to the extended Stein loss (7.3).

**Proof** This dominance result is proved along the same arguments as in Dey and Srinivasan (1985). For details, see Tsukuma (2014a).  $\square$

Proposition 7.3 implies that the best invariant estimator  $\widehat{\Sigma}^{JS}$  under  $\mathbb{L}_p^{(+)}$  is dominated by the invariant estimator  $\widehat{\Sigma}^M$  under the commutator subgroup of  $\mathbb{L}_p^{(+)}$ , namely, under  $\mathbb{L}_p^{(1)}$ . Since the lower triangular group  $\mathbb{L}_p^{(1)}$  is solvable,  $\mathbb{L}_p^{(1)}$  also has a commutator subgroup. It is still not known whether there exists an invariant estimator under the commutator subgroup of  $\mathbb{L}_p^{(1)}$  which dominates  $\widehat{\Sigma}^M$ .

## 7.4 Orthogonally Invariant Estimators

### 7.4.1 Class of Orthogonally Invariant Estimators

The triangular invariant estimator  $\widehat{\Sigma}^T$  in (7.6) depends on the coordinate system. This fact causes inadmissibility of the James-Stein estimator  $\widehat{\Sigma}^{JS}$  in (7.12). For example, the inadmissibility can be shown by using the same arguments as in Stein (1956): Let  $\mathbf{P}$  be a  $p \times p$  unit anti-diagonal matrix of the form

$$\mathbf{P} = \begin{pmatrix} 0 & & 1 \\ & \ddots & \\ 1 & & 0 \end{pmatrix},$$

which is symmetric and orthogonal. Denote  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^\top \in \mathbb{R}^{n \times p}$ . Note that, for each  $i = 1, \dots, n$ ,  $\mathbf{P}\mathbf{y}_i$  is a  $p$ -dimensional vector obtained from  $\mathbf{y}_i$  by reversing the order of coordinates. Define the Cholesky decomposition of  $\mathbf{PSP} = \mathbf{P}\mathbf{Y}^\top \mathbf{Y} \mathbf{P}$  as  $\mathbf{T}_* \mathbf{T}_*^\top$ , where  $\mathbf{T}_* \in \mathbb{L}_{p,v}^{(+)}$  and  $\mathbf{T}_* \mathbf{T}_*^\top$  is not the same as the Cholesky decomposition of  $\mathbf{S}$  because of its uniqueness. Let  $\widehat{\Sigma}^U = \mathbf{T}_* \mathbf{D}_v^{JS} \mathbf{T}_*$ , which is the best triangular invariant estimator of  $\mathbf{PSP}$ . Here  $\mathbf{P}\widehat{\Sigma}^U \mathbf{P}$  becomes an estimator of  $\Sigma$  and has the same risk as  $\widehat{\Sigma}^{JS}$ . From the convexity of the extended Stein loss (7.3), it is easily proved that  $\widehat{\Sigma}^{JS}$  is dominated by a combination estimator  $(\widehat{\Sigma}^{JS} + \mathbf{P}\widehat{\Sigma}^U \mathbf{P})/2$ .

In this section, we will consider a general class of estimators not depending on the coordinate system and aim to find better estimators dominating  $\widehat{\Sigma}^{JS}$  and  $\widehat{\Sigma}^{BS}$  relative to the extended Stein loss (7.3). Write the eigenvalue decomposition of  $\mathbf{S}$  as

$$\mathbf{S} = \mathbf{H} \mathbf{L} \mathbf{H}^\top,$$

where  $\mathbf{L} = \text{diag}(\ell_1, \dots, \ell_v) \in \mathbb{D}_v^{(\geq 0)}$  and  $\mathbf{H} \in \mathbb{V}_{p,v}$ . The general class of estimators is defined by

$$\widehat{\Sigma}^O = \widehat{\Sigma}^O(\mathbf{S}) = \mathbf{H} \text{diag}(\ell_1 \phi_1(\mathbf{L}), \dots, \ell_v \phi_v(\mathbf{L})) \mathbf{H}^\top = \mathbf{H} \mathbf{L} \Phi(\mathbf{L}) \mathbf{H}^\top,$$

where  $\Phi(\mathbf{L}) = \text{diag}(\phi_1(\mathbf{L}), \dots, \phi_v(\mathbf{L}))$  and the  $\phi_i(\mathbf{L})$ 's are absolutely continuous functions of  $\mathbf{L}$ . The class  $\widehat{\Sigma}^O$  is not only invariant with respect to exchanging coordinates, but also, more generally, orthogonally invariant in the sense that it satisfies  $\mathbf{O} \widehat{\Sigma}^O(\mathbf{S}) \mathbf{O}^\top = \widehat{\Sigma}^O(\mathbf{O} \mathbf{S} \mathbf{O}^\top)$  for any  $\mathbf{O} \in \mathbb{O}_p$ .

### 7.4.2 Unbiased Risk Estimate

Here, we will derive an unbiased risk estimate for orthogonally invariant estimators  $\widehat{\Sigma}^O$ . Abbreviate  $\phi_i(\mathbf{L})$  by  $\phi_i$ . Since

$$\mathbf{H} = \mathbf{H}\mathbf{H}^\top \mathbf{H} = \mathbf{S}\mathbf{S}^+ \mathbf{H} = \mathbf{Y}^\top \mathbf{Y} \mathbf{S}^+ \mathbf{H}$$

with  $\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mathbf{0}_{n \times p}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ , it follows from Theorem 5.1 and Lemma 5.3 that

$$\begin{aligned} E[\text{tr}[\text{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^O)]] &= E[\text{tr} \boldsymbol{\Sigma}^{-1} \mathbf{H} \mathbf{L} \boldsymbol{\Phi}(\mathbf{L}) \mathbf{H}^\top] \\ &= E[\text{tr} \nabla_Y^\top \mathbf{Y} \mathbf{S}^+ \mathbf{H} \mathbf{L} \boldsymbol{\Phi}(\mathbf{L}) \mathbf{H}^\top] \\ &= E \left[ \sum_{i=1}^v \left\{ (|n-p|+1)\phi_i + 2\ell_i \frac{\partial \phi_i}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i - \ell_j \phi_j}{\ell_i - \ell_j} \right\} \right]. \end{aligned} \quad (7.17)$$

Note that  $|\text{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^O)| = |\mathbf{H}^\top \boldsymbol{\Sigma}^{-1} \mathbf{H} \mathbf{L}| \cdot |\boldsymbol{\Phi}(\mathbf{L})| = |\text{Ch}(\boldsymbol{\Sigma}^{-1} \mathbf{S})| \prod_{i=1}^v \phi_i$ . The risk of  $\widehat{\boldsymbol{\Sigma}}^O$  becomes

$$\begin{aligned} R_{ES}(\widehat{\boldsymbol{\Sigma}}^O, \boldsymbol{\Sigma}) &= E[\text{tr} \boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^O - \log |\text{Ch}(\boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\Sigma}}^O)| - v] \\ &= E \left[ \sum_{i=1}^v \left\{ (|n-p|+1)\phi_i + 2\ell_i \frac{\partial \phi_i}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i - \ell_j \phi_j}{\ell_i - \ell_j} - \log \phi_i \right\} \right] \\ &\quad - r_{K,v} - v. \end{aligned}$$

Hence, we obtain the unbiased risk estimate for  $\widehat{\boldsymbol{\Sigma}}^O$ ,

$$\widehat{R}_{ES}(\widehat{\boldsymbol{\Sigma}}^O) = \sum_{i=1}^v \left\{ (|n-p|+1)\phi_i + 2\ell_i \frac{\partial \phi_i}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i - \ell_j \phi_j}{\ell_i - \ell_j} - \log \phi_i \right\} - r_{K,v} - v. \quad (7.18)$$

Comparing (7.5), or (7.13), with (7.18) gives a sufficient condition that  $\widehat{\boldsymbol{\Sigma}}^O$  dominates  $\widehat{\boldsymbol{\Sigma}}^{BS}$ , or  $\widehat{\boldsymbol{\Sigma}}^{JS}$ , relative to the extended Stein loss (7.3). For example, we denote the difference between  $\widehat{R}_{ES}(\widehat{\boldsymbol{\Sigma}}^O)$  and (7.13) by

$$\begin{aligned} \widehat{\Delta}(\widehat{\boldsymbol{\Sigma}}^O) &= \widehat{R}_{ES}(\widehat{\boldsymbol{\Sigma}}^O) - R_{ES}(\widehat{\boldsymbol{\Sigma}}^{JS}, \boldsymbol{\Sigma}) \\ &= \sum_{i=1}^v \left\{ (|n-p|+1)\phi_i + 2\ell_i \frac{\partial \phi_i}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i - \ell_j \phi_j}{\ell_i - \ell_j} - \log \phi_i \right\} \\ &\quad + \sum_{i=1}^v \log d_i^{JS} - v, \end{aligned} \quad (7.19)$$

implying that if  $\widehat{\Delta}(\widehat{\boldsymbol{\Sigma}}^O) \leq 0$  for every  $\mathbf{L} \in \mathbb{D}_v^{(\geq 0)}$  then  $\widehat{\boldsymbol{\Sigma}}^O$  dominates  $\widehat{\boldsymbol{\Sigma}}^{JS}$ .

### 7.4.3 Examples

#### 7.4.3.1 Haff's Empirical Bayes Estimator

When  $n \geq p$ , Haff (1980) considered the empirical Bayes estimation of  $\Sigma$ . For the sake of simplicity, let a prior distribution of  $\Sigma^{-1}$  be  $\mathcal{W}_p(p+1, \gamma^{-1} \mathbf{I}_p)$ , where  $\gamma$  is an unknown hyperparameter. The resulting posterior distribution of  $\Sigma^{-1}$  given  $\mathbf{S}$  is  $\Sigma^{-1} | \mathbf{S} \sim \mathcal{W}_p(n+p+1, (\mathbf{S} + \gamma \mathbf{I}_p)^{-1})$ , so that the posterior mean of  $\Sigma$  is  $E[\Sigma | \mathbf{S}] = n^{-1}(\mathbf{S} + \gamma \mathbf{I}_p)$ , where this expectation will be explained in Sect. 7.6.4. The hyperparameter  $\gamma$  is estimated from the marginal density of  $\mathbf{S}$  proportional to

$$\gamma^{p(p+1)/2} |\mathbf{S}|^{(n-p-1)/2} |\mathbf{S} + \gamma \mathbf{I}_p|^{-(n+p+1)/2}.$$

The log-likelihood function with ignoring a constant has the form

$$l(\gamma | \mathbf{S}) = \frac{p(p+1)}{2} \log \gamma - \frac{n+p+1}{2} \log |\mathbf{I}_p + \gamma \mathbf{S}^{-1}|.$$

The first order approximation of  $\log |\mathbf{I}_p + \gamma \mathbf{S}^{-1}|$  is  $\gamma \operatorname{tr} \mathbf{S}^{-1}$ , so that the log-likelihood function can be approximated as

$$l(\gamma | \mathbf{S}) \approx \frac{p(p+1)}{2} \log \gamma - \frac{n+p+1}{2} \gamma \operatorname{tr} \mathbf{S}^{-1},$$

which attains a maximum at

$$\hat{\gamma} = \frac{p(p+1)}{n+p+1} \frac{1}{\operatorname{tr} \mathbf{S}^{-1}}.$$

The estimate  $\hat{\gamma}$  is an approximated maximum likelihood estimate for  $\gamma$ . Thus we obtain an empirical Bayes estimator of the form

$$\hat{\Sigma}^{EB} = \frac{1}{n}(\mathbf{S} + \hat{\gamma} \mathbf{I}_p) = \frac{1}{n} \left( \mathbf{S} + \frac{p(p+1)}{n+p+1} \frac{1}{\operatorname{tr} \mathbf{S}^{-1}} \mathbf{I}_p \right).$$

Haff (1980) defined a general class of empirical Bayes estimators and gave some dominance results.

