

EMPIRICAL BAYES ESTIMATION OF THE MULTIVARIATE NORMAL COVARIANCE MATRIX

BY L. R. HAFF

University of California, San Diego, La Jolla

Let $S_{p \times p}$ have a Wishart distribution with scale matrix Σ and k degrees of freedom. Estimators of Σ are given for each of the loss functions $L_1(\hat{\Sigma}, \Sigma) = \text{tr}(\hat{\Sigma}\Sigma^{-1}) - \log \det(\hat{\Sigma}\Sigma^{-1}) - p$ and $L_2(\hat{\Sigma}, \Sigma) = \text{tr}(\hat{\Sigma}\Sigma^{-1} - I)^2$. The obvious estimators of Σ are the scalar multiples of S , i.e., aS where $0 < a < 1/k$. (Recall that $(1/k)S$ is unbiased.) For each problem $(\Sigma, \hat{\Sigma}, L_i)$, $i = 1, 2$, we provide empirical Bayes estimators which dominate aS by a substantial amount. It is seen that the uniform reduction in the risk function determined by L_2 is at least $100(p+1)/(k+p+1)\%$. Dominance results for L_1 and L_2 were first given by James and Stein.

1. Introduction and summary. The problem under consideration here is that of estimating the multi-normal covariance matrix Σ . Empirical Bayes (EB) alternatives are derived which dominate *all* scalar multiples of the unbiased estimator. The risk of the unbiased estimator is reduced by a significant amount in each case.

In [3] and [4], the author treated the problem of estimating Σ^{-1} . An identity for the Wishart distribution was derived and used to compute an unbiased estimator for the risk function. Dominance results were then obtained by working with the unbiased estimator. The same technique is widely used in the present paper. (The Wishart identity is stated in line (2.4) below. Unknown to the author, Charles Stein essentially derived it several years ago. His approach was different, however, and his work remains unpublished—see [12].)

Our present results have the flavor of those given by James and Stein [6], pages 376, 377. We provide estimators which dominate the usual unbiased estimator for each of two invariant loss functions. Our present methods, however, supersede those of [6], and the implications are deeper. Whereas the estimator in [6] is only slightly better than the usual one, the EB estimators are substantially better.

Let $S_{p \times p}$ have a Wishart distribution with unknown matrix Σ and k degrees of freedom, i.e.,

$$(1.1) \quad S \sim W(\Sigma, k) \text{ for } k - p - 1 > 0.$$

Also, let $\hat{\Sigma}$ be an estimator of Σ . We assume that the loss function is

$$(1.2) \quad L_1(\hat{\Sigma}, \Sigma) = \text{tr}(\hat{\Sigma}\Sigma^{-1}) - \log \det(\hat{\Sigma}\Sigma^{-1}) - p \text{ or } L_2(\hat{\Sigma}, \Sigma) = \text{tr}(\hat{\Sigma}\Sigma^{-1} - I)^2$$

(from [6]), and define the risk function by

$$R_i(\hat{\Sigma}, \Sigma) \equiv E[L_i(\hat{\Sigma}, \Sigma)|\Sigma, k], \quad i = 1 \text{ or } 2.$$

Received January 1978; revised February 1979.

AMS 1970 subject classifications. Primary 62F10; secondary 62C99.

Key words and phrases. Covariance matrix, empirical Bayes estimators, unbiased estimation of risk function.

(The latter is an average with respect to the $W(\Sigma, k)$ distribution.) If $\hat{\Sigma}$ and $\hat{\Sigma}_*$ are competing estimators of Σ , then " $\hat{\Sigma}$ dominates $\hat{\Sigma}_*$ (mod L_i)" will mean $R_i(\hat{\Sigma}, \Sigma) < R_i(\hat{\Sigma}_*, \Sigma)$ ($\forall \Sigma$).

Our EB estimators have the form

$$(1.3) \quad \hat{\Sigma} = a[\mathbf{S} + \mathbf{u}t(\mathbf{u})C]$$

with $0 < a < 1/k$, $\mathbf{u} = 1/\text{tr}(\mathbf{S}^{-1}C)$, $t(\cdot)$ nonincreasing, and C an arbitrary positive definite matrix. For $t(\cdot) \equiv 0$, we have the obvious estimators, the scalar multiples of \mathbf{S} . In particular, note that $E(1/k)\mathbf{S} = \Sigma$.

For loss function L_1 , Stein [6] derived a minimax estimator and proved that for any $\beta > 0$, $\beta(\mathbf{S}/k)$ is not minimax. He obtained similar results for L_2 , but was unable to get an explicit formula for a minimax estimator in this case. Selliah [10] gave additional results for L_2 . Later, Perlman [8] showed that if $(kp - 2)/(kp + 2) < \beta < 1$, then $\beta(\mathbf{S}/k)$ dominates \mathbf{S}/k with respect to a general quadratic loss function. All these estimators (from [6], [8], and [10]) are only slightly better than the unbiased estimator. Finally, in the 1975 Rietz lecture [12], Stein described one (for L_1) which is substantially better. His methods were similar to those used in the following.

A summary of the present work. Among the scalar multiples of \mathbf{S} , we prove that the best estimator (mod L_1) is the unbiased estimator

$$(1.4) \quad \hat{\Sigma}_1 = (1/k)\mathbf{S},$$

and the best estimator (mod L_2) is

$$\hat{\Sigma}_2 = [1/(k + p + 1)]\mathbf{S}.$$

Our main results concern the estimators in line (1.3). These, (1.3), are derived as empirical Bayes estimators. Then, for each loss function, conditions are given under which they dominate the best scalar multiple of \mathbf{S} . Under L_2 , we have the following numerical comparisons: it is seen that $R_2(\hat{\Sigma}_1, \Sigma) = p(p + 1)/k$ and $R_2(\hat{\Sigma}, \Sigma) < R_2(\hat{\Sigma}_2, \Sigma) = p(p + 1)/(k + p + 1)$ ($\forall \Sigma$). Thus $R_2(\hat{\Sigma}, \Sigma)$ is less than $R_2(\hat{\Sigma}_1, \Sigma)$ by at least $100(p + 1)/(k + p + 1)\%$ ($\forall \Sigma$). For $p = 3$ and $k = 8$, say, the uniform improvement is at least 33%. In addition, our Monte Carlo results (unpublished) indicate that $R_2(\hat{\Sigma}, \Sigma)$ is less than $R_2(\hat{\Sigma}_2, \Sigma)$ by as much as 16% ($p = 3$ and $k = 8$).