Here, taking into account both the cases of  $n \geq p$  and  $p > n$ , we define a Haff type estimator as

$$\hat{\Sigma}^{HF} = \hat{\Sigma}^{BS} + \frac{a}{\kappa \operatorname{tr} \mathbf{S}^+} \mathbf{S} \mathbf{S}^+,$$

where  $a$  is a positive constant. Since  $\operatorname{tr} \mathbf{S}^+ = \operatorname{tr} \mathbf{L}^{-1}$  and  $\mathbf{S} \mathbf{S}^+ = \mathbf{H} \mathbf{H}^\top$  for  $\mathbf{S} = \mathbf{H} \mathbf{L} \mathbf{H}^\top$ , the Haff estimator can be expressed as

$$\widehat{\Sigma}^{HF} = \mathbf{H} \mathbf{L} \Phi^{HF}(\mathbf{L}) \mathbf{H}^\top, \quad \Phi^{HF}(\mathbf{L}) = \frac{1}{\kappa} \left( \mathbf{I}_v + \frac{a}{\text{tr } \mathbf{L}^{-1}} \mathbf{L}^{-1} \right).$$

**Proposition 7.4** *If  $0 < a \leq 2(v-1)/(|n-p|+1)$ , then  $\widehat{\Sigma}^{HF}$  dominates  $\widehat{\Sigma}^{BS}$  relative to the extended Stein loss (7.3).*

**Proof** For  $i = 1, \dots, v$ , let  $\phi_i^{HF} = \kappa^{-1}(1 + a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1})$ . Note that

$$\sum_{i=1}^v \ell_i \frac{\partial \phi_i^{HF}}{\partial \ell_i} = \frac{a}{\kappa} \left( -1 + \frac{\text{tr } \mathbf{L}^{-2}}{(\text{tr } \mathbf{L}^{-1})^2} \right) \leq 0,$$

so that, by (7.18),

$$\begin{aligned} \widehat{R}_{ES}(\widehat{\Sigma}^{HF}) &= \sum_{i=1}^v \left\{ (|n-p|+1) \phi_i^{HF} + 2\ell_i \frac{\partial \phi_i^{HF}}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i^{HF} - \ell_j \phi_j^{HF}}{\ell_i - \ell_j} \right. \\ &\quad \left. - \log \phi_i^{HF} \right\} - r_{\kappa, v} - v \\ &\leq \widehat{R}_{ES}(\widehat{\Sigma}^{BS}) + (|n-p|+1) \frac{a}{\kappa} - \sum_{i=1}^v \log(1 + a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1}). \end{aligned}$$

Since  $\log(1+x) \geq 2x/(2+x)$  for  $x \geq 0$ , it holds that

$$\begin{aligned} \sum_{i=1}^v \log(1 + a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1}) &\geq \sum_{i=1}^v \frac{2a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1}}{2 + a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1}} \\ &\geq \sum_{i=1}^v \frac{2a\ell_i^{-1}/\text{tr } \mathbf{L}^{-1}}{2+a} = \frac{2a}{2+a}, \end{aligned}$$

which yields

$$\widehat{R}_{ES}(\widehat{\Sigma}^{HF}) - \widehat{R}_{ES}(\widehat{\Sigma}^{BS}) \leq (|n-p|+1) \frac{a}{\kappa} - \frac{2a}{2+a} = \frac{a}{2+a} \{(|n-p|+1)a - 2(v-1)\}.$$

Hence we complete the proof.  $\square$

#### 7.4.3.2 The Efron-Morris-Dey Shrinkage Estimator

Since  $\widehat{\Sigma}^{HF} \succeq \widehat{\Sigma}^{BS}$ , the Haff estimator  $\widehat{\Sigma}^{HF}$  is an expansion estimator in the Löwner sense. Next, we give an improved estimator shrinking  $\widehat{\Sigma}^{BS}$ .

For positive constants  $b$  and  $c$ , define

$$\widehat{\Sigma}^{SH} = \mathbf{H} \mathbf{L} \Phi^{SH}(\mathbf{L}) \mathbf{H}^\top, \quad \Phi^{SH}(\mathbf{L}) = \frac{1}{\kappa} \left( \mathbf{I}_v + \frac{b}{\text{tr } \mathbf{L}^c} \mathbf{L}^c \right)^{-1},$$

where  $\mathbf{L}^c = \text{diag}(\ell_1^c, \dots, \ell_v^c)$ . This estimator is inspired from Efron and Morris (1976) and Dey (1987) for estimating  $\Sigma^{-1}$  under certain quadratic losses. Indeed, when  $n \geq p$ , it corresponds to Efron and Morris (1976) for  $c = 1$  and to Dey (1987) for  $c = 2$ . For  $b > 0$ , it holds that  $\Phi^{SH}(\mathbf{L}) \preceq \kappa^{-1} \mathbf{I}_v$ , so that  $\widehat{\Sigma}^{SH} \preceq \widehat{\Sigma}^{BS}$  in the Löwner sense. We can obtain a dominance result on the shrinkage estimator  $\widehat{\Sigma}^{SH}$  as follows.

**Proposition 7.5** *For  $0 < b \leq (v - 1)/\kappa$  and  $c \geq 1$ ,  $\widehat{\Sigma}^{SH}$  dominates  $\widehat{\Sigma}^{BS}$  relative to the extended Stein loss (7.3).*

**Proof** Let  $\phi_i^{SH} = \kappa^{-1}(1 + b\ell_i^c / \text{tr } \mathbf{L}^c)^{-1}$  and  $\phi_i^{SH*} = b\ell_i^c \{(1 + b)\kappa \text{tr } \mathbf{L}^c\}^{-1}$  for  $i = 1, \dots, v$ . Note that  $\Phi^{SH}(\mathbf{L}) = \text{diag}(\phi_1^{SH}, \dots, \phi_v^{SH})$  and, for  $i = 1, \dots, v$ ,

$$\begin{aligned} \phi_i^{SH} &= \kappa^{-1} - \frac{b\ell_i^c}{\kappa \text{tr } \mathbf{L}^c} \left(1 + \frac{b\ell_i^c}{\text{tr } \mathbf{L}^c}\right)^{-1} \\ &\leq \kappa^{-1} - \frac{b\ell_i^c}{\kappa \text{tr } \mathbf{L}^c} (1 + b)^{-1} = \kappa^{-1} - \phi_i^{SH*}. \end{aligned}$$

Thus,

$$E[\text{tr}[\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{SH})]] \leq E[\text{tr} \Sigma^{-1} \widehat{\Sigma}^{BS}] - E[\text{tr} \Sigma^{-1} \mathbf{H} \mathbf{L} \Phi^{SH*} \mathbf{H}^\top],$$

where  $\Phi^{SH*} = \text{diag}(\phi_1^{SH*}, \dots, \phi_v^{SH*})$ . It follows that

$$\begin{aligned} \log |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{SH})| &= \log |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{BS})| - \log \left| \mathbf{I}_v + \frac{b}{\text{tr } \mathbf{L}^c} \mathbf{L}^c \right| \\ &\geq \log |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{BS})| - \text{tr} \left( \frac{b}{\text{tr } \mathbf{L}^c} \mathbf{L}^c \right) = \log |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{BS})| - b, \end{aligned}$$

so that

$$R_{ES}(\widehat{\Sigma}^{SH}, \Sigma) - R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) \leq -E[\text{tr} \Sigma^{-1} \mathbf{H} \mathbf{L} \Phi^{SH*} \mathbf{H}^\top] + b. \quad (7.20)$$

Using the identity (7.17) gives

$$\begin{aligned} &E[\text{tr} \Sigma^{-1} \mathbf{H} \mathbf{L} \Phi^{SH*} \mathbf{H}^\top] \\ &= E \left[ \sum_{i=1}^v \left\{ (|n - p| + 1) \phi_i^{SH*} + 2\ell_i \frac{\partial \phi_i^{SH*}}{\partial \ell_i} + 2 \sum_{j>i}^v \frac{\ell_i \phi_i^{SH*} - \ell_j \phi_j^{SH*}}{\ell_i - \ell_j} \right\} \right]. \end{aligned} \quad (7.21)$$

For  $b > 0$  and  $c \geq 1$ ,

$$\sum_{i=1}^v \ell_i \frac{\partial \phi_i^{SH*}}{\partial \ell_i} = \frac{bc}{(1 + b)\kappa} \left( 1 - \frac{\text{tr } \mathbf{L}^{2c}}{(\text{tr } \mathbf{L}^c)^2} \right) \geq 0. \quad (7.22)$$

Since, for  $1 \leq i < j \leq v$  and  $c \geq 1$ ,

$$\frac{\ell_i^{c+1} - \ell_j^{c+1}}{\ell_i - \ell_j} \geq \ell_i^c + \ell_j^c,$$

and  $\sum_{i=1}^v \sum_{j>i}^v (\ell_i^c + \ell_j^c) = (v-1) \operatorname{tr} \mathbf{L}^c$ , we obtain

$$\sum_{i=1}^v \sum_{j>i}^v \frac{\ell_i \phi_i^{SH*} - \ell_j \phi_j^{SH*}}{\ell_i - \ell_j} \geq \frac{b}{(1+b)\kappa \operatorname{tr} \mathbf{L}^c} \sum_{i=1}^v \sum_{j>i}^v (\ell_i^c + \ell_j^c) = \frac{(v-1)b}{(1+b)\kappa}. \quad (7.23)$$

Combining (7.20)–(7.23) provides

$$R_{ES}(\widehat{\Sigma}^{SH}, \Sigma) - R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) \leq -\frac{(n+p-1)b}{(1+b)\kappa} + b = \frac{\kappa b^2 - (v-1)b}{(1+b)\kappa},$$

which is not positive if  $0 < b \leq (v-1)/\kappa$ .  $\square$

### 7.4.3.3 Stein's Simple Estimator and Risk Minimization Method

In the nonsingular case, namely,  $n \geq p$ , the unbiased estimator  $\widehat{\Sigma}^{UB}$  can be expressed as  $\widehat{\Sigma}^{UB} = \mathbf{H} \mathbf{L} \mathbf{D}_p^{UB} \mathbf{H}$  with  $\mathbf{L} = \operatorname{diag}(\ell_1, \dots, \ell_p)$  and  $\mathbf{D}_p^{UB} = \operatorname{diag}(n^{-1}, \dots, n^{-1})$ . Then the first and last diagonal elements of  $\mathbf{L} \mathbf{D}_p^{UB}$  are, respectively,  $\ell_1/n$  and  $\ell_p/n$ , which are the largest and smallest eigenvalues of  $\widehat{\Sigma}^{UB}$ . Let  $\lambda_1$  be the largest eigenvalue of  $\Sigma$  and let  $\eta$  be the corresponding normalized eigenvector. Now,

$$\lambda_1 = \eta^\top \Sigma \eta = E[\eta^\top \widehat{\Sigma}^{UB} \eta] \leq E \left[ \max_{\xi \in \mathbb{R}^p: \|\xi\|=1} \xi^\top \widehat{\Sigma}^{UB} \xi \right] = E[\ell_1/n].$$

Similarly, it can be shown that  $E[\ell_p/n] \leq \lambda_p$ , where  $\lambda_p$  is the smallest eigenvalue of  $\Sigma$ . These imply that  $\ell_1/n$  overestimates  $\lambda_1$  and  $\ell_p/n$  underestimates  $\lambda_p$ , so that  $\widehat{\Sigma}^{UB}$  should probably be modified by shrinking its largest eigenvalue  $\ell_1/n$  and expanding its smallest eigenvalue  $\ell_p/n$ . In other words, an orthogonally invariant estimator  $\widehat{\Sigma}^O = \mathbf{H} \mathbf{L} \Phi(\mathbf{L}) \mathbf{H}^\top$  with  $\Phi(\mathbf{L}) = \operatorname{diag}(\phi_1(\mathbf{L}), \dots, \phi_p(\mathbf{L}))$  should satisfy  $\phi_1(\mathbf{L}) \leq \dots \leq \phi_p(\mathbf{L})$ . Moreover from the eigenvalue decomposition  $\mathbf{S} = \mathbf{H} \mathbf{L} \mathbf{H}^\top$ ,  $\mathbf{L}$  is defined by the diagonal matrix consisting of ordered eigenvalues  $\ell_1 \geq \dots \geq \ell_p$ . Thus, for the orthogonally invariant estimators  $\widehat{\Sigma}^O$ , it seems reasonable that the  $\ell_i \phi_i(\mathbf{L})$ 's are required to have the property  $\ell_1 \phi_1(\mathbf{L}) \geq \dots \geq \ell_p \phi_p(\mathbf{L})$ . As in Perron (1992), the properties  $\phi_1(\mathbf{L}) \leq \dots \leq \phi_p(\mathbf{L})$  and  $\ell_1 \phi_1(\mathbf{L}) \geq \dots \geq \ell_p \phi_p(\mathbf{L})$  shall be called, respectively, **the shrinkage and ordering properties**. Sheena and Takemura (1992) called  $\widehat{\Sigma}^O$  order-preserving when  $\widehat{\Sigma}^O$  has the ordering property.

Here we will provide well known, two orthogonally invariant estimators with the shrinkage and ordering properties given by Stein. For details, see Stein (1975, 1977) and also Dey and Srinivasan (1985).



The first estimator is given by

$$\widehat{\Sigma}^{ST} = \mathbf{H} \mathbf{L} \mathbf{D}_v^{JS} \mathbf{H}^\top, \quad (7.24)$$

where  $\mathbf{D}_v^{JS}$  is given in (7.12). Since  $d_1^{JS} \leq \dots \leq d_v^{JS}$ ,  $\widehat{\Sigma}^{ST}$  has the shrinkage property, but lacks the ordering property. As in the following proposition, the simple estimator  $\widehat{\Sigma}^{ST}$  improves  $\widehat{\Sigma}^{JS}$ .

**Proposition 7.6**  $\widehat{\Sigma}^{ST}$  dominates  $\widehat{\Sigma}^{JS}$  relative to the extended Stein loss (7.3).

**Proof** This proposition can be proved in the same lines as in Dey and Srinivasan (1985, Theorem 3.1). Using (7.19), we can write the difference between the unbiased risk estimate of  $\widehat{\Sigma}^{ST}$  and the constant risk of  $\widehat{\Sigma}^{JS}$  as

$$\begin{aligned} \widehat{\Delta}(\widehat{\Sigma}^{ST}) &= \widehat{R}_{ES}(\widehat{\Sigma}^{ST}) - R_{ES}(\widehat{\Sigma}^{JS}, \Sigma) \\ &= \sum_{i=1}^v \left\{ (|n-p|+1)d_i^{JS} + 2 \sum_{j>i}^v \frac{d_i^{JS}\ell_i - d_j^{JS}\ell_j}{\ell_i - \ell_j} \right\} - v. \end{aligned}$$

It is observed that

$$\begin{aligned} \sum_{i=1}^v \sum_{j>i}^v \frac{d_i^{JS}\ell_i - d_j^{JS}\ell_j}{\ell_i - \ell_j} &= \sum_{i=1}^v \sum_{j>i}^v \frac{d_i^{JS}(\ell_i - \ell_j) + (d_i^{JS} - d_j^{JS})\ell_j}{\ell_i - \ell_j} \\ &< \sum_{i=1}^v \sum_{j>i}^v d_i^{JS} = \sum_{i=1}^v (v-i)d_i^{JS}, \end{aligned}$$

where the inequality is verified by the facts that  $\ell_1 > \dots > \ell_v$  and that  $d_1^{JS} < \dots < d_v^{JS}$ . Thus,

$$\widehat{\Delta}(\widehat{\Sigma}^{ST}) < \sum_{i=1}^v (|n-p|+1+2v-2i)d_i^{JS} - v = \sum_{i=1}^v (n+p-2i+1)d_i^{JS} - v = 0,$$

which completes the proof.  $\square$

The other well-known estimator due to Stein (1975, 1977) is obtained from minimizing the unbiased risk estimate (7.18) subject to  $\phi_i$ 's with ignoring differential terms. The unbiased risk estimate (7.18) with ignoring differential terms can be rewritten as

$$\widehat{R}_*(\widehat{\Sigma}^O) = \sum_{i=1}^v \left\{ (|n-p|+1)\phi_i + 2 \sum_{j \neq i}^v \frac{\ell_i \phi_i}{\ell_i - \ell_j} - \log \phi_i \right\} + \text{const.}$$

which is minimized by, for  $i = 1, \dots, v$ ,

$$\phi_i^{RM} = 1/\omega_i(\mathbf{L}), \quad \omega_i(\mathbf{L}) = |n - p| + 1 + 2\ell_i \sum_{j \neq i}^v \frac{1}{\ell_i - \ell_j}.$$

The  $\phi_i^{RM}$ 's are sometimes negative and not satisfying the ordering property  $\ell_1 \phi_1^{RM} \geq \dots \geq \ell_v \phi_v^{RM}$ . To modify them, Stein (1977) suggested applying the isotonic regression to  $\phi_i^{RM}$ 's. For a detailed example of the isotonic regression, see Lin and Perlman (1985). No exact dominance result exists for the resulting modified estimator with the ordering property, but it has been much used for numerical comparison in the literature including Lin and Perlman (1985), Haff (1991), Yang and Berger (1994) and Ledoit and Wolf (2004).

#### 7.4.3.4 The Dey-Srinivasan Estimator

Next, we introduce an improved estimator, which is based on Dey and Srinivasan (1985). Let  $\Phi^{DS}(\mathbf{L}) = \text{diag}(\phi_1^{DS}, \dots, \phi_v^{DS})$  with

$$\phi_i^{DS} = d_i^{JS} - \frac{g(u) \log \ell_i}{b + u}, \quad u = \sum_{i=1}^v (\log \ell_i)^2,$$

where  $b$  is a suitable constant and  $g(u)$  is a differentiable function of  $u$ . Define

$$\widehat{\Sigma}^{DS} = \mathbf{H} \mathbf{L} \Phi^{DS}(\mathbf{L}) \mathbf{H}^\top.$$

Of course,  $\phi_1^{DS} \leq \dots \leq \phi_v^{DS}$ , so  $\widehat{\Sigma}^{DS}$  has the shrinkage property.

**Proposition 7.7** Suppose  $v \geq 3$  and  $b \geq 144(v-2)^2/\{25(n+p-1)^2\}$ . If  $g(u)$  is nondecreasing in  $u$  and  $0 < g(u) \leq 12(v-2)/\{5(n+p-1)^2\}$ , then  $\widehat{\Sigma}^{DS}$  dominates  $\widehat{\Sigma}^{JS}$  and  $\widehat{\Sigma}^{ST}$  relative to the extended Stein loss (7.3).

**Proof** From a straightforward calculation after substituting  $\phi_i = \phi_i^{DS}$  into (7.18), the unbiased risk estimate of  $\widehat{\Sigma}^{DS}$  can be expressed as

$$\widehat{R}_{ES}(\widehat{\Sigma}^{DS}) = \widehat{R}_{ES}(\widehat{\Sigma}^{ST}) + \widehat{\Delta}^{DS},$$

where

$$\begin{aligned} \widehat{\Delta}^{DS} = \sum_{i=1}^v \left\{ -(|n-p|+1) \frac{g(u) \log \ell_i}{b+u} - 2\ell_i \frac{\partial}{\partial \ell_i} \frac{g(u) \log \ell_i}{b+u} \right. \\ \left. - 2 \frac{g(u)}{b+u} \sum_{j>i}^v \frac{\ell_i \log \ell_i - \ell_j \log \ell_j}{\ell_i - \ell_j} - \log \left( 1 - \frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} \right) \right\}. \end{aligned}$$

If  $\widehat{\Delta}^{DS} \leq 0$  then the proposition is verified.

Note that

$$\sum_{j>i}^v \frac{\ell_i \log \ell_i - \ell_j \log \ell_j}{\ell_i - \ell_j} = (v-i) \log \ell_i + \sum_{j>i}^v \frac{\log \ell_i - \log \ell_j}{\ell_i - \ell_j} \ell_j \geq (v-i) \log \ell_i$$

because  $\ell_1 \geq \dots \geq \ell_v$ . Further noting that

$$\ell_i \frac{\partial}{\partial \ell_i} \frac{g(u) \log \ell_i}{b+u} = \frac{g(u)}{b+u} + 2 \frac{g'(u)(\log \ell_i)^2}{b+u} - 2 \frac{g(u)(\log \ell_i)^2}{(b+u)^2},$$

we obtain

$$\begin{aligned} \widehat{\Delta}^{DS} \leq \sum_{i=1}^v \left\{ -\frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} - 2 \frac{g(u)}{b+u} - 4 \frac{g'(u)(\log \ell_i)^2}{b+u} \right. \\ \left. + 4 \frac{g(u)(\log \ell_i)^2}{(b+u)^2} - \log \left( 1 - \frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} \right) \right\}. \end{aligned} \quad (7.25)$$

It follows that  $|x|/(b+x^2) \leq \{2\sqrt{b}\}^{-1}$  for  $b > 0$ , implying that for each  $i$

$$\frac{|\log \ell_i|}{b+u} < \frac{|\log \ell_i|}{b + (\log \ell_i)^2} \leq \frac{1}{2\sqrt{b}}.$$

Combining this inequality and the given conditions on  $b$  and  $g(u)$  gives

$$\frac{g(u)|\log \ell_i|}{d_i^{JS}(b+u)} \leq (n+p-1) \frac{g(u)|\log \ell_i|}{b+u} < \frac{1}{2\sqrt{b}} \times \frac{12(v-2)}{5(n+p-1)} < \frac{1}{2}.$$

Since

$$\log(1+x) \geq x - \frac{5}{6}x^2$$

for  $|x| \leq 1/2$  (see Dey and Srinivasan 1985, Lemma 2.2), using the given condition on  $g(u)$  yields

$$\begin{aligned} \log \left( 1 - \frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} \right) &\geq -\frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} - \frac{5}{6} \left\{ \frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} \right\}^2 \\ &> -\frac{g(u) \log \ell_i}{d_i^{JS}(b+u)} - 2(v-2) \frac{g(u)(\log \ell_i)^2}{(b+u)^2}. \end{aligned} \quad (7.26)$$

Combining (7.25) and (7.26) gives

$$\begin{aligned}\widehat{\Delta}^{DS} &< \sum_{i=1}^v \left\{ -2 \frac{g(u)}{b+u} - 4 \frac{g'(u)(\log \ell_i)^2}{b+u} + 2v \frac{g(u)(\log \ell_i)^2}{(b+u)^2} \right\} \\ &= -2v \frac{g(u)}{b+u} - 4 \frac{ug'(u)}{b+u} + 2v \frac{ug(u)}{(b+u)^2} < -4 \frac{ug'(u)}{b+u} \leq 0,\end{aligned}$$

which completes the proof.  $\square$

### 7.4.3.5 Sheena and Takemura's Methods for Improving Non-order-preserving Estimators

Let  $\boldsymbol{\varphi} = (\varphi_1, \dots, \varphi_v)$  and  $\text{diag}(\boldsymbol{\varphi}) = \text{diag}(\varphi_1, \dots, \varphi_v)$ , where  $\varphi_i = \ell_i \phi_i(\mathbf{L})$  for  $i = 1, \dots, v$ . Here as in Sheena and Takemura (1992), we call  $\widehat{\boldsymbol{\Sigma}}^O = \mathbf{H} \text{diag}(\boldsymbol{\varphi}) \mathbf{H}^\top$  order-preserving if  $\widehat{\boldsymbol{\Sigma}}^O$  has the ordering property  $\varphi_1 \geq \dots \geq \varphi_v$ . Note that  $\widehat{\boldsymbol{\Sigma}}^{ST}$  and  $\widehat{\boldsymbol{\Sigma}}^{DS}$  are not always order-preserving since  $\ell_1 \geq \dots \geq \ell_v$ ,  $d_1^{JS} \leq \dots \leq d_v^{JS}$  and  $\phi_1^{DS} \leq \dots \leq \phi_v^{DS}$ . Sheena and Takemura (1992) give two methods of improving such a non-order-preserving estimator in the nonsingular case. Now, the methods are unifiedly treated for the nonsingular and singular cases:

- (i) For  $i = 1, \dots, v$ , let  $\varphi_{(i)}$  be the  $i$ -th largest element in  $\boldsymbol{\varphi} = (\varphi_1, \dots, \varphi_v)$ . Define  $\widehat{\boldsymbol{\Sigma}}^{OS} = \mathbf{H} \text{diag}(\boldsymbol{\varphi}^{OS}) \mathbf{H}^\top$  with  $\boldsymbol{\varphi}^{OS} = (\varphi_{(1)}, \dots, \varphi_{(v)})$ .
- (ii) Denote by  $(\varphi_1^{IR}, \dots, \varphi_v^{IR})$  the isotonic regression of  $(\varphi_1, \dots, \varphi_v)$ , satisfying

$$\min_{c_1 \geq \dots \geq c_v} \sum_{i=1}^v (c_i - \varphi_i)^2 = \sum_{i=1}^v (\varphi_i^{IR} - \varphi_i)^2.$$

Define  $\widehat{\boldsymbol{\Sigma}}^{IR} = \mathbf{H} \text{diag}(\boldsymbol{\varphi}^{IR}) \mathbf{H}^\top$  with  $\boldsymbol{\varphi}^{IR} = (\varphi_1^{IR}, \dots, \varphi_v^{IR})$ .

The  $\varphi_i^{IR}$ 's are given by

$$\varphi_i^{IR} = \min_{s \leq i} \max_{t \geq i} \frac{\sum_{r=s}^t \varphi_r}{t - s + 1},$$

so that  $\sum_{i=1}^v g(\varphi_i^{IR}) \leq \sum_{i=1}^v g(\varphi_i)$  for any convex function  $g$ . For computation algorithm and mathematical properties of the isotonic regression, see Robertson et al. (1988, Chap. 1).

We first show the following theorem.

**Theorem 7.2** *Let  $\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}^*) = \mathbf{H} \text{diag}(\boldsymbol{\varphi}^*) \mathbf{H}^\top$  be an orthogonally invariant estimator of  $\boldsymbol{\Sigma}$ , where  $\boldsymbol{\varphi}^* = (\varphi_1^*, \dots, \varphi_v^*)$  and the  $\varphi_i^*$ 's are functions of  $\mathbf{L}$ . Assume that*

$$\sum_{i=1}^j \varphi_i^* \geq \sum_{i=1}^j \varphi_i \text{ for } 1 \leq j \leq v-1, \text{ and } \sum_{i=1}^v \varphi_i^* = \sum_{i=1}^v \varphi_i.$$

If  $\Pr(\boldsymbol{\varphi}^* = \boldsymbol{\varphi}) \neq 1$  and  $\sum_{i=1}^v \log \varphi_i^* \geq \sum_{i=1}^v \log \varphi_i$ , then  $\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}^*)$  dominates  $\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi})$  relative to the extended Stein loss (7.3).