Finally, there is a parallel between our results and those of Stein [6]. Roughly speaking, the EB estimators which perform well (mod L_1) also perform well (mod L_2). Although we give explicit estimators under L_2 , the calculations are difficult, and the results lack the generality of the L_1 results.

2. A brief description of the mathematics. Let δ be a real number, $V_{p \times p}$ a symmetric matrix. If we apply the expansion

$$(2.1) \quad \log \det(I + \delta V) = \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} \delta^n \text{tr}(V^n)$$

to L_1 , then a relationship is obtained between L_1 and L_2 . (The series converges if the spectral radius of V is less than unity and $0 < \delta < 1$.) In particular, set $\delta = 1$, factor Σ^{-1} as $\Sigma^{-1} = \Omega^2$, and expand

$$\begin{aligned}
 \log \det(\hat{\Sigma}\Sigma^{-1}) &= \log \det(\Omega\hat{\Sigma}\Omega) \\
 &= \log \det(I + V) \quad (V = \Omega\hat{\Sigma}\Omega - I) \\
 (2.2) \quad &= \text{tr}(V) - \left(\frac{1}{2}\right)\text{tr}(V^2) + \left(\frac{1}{3}\right)\text{tr}(V^3) - \dots \\
 &= \text{tr}(\hat{\Sigma}\Sigma^{-1} - I) - \left(\frac{1}{2}\right)\text{tr}(\hat{\Sigma}\Sigma^{-1} - I)^2 + \dots
 \end{aligned}$$

From (2.2) loss function L_1 can be written as

$$(2.3) \quad L_1(\hat{\Sigma}, \Sigma) = \left(\frac{1}{2}\right)L_2(\hat{\Sigma}, \Sigma) - \left(\frac{1}{3}\right)\text{tr}(\hat{\Sigma}\Sigma^{-1} - I)^3 + \dots,$$

so it is plausible that estimators which perform well (mod L_1) also perform well (mod L_2).

The series (2.1) is used to obtain an approximation for $\alpha_1(\Sigma) \equiv R_1(\hat{\Sigma}, \Sigma) - R_1(\hat{\Sigma}_1, \Sigma)$. This approximation, call it $\underline{\alpha}_1(\Sigma)$, has the property that $\underline{\alpha}(\Sigma) < 0$ ($\forall \Sigma$) implies $\alpha_1(\Sigma) < 0$ ($\forall \Sigma$). We shall force $\underline{\alpha}_1(\Sigma) < 0$ ($\forall \Sigma$) by using an identity for the Wishart distribution.

The identity and its application. For suitable choices of a matrix $T = T(\mathbf{S}, \Sigma)_{p \times p}$, a scalar $h(\mathbf{S})$, and a constant $r \neq 0$; the identity is stated in terms of the following:

- (i) $\text{diag}(T)$, a diagonal matrix with diagonal elements equal to those of T ;
- (ii) $T_{(r)} \equiv rT + (1 - r) \text{diag}(T)$;
- (iii) $D^*T_{(r)} \equiv \sum_{i=1}^p \partial t_{ii} / \partial s_{ii} + r \sum_{i \neq j} \partial t_{ij} / \partial s_{ij}$; and
- (iv) $\partial h / \partial \mathbf{S} \equiv (\partial h / \partial s_{ij})_{p \times p}$.

The identity is given by

$$\begin{aligned}
 (2.4) \quad E[h(\mathbf{S})\text{tr}(T\Sigma^{-1})] &= 2E\left[h(\mathbf{S})D^*T_{\left(\frac{1}{2}\right)}\right] + 2E\text{tr}\left[\frac{\partial h(\mathbf{S})}{\partial \mathbf{S}} \cdot T_{\left(\frac{1}{2}\right)}\right] \\
 &\quad + (k - p - 1)E[h(\mathbf{S})\text{tr}(\mathbf{S}^{-1}T)].
 \end{aligned}$$

First, we apply (2.4) to risk function R_1 . It happens that $\underline{\alpha}_1(\Sigma)$ has terms under the expectation of the form $h(\mathbf{S})\text{tr}(T\Sigma^{-1})$ with $T = T(\mathbf{S})$ (see the left side of (2.4)). Consequently, we readily obtain $\hat{\underline{\alpha}}_1(\Sigma)$, an unbiased estimator of $\underline{\alpha}_1(\Sigma)$. Theorem 4.3 gives conditions under which $\hat{\underline{\alpha}}_1(\Sigma) < 0$ ($\forall \mathbf{S}$). The inequality implies $\underline{\alpha}_1(\Sigma) < 0$ ($\forall \Sigma$). Risk function R_2 introduces additional difficulties. In particular, $\alpha_2(\Sigma) \equiv R_2(\hat{\Sigma}, \Sigma) - R_2(\hat{\Sigma}_2, \Sigma)$ has terms which are quadratic in Σ^{-1} (not merely linear like $\underline{\alpha}_1$). We obtain an unbiased estimator for $\alpha_2(\Sigma)$ by a method involving repeated application of (2.4).

Various combinations of $h(\mathbf{S})$ and $T(\mathbf{S}, \Sigma)$ are used in the following. For each, the identity (2.4) is verified by straightforward modification of [3], pages 377-379. Most of these details are omitted. However, some general conditions for the validity of (2.4) are established in Haff [5].

3. Empirical Bayes estimators. In this section, we describe the empirical Bayes character of the estimators

$$(3.1) \quad \hat{\Sigma} = a[\mathbf{S} + \mathbf{u}t(\mathbf{u})C]$$

where $0 < a < 1/k$, $\mathbf{u} = 1/\text{tr}(\mathbf{S}^{-1}C)$, $t(\mathbf{u}) \searrow$, and $C_{p \times p} > 0$. Our dominance results are stated in Section 4. Together, these sections provide a basis for choosing a , $t(\mathbf{u})$, and C .

Assume that

$$(3.2) \quad \mathbf{S} \sim W_{p \times p}(\Sigma, k) \text{ and } \Sigma^{-1} \sim W_{p \times p}[(1/\gamma)C^{-1}, k']$$

with γ unknown, C p.d. and known, and $k' > 0$ a known integer. It is seen that

$$(3.3) \quad \Sigma^{-1}|\mathbf{S} \sim W_{p \times p}[(\mathbf{S} + \gamma C)^{-1}, k + k']$$

$$\text{and } \Sigma|\mathbf{S} \sim W_{p \times p}^{-1}[\mathbf{S} + \gamma C, (k + k') + p + 1].$$

Here $W_{p \times p}^{-1}$ indicates the inverse Wishart distribution—see Press [9], pages 109–112 for a useful description. From (3.3), the posterior mean of Σ is

$$(3.4) \quad E(\Sigma|\mathbf{S}, \gamma) = a(\mathbf{S} + \gamma C), \quad a = 1/(k + k' - p - 1).$$

Our estimator in (3.1) is a modification of the posterior mean—being “empirical Bayes” in the sense that $\mathbf{u}t(\mathbf{u})$ is an estimator of γ . We now derive $\hat{\gamma} = \mathbf{u}t(\mathbf{u})$ as a generalized maximum likelihood estimator. From (3.2), the marginal density of \mathbf{S} is

$$(3.5) \quad f(\mathbf{S}|\gamma) \propto \gamma^{pk'/2} |\mathbf{S}|^{(k-p-1)/2} |\mathbf{S} + \gamma C|^{-(k+k')/2} \text{ for } \mathbf{S} > 0;$$

hence, the likelihood function is

$$l(\gamma|\mathbf{S}) = \gamma^{pk'/2} [\det(\mathbf{I} + \gamma C \mathbf{S}^{-1})]^{-(k+k')/2}.$$

Let

$$\begin{aligned} l^* &= \log l(\gamma|\mathbf{S}) \\ &= \frac{pk'}{2} \log \gamma - \frac{(k + k')}{2} \log \det(\mathbf{I} + \gamma C \mathbf{S}^{-1}) \\ &= \frac{pk'}{2} \log \gamma - \frac{(k + k')}{2} \log \det(\mathbf{I} + \gamma C^{\frac{1}{2}} \mathbf{S}^{-1} C^{\frac{1}{2}}) \\ &= \frac{pk'}{2} \log \gamma - \frac{(k + k')}{2} \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} \gamma^n \text{tr}(C^{\frac{1}{2}} \mathbf{S}^{-1} C^{\frac{1}{2}})^n \end{aligned}$$

[see (2.1)]. We shall proceed on a formal basis and discuss convergence later on. The first order approximation

$$(3.6) \quad l^* \doteq \frac{pk'}{2} \log \gamma - \frac{(k + k')}{2} \gamma \text{tr}(C^{\frac{1}{2}} \mathbf{S}^{-1} C^{\frac{1}{2}})$$

has a maximum at

$$(3.7) \quad \begin{aligned} \hat{\gamma} &= c^*/\text{tr}(C^{\frac{1}{2}} \mathbf{S}^{-1} C^{\frac{1}{2}}) \\ &= c^* \mathbf{u} \end{aligned}$$

where $c^* = pk/(k + k')$ and $u = 1/\text{tr}(S^{-1}C)$. Because of the approximation in (3.6), it is plausible that the constant c^* may be improved upon. Actually, functions $t(u)$ are given for which $\hat{\Sigma}$ dominates the best scalar multiple of S . Among these are constant functions which outperform $t(u) \equiv c^*$. (See Theorem 4.3 and 4.4.)

We note that the expansion of I^* converges at any $\hat{\gamma} = ut(u)$, $0 \leq t(u) < 1$; i.e.,

$$\sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} (\hat{\gamma})^n \text{tr} (C^{\frac{1}{2}} S^{-1} C^{\frac{1}{2}})^n = \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} t^n(u) \text{tr} \left[\frac{C^{\frac{1}{2}} S^{-1} C^{\frac{1}{2}}}{\text{tr} (C^{\frac{1}{2}} S^{-1} C^{\frac{1}{2}})} \right]^n,$$

and the series converges since $t^n(u) \searrow 0$ and the matrix inside brackets has spectral radius less than 1.

4. The main results. We now state the main results of this paper. The proofs are given in Section 5.

THEOREM 4.1. *Under loss function L_1 , the best estimator of the form aS is $\hat{\Sigma}_1 = (1/k)S$.*

THEOREM 4.2. *Under loss function L_2 , the best estimator of the form aS is $\hat{\Sigma}_2 = [1/(k + p + 1)]S$.*

Our main results concern the estimators

$$(4.1) \quad \hat{\Sigma} = a[S + ut(u)C], \quad a = 1/(k + k' - p - 1).$$

For the purpose of comparing $\hat{\Sigma}$ with $\hat{\Sigma}_i$, $i = 1, 2$, we shall set k' equal to $p + 1$ and $2(p + 1)$. Note that each problem $(\Sigma, \hat{\Sigma}, L_i)$, $i = 1, 2$, is invariant under the transformations $S \rightarrow ASA'$, $\Sigma \rightarrow A\Sigma A'$, $\hat{\Sigma} \rightarrow A\hat{\Sigma}A'$, where A is an arbitrary non-singular matrix. In particular, let $A = C^{-\frac{1}{2}}$. Since $R_i(\hat{\Sigma}, \Sigma) = R_i(A\hat{\Sigma}A', A\Sigma A)$, we assume without loss in generality that $C = I$.