**Proof** Since  $|\text{Ch}(\boldsymbol{\Sigma}^{-1} \mathbf{H} \text{diag}(\boldsymbol{\varphi}^*) \mathbf{H}^\top)| = |\mathbf{H}^\top \boldsymbol{\Sigma}^{-1} \mathbf{H}| \prod_{i=1}^v \varphi_i^*$ , we obtain

$$\begin{aligned} \log |\text{Ch}(\boldsymbol{\Sigma}^{-1} \mathbf{H} \text{diag}(\boldsymbol{\varphi}^*) \mathbf{H}^\top)| &= \log |\mathbf{H}^\top \boldsymbol{\Sigma}^{-1} \mathbf{H}| + \sum_{i=1}^v \log \varphi_i^* \\ &\geq \log |\mathbf{H}^\top \boldsymbol{\Sigma}^{-1} \mathbf{H}| + \sum_{i=1}^v \log \varphi_i \\ &= \log |\text{Ch}(\boldsymbol{\Sigma}^{-1} \mathbf{H} \text{diag}(\boldsymbol{\varphi}) \mathbf{H}^\top)|, \end{aligned}$$

so that

$$R_{ES}(\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}^*), \boldsymbol{\Sigma}) - R_{ES}(\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}), \boldsymbol{\Sigma}) \leq E[\text{tr} \boldsymbol{\Sigma}^{-1} \mathbf{H} \{ \text{diag}(\boldsymbol{\varphi}^*) - \text{diag}(\boldsymbol{\varphi}) \} \mathbf{H}^\top].$$

For  $i = 1, \dots, v$ , let  $a_i = \{\mathbf{H}^\top \boldsymbol{\Sigma}^{-1} \mathbf{H}\}_{ii}$ . From (3.3),

$$\begin{aligned} &E[\text{tr} \boldsymbol{\Sigma}^{-1} \mathbf{H} \{ \text{diag}(\boldsymbol{\varphi}^*) - \text{diag}(\boldsymbol{\varphi}) \} \mathbf{H}^\top] \\ &= E\left[\sum_{i=1}^v (\varphi_i^* - \varphi_i) a_i\right] \\ &= c \int_{\mathbb{D}_v^{(\geq 0)}} \sum_{i=1}^v (\varphi_i^* - \varphi_i) E^*(a_i | \mathbf{L}) |\mathbf{L}|^{(ln-p-1)/2} \left( \prod_{1 \leq i < j \leq v} (\ell_i - \ell_j) \right) (d\mathbf{L}), \end{aligned}$$

where  $c$  is a constant and

$$E^*(a_i | \mathbf{L}) = \int_{\mathbb{V}_{p,v}} a_i \exp\left(-\frac{1}{2} \sum_{k=1}^v a_k \ell_k\right) (\mathbf{H}^\top d\mathbf{H}).$$

Note that

$$\begin{aligned} &\sum_{i=1}^v (\varphi_i^* - \varphi_i) E^*(a_i | \mathbf{L}) \\ &= (\varphi_1^* - \varphi_1) \{E^*(a_1 | \mathbf{L}) - E^*(a_2 | \mathbf{L})\} + (\varphi_1^* + \varphi_2^* - \varphi_1 - \varphi_2) \{E^*(a_2 | \mathbf{L}) - E^*(a_3 | \mathbf{L})\} \\ &\quad + \dots + (\varphi_1^* + \dots + \varphi_{v-1}^* - \varphi_1 - \dots - \varphi_{v-1}) \{E^*(a_{v-1} | \mathbf{L}) - E^*(a_v | \mathbf{L})\}. \end{aligned}$$

Hence, if

$$\Delta_i(\mathbf{L}) \equiv E^*(a_i | \mathbf{L}) - E^*(a_{i+1} | \mathbf{L}) \leq 0$$

for  $i = 1, \dots, v-1$  then  $R_{ES}(\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}^*), \boldsymbol{\Sigma}) \leq R_{ES}(\widehat{\boldsymbol{\Sigma}}(\boldsymbol{\varphi}), \boldsymbol{\Sigma})$ .

Since  $(\mathbf{H}^\top d\mathbf{H})$  is invariant under any orthogonal transformation, it is invariant under permutation of columns of  $\mathbf{H}$ . Exchanging the  $i$ -th and  $(i+1)$ -th columns of

$\mathbf{H}$  gives

$$\Delta_i(\mathbf{L}) = \int_{\mathbb{V}_{p,v}} (a_{i+1} - a_i) \exp \left( -\frac{1}{2}a_i\ell_{i+1} - \frac{1}{2}a_{i+1}\ell_i - \frac{1}{2} \sum_{k \neq i, i+1}^v a_k \ell_k \right) (\mathbf{H}^\top d\mathbf{H}),$$

so that

$$\begin{aligned} 2\Delta_i(\mathbf{L}) &= \int_{\mathbb{V}_{p,v}} (a_i - a_{i+1}) \exp \left( -\frac{1}{2} \sum_{k=1}^v a_k \ell_k \right) (\mathbf{H}^\top d\mathbf{H}) \\ &\quad + \int_{\mathbb{V}_{p,v}} (a_{i+1} - a_i) \exp \left( -\frac{1}{2}a_i\ell_{i+1} - \frac{1}{2}a_{i+1}\ell_i - \frac{1}{2} \sum_{k \neq i, i+1}^v a_k \ell_k \right) (\mathbf{H}^\top d\mathbf{H}) \\ &= \int_{\mathbb{V}_{p,v}} (a_i - a_{i+1}) \left\{ 1 - \exp \left( \frac{1}{2}(a_i - a_{i+1})(\ell_i - \ell_{i+1}) \right) \right\} \\ &\quad \times \exp \left( -\frac{1}{2} \sum_{k=1}^v a_k \ell_k \right) (\mathbf{H}^\top d\mathbf{H}). \end{aligned}$$

For both when  $a_i - a_{i+1} \geq 0$  and when  $a_i - a_{i+1} < 0$ , we can verify  $\Delta_i(\mathbf{L}) \leq 0$ . Thus the proof is complete.  $\square$

It is easy to check that both estimators  $\widehat{\Sigma}^{OS}$  and  $\widehat{\Sigma}^{IR}$  satisfy conditions of Theorem 7.2. Hence we obtain the following proposition.

**Proposition 7.8**  $\widehat{\Sigma}^{OS}$  and  $\widehat{\Sigma}^{IR}$  are better than  $\widehat{\Sigma}^O$  relative to the extended Stein loss (7.3).

#### 7.4.3.6 The Perron Estimator

Consider the case of  $n \geq p$ . For every  $\mathbf{O} \in \mathbb{O}_p$ , let  $\mathbf{T}_O \mathbf{T}_O^\top$  be the Cholesky decomposition of  $\mathbf{O}^\top \mathbf{S} \mathbf{O}$ , where  $\mathbf{T}_O \in \mathbb{L}_p^{(+)}$ . Using the James-Stein estimator  $\widehat{\Sigma}^{JS} = \widehat{\Sigma}^{JS}(\mathbf{S})$ , we define an estimator of  $\Sigma$  as

$$\widehat{\Sigma}^E = E[\mathbf{O} \widehat{\Sigma}^{JS}(\mathbf{O}^\top \mathbf{S} \mathbf{O}) \mathbf{O}^\top | \mathbf{S}] = E[\mathbf{O} \mathbf{T}_O \mathbf{D}_p^{JS} \mathbf{T}_O^\top \mathbf{O}^\top | \mathbf{S}],$$

where  $E[\cdot | \mathbf{S}]$  stands for conditional expectation with respect to the uniform distribution on  $\mathbb{O}_p$  given  $\mathbf{S}$ . Eaton (1970) suggested that  $\widehat{\Sigma}^E$  is an orthogonally invariant estimator improving  $\widehat{\Sigma}^{JS}$  relative to the ordinary Stein loss (7.2). Denote  $\widehat{\Sigma}^E = \mathbf{H} \mathbf{L} \Phi^E(\mathbf{L}) \mathbf{H}^\top$  with  $\Phi^E(\mathbf{L}) = \text{diag}(\phi_1^E, \dots, \phi_p^E)$ . The computation of  $\Phi^E(\mathbf{L})$  was done by Sharma and Krishnamoorthy (1983) for  $p = 2$  and by Takemura (1984) for  $p = 3$ . Takemura (1984) also provided a detailed discussion on  $\widehat{\Sigma}^E$  and showed that, for  $i = 1, \dots, p$ ,  $\phi_i^E$  can be expressed by

$$\phi_i^E = \sum_{j=1}^p w_{ij}(\mathbf{L}) d_j^{JS},$$

where  $\mathbf{W}(\mathbf{L}) = (w_{ij}(\mathbf{L}))$  is a doubly stochastic matrix. A closed form of  $\mathbf{W}(\mathbf{L})$  is hard to clarify. Perron (1992) proposed an approximation method for  $\mathbf{W}(\mathbf{L})$  and established a dominance result of the resulting estimator. Here, we give a unified Perron (1992) type estimator for the  $n \geq p$  and the  $p > n$  cases.

For  $i, j \in \{1, \dots, v\}$ , let

$$w_{ij}^{PR}(\mathbf{L}) = \frac{\text{tr}_{j-1}(\mathbf{L}_i)}{\text{tr}_{j-1}(\mathbf{L})} - \frac{\text{tr}_j(\mathbf{L}_i)}{\text{tr}_j(\mathbf{L})},$$

where  $\mathbf{L}_i = \text{diag}(\ell_1, \dots, \ell_{i-1}, 0, \ell_{i+1}, \dots, \ell_v)$  and

$$\text{tr}_j(\mathbf{L}) = \begin{cases} 1, & \text{if } j = 0, \\ \sum_{1 \leq i_1 < \dots < i_j \leq v} \prod_{k=1}^j \ell_{i_k}, & \text{if } j \in \{1, \dots, v\}, \\ 0, & \text{otherwise.} \end{cases}$$

Let  $\phi_i^{PR} = \sum_{j=1}^v w_{ij}^{PR}(\mathbf{L}) d_j^{JS}$  for  $i = 1, \dots, v$  and let  $\Phi^{PR}(\mathbf{L}) = \text{diag}(\phi_1^{PR}, \dots, \phi_v^{PR})$ . The Perron (1992) type estimator is defined by

$$\hat{\Sigma}^{PR} = \mathbf{H} \mathbf{L} \Phi^{PR}(\mathbf{L}) \mathbf{H}^\top.$$

Then we obtain the following proposition.

**Proposition 7.9** *The Perron type estimator  $\hat{\Sigma}^{PR}$  is better than  $\hat{\Sigma}^{JS}$  relative to the extended Stein loss (7.3).*

The proof of Proposition 7.9 is omitted since it can be done along the same lines as in Perron (1992). A noteworthy fact is that  $\hat{\Sigma}^{PR}$  has the ordering and shrinkage properties

$$\ell_1 \phi_1^{PR}(\mathbf{L}) \geq \dots \geq \ell_v \phi_v^{PR}(\mathbf{L}), \quad \phi_1^{PR}(\mathbf{L}) \leq \dots \leq \phi_v^{PR}(\mathbf{L}),$$

which can be shown in the same arguments as in Perron (1992). Hence  $\hat{\Sigma}^{PR}$  is order-preserving as well.

## 7.5 Improvement Using Information on Mean Statistic

Kubokawa and Tsai (2006) suggested a truncation method for improving the existing estimators of  $\Sigma$  based on both  $\mathbf{Y}$  and  $\mathbf{X}$  in (7.1) with  $n \geq p$ . The truncation method was generalized by Tsukuma and Kubokawa (2016) for any possible ordering among  $m, n$  and  $p$ . This section will introduce the generalized truncation method.

### 7.5.1 A Class of Estimators and Its Risk Function

This section uses the same notation as in Sect. 6.3. Recall that  $\nu = n \wedge p$ ,  $\tau = m \wedge n \wedge p$  and  $\mathbf{XHL}^{-1/2} = \mathbf{RF}^{1/2}\mathbf{V}^\top$ , where  $\mathbf{R} \in \mathbb{V}_{m,\tau}$ ,  $\mathbf{V} \in \mathbb{V}_{\nu,\tau}$  and  $\mathbf{F}^{1/2} = \text{diag}(\sqrt{f_1}, \dots, \sqrt{f_\tau}) \in \mathbb{D}_\tau^{(\geq 0)}$ . As seen in Lemma 6.3,  $\mathbf{Q}^- = \mathbf{V}^\top \mathbf{L}^{1/2} \mathbf{H}^\top$  is the generalized inverse of  $\mathbf{Q} = \mathbf{HL}^{-1/2} \mathbf{V}$ .

Let  $c_0 = \kappa^{-1} = (n \vee p)^{-1}$ . A class of estimators treated in this section is of the form

$$\widehat{\Sigma}(\Psi) = \widehat{\Sigma}^{BS} + c_0(\mathbf{Q}^-)^\top \Psi(\mathbf{F}) \mathbf{Q}^- = c_0\{S + (\mathbf{Q}^-)^\top \Psi(\mathbf{F}) \mathbf{Q}^-\}, \quad (7.27)$$

where  $\Psi(\mathbf{F}) \in \mathbb{D}_\tau$  satisfies  $\Psi(\mathbf{F}) + \mathbf{I}_\tau$  is positive definite and the diagonal elements of  $\Psi(\mathbf{F})$  are absolutely continuous functions of  $\mathbf{F}$ . The class (7.27) can be rewritten by

$$\widehat{\Sigma}(\Psi) = \widehat{\Sigma}^{BS} + c_0 \mathbf{S} \mathbf{S}^+ \mathbf{X}^\top \mathbf{R} \mathbf{F}^{-1} \Psi(\mathbf{F}) \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+$$

because  $\mathbf{Q}^- = \mathbf{F}^{-1/2} \mathbf{R}^\top \mathbf{X} \mathbf{S} \mathbf{S}^+$  from (6.14). Also, the definition of  $\mathbf{Q}^-$  implies that

$$\widehat{\Sigma}(\Psi) = c_0 \mathbf{H} \mathbf{L}^{1/2} (\mathbf{I}_\nu + \mathbf{V} \Psi(\mathbf{F}) \mathbf{V}^\top) \mathbf{L}^{1/2} \mathbf{H}^\top.$$

When  $\Psi(\mathbf{F}) + \mathbf{I}_\tau$  is positive definite,  $\mathbf{I}_\nu + \mathbf{V} \Psi(\mathbf{F}) \mathbf{V}^\top$  is positive definite and thus both  $\widehat{\Sigma}(\Psi)$  and  $\Sigma^{-1} \widehat{\Sigma}(\Psi)$  are of rank  $\nu$  with probability one. Here, we will give a useful risk expression for  $\widehat{\Sigma}(\Psi)$  in the following theorem.