THEOREM 4.3. *In (4.1), let $\hat{\Sigma}$ be given by*

- (i) $k' = p + 1$;
- (ii) C an arbitrary p.d. matrix; and
- (iii) $t(u)$ an absolutely continuous and nonincreasing function, $0 \leq t(u) \leq 2(p - 1)/k$.

Then $\hat{\Sigma}$ dominates $\hat{\Sigma}_1$ (mod L_1), i.e., $R_1(\hat{\Sigma}, \Sigma) \leq R_1(\hat{\Sigma}_1, \Sigma) (\forall \Sigma)$.

If t is a constant, then an optimal value of t is $(p - 1)/k$ (as seen from the proof).

Corollary 4.4 is useful when we have an a priori upper bound on γ , say $\bar{\gamma}$.

COROLLARY 4.4. *Assume the conditions of Theorem 4.3 with*

$$(4.2) \quad \begin{aligned} t(u) &= \bar{\gamma}/u \text{ for } u \geq k\bar{\gamma}/(p - 1) \\ &= (p - 1)/k \text{ otherwise.} \end{aligned}$$

Then $R_1(\hat{\Sigma}, \Sigma) \leq R_1(\hat{\Sigma}_1, \Sigma) (\forall \Sigma)$.

According to (4.2), the estimator $ut(u)$ never exceeds $\bar{\gamma}$. Theorems 4.5 and 4.6 concern loss function L_2 .

THEOREM 4.5. Let $\hat{\Sigma}$ (see equation 4.1) be given by

- (i) $k' = 2(p + 1)$;
- (ii) C an arbitrary p.d. matrix; and
- (iii) $0 \leq t \leq 2(p - 1)/(k - p + 3)$, t a constant.

Then $\hat{\Sigma}$ dominates $\hat{\Sigma}_2 \pmod{L_2}$, i.e., $R_2(\hat{\Sigma}, \Sigma) \leq R_2(\hat{\Sigma}_2, \Sigma) \ (\forall \Sigma)$.

The choice $t = (p - 1)/(k - p + 3)$ is optimal (again, see the proof). As mentioned before, the R_2 calculations are relatively difficult, so our present Theorem 4.5 lacks the generality of Theorem 4.3. The univariate case of Theorem 4.5 is interesting. For $p = 1$, our estimator becomes $\hat{\Sigma} = \hat{\Sigma}_2 = \mathbf{S}/(k + 2)$, the unique admissible minimax estimator of Σ (see Lehmann [7], pages 4–15).

An analogue of Corollary 4.4 is

THEOREM 4.6. Let $\hat{\Sigma}$ be given by

- (i) $k' = 2(p + 1)$, $k - p - 3 > 0$;
- (ii) C an arbitrary p.d. matrix; and
- (iii) $t(\mathbf{u}) = \bar{\gamma}/\mathbf{u}$ for $\mathbf{u} \geq \bar{\gamma}(k - p + 3)/(p - 1) = (p - 1)/(k - p + 3)$ otherwise.

Then $R_2(\hat{\Sigma}, \Sigma) \leq R_2(\bar{\Sigma}_2, \Sigma) \ (\forall \Sigma)$.

5. Mathematical details.

PROOF OF THEOREM 4.1. Let $\hat{\Sigma}_{1c} = (1/k)(1 + c)\mathbf{S}$, $|c| < 1$, and $\hat{\Sigma}_1 = (1/k)\mathbf{S}$. We show that

$$R_1(\hat{\Sigma}_{1c}, \Sigma) - R_1(\hat{\Sigma}_1, \Sigma) = E[(c/k) \text{tr}(\mathbf{S}\Sigma^{-1}) - p \log(1 + c)] \geq 0 \ (\forall \Sigma).$$

This inequality is true if $c - \sum_{n=1}^{\infty} (-1)^{n+1} (c^n/n) = \sum_{n=2}^{\infty} (-1)^n (c^n/n) > 0$. The latter holds because each positive term dominates its successor; and the proof is complete. \square

PROOF OF THEOREM 4.2. A scalar multiple of \mathbf{S} has risk $R_2(a\mathbf{S}, \Sigma) = E \text{tr}(a\mathbf{S}\Sigma^{-1} - I)^2 = a^2 \text{tr}E(\mathbf{V}^2) - 2kpa + p$ in which $\mathbf{V} = \Sigma^{-\frac{1}{2}}\mathbf{S}\Sigma^{-\frac{1}{2}} \sim W(I, k)$. Thus $R_2(a\mathbf{S}, \Sigma) = kp(k + p + 1)a^2 - 2kpa + 2$, and the latter is minimized at $a = 1/(k + p + 1)$. \square

The remaining calculations depend on special cases of (2.4). In particular, we shall need

- (i) $E[h(\mathbf{S}) \text{tr}(Q\Sigma^{-1})] = \text{tr} E[(k - p - 1)h(\mathbf{S})\mathbf{S}^{-1}Q + 2(\partial h(\mathbf{S})/\partial \mathbf{S}) \cdot Q_{(\frac{1}{2})}]$ where $Q_{p \times p}$ is a matrix of constants,
- (5.1)
- (ii) $E[h(\mathbf{S}) \text{tr}(\mathbf{S}\Sigma^{-1})] = E[pkh(\mathbf{S}) + 2 \text{tr}(\partial h(\mathbf{S})/\partial \mathbf{S}) \cdot \mathbf{S}_{(\frac{1}{2})}]$ and
 - (iii) $E[h(\mathbf{S}) \text{tr}(\mathbf{S}^{-1}Q\Sigma^{-1})] = E\{(k - p - 2)h(\mathbf{S}) \text{tr}(\mathbf{S}^{-2}Q) - h(\mathbf{S}) (\text{tr} \mathbf{S}^{-1})(\text{tr} \mathbf{S}^{-1}Q) + 2 \text{tr}[\partial h(\mathbf{S})/\partial \mathbf{S}) \cdot (\mathbf{S}^{-1}Q)_{(\frac{1}{2})}]\}$ with Q as in part (i).