**Theorem 7.3** *Let  $\tau = m \wedge n \wedge p$  and  $\Psi = \Psi(\mathbf{F}) = \text{diag}(\psi_1, \dots, \psi_\tau)$ . For any order among  $m, n$  and  $p$ , the risk function of  $\widehat{\Sigma}(\Psi)$  in (7.27) relative to the extended Stein loss (7.3) is expressed as*

$$\begin{aligned} R_{ES}(\widehat{\Sigma}(\Psi), \Sigma) &= R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) \\ &\quad + c_0 E \left[ \sum_{i=1}^{\tau} \alpha_i \psi_i - 2g_1(\Psi) - 2g_2(\Psi) - c_0^{-1} \log |\mathbf{I}_\tau + \Psi| \right], \end{aligned} \quad (7.28)$$

where  $\alpha_i = |n - p| + 2i - 1$  for  $i = 1, \dots, \tau$  and

$$g_1(\Psi) = \sum_{i=1}^{\tau} f_i \frac{\partial \psi_i}{\partial f_i}, \quad g_2(\Psi) = \sum_{i=1}^{\tau} \sum_{j>i}^{\tau} \frac{\psi_i - \psi_j}{f_i - f_j} f_j.$$

**Proof** Both  $\mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2}$  and  $\mathbf{I}_\nu + \mathbf{V} \Psi \mathbf{V}^\top$  are of rank  $\nu$ , so that

$$\begin{aligned} |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}(\Psi))| &= |\text{Ch}(c_0 \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2} (\mathbf{I}_\nu + \mathbf{V} \Psi \mathbf{V}^\top))| \\ &= |c_0 \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2} (\mathbf{I}_\nu + \mathbf{V} \Psi \mathbf{V}^\top)| \\ &= |c_0 \mathbf{L}^{1/2} \mathbf{H}^\top \Sigma^{-1} \mathbf{H} \mathbf{L}^{1/2}| \times |\mathbf{I}_\nu + \mathbf{V} \Psi \mathbf{V}^\top| \\ &= |\text{Ch}(\Sigma^{-1} \widehat{\Sigma}^{BS})| \times |\mathbf{I}_\tau + \Psi|. \end{aligned}$$



Since  $\text{tr}[\text{Ch}(\Sigma^{-1}\widehat{\Sigma}^{BS})] = \text{tr}\Sigma^{-1}\widehat{\Sigma}^{BS}$  and  $\text{tr}[\text{Ch}(\Sigma^{-1}\widehat{\Sigma}(\Psi))] = \text{tr}\Sigma^{-1}\widehat{\Sigma}(\Psi)$ , the risk of  $\widehat{\Sigma}(\Psi)$  with respect to the extended Stein loss (7.3) is written as

$$\begin{aligned}
 R_{ES}(\widehat{\Sigma}(\Psi), \Sigma) &= E[L_{ES}(\widehat{\Sigma}(\Psi), \Sigma)] \\
 &= E[\text{tr}[\text{Ch}(\Sigma^{-1}\widehat{\Sigma}(\Psi))] - \log |\text{Ch}(\Sigma^{-1}\widehat{\Sigma}(\Psi))| - \nu] \\
 &= E[\text{tr}\Sigma^{-1}\widehat{\Sigma}^{BS} - \log |\text{Ch}(\Sigma^{-1}\widehat{\Sigma}^{BS})| - \nu] \\
 &\quad + c_0 E[\text{tr}\Sigma^{-1}(\mathbf{Q}^-)^\top \Psi \mathbf{Q}^- - c_0^{-1} \log |\mathbf{I}_\tau + \Psi|] \\
 &= R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) + c_0 E[\text{tr}\Sigma^{-1}(\mathbf{Q}^-)^\top \Psi \mathbf{Q}^- - c_0^{-1} \log |\mathbf{I}_\tau + \Psi|].
 \end{aligned} \tag{7.29}$$

Applying the Stein identity (5.3) and Lemma 6.7 to  $E[\text{tr}\Sigma^{-1}(\mathbf{Q}^-)^\top \Psi \mathbf{Q}^-]$  leads to

$$\begin{aligned}
 E[\text{tr}\Sigma^{-1}(\mathbf{Q}^-)^\top \Psi \mathbf{Q}^-] &= E[\text{tr}\Sigma^{-1}SS^\top X^\top \mathbf{R}F^{-1}\Psi \mathbf{R}^\top XSS^\top] \\
 &= E\left[\sum_{i=1}^{\tau}\left\{a\psi_i - 2f_i\frac{\partial\psi_i}{\partial f_i} - 2\sum_{j>i}^{\tau}\frac{f_i\psi_i - f_j\psi_j}{f_i - f_j}\right\}\right],
 \end{aligned}$$

where  $a = n + p - 2\nu + 2\tau - 1$ . It follows that

$$\sum_{i=1}^{\tau}\sum_{j>i}^{\tau}\frac{f_i\psi_i - f_j\psi_j}{f_i - f_j} = \sum_{i=1}^{\tau}\left\{(\tau - i)\psi_i + \sum_{j>i}^{\tau}\frac{\psi_i - \psi_j}{f_i - f_j}f_j\right\},$$

implying that

$$E[\text{tr}\Sigma^{-1}(\mathbf{Q}^-)^\top \Psi \mathbf{Q}^-] = E\left[\sum_{i=1}^{\tau}\left\{\alpha_i\psi_i - 2f_i\frac{\partial\psi_i}{\partial f_i} - 2\sum_{j>i}^{\tau}\frac{\psi_i - \psi_j}{f_i - f_j}f_j\right\}\right]. \tag{7.30}$$

where  $\alpha_i = a - 2(\tau - i) = |n - p| + 2i - 1$ . Combining (7.29) and (7.30), we obtain (7.28). Thus the proof is complete.  $\square$

## 7.5.2 Examples of Improved Estimators

### 7.5.2.1 The Stein-Type Estimator

A Stein-type estimator similar to (7.24) is described by

$$\widehat{\Sigma}(\Psi^{ST}) = c_0\{S + (\mathbf{Q}^-)^\top \Psi^{ST}(\mathbf{F})\mathbf{Q}^-\}, \quad \Psi^{ST}(\mathbf{F}) = \text{diag}(\psi_1^{ST}, \dots, \psi_\tau^{ST}),$$

where for  $i = 1, \dots, \tau$ ,

$$\psi_i^{ST} = \frac{1}{c_0 \alpha_i} - 1 = \frac{\nu - 2i + 1}{|n - p| + 2i - 1}.$$

The risk function of  $\widehat{\Sigma}(\Psi^{ST})$  under the extended Stein loss (7.3) is expressed as

$$R_{ES}(\widehat{\Sigma}(\Psi^{ST}), \Sigma) = R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) + \sum_{i=1}^{\tau} \{c_0 \alpha_i \psi_i^{ST} - \log(1 + \psi_i^{ST})\} - 2c_0 E[g_2(\Psi^{ST})].$$

It is observed that  $\psi_i^{ST} - \psi_j^{ST} > 0$  for  $j > i$ , so that  $g_2(\Psi^{ST}) > 0$ . Also, we see that for  $i = 1, \dots, \tau$

$$c_0 \alpha_i \psi_i^{ST} - \log(1 + \psi_i^{ST}) = -\{c_0 \alpha_i - \log(c_0 \alpha_i) - 1\} \leq 0$$

because  $x - \log x - 1 \geq 0$  for  $x > 0$ . Thus from Theorem 7.3,  $\widehat{\Sigma}(\Psi^{ST})$  dominates  $\widehat{\Sigma}^{BS}$  for any order of  $m$ ,  $n$  and  $p$ .

Further, if  $\tau = \nu$ , namely,  $m > n \wedge p$ , then  $\widehat{\Sigma}(\Psi^{ST})$  dominates the James-Stein estimator  $\widehat{\Sigma}^{JS}$  in (7.12). In fact, since  $\sum_{i=1}^{\nu} (c_0 \alpha_i - 1) = 0$  and

$$\sum_{i=1}^{\nu} \log(c_0 \alpha_i) = -\nu \log \kappa + \sum_{i=1}^{\nu} \log \alpha_i = -\nu \log \kappa + \sum_{i=1}^{\nu} \log(n + p - 2i + 1),$$

it follows that, by (7.5) and (7.13),

$$\begin{aligned} R_{ES}(\widehat{\Sigma}(\Psi^{ST}), \Sigma) &< R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) - \sum_{i=1}^{\nu} \{c_0 \alpha_i - \log(c_0 \alpha_i) - 1\} \\ &= \sum_{i=1}^{\nu} \log(n + p - 2i + 1) - r_{\kappa, \nu} = R_{ES}(\widehat{\Sigma}^{JS}, \Sigma). \end{aligned}$$

This shows that if  $\tau = \nu$  then  $\widehat{\Sigma}(\Psi^{ST})$  dominates  $\widehat{\Sigma}^{JS}$  relative to the extended Stein loss (7.3).

### 7.5.2.2 The Haff Type Estimator

As a reasonable estimator, we define the Haff (1980) type estimator as

$$\begin{aligned} \widehat{\Sigma}(\Psi^{HF}) &= c_0 \{S + (Q^-)^T \Psi^{HF} Q^-\}, \quad \Psi^{HF}(F) = \text{diag}(\psi_1^{HF}, \dots, \psi_{\tau}^{HF}), \\ \psi_i^{HF} &= \frac{a}{\text{tr } F} f_i \quad (i = 1, \dots, \tau), \quad a > 0. \end{aligned}$$

Using Theorem 7.3, we can show that the Haff type estimator  $\widehat{\Sigma}(\Psi^{HF})$  dominates  $\widehat{\Sigma}^{BS}$  if constant  $a$  satisfies the inequality  $0 < a \leq 2(\nu - 1)/(|n - p| + 1)$  for  $\nu > 1$ .

In fact, it is noted that

$$g_2(\Psi^{HF}) = \sum_{i=1}^{\tau} \sum_{j>i}^{\tau} \frac{\psi_i^{HF} - \psi_j^{HF}}{f_i - f_j} f_j = \frac{a}{\text{tr } \mathbf{F}} \sum_{i=1}^{\tau} \sum_{j>i}^{\tau} f_j = \frac{a}{\text{tr } \mathbf{F}} \sum_{i=1}^{\tau} (i-1) f_i.$$

The difference in risk of  $\widehat{\Sigma}(\Psi^{HF})$  and  $\widehat{\Sigma}^{BS}$  is written as

$$\begin{aligned} R_{ES}(\widehat{\Sigma}(\Psi^{HF}), \Sigma) - R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) \\ = c_0(|n-p|+1)a - 2c_0 E[g_1(\Psi^{HF})] - E \left[ \sum_{i=1}^{\tau} \log(1 + \psi_i^{HF}) \right]. \end{aligned}$$

It follows that for any  $\mathbf{F} \in \mathbb{D}_{\tau}^{(\geq 0)}$

$$g_1(\Psi^{HF}) = a \left( 1 - \frac{\text{tr } \mathbf{F}^2}{(\text{tr } \mathbf{F})^2} \right) \geq 0.$$

Since  $\log(1+x) \geq 2x/(2+x)$  for  $x \geq 0$  and  $\sum_{i=1}^{\tau} \psi_i^{HF} = a$ , we observe

$$\sum_{i=1}^{\tau} \log(1 + \psi_i^{HF}) \geq \sum_{i=1}^{\tau} \frac{2\psi_i^{HF}}{2 + \psi_i^{HF}} \geq \sum_{i=1}^{\tau} \frac{2\psi_i^{HF}}{2 + a} = \frac{2a}{2 + a}.$$

Thus,

$$\begin{aligned} R_{ES}(\widehat{\Sigma}(\Psi^{HF}), \Sigma) - R_{ES}(\widehat{\Sigma}^{BS}, \Sigma) &\leq c_0 \left\{ (|n-p|+1)a - c_0^{-1} \frac{2a}{2+a} \right\} \\ &= c_0(|n-p|+1) \frac{a}{2+a} \left[ a - \frac{2(v-1)}{|n-p|+1} \right], \end{aligned}$$

which shows the dominance result.

### 7.5.3 Further Improvements with a Truncation Rule

First, we provide a useful lemma which will be a key tool to show further dominance results.

**Lemma 7.1** *Let  $\Phi(\mathbf{F}) \in \mathbb{D}_{\tau}$  such that the diagonal elements are absolutely continuous and nonnegative functions of  $\mathbf{F}$ . Then we have*

$$E[\text{tr } \Sigma^{-1}(\mathbf{Q}^{-})^{\top} (\mathbf{I}_{\tau} + \mathbf{F}) \Phi(\mathbf{F}) \mathbf{Q}^{-}] \geq E[(\kappa + m) \text{tr } \Phi(\mathbf{F})].$$

**Proof** See Tsukuma and Kubokawa (2016). □

By using Lemma 7.1, we will improve on  $\widehat{\Sigma}(\Psi^{ST})$  and  $\widehat{\Sigma}(\Psi^{HF})$ . Let  $[\Psi]^{TR} = \text{diag}(\psi_1^{TR}(\mathbf{F}), \dots, \psi_\tau^{TR}(\mathbf{F})) \in \mathbb{D}_\tau$  such that the  $i$ -th diagonal element is given by

$$\psi_i^{TR}(\mathbf{F}) = \min\left\{\psi_i(\mathbf{F}), \frac{1 + f_i}{c_0(\kappa + m)} - 1\right\},$$

where  $\Psi(\mathbf{F}) = \text{diag}(\psi_1(\mathbf{F}), \dots, \psi_\tau(\mathbf{F}))$ . Then we obtain a general dominance result for improvement on the class (7.27).