PROOF OF THEOREM 4.3. Write (4.1) as $\hat{\Sigma} = (1/k)\mathbf{S} + g(\mathbf{S})I$, $g(\mathbf{S}) = (1/k)\mathbf{u}t(\mathbf{u})$, and set $\alpha_1(\Sigma) = R_1(\hat{\Sigma}, \Sigma) - R_1(\hat{\Sigma}_1, \Sigma) = E[g(\mathbf{S})\text{tr}\Sigma^{-1} - \log \det (I + kg(\mathbf{S})\mathbf{S}^{-1})]$. We must show that $\alpha_1 \leq 0$ ($\forall \Sigma$). From (5.1(i)) with $Q = I$, we have

$$\alpha_1(\Sigma) = E \left[(k - p - 1)g(\mathbf{S})\text{tr} \mathbf{S}^{-1} + 2 \text{tr} (\partial g(\mathbf{S})/\partial S) - \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} t^n(\mathbf{u})\text{tr} (\mathbf{S}^{-1}/\text{tr} \mathbf{S}^{-1})^n \right].$$

This series is bounded below by a quadratic; i.e.,

$$\sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} t^n(\mathbf{u})\text{tr} (\mathbf{S}^{-1}/\text{tr} \mathbf{S}^{-1})^n > t(\mathbf{u}) - 1/2t^2(\mathbf{u}),$$

because the terms are in decreasing magnitude. A sufficient condition for $\alpha_1(\Sigma) < 0$ ($\forall \Sigma$) is

$$(5.2) \quad (k - p - 1)g(\mathbf{S})\text{tr} \mathbf{S}^{-1} + 2\text{tr} (\partial g(\mathbf{S})/\partial S) - t(\mathbf{u}) + 1/2 t^2(\mathbf{u}) = - (p + 1)t(\mathbf{u}) + 2 \left[\text{tr} \frac{\partial}{\partial S} \{ \mathbf{u}t(\mathbf{u}) \} \right] + (k/2)t^2(\mathbf{u}) < 0.$$

From [4], it can be seen that

$$\text{tr} \frac{\partial}{\partial S} \{ \mathbf{u}t(\mathbf{u}) \} = [\mathbf{u}t'(\mathbf{u}) + t(\mathbf{u})] \left[\text{tr} \mathbf{S}^{-2}/(\text{tr} \mathbf{S}^{-1})^2 \right].$$

Finally, since $t'(\mathbf{u}) \leq 0$ and $\text{tr} \mathbf{S}^{-2}/(\text{tr} \mathbf{S}^{-1})^2 \leq 1$, it is easily seen that $0 < t(\mathbf{u}) < 2(p - 1)/k$ is sufficient for (5.2). \square

If t is a constant, then (5.2) is bounded above by

$$- (p - 1)t + (k/2)t^2,$$

and the latter is minimized at $t = (p - 1)/k$.

Theorems 4.5 and 4.6 concern comparisons between $\hat{\Sigma} = a[\mathbf{S} + \mathbf{u}t(\mathbf{u})I]$ and $\hat{\Sigma}_2 = a\mathbf{S}$, $a = 1/(k + p + 1)$. In particular, for

$$(5.3) \quad \text{and} \quad R_2(\hat{\Sigma}, \Sigma) = R_2(\hat{\Sigma}_2, \Sigma) + \alpha_2(\Sigma)$$

$$\alpha_2(\Sigma) \equiv E[2ag(\mathbf{S})\text{tr} (\mathbf{S}\Sigma^{-2}) - 2g(\mathbf{S})\text{tr}\Sigma^{-1} + g^2(\mathbf{S})\text{tr}\Sigma^{-2}],$$

they give conditions under which $\alpha_2(\Sigma) \leq 0$ ($\forall \Sigma$).

We need the following properties of D^* (see the definition above line (2.4)):

- (i) For a matrix $F(S)_{p \times p}$ and a scalar $\varphi(S)$,
 $D^*(\varphi F) = \varphi D^*F + \text{tr} [(\partial \varphi(S)/\partial S) \cdot F]$;
- (ii) For a matrix Q of constants,
 $D^*(QS)_{(\frac{1}{2})} = \left(\frac{p+1}{2} \right) \text{tr} Q$;
- (iii) $D^*S_{(\frac{1}{2})}^{-1} = -(\frac{1}{2}) \text{tr} S^{-2} - (\frac{1}{2})(\text{tr} S^{-1})^2$;
- (iv) $D^*S_{(\frac{1}{2})}^{-2} = -\text{tr} S^{-3} - (\text{tr} S^{-1})(\text{tr} S^{-2})$.

The calculation (5.4(i)) is routine, hence omitted. See the Appendix for verification of (ii) and (iv).

Identities for the terms in (5.3) are given by

LEMMA 5.1. Assume that $g_{ij}(\mathbf{S}) = \partial g(\mathbf{S})/\partial s_{ij}$ (i, j, \dots, p) satisfy the conditions of Stoke's theorem—see Haff [3] and [4].