**Theorem 7.4** *For any possible ordering among  $m$ ,  $n$  and  $p$ , the truncated estimator  $\widehat{\Sigma}([\Psi]^{TR})$  dominates  $\widehat{\Sigma}(\Psi)$  relative to the extended Stein loss (7.3) if  $\Pr([\Psi]^{TR} \neq \Psi) > 0$ .*

**Proof** Abbreviate  $\Psi(\mathbf{F})$  to  $\Psi$ . The difference in risk of  $\widehat{\Sigma}(\Psi)$  and  $\widehat{\Sigma}([\Psi]^{TR})$  can be expressed as

$$\begin{aligned} & R_{ES}(\widehat{\Sigma}(\Psi), \Sigma) - R_{ES}(\widehat{\Sigma}([\Psi]^{TR}), \Sigma) \\ &= E[c_0 \text{tr} \Sigma^{-1}(\mathbf{Q}^-)^\top (\Psi - [\Psi]^{TR}) \mathbf{Q}^- - \log |\mathbf{I}_\tau + \Psi| + \log |\mathbf{I}_\tau + [\Psi]^{TR}|] \\ &\geq E[c_0(\kappa + m) \text{tr}(\mathbf{I}_\tau + \mathbf{F})^{-1} (\Psi - [\Psi]^{TR}) - \log |\mathbf{I}_\tau + \Psi| + \log |\mathbf{I}_\tau + [\Psi]^{TR}|], \end{aligned}$$

where the inequality follows directly from Lemma 7.1. The last r.h.s. can be written by  $E[\sum_{i=1}^\tau \Delta_i]$  with

$$\Delta_i = c_0(\kappa + m) \cdot \frac{\psi_i - \psi_i^{TR}}{1 + f_i} - \log(1 + \psi_i) + \log(1 + \psi_i^{TR}).$$

When  $\psi_i^{TR} = \psi_i$ , we get  $\Delta_i = 0$ . When  $\psi_i^{TR} = c_0^{-1}(1 + f_i)/(\kappa + m) - 1 < \psi_i$ , it is observed that

$$\Delta_i = c_0(\kappa + m) \cdot \frac{1 + \psi_i}{1 + f_i} - \log \left[ c_0(\kappa + m) \cdot \frac{1 + \psi_i}{1 + f_i} \right] - 1 \geq 0,$$

which completes the proof.  $\square$

The following proposition is derived immediately from Theorem 7.4.

**Proposition 7.10** *The truncated estimator  $\widehat{\Sigma}([\Psi^{ST}]^{TR})$  dominates  $\widehat{\Sigma}(\Psi^{ST})$  relative to the extended Stein loss (7.3). Also,  $\widehat{\Sigma}([\Psi^{HF}]^{TR})$  dominates  $\widehat{\Sigma}(\Psi^{HF})$  relative to the extended Stein loss (7.3).*

## 7.6 Related Topics

### 7.6.1 Decomposition of the Estimation Problem

When  $n \geq p$ , covariance estimation under the ordinary Stein loss (7.2) is closely related to simultaneous estimation of mean vectors and variances. Here we will briefly introduce the relationship and then give a simple improved procedure on  $\widehat{\Sigma}^{JS}$  via the James-Stein (1961) shrinkage estimators of multivariate normal mean vectors.

We use the same notation as in Proposition 3.10 for  $n \geq p$ . Recall that  $\mathbf{T}\mathbf{T}^\top$  and  $\Xi\Xi^\top$  are the Cholesky decompositions of the Wishart matrix  $\mathbf{S}$  and the covariance matrix  $\Sigma$ , respectively, where  $\mathbf{T} = (t_{i,j}) \in \mathbb{L}_p^{(+)}$  and  $\Xi = (\xi_{i,j}) \in \mathbb{L}_p^{(+)}$ . Denote  $\mathbf{T}_{(p)} = t_{p,p}$  and  $\Xi_{(p)} = \xi_{p,p}$  and, for  $i = p-1, \dots, 1$ , define  $\mathbf{T}_{(i)}$  and  $\Xi_{(i)}$  inductively as, respectively,

$$\mathbf{T}_{(i)} = \begin{pmatrix} t_{i,i} & \mathbf{0}_{p-i}^\top \\ \mathbf{t}_{(i)} & \mathbf{T}_{(i+1)} \end{pmatrix} \in \mathbb{L}_{p-i+1}^{(+)}, \quad \Xi_{(i)} = \begin{pmatrix} \xi_{i,i} & \mathbf{0}_{p-i}^\top \\ \boldsymbol{\xi}_{(i)} & \Xi_{(i+1)} \end{pmatrix} \in \mathbb{L}_{p-i+1}^{(+)}$$

with  $\mathbf{t}_{(i)} = (t_{i+1,i}, \dots, t_{p,i})^\top$  and  $\boldsymbol{\xi}_{(i)} = (\xi_{i+1,i}, \dots, \xi_{p,i})^\top$ . Note that  $\mathbf{T}_{(1)} = \mathbf{T}$  and  $\Xi_{(1)} = \Xi$ . By Proposition 3.10, the columns of  $\mathbf{T}$  are mutually independent and

$$\begin{cases} t_{i,i}^2 \sim \sigma_i^2 \chi_{n-i+1}^2 & \text{for } i = 1, \dots, p, \\ \mathbf{t}_{(i)} | t_{i,i} \sim N_{p-i}(t_{i,i} \boldsymbol{\gamma}_{(i)}, \Sigma_{(i+1)}) & \text{for } i = 1, \dots, p-1, \end{cases} \quad (7.31)$$

where  $\boldsymbol{\gamma}_{(i)} = \boldsymbol{\xi}_{(i)}/\xi_{i,i}$  for  $i = 1, \dots, p-1$ ,  $\sigma_i^2 = \xi_{i,i}^2$  for  $i = 1, \dots, p$ , and  $\Sigma_{(i)} = \Xi_{(i)}\Xi_{(i)}^\top$  for  $i = 1, \dots, p$ . Recall also that for  $i = 1, \dots, p-1$

$$\Sigma_{(i)} = \begin{pmatrix} 1 & \mathbf{0}_{p-i}^\top \\ \boldsymbol{\gamma}_{(i)} & \mathbf{I}_{p-i} \end{pmatrix} \begin{pmatrix} \sigma_i^2 & \mathbf{0}_{p-i}^\top \\ \mathbf{0}_{p-i} & \Sigma_{(i+1)} \end{pmatrix} \begin{pmatrix} 1 & \boldsymbol{\gamma}_{(i)}^\top \\ \mathbf{0}_{p-i} & \mathbf{I}_{p-i} \end{pmatrix}. \quad (7.32)$$

Let the  $\widehat{\sigma}_i^2$ 's and the  $\widehat{\boldsymbol{\gamma}}_{(i)}$ 's be certain estimators of the  $\sigma_i^2$ 's and the  $\boldsymbol{\gamma}_{(i)}$ 's, respectively. Set  $\widehat{\Sigma}_{(p)} = \widehat{\sigma}_p^2$ . For  $i = p-1, \dots, 1$ , we define  $\widehat{\Sigma}_{(i)} (\in \mathbb{S}_{p-i+1}^{(+)})$  inductively as

$$\widehat{\Sigma}_{(i)} = \begin{pmatrix} 1 & \mathbf{0}_{p-i}^\top \\ \widehat{\boldsymbol{\gamma}}_{(i)} & \mathbf{I}_{p-i} \end{pmatrix} \begin{pmatrix} \widehat{\sigma}_i^2 & \mathbf{0}_{p-i}^\top \\ \mathbf{0}_{p-i} & \widehat{\Sigma}_{(i+1)} \end{pmatrix} \begin{pmatrix} 1 & \widehat{\boldsymbol{\gamma}}_{(i)}^\top \\ \mathbf{0}_{p-i} & \mathbf{I}_{p-i} \end{pmatrix}. \quad (7.33)$$

Then  $\widehat{\Sigma}^A = \widehat{\Sigma}_{(1)}$  is an estimator of  $\Sigma$ . Conversely, for any estimator  $\widehat{\Sigma}$ , the LDL<sup>⊤</sup> decomposition of  $\widehat{\Sigma}$  can be obtained uniquely from (7.33).

Combining (7.32) and (7.33) yields

$$\text{tr } \Sigma_{(i)}^{-1} \widehat{\Sigma}_{(i)} = \widehat{\sigma}_i^2 / \sigma_i^2 + \widehat{\sigma}_i^2 (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)})^\top \Sigma_{(i+1)}^{-1} (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)}) + \text{tr } \Sigma_{(i+1)}^{-1} \widehat{\Sigma}_{(i+1)},$$

which is used again and again to obtain

$$\text{tr } \Sigma^{-1} \widehat{\Sigma}^A = \sum_{i=1}^p \frac{\widehat{\sigma}_i^2}{\sigma_i^2} + \sum_{i=1}^{p-1} \widehat{\sigma}_i^2 (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)})^\top \Sigma_{(i+1)}^{-1} (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)}).$$

Since  $|\Sigma^{-1} \Sigma^A| = \sum_{i=1}^p \widehat{\sigma}_i^2 / \sigma_i^2$ , the ordinary Stein loss (7.2) of  $\widehat{\Sigma}^A$  derived from (7.33) can alternatively be written as

$$L_S(\widehat{\Sigma}^A, \Sigma) = \sum_{i=1}^p \left( \frac{\widehat{\sigma}_i^2}{\sigma_i^2} - \log \frac{\widehat{\sigma}_i^2}{\sigma_i^2} - 1 \right) + \sum_{i=1}^{p-1} \widehat{\sigma}_i^2 (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)})^\top \Sigma_{(i+1)}^{-1} (\widehat{\boldsymbol{\gamma}}_{(i)} - \boldsymbol{\gamma}_{(i)}). \quad (7.34)$$

This suggests that the covariance estimation problem with the ordinary Stein loss (7.2) is considered as the problem of simultaneously estimating the  $\sigma_i^2$ 's and the  $\boldsymbol{\gamma}_{(i)}$ 's under the decomposed loss (7.34) in the decomposed model (7.31).

If  $\widehat{\Sigma}^A = \widehat{\Sigma}^{JS}$ , then estimators of the  $\sigma_i^2$ 's and the  $\boldsymbol{\gamma}_{(i)}$ 's can be written as, respectively,

$$\begin{cases} \widehat{\sigma}_i^{2JS} = d_i^{JS} t_{i,i}^2 & \text{for } i = 1, \dots, p, \\ \widehat{\boldsymbol{\gamma}}_{(i)}^{JS} = \mathbf{t}_{(i)} / t_{i,i} & \text{for } i = 1, \dots, p-1. \end{cases}$$

Hence the risk of  $\widehat{\Sigma}^{JS}$  is expressed by  $R_S(\widehat{\Sigma}^{JS}, \Sigma) = E[L_S(\widehat{\Sigma}^{JS}, \Sigma)] = R_1^{JS} + R_2^{JS}$ , where

$$\begin{aligned} R_1^{JS} &= E \left[ \sum_{i=1}^p \left( \frac{\widehat{\sigma}_i^{2JS}}{\sigma_i^2} - \log \frac{\widehat{\sigma}_i^{2JS}}{\sigma_i^2} - 1 \right) \right], \\ R_2^{JS} &= E \left[ \sum_{i=1}^{p-1} \widehat{\sigma}_i^{2JS} (\widehat{\boldsymbol{\gamma}}_{(i)}^{JS} - \boldsymbol{\gamma}_{(i)})^\top \Sigma_{(i+1)}^{-1} (\widehat{\boldsymbol{\gamma}}_{(i)}^{JS} - \boldsymbol{\gamma}_{(i)}) \right], \end{aligned}$$

where the expectations are taken with respect to (7.31). Here, when  $p \geq 4$ , we consider improvement on

$$R_2^{JS} = E \left[ \sum_{i=1}^{p-1} \widehat{\sigma}_i^{2JS} (p-i) / t_{i,i}^2 \right] = \sum_{i=1}^{p-1} (p-i) d_i^{JS}.$$

For  $i = 1, \dots, p-1$ , denote  $\mathbf{x}_{(i)} = \mathbf{t}_{(i)} / t_{i,i}$  and  $\mathbf{S}_{(i)} = \mathbf{T}_{(i)} \mathbf{T}_{(i)}^\top$ . Define

$$\widehat{\boldsymbol{\gamma}}_{(i)}^{SH} = \begin{cases} \left\{ 1 - \frac{p-i-2}{(n-p+3)t_{i,i}^2 \mathbf{x}_{(i)}^\top \mathbf{S}_{(i+1)}^{-1} \mathbf{x}_{(i)}} \right\} \mathbf{x}_{(i)} & \text{for } i = 1, \dots, p-3, \\ \mathbf{x}_{(i)} & \text{for } i = p-2 \text{ and } p-1. \end{cases}$$

For  $i = 1, \dots, p-3$ ,  $\widehat{\boldsymbol{\gamma}}_{(i)}^{SH}$  is the James-Stein (1961, Eq. 23) shrinkage estimator in estimation of the multivariate normal mean vector with unknown covariance matrix.

Note that  $\mathbf{S}_{(i+1)} \sim \mathcal{W}_{p-i}(n-i, \mathbf{\Sigma}_{(i+1)})$  independent of  $\mathbf{x}_{(i)}$  for each  $i$ . In a similar way to James and Stein (1961), we can show that for  $i = 1, \dots, p-3$

$$E \left[ \hat{\sigma}_i^{2JS} (\hat{\mathbf{y}}_{(i)}^{SH} - \mathbf{y}_{(i)})^\top \mathbf{\Sigma}_{(i+1)}^{-1} (\hat{\mathbf{y}}_{(i)}^{SH} - \mathbf{y}_{(i)}) \right] \leq (p-i) d_i^{JS},$$

which implies that  $\hat{\mathbf{\Sigma}}^A$  obtained from using the  $\hat{\sigma}_i^{2JS}$ 's and the  $\hat{\mathbf{y}}_{(i)}^{SH}$ 's dominates  $\hat{\mathbf{\Sigma}}^{JS}$  relative to the ordinary Stein loss (7.2).