(i) For $h(\mathbf{S}) = 2ag(\mathbf{S})$, we have

$$E[h(\mathbf{S}) \operatorname{tr}(\mathbf{S}\Sigma^{-2})] = E\left\{k(k-p-1)h(\mathbf{S}) \operatorname{tr} \mathbf{S}^{-1} + 2(2k-p-1) \operatorname{tr}(\partial h(\mathbf{S})/\partial S) + 4D^*[(\partial h(\mathbf{S})/\partial S)_{(\frac{1}{2})} \cdot \mathbf{S}]_{(\frac{1}{2})}\right\}.$$

(ii) For $h(\mathbf{S}) = -2g(\mathbf{S})$, we have

$$E[h(\mathbf{S}) \operatorname{tr} \Sigma^{-1}] = E[(k-p-1)h(\mathbf{S}) \operatorname{tr} \mathbf{S}^{-1} + 2\operatorname{tr} \partial h(\mathbf{S})/\partial S].$$

(iii) For $h(\mathbf{S}) = g^2(\mathbf{S})$, we have

$$E[h(\mathbf{S}) \operatorname{tr} \Sigma^{-2}] = E\left\{4D^*[\partial h(\mathbf{S})/\partial S]_{(\frac{1}{4})} + 4(k-p-1) \operatorname{tr}[\partial h(\mathbf{S})/\partial S \cdot \mathbf{S}^{-1}] + 2(k-p-1)h(\mathbf{S})D^*\mathbf{S}^{-1}_{(\frac{1}{2})} + (k-p-1)^2h(\mathbf{S}) \operatorname{tr} \Sigma^{-2}\right\}.$$

PROOF OF (i). Set $\mathbf{T} = \mathbf{S}\Sigma^{-1}$. From (2.4) we have

$$\begin{aligned} E[h(\mathbf{S}) \operatorname{tr}(\mathbf{T}\Sigma^{-1})] &= E\left\{2h(\mathbf{S})D^*T_{(\frac{1}{2})} + 2 \operatorname{tr}[\partial h(\mathbf{S})/\partial S \cdot \mathbf{T}_{(\frac{1}{2})}] + (k-p-1)h(\mathbf{S}) \operatorname{tr} \Sigma^{-1}\right\} \\ &= E\left\{kh(\mathbf{S}) \operatorname{tr} \Sigma^{-1} + 2 \operatorname{tr}[(\partial h(\mathbf{S})/\partial S)_{(\frac{1}{2})}\mathbf{S}\Sigma^{-1}]\right\}. \end{aligned}$$

(Use (5.4(ii)) and combine terms.) From the last equation, the result follows by applying (5.1(i)) to the first term and then Lemma 5.1 to the second ($\mathbf{T} = (\partial h(\mathbf{S})/\partial S)_{(\frac{1}{2})} \cdot \mathbf{S}$).

PROOF OF (ii). In (5.1(i)), set $Q = I$.

PROOF OF (iii). Set $T = \Sigma^{-1}$. From (2.4) we have

$$E[h(\mathbf{S}) \operatorname{tr}(T\Sigma^{-1})] = E\left\{2 \operatorname{tr}[\partial h(\mathbf{S})/\partial S]_{(\frac{1}{2})} \Sigma^{-1} + (k-p-1)h(\mathbf{S}) \operatorname{tr}(\mathbf{S}^{-1}\Sigma^{-1})\right\}.$$

Our result follows by applying (2.4) to the first term under the expectation with $\mathbf{T} = [\partial h(\mathbf{S})/\partial S]_{(\frac{1}{2})}$ and applying it to the second with $\mathbf{T} = \mathbf{S}^{-1}$.

Now we specialize Lemma 5.1 by taking $g(\mathbf{S}) = a\mathbf{u}t(\mathbf{u})$, $\mathbf{u} = 1/\operatorname{tr}(\mathbf{S}^{-1})$. The matrix

$$\partial \mathbf{u}/\partial S = \mathbf{u}^2 S_{(2)}^{-2}$$

follows from standard perturbation results—see [4]. Thus we obtain

$$(5.5) \quad \partial g(S)/\partial S = a[\mathbf{u}^3 t'(\mathbf{u}) + \mathbf{u}^2 t(\mathbf{u})] S_{(2)}^{-2}.$$

Let us introduce the notation

$$(5.6) \quad \rho_m \equiv (\text{tr } S^{-m}) / (\text{tr } S^{-1})^m, \quad m = 1, 2, 3, \dots$$

Our specialization of Lemma 5.1 is

LEMMA 5.1*. For $g(\mathbf{S}) = a\mathbf{u}t(\mathbf{u})$ and $\tau(\mathbf{u}) = \mathbf{u}^3t'(\mathbf{u}) + \mathbf{u}^2t(\mathbf{u})$,

$$(i) \quad E[2a^2\mathbf{u}t(\mathbf{u}) \text{tr}(\mathbf{S}\Sigma^{-2})] = E[2a^2k(k-p-1)t - 4a^2(\tau/\mathbf{u}^2) + 4a^2(2k-p-2)(\tau/\mathbf{u}^2)\rho_2 + 8a^2(\tau'/\mathbf{u})\rho_3,$$

$$(ii) \quad E[-2a\mathbf{u}t(\mathbf{u}) \text{tr}\Sigma^{-1}] = E[-2a(k-p-1)t - 4a(\tau/\mathbf{u}^2)\rho_2,$$

and

$$(iii) \quad E a^2\mathbf{u}^2t^2(\mathbf{u}) = E\{[(k-p-1)t^2 - (k-p-1)t - 8(\tau/\mathbf{u}^2)]a^2t\rho_2 + 8(k-p-2)a^2t(\tau/\mathbf{u}^2)\rho_3 + 8[(\tau/\mathbf{u}^2)^2 + t(\tau'/\mathbf{u})]a^2\rho_4 - (k-p-1)a^2t^2\}.$$

PROOF. In part (i), $h(S) = 2ag(S) = 2a^2\mathbf{u}t(\mathbf{u})$. We shall compute $D^*[(\partial h(S)/\partial S)_{(\frac{1}{2})} \cdot S]_{(\frac{1}{2})}$ and leave the remaining details for the reader. Part (ii) is immediate. Also, we omit the details of part (iii). The latter calculation is lengthy; nevertheless, it is straightforward from (5.4), (5.5), and (5.6). We have

$$\begin{aligned} D^*[(\partial h(S)/\partial S)_{(\frac{1}{2})} \cdot S]_{(\frac{1}{2})} &= 2a^2D^*[(\partial \mathbf{u}t(\mathbf{u})/\partial S)_{(\frac{1}{2})} \cdot S]_{(\frac{1}{2})} \\ &= 2a^2D^*\left[\left\{\tau(\mathbf{u})S_{(\frac{1}{2})}^{-1}\right\} \quad (\text{from(5.5)})\right] \\ &= 2a^2\left\{\tau(\mathbf{u})D^*S_{(\frac{1}{2})}^{-1} + \text{tr}\left[\partial\tau(\mathbf{u})/\partial S \cdot S_{(\frac{1}{2})}^{-1}\right]\right\} \quad (\text{from (5.4(i))}) \\ &= 2a^2\left\{\tau(\mathbf{u})\left[-\left(\frac{1}{2}\right)\text{tr}S^{-2} - \left(\frac{1}{2}\right)(\text{tr}S^{-1})^2\right.\right. \\ &\quad \left.\left.+ \tau'(\mathbf{u})\text{tr}\left[\left(\partial\mathbf{u}/\partial S\right) \cdot S_{(\frac{1}{2})}^{-1}\right]\right]\right\} \quad (\text{from (5.4(iii))}) \\ &= 2a^2\left\{\left(-\frac{1}{2}\right)\left[\tau(\mathbf{u})/u^2\right]\left[\rho_2 + 1\right]\right. \\ &\quad \left.+ \tau'(\mathbf{u})\text{tr}\left[u^2S_{(2)}^{-2} \cdot S_{(\frac{1}{2})}^{-1}\right]\right\} \\ &= 2a^2\left\{\left(-\frac{1}{2}\right)\left[\tau(\mathbf{u})/u^2\right]\left[\rho_2 + 1\right] + \left[\tau'(\mathbf{u})/u\right]\rho_3\right\}. \quad \square \end{aligned}$$

Finally, the proofs of Theorems 4.5 and 4.6 depend on a result from Bellman [1], page 137, namely,

LEMMA 5.2. We have the inequality

$$(1/p)^{m-1} \leq \rho_m \leq 1, \quad m = 0, 1, 2, \dots$$

PROOF. Omitted. \square

In addition, it is convenient to note that $\rho_m = \text{tr} [S^{-1}/(\text{tr } S^{-1})]^m$ decreases in m .

PROOF OF THEOREM 4.5. The function $\alpha_2(\Sigma)$ is given by Lemma 5.1*, being the sum of (i), (ii) and (iii). Let $t(u)$ be a constant, say t . Then $\tau(u)/u^2 = t$ and

$\tau'(u)/u = 2t$. The unbiased estimator of $\alpha_2(\Sigma)$ is

$$\begin{aligned} \hat{\alpha}_2(\Sigma) = & \{2ak[k - p - 1 + 2\rho_2] \\ & + 4a[(k - p - 2)\rho_2 + 4\rho_3 - 1] - 2[k - p - 1 + 2\rho_2]\}at \\ & + \{[(k - p - 1)^2 - (k - p - 1) - 8]\rho_2 + 8[k - p - 2]\rho_3 + 24\rho_4 \\ & - (k - p - 1)\}a^2t^2. \end{aligned}$$

Recall that $a = 1/(k + p + 1)$. The above coefficient of $at(u)$ can be written as

$$(5.7) \quad -2a[2 + k(p + 1) - (p + 1)^2] + 4a(k - 2p - 3)\rho_2 + 16a\rho_3.$$

In (5.7), the first term is negative because $k > p + 1$. Thus the entire quantity is bounded above by

$$\begin{aligned} (5.8) \quad & -2a\rho_2[2 + k(p + 1) - (p + 1)^2] + 4a\rho_2(k - 2p + 3) + 16a\rho_2 \\ & = -2a\rho_2(k - p + 1)(p - 1). \quad (\text{since } 0 \leq \rho_m \leq 1 \text{ and } \rho_m \searrow) \end{aligned}$$

The coefficient of a^2t^2 is bounded above by

$$\begin{aligned} (5.9) \quad & [(k - p - 1)^2 - (k - p - 1) - 8]\rho_2 + 8(k - p - 2)\rho_2 + 24\rho_2 \\ & = \rho_2(k - p + 3)(k - p + 1). \end{aligned}$$

Finally, from (5.8) and (5.9), a sufficient condition for $\alpha_2(\Sigma) < 0$ ($\forall \Sigma$) is

$$(5.10) \quad -2a^2\rho_2(k - p + 1)(p - 1)t + a^2\rho_2(k - p + 3)(k - p + 1)t^2 < 0$$

which is equivalent to

$$0 \leq t \leq 2(p - 1)/(k - p + 3). \quad \square$$

With respect to (5.10), the optimal t is $t = (p - 1)/(k - p + 3)$.