For more details and improvement on  $R_1^{JS}$ , see Tsukuma (2014a, 2016b). The improvement on  $\hat{\mathbf{\Sigma}}^{JS}$  can also be done by using matricial shrinkage estimators of the mean matrix and this was discussed in Ma et al. (2012).

### 7.6.2 Decision-Theoretic Studies Under Quadratic Losses

Instead of the ordinary Stein loss (7.2) or the extended Stein loss (7.3), some quadratic-type loss functions have often been used for obtaining decision-theoretic results on covariance estimation. A typical quadratic loss is

$$L_1(\hat{\mathbf{\Sigma}}, \mathbf{\Sigma}) = \text{tr } \mathbf{\Sigma}^{-1}(\hat{\mathbf{\Sigma}} - \mathbf{\Sigma})\mathbf{\Sigma}^{-1}(\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}) = \text{tr } (\mathbf{\Sigma}^{-1}\hat{\mathbf{\Sigma}} - \mathbf{I}_p)^2.$$

The  $L_1$ -loss is invariant under the general scale transformation  $\mathbf{\Sigma} \rightarrow \mathbf{U}^\top \mathbf{\Sigma} \mathbf{U}$  and  $\hat{\mathbf{\Sigma}} \rightarrow \mathbf{U}^\top \hat{\mathbf{\Sigma}} \mathbf{U}$  for any  $\mathbf{U} \in \mathbb{U}_p$ . Selliah (1964) addressed the  $n \geq p$  case of covariance estimation under the  $L_1$ -loss and obtained a minimax estimator based on the Cholesky decomposition of the Wishart matrix. For other approaches, see Haff (1979b, 1980, 1991), Yang and Berger (1994) and Tsukuma (2014b). See also Konno (2009), who discussed the  $p > n$  case under the  $L_1$ -loss.

A multivariate extension of squared error loss to covariance estimation may be defined by

$$L_2(\hat{\mathbf{\Sigma}}, \mathbf{\Sigma}) = \text{tr } (\hat{\mathbf{\Sigma}} - \mathbf{\Sigma})^2,$$

namely, the squared Frobenius norm of  $\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}$ . The  $L_2$ -loss has orthogonal invariance under the orthogonal transformation  $\mathbf{\Sigma} \rightarrow \mathbf{O}^\top \mathbf{\Sigma} \mathbf{O}$  and  $\hat{\mathbf{\Sigma}} \rightarrow \mathbf{O}^\top \hat{\mathbf{\Sigma}} \mathbf{O}$  for any  $\mathbf{O} \in \mathbb{O}_p$ . In the literature, a much-discussed estimator is a linear shrinkage estimator  $\hat{\mathbf{\Sigma}}^{LS} = \alpha \hat{\mathbf{\Sigma}}^{UB} + (1 - \alpha)(\text{tr } \hat{\mathbf{\Sigma}}^{UB}/p)\mathbf{I}_p$ , where  $0 \leq \alpha \leq 1$  and  $\hat{\mathbf{\Sigma}}^{UB} = \mathbf{S}/n$  is the unbiased estimator of  $\mathbf{\Sigma}$  under the normality assumption on error distribution. Leung and Chan (1998) suggested using  $\alpha = n/(n+2)$  from a decision-theoretic point of view. This suggestion of Leung and Chan (1998) was extended to an elliptically contoured distribution model by Leung and Ng (2004). Ledoit and Wolf (2004) took an asymptotic approach to estimating an optimal  $\alpha$  from sample under a general error distribution.

For  $\Sigma = (\sigma_{ij})$  and  $\hat{\Sigma} = (\hat{\sigma}_{ij})$ , the  $L_2$ -loss can be generalized as

$$L_3(\hat{\Sigma}, \Sigma) = \sum_{1 \leq i \leq j \leq p} w_{ij} (\hat{\sigma}_{ij} - \sigma_{ij})^2,$$

where  $w_{ij} \geq 0$  for  $1 \leq i \leq j \leq p$ . When  $w_{ii} = 1$  for  $i = 1, \dots, p$  and  $w_{ij} = 2$  for  $1 \leq i < j \leq p$ , the  $L_3$ -loss coincides with the  $L_2$ -loss. However, the  $L_3$ -loss does not include the  $L_1$ -loss. For decision-theoretic results under the  $L_3$ -loss, see Perlman (1972) and Haff (1979b).

For  $n \geq p$ , as a variant type of the  $L_1$ -loss, we define

$$L_4(\hat{\Sigma}, \Sigma) = \text{tr } \hat{\Sigma}^{-1}(\hat{\Sigma} - \Sigma)\Sigma^{-1}(\hat{\Sigma} - \Sigma) = \text{tr } \Sigma^{-1}\hat{\Sigma} + \text{tr } (\Sigma^{-1}\hat{\Sigma})^{-1} - 2p.$$

The  $L_4$ -loss can also be obtained from the sum of the ordinary Stein loss (7.2) and its different entropy-type loss  $L_P(\hat{\Sigma}, \Sigma) = \text{tr } (\Sigma^{-1}\hat{\Sigma})^{-1} - \log |(\Sigma^{-1}\hat{\Sigma})^{-1}| - p$ . The invariance of the  $L_4$ -loss can easily be verified under a general scale transformation. Improved estimation under the  $L_4$ -loss was studied by Kubokawa and Konno (1990), Gupta and Ofori-Nyarko (1995) and Sun and Sun (2005).

### 7.6.3 Estimation of the Generalized Variance

Some statistical measures are formulated as functions of the covariance matrix. Generalized variance is the determinant of the covariance matrix and interpreted as a scalar measure of uncertainty.

Consider the case of  $n \geq p$  in the model (7.1). The generalized variance is defined by  $|\Sigma|$ . Denote an estimator of  $|\Sigma|$  by  $|\hat{\Sigma}|$  and we now treat decision-theoretic estimation of  $|\Sigma|$  relative to a quadratic-type loss

$$L_G(|\hat{\Sigma}|, |\Sigma|) = |\Sigma|^{-2}(|\hat{\Sigma}| - |\Sigma|)^2.$$

The Cholesky decomposition of  $S$  is denoted by  $S = TT^\top$ , where  $T = (t_{ij}) \in \mathbb{L}_p^{(+)}$ . Due to Proposition 3.10, it turns out that

$$E[|S|] = \prod_{i=1}^p E[t_{ii}^2] = \prod_{i=1}^p (n - i + 1)\sigma_i^2 = |\Sigma| \prod_{i=1}^p (n - i + 1),$$

so that

$$|\hat{\Sigma}^{UB}| = \left\{ \prod_{i=1}^p (n - i + 1)^{-1} \right\} |S| = \frac{(n - p + 2)!}{n!} |S|$$



is the unbiased estimator of  $|\Sigma|$ . However under the  $L_G$ -loss,  $|\widehat{\Sigma}^{UB}|$  is not the best among constant-multiple estimators of the form  $c|\mathbf{S}|$  with positive constant  $c$ . The best constant  $c$  that minimizes risk of estimator  $c|\mathbf{S}|$  is

$$c_0 = \frac{(n - p + 2)!}{(n + 2)!},$$

which can easily be verified by Proposition 3.10.

Here, any constant-multiple estimator  $c|\mathbf{S}|$  is invariant under an affine transformation. Shorrock and Zidek (1976) discussed improvement on the best affine invariant estimator  $|\widehat{\Sigma}^{BC}| = c_0|\mathbf{S}|$  by using the information on  $\mathbf{X}$ . They showed that  $|\widehat{\Sigma}^{BC}|$  is dominated by

$$|\widehat{\Sigma}^{SZ}| = \min \left\{ \frac{(n - p + 2)!}{(n + 2)!} |\mathbf{S}|, \frac{(n + m - p + 2)!}{(n + m + 2)!} |\mathbf{S} + \mathbf{X}^\top \mathbf{X}| \right\}$$

relative to the  $L_G$ -loss. Clearly, the probability of  $|\widehat{\Sigma}^{SZ}| \leq |\widehat{\Sigma}^{BC}|$  is one, and hence  $|\widehat{\Sigma}^{SZ}|$  is shrinking  $|\widehat{\Sigma}^{BC}|$  toward the zero.

A different approach to proving the above dominance result is given by Sinha (1976). Rukhin and Sinha (1990) provided another dominance result without using the information on  $\mathbf{X}$ . Some results under an entropy-type loss are obtained by Sinha and Ghosh (1987) and Kubokawa and Srivastava (2003). On the other hand, a dominance result in the case of  $p > n$  is still not known.

#### 7.6.4 Estimation of the Precision Matrix

Assume now that  $n \geq p$ . Recall that  $\mathbf{S} = \mathbf{Y}^\top \mathbf{Y} \sim \mathcal{W}_p(n, \Sigma)$ . For any constant matrix  $\mathbf{A} \in \mathbb{S}_p$ , an application of the Haff identity (5.5) to  $\text{tr } \Sigma^{-1} \mathbf{A}$  yields

$$\text{tr } \Sigma^{-1} \mathbf{A} = E[(n - p - 1) \text{tr } \mathbf{S}^{-1} \mathbf{A} + 2 \text{tr } \mathbf{D}_S \mathbf{A}] = E[(n - p - 1) \text{tr } \mathbf{S}^{-1} \mathbf{A}].$$

From the arbitrariness of  $\mathbf{A}$ , we obtain  $\Sigma^{-1} = E[(n - p - 1) \mathbf{S}^{-1}]$ , so that

$$\widehat{\Sigma}_{UB}^{-1} = (n - p - 1) \mathbf{S}^{-1}$$

is the unbiased estimator of  $\Sigma^{-1}$ .

The inverse of the covariance matrix  $\Sigma$  is commonly called the precision matrix. The estimation problem of the precision matrix  $\Sigma^{-1}$  has been studied since Efron and Morris (1976). They pointed out that certain empirical Bayes estimation for a normal mean matrix is closely related to the problem of estimating  $\Sigma^{-1}$  under a quadratic-type loss

$$L_{EM}(\widehat{\Sigma}^{-1}, \Sigma^{-1} | \mathbf{S}) = \text{tr} (\widehat{\Sigma}^{-1} - \Sigma^{-1})^2 \mathbf{S}.$$

Here, we briefly introduce a unified approach to the  $n \geq p$  and the  $p > n$  cases based on the Efron-Morris (1976) estimator.

Let

$$\widehat{\Sigma}_{BS}^{-1} = a_0 \mathbf{S}^+, \quad a_0 = |n - p| - 1.$$

Among estimators of the form  $a \mathbf{S}^+$  with positive constant  $a$ ,  $\widehat{\Sigma}_{BS}^{-1}$  is the best estimator relative to the  $L_{EM}$ -loss, and, when  $n \geq p$ ,  $\widehat{\Sigma}_{BS}^{-1}$  is equivalent to  $\widehat{\Sigma}_{UB}^{-1}$ . To improve  $\widehat{\Sigma}_{BS}^{-1}$ , we define the Efron-Morris (1976) type estimator as

$$\widehat{\Sigma}_{EM}^{-1} = \widehat{\Sigma}_{BS}^{-1} + \frac{(\nu - 1)(\nu + 2)}{\text{tr } \mathbf{S}} \mathbf{S}^+$$

with  $\nu = n \wedge p$ . If  $n \geq p$ , then  $\widehat{\Sigma}_{EM}^{-1}$  is the same as in Efron and Morris (1976). Denote by  $\mathbf{S} = \mathbf{H} \mathbf{L} \mathbf{H}^\top$  the eigenvalue decomposition of  $\mathbf{S}$ , where  $\mathbf{L} \in \mathbb{D}_\nu^{(\geq 0)}$  and  $\mathbf{H} \in \mathbb{V}_{p, \nu}$ , and then note that

$$\widehat{\Sigma}_{EM}^{-1} = \mathbf{H} \Phi^{EM} \mathbf{H}^\top, \quad \Phi^{EM} = a_0 \mathbf{L}^{-1} + \frac{(\nu - 1)(\nu + 2)}{\text{tr } \mathbf{L}} \mathbf{I}_\nu.$$

**Proposition 7.11**  $\widehat{\Sigma}_{EM}^{-1}$  dominates  $\widehat{\Sigma}_{BS}^{-1}$  relative to the  $L_{EM}$ -loss.

The proof of Proposition 7.11 can be provided by an unbiased risk estimate method and it is omitted.

For decision-theoretic estimation of the precision matrix with  $n \geq p$ , other procedures for improving the unbiased estimator can be found in Haff (1977, 1979a, b), Dey (1987) and Dey et al. (1990). An improving method via using information on means was considered by Sinha and Ghosh (1987). Eaton and Olkin (1987) and Krishnamoorthy and Gupta (1989) provided minimax estimators based on the Cholesky decomposition of the Wishart matrix  $\mathbf{S}$  relative to the Stein-type loss

$$L_P(\widehat{\Sigma}^{-1}, \Sigma^{-1}) = \text{tr } \Sigma \widehat{\Sigma}^{-1} - \log |\Sigma \widehat{\Sigma}^{-1}| - p.$$

See also Zhou et al. (2001) and Tsukuma (2014b) for related works to improved minimax estimation. Orthogonally invariant minimax estimators are obtained by Perron (1997) and Kubokawa (2005) for  $p = 2$  and by Sheena (2003) for  $p = 3$ . However, for  $p \geq 4$ , orthogonally invariant minimax estimators are still not provided.

The  $p > n$  case is treated by Kubokawa and Srivastava (2008), who propose some improved estimators under quadratic-type losses.

## References

- D.K. Dey, Improved estimation of a multinormal precision matrix. *Stat. Probab. Lett.* **6**, 125–128 (1987)
- D.K. Dey, M. Ghosh, C. Srinivasan, A new class of improved estimators of a multinormal precision matrix. *Stat. Decisions* **8**, 141–151 (1990)
- D.K. Dey, C. Srinivasan, Estimation of a covariance matrix under Stein's loss. *Ann. Stat.* **13**, 1581–1591 (1985)
- M.L. Eaton, Some problems in covariance estimation. Technical Reports No. 49, (Department of Statistics, Stanford University, 1970)
- M.L. Eaton, I. Olkin, Best equivariant estimators of a Cholesky decomposition. *Ann. Stat.* **15**, 1639–1650 (1987)
- B. Efron, C. Morris, Multivariate empirical Bayes and estimation of covariance matrices. *Ann. Stat.* **4**, 22–32 (1976)
- A.K. Gupta, S. Ofori-Nyarko, Improved minimax estimators of normal covariance and precision matrices. *Statistics* **26**, 19–25 (1995)
- L.R. Haff, Minimax estimators for a multinormal precision matrix. *J. Multivar. Anal.* **7**, 374–385 (1977)
- L.R. Haff, Estimation of the inverse covariance matrix: Random mixtures of the inverse Wishart matrix and the identity. *Ann. Stat.* **7**, 1264–1276 (1979a)
- L.R. Haff, An identity for the Wishart distribution with applications. *J. Multivar. Anal.* **9**, 531–544 (1979b)
- L.R. Haff, Empirical Bayes estimation of the multivariate normal covariance matrix. *Ann. Stat.* **8**, 586–597 (1980)
- L.R. Haff, The variational form of certain Bayes estimators. *Ann. Stat.* **19**, 1163–1190 (1991)
- W. James, C. Stein, Estimation with quadratic loss, in *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, ed. by J. Neyman, (University of California Press, Berkeley, 1961), pp. 361–379
- J. Kiefer, Invariance, minimax sequential estimation, and continuous time processes. *Ann. Math. Stat.* **28**, 573–601 (1957)
- Y. Konno, Shrinkage estimators for large covariance matrices in multivariate real and complex normal distributions under an invariant quadratic loss. *J. Multivar. Anal.* **100**, 2237–2253 (2009)
- K. Krishnamoorthy, A.K. Gupta, Improved minimax estimation of a normal precision matrix. *Can. J. Stat.* **17**, 91–102 (1989)
- T. Kubokawa, A revisit to estimating of the precision matrix of the Wishart distribution. *J. Stat. Res.* **39**, 91–114 (2005)
- T. Kubokawa, Y. Konno, Estimating the covariance matrix and the generalized variance under a symmetric loss. *Ann. Inst. Stat. Math.* **42**, 331–343 (1990)
- T. Kubokawa, C. Robert, AKMdE Saleh, Empirical Bayes estimation of the variance parameter of a normal distribution with unknown mean under an entropy loss. *Sankhyā Ser. A* **54**, 402–410 (1992)
- T. Kubokawa, M.S. Srivastava, Estimating the covariance matrix: a new approach. *J. Multivar. Anal.* **86**, 28–47 (2003)
- T. Kubokawa, M.S. Srivastava, Estimation of the precision matrix of a singular Wishart distribution and its application in high-dimensional data. *J. Multivar. Anal.* **99**, 1906–1928 (2008)
- T. Kubokawa, M.-T. Tsai, Estimation of covariance matrices in fixed and mixed effects linear models. *J. Multivar. Anal.* **97**, 2242–2261 (2006)
- O. Ledoit, M. Wolf, A well-conditioned estimator for large-dimensional covariance matrices. *J. Multivar. Anal.* **88**, 365–411 (2004)
- P.L. Leung, W.Y. Chan, Estimation of the scale matrix and its eigenvalues in the Wishart and the multivariate F distributions. *Ann. Inst. Stat. Math.* **50**, 523–530 (1998)

- P.L. Leung, F.Y. Ng, Improved estimation of a covariance matrix in an elliptically contoured matrix distribution. *J. Multivar. Anal.* **88**, 131–137 (2004)
- S.P. Lin, M.D. Perlman, A monte carlo comparison of four estimators for a covariance matrix, in *Multivar. Anal. VI*, ed. by P.R. Krishnaiah (North-Holland, Amsterdam, 1985), pp. 411–429
- T. Ma, L. Jia, Y. Su, A new estimator of covariance matrix. *J. Stat. Plan. Infer.* **142**, 529–536 (2012)
- M.D. Perlman, Reduced mean square error estimation for several parameters. *Sankhyā Ser. B* **34**, 89–92 (1972)
- F. Perron, Equivariant estimators of the covariance matrix. *Can. J. Stat.* **18**, 179–182 (1990)
- F. Perron, Minimax estimators of a covariance matrix. *J. Multivar. Anal.* **43**, 16–28 (1992)
- F. Perron, On a conjecture of Krishnamoorthy and Gupta. *J. Multivar. Anal.* **62**, 110–120 (1997)
- T. Robertson, F.T. Wright, R.L. Dykstra, *Order Restricted Statistical Inference* (Wiley, New York, 1988)
- A.L. Rukhin, B.K. Sinha, Decision-theoretic estimation of the product of gamma scales and generalized variance. *Calcutta Stat. Assoc. Bull.* **40**, 257–265 (1990)
- D. Sharma, K. Krishnamoorthy, Orthogonal equivariant minimax estimators of bivariate normal covariance matrix and precision matrix. *Calcutta Stat. Assoc. Bull.* **32**, 23–46 (1983)
- J.B. Selliah, Estimation and testing problems in a Wishart distribution. Technical reports No.10 (Department of Statistics, Stanford University, 1964)
- Y. Sheena, On minimaxity of the normal precision matrix estimator of Krishnamoorthy and Gupta. *Statistics* **37**, 387–399 (2003)
- Y. Sheena, A. Takemura, Inadmissibility of non-order-preserving orthogonally invariant estimators of the covariance matrix in the case of Stein's loss. *J. Multivar. Anal.* **41**, 117–131 (1992)
- R.B. Shorrock, J.V. Zidek, An improved estimator of the generalized variance. *Ann. Stat.* **4**, 629–638 (1976)
- B.K. Sinha, On improved estimators of the generalized variance. *J. Multivar. Anal.* **6**, 617–626 (1976)
- B.K. Sinha, M. Ghosh, Inadmissibility of the best equivariant estimators of the variance-covariance matrix, the precision matrix, and the generalized variance under entropy loss. *Stat. Dec.* **5**, 201–227 (1987)
- C. Stein, Some problems in multivariate analysis, Part I. Technical Reports No.6 (Department of Statistics, Stanford University, 1956)
- C. Stein, Inadmissibility of the usual estimator for the variance of a normal distribution with unknown mean. *Ann. Inst. Stat. Math.* **16**, 155–160 (1964)
- C. Stein, *Estimation of a Covariance Matrix, Rietz Lecture, 39th Annual Meeting IMS* (Atlanta, GA, 1975)
- C. Stein, Lectures on the theory of estimation of many parameters, in *Proceedings of Scientific Seminars of the Steklov Institute Studies in the Statistical Theory of Estimation, Part I*, vol. 74, eds. by I.A. Ibragimov, M.S. Nikulin (Leningrad Division, 1977), pp. 4–65
- W.E. Strawderman, Minimaxity. *J. Am. Stat. Assoc.* **95**, 1364–1368 (2000)
- D. Sun, X. Sun, Estimation of the multivariate normal precision and covariance matrices in a star-shape model. *Ann. Inst. Stat. Math.* **57**, 455–484 (2005)
- A. Takemura, An orthogonally invariant minimax estimator of the covariance matrix of a multivariate normal population. *Tsukuba J. Math.* **8**, 367–376 (1984)
- H. Tsukuma, Minimax covariance estimation using commutator subgroup of lower triangular matrices. *J. Multivar. Anal.* **124**, 333–344 (2014a)
- H. Tsukuma, Improvement on the best invariant estimators of the normal covariance and precision matrices via a lower triangular subgroup. *J. Jpn. Stat. Soc.* **44**, 195–218 (2014b)
- H. Tsukuma, Estimation of a high-dimensional covariance matrix with the Stein loss. *J. Multivar. Anal.* **148**, 1–17 (2016a)
- H. Tsukuma, Minimax estimation of a normal covariance matrix with the partial Iwasawa decomposition. *J. Multivar. Anal.* **145**, 190–207 (2016b)
- H. Tsukuma, T. Kubokawa, Minimaxity in estimation of restricted and non-restricted scale parameter matrices. *Ann. Inst. Stat. Math.* **67**, 261–285 (2015)

- H. Tsukuma, T. Kubokawa, Unified improvements in estimation of a normal covariance matrix in high and low dimensions. *J. Multivar. Anal.* **143**, 233–248 (2016)
- R. Yang, J.O. Berger, Estimation of a covariance matrix using the reference prior. *Ann. Stat.* **22**, 1195–1211 (1994)
- X. Zhou, X. Sun, J. Wang, Estimation of the multivariate normal precision matrix under the entropy loss. *Ann. Inst. Stat. Math.* **53**, 760–768 (2001)

# Index

## A

Admissibility, 2  
Affine transformation, 3, 106  
Anti-diagonal matrix, 84

## B

Bartlett decomposition, 23

## C

Chi-square identity, 37, 83  
Cholesky decomposition, 11, 15, 16, 23, 79, 104, 107  
Commutator, 81  
Commutator subgroup, 81

## D

Decision space, 31  
Descending wedge symbol, 12  
Digamma function, 78

## E

Eigenvalue, 12, 89  
Eigenvalue decomposition, 11, 15, 39, 50, 60, 68, 84, 89, 107  
Eigenvector, 12  
Elliptically contoured distribution, 67  
Empirical Bayes method, 3, 45, 48, 86  
Error covariance, 29, 75  
Error matrix, 27, 64, 67  
Exterior product, 14

## G

Gauss divergence theorem, 37, 40  
Generalized inverse, 9, 52  
Generalized variance, 105  
GMANOVA model, 63  
    canonical form, 64  
Group invariance, 31  
Growth curve model, 28, 64

## H

Haff identity, 37, 106

## I

Inadmissibility, 2  
Invariance, 31  
    of estimation problem, 31  
    of estimator, 31  
    of loss function, 31  
    of multivariate linear model, 31  
    of statistical model, 31  
Isotonic regression, 91, 93

## J

Jacobian, 14  
Jacobian matrix, 14  
James-Stein estimator  
    of covariance matrix, 79  
    of mean vector, 2, 103

## K

Kronecker delta, 37, 38

Kronecker product, 10

## L

Landau symbol, 40

LDL<sup>T</sup> decomposition, 11, 81, 82, 102

Least squares estimator, 28

Least squares method, 28

Linear shrinkage estimator, 104

Loss function

definition, 1

loss matrix

of mean matrix estimation, 62

quadratic loss

of covariance estimation, 104

of generalized variance estimation,  
105

of mean matrix estimation, 31, 47, 64

of mean vector estimation, 2

of precision matrix estimation, 106

Stein loss, 76, 102, 104

extended, 76

Lower triangular group, 8, 79

Löwner order, 8, 62, 87

LQ decomposition, 11, 14, 23

LU decomposition, 11

## M

MANOVA model, 28

Matricial shrinkage estimator, 46

Matrix decomposition, 11

Matrix differential operator, 35, 68

Matrix factorization, 11

Matrix square root, 8, 53

Matrix transformation, 14

Matrix-variate normal distribution, 17

Maximum likelihood estimator, 3, 28, 45

Minimaxity, 2, 45, 47, 54, 76, 81, 104, 107

invariance approach, 81

least favorable prior approach, 81

Moore-Penrose inverse, 9, 49

differential, 37

Multivariate gamma function, 16

Multivariate linear model, 27

canonical form, 30, 46, 75

Multivariate normal distribution, 13

## O

Ordering property, 89, 93, 96

Order-preserving, 89, 93, 96

Orthogonal group, 8

Orthogonally invariant, 84, 89, 104, 107

Orthogonal projection matrix, 29, 52

## P

Positive-part rule, 46, 59

Precision matrix, 46, 106

## Q

QR decomposition, 11, 29, 30

## R

Regression coefficient matrix, 27, 64, 75

Residual sum of squares matrix, 29

Risk function, 2

## S

Scalar shrinkage estimator, 46

Scale invariant, 77, 79, 104

Shrinkage estimator, 3

Shrinkage property, 89, 96

Singular value, 12

Singular value decomposition, 12, 15, 21, 50,  
61

Skew-symmetric, 69

Spectral decomposition, 12

Stein identity, 4, 35, 63, 66, 67, 98

Stein's lemma, 4

Stiefel manifold, 8

## T

Transformation group, 31

Triangular invariant, 79, 84

## U

Unbiased risk estimate, 4, 54, 55, 84, 90, 91

## V

Vec operator, 11

## W

Wedge symbol, 12

Wishart distribution, 21, 37

pseudo-Wishart distribution, 21

Wishart matrix, 21, 23, 30, 47, 75, 102, 104,  
107