PROOF OF THEOREM 4.6. We have assumed that $0 \leq t \leq (p - 1)/(k - p + 3)$ and

$$\begin{aligned} t(u) = & \bar{\gamma}/u \text{ for } u \geq \bar{\gamma}(k - p + 3)/(p - 1) \\ = & (p - 1)/(k - p + 3) \text{ otherwise.} \end{aligned}$$

Recall that $\tau(u) = u^3t'(u) + u^2t(u)$. (Again, see Lemma 5.1*.) Let $\mathfrak{N} = \{u: u \geq \bar{\gamma}(k - p + 3)/(p - 1)\}$. For $u \in \mathfrak{N}$, $\tau(u) = \tau'(u) = 0$. In this case, simple algebra shows that $\hat{\alpha}_2(\Sigma) \leq 0$ if $0 < t \leq (p + 1)/(k - p - 3)$. The latter inequality is satisfied by hypothesis. For $u \notin \mathfrak{N}$, $t = (p - 1)/(k - p + 3)$, and again $\hat{\alpha}_2(\Sigma) \leq 0$. The proof is complete. \square

APPENDIX

THE CALCULATIONS IN LINE (5.4).

Calculation (5.4(ii)). We show that

$$D^*(QS)_{(\frac{1}{2})} = \frac{p + 1}{2} \text{tr} Q.$$

Let $S = (S_1, \dots, S_p)$ and $Q = (Q_1, \dots, Q_p)$; thus $QS = [Q_i' S_j]$. We have

$$\begin{aligned} D^*(QS)_{(\frac{1}{2})} &= \sum_i (\partial/\partial s_{ii}) Q_i' S_i + (\frac{1}{2}) \sum_{i \neq j} (\partial/\partial s_{ij}) Q_i' S_j \\ &= \sum_i q_{ii} + (\frac{1}{2}) \sum_{i \neq j} q_{ij} \\ &= (\frac{1}{2}) \sum_i q_{ii} + (\frac{1}{2}) \sum_{(i,j)} q_{ii} \\ &= (p + 1)/2 \text{tr} Q. \end{aligned}$$

Calculation (5.4(iv)). We show that

$$D^* S_{(\frac{1}{2})}^{-2} = -\text{tr} (S^{-3}) - (\text{tr} S^{-1})(\text{tr} S^{-2}).$$

Let $S^{-1} = A = (a_{ij})$. We have

$$\begin{aligned} D^* A_{(\frac{1}{2})}^2 &= \sum_i \sum_t \frac{\partial}{\partial s_{ii}} a_{it} a_{it} + (\frac{1}{2}) \sum_{i \neq j} \sum_t \frac{\partial}{\partial s_{ij}} a_{it} a_{ij} \\ &= \sum_i \sum_t 2a_{it} (\partial a_{it} / \partial s_{ii}) + (\frac{1}{2}) \sum_{i \neq j} \sum_t [a_{it} (\partial a_{ij} / \partial s_{ij}) + a_{ij} (\partial a_{it} / \partial s_{ij})] \\ &= -\sum_i \sum_t 2a_{it} a_{it} a_{ii} - (\frac{1}{2}) \sum_{i \neq j} \sum_t a_{it} (a_{ii} a_{ij} + a_{ij} a_{ij}) \\ &\quad - (\frac{1}{2}) \sum_{i \neq j} \sum_t a_{ij} (a_{ii} a_{ij} + a_{ij} a_{ii}) \\ &= -(\sum_i \sum_t a_{it} a_{it}^2 + \sum_{i \neq j} \sum_t a_{ij} a_{it} a_{ij}) \\ &\quad - (\sum_i \sum_t a_{it} a_{it}^2 + (\frac{1}{2}) \sum_{i \neq j} \sum_t a_{it}^2 a_{ij} + (\frac{1}{2}) \sum_{i \neq j} \sum_t a_{ij}^2 a_{ii}) \\ &= -\text{tr}(S^{-3}) - (\text{tr} S^{-1})(\text{tr} S^{-2}). \end{aligned}$$

□

REFERENCES

- [1] BELLMAN, R. (1970). *Introduction to Matrix Analysis*. McGraw-Hill, New York.
- [2] EFRON, B. and MORRIS, C. (1976). Multivariate Empirical Bayes and estimation of covariance matrices. *Ann. Statist.* 4 22-32.
- [3] HAFF, L. R. (1977a). Minimax estimators for a multinormal precision matrix. *J. Multivariate Anal.* 7 374-385.
- [4] HAFF, L. R. (1977b). Estimation of the inverse covariance matrix: Random mixtures of the inverse Wishart matrix and the identity. *Ann. Statist.* 7 1264-1276.
- [5] HAFF, L. R. (1978). An identity for the Wishart distribution with applications. *J. Multivariate Anal.* 9 (in press).
- [6] JAMES W. and STEIN, C. (1961). Estimation with quadratic loss. In *Fourth Berkeley Symp. Math. Statist. Probability* Univ. California Press, Berkeley.
- [7] LEHMANN, E. L. (1951). Notes on estimation. Unpublished manuscript.

- [8] PERLMAN, M. D. (1972). Reduced mean square error estimation for several parameters. *Sankhyā* 34 89–92.
- [9] PRESS, S. J. (1972). *Applied Multivariate Analysis*. Holt, Rinehart and Winston, New York.
- [10] SELLIAH, J. (1964). Estimation and testing problems in a Wishart distribution. Ph.D. thesis, Depart. Statist. Stanford, Univ.
- [11] STEIN, C. (1973). Estimation of the mean of a multivariate normal distribution. In *Proc. Prague Symp. Asymptotic Statist.* 345–381.
- [12] STEIN, C. (1975). Reitz lecture. 38th annual meeting IMS. Atlanta, Georgia.
- [13] STEIN, C., EFRON, B. and MORRIS, C. (1972). Improving the usual estimator of a normal covariance matrix. Technical Report No. 37, Depart. Statist., Stanford Univ.

DEPARTMENT OF MATHEMATICS
UNIVERSITY OF CALIFORNIA, SAN DIEGO
LA JOLLA, CALIFORNIA 92093