

MULTICOLLINEARITY AND REDUCED-RANK ESTIMATION

3.1 Introduction

Multicollinearity has been a topic of concern in econometrics ever since the publication of Frisch's monograph (1934). Two approaches have already been discussed in the previous chapter (sections 2.7 and 2.8). In this chapter attention will revolve around the singular-value decomposition of a matrix; in the case of an $n \times k$ observation matrix X , its singular values are the positive square roots of the eigenvalues of $X'X$. If some of these are very small (and because of rounding error, a computer cannot easily distinguish between "small" and zero), classical methods of computing least-squares estimates (e.g., the Gauss-Seidel procedure) tend to be highly inaccurate. The method of computing the singular-value decomposition and replacing small singular values by zeros produces much more reliable results. Interestingly enough, statistical theory reaches a similar conclusion: replacing small singular values by zeros (which amounts to approximating X by an $n \times k$ matrix $X_{(l)}$ of reduced rank, l) leads to estimators with lower mean-square error. This theory is the subject of the present chapter.

3.2 Singular-value decomposition of a matrix

The concept of a singular-value decomposition of a matrix was introduced (independently) by Eckart & Young (1939, pp. 118–121) and Mirsky (1960).

DEFINITION 3.2.2. *Let X be any $n \times k$ matrix. If a triple (s, p, q) exists, where s is a nonnegative scalar, p is a $k \times 1$ vector, and q is an $n \times 1$ vector, such that*

$$Xp = sq \quad \text{and} \quad q'X = sp',$$

then s is called a singular value of X , and p and q are respectively called right and left singular vectors of X .

DEFINITION 3.2.2. *An $n \times k$ matrix $D = [d_{ij}]$ is called a (rectangular) diagonal matrix if $d_{ij} = 0$ for $i \neq j$.*

THEOREM 3.2.1. *Let X be any $n \times k$ matrix of rank r , and define $m = \min(n, k)$. Then there exists an $n \times n$ orthogonal matrix Q , a $k \times k$ orthogonal matrix P , and an $n \times k$ diagonal matrix D , with diagonal elements $s_1 \geq s_2 \geq \dots \geq s_m \geq 0$ (the singular values of X), such that*

$$(3.2.1) \quad X = QDP',$$

where the columns of Q and P are respectively the corresponding unit-length left and right singular vectors of X . These singular values are unique.

PROOF. Since the $k \times k$ matrix $X'X$ is symmetric, there exists an orthogonal $k \times k$ matrix P such that

$$(3.2.2) \quad P'X'XP = \Lambda = \text{diag}\{\lambda_i\}$$

where the λ_i 's are the k (nonnegative) eigenvalues of $X'X$, and the $k \times 1$ columns p^i of P (the eigenvectors of $X'X$) may be ordered so that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$. We define $s_i = +\sqrt{\lambda_i}$ for $i = 1, 2, \dots, m$.

Since $r = \text{rank } X$, $s_r > 0$ and $s_{r+1} = \dots = s_m = 0$. Define

$$(3.2.3) \quad q^i = Xp^i/s_i \quad \text{for } i = 1, 2, \dots, r.$$

Then for $i, j = 1, 2, \dots, r$, we have from (3.2.3) and (3.2.2)

$$q^i \cdot q^j = \frac{p^{i'}X'Xp^j}{s_i s_j} = \delta_{ij}$$

(where δ_{ij} is the Kronecker delta). For $i = r+1, \dots, n$, choose q^i so that the columns of

$$(3.2.4) \quad Q = [q^1, q^2, \dots, q^r, q^{r+1}, \dots, q^n]$$

form an orthonormal set; then $Q'Q = QQ' = I_n$. From (3.2.2) we have clearly

$$(3.2.5) \quad Xp^i = 0 = s_i q^i \quad \text{for } i = r+1, \dots, m.$$

Thus, from (3.2.3) and (3.2.5) we have

$$(3.2.6) \quad \begin{aligned} X &= XPP' = X[p^1, p^2, \dots, p^k] \begin{bmatrix} p^{1'} \\ p^{2'} \\ \vdots \\ p^{k'} \end{bmatrix} \\ &= \sum_{i=1}^k Xp^i p^{i'} = \sum_{i=1}^m q^i s_i p^{i'} \quad (\text{since } Xp^i = 0 \text{ for } m < i < k) \end{aligned}$$

$$= [q^1, q^2, \dots, q^m] \begin{bmatrix} s_1 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_m \end{bmatrix} \begin{bmatrix} p^{1'} \\ p^{2'} \\ \vdots \\ p^{m'} \end{bmatrix}$$

Denoting the square diagonal matrix in (3.2.6) as $S = \text{diag}\{s_i\}$, (3.2.6) is equivalent to (3.2.1), where $D' = [S, 0]$ for $n \geq k = m$ and $D = [S, 0]$ for $m = n \geq k$.

Since, by (3.2.2), the singular values s_i of X coincide with the nonnegative square roots of the eigenvalues λ_i of $X'X$, for $i = 1, 2, \dots, m$ (the remaining eigenvalues in the case $n < k$ being all zero), and these eigenvalues are unique, therefore, the singular values are unique. \square

Note that if $n > k$ (hence $m = k$) and we partition Q as

$$(3.2.7) \quad Q = [Q_1; Q_2] = [q^1, q^2, \dots, q^m; q^{m+1}, \dots, q^n],$$

we may write (3.2.6) or (3.2.1) in the form

$$(3.2.8) \quad X = Q_1 S P'.$$

Likewise, if $n < k$ (hence $m = n$) and we partition P as

$$(3.2.9) \quad P = [P_1; P_2] = [p^1, p^2, \dots, p^m; p^{m+1}, \dots, p^k],$$

we may write (3.2.6) or (3.2.1) in the form

$$(3.2.10) \quad X = Q S P'_1.$$

Since one normally has $n > k$, the form (3.2.8) will frequently be used in applications.

We may use Theorem 3.2.1 as a way to define the Moore-Penrose generalized inverse of a matrix (Definition 2.3.1). First, we define the Moore-Penrose generalized inverse of a rectangular diagonal matrix, as follows. If D is an $n \times k$ diagonal matrix with diagonal elements $d_{ii} = s_i$, then defining

$$(3.2.11) \quad s_i^\dagger = \begin{cases} 1/s_i & \text{if } s_i \neq 0, \\ 0 & \text{if } s_i = 0, \end{cases}$$

the generalized inverse D^\dagger is the $k \times n$ diagonal matrix in which each diagonal element in the transpose of D is replaced by its generalized inverse as defined by (3.2.11). Thus, denoting $S^\dagger = \text{diag}\{s_i^\dagger\}$, if $n > k$ we have $D' = [S', 0]$ and $D^\dagger = [S^\dagger, 0]$. The Moore-Penrose generalized inverse of X is then given by

$$(3.2.12) \quad X^\dagger = P D^\dagger Q'.$$

It is easily verified that X^\dagger and D^\dagger satisfy the four properties of Definition 2.3.1. From (3.2.12) it is clear that once the singular-value decomposition of a matrix X has been computed, computation of its Moore-Penrose generalized inverse X^\dagger is trivial. Computation of the singular-value decomposition, however, is far from trivial (see Golub & Kahan 1965, Golub & Reinsch 1970, Noble 1976, and Golub & Van Loan 1983).

Since one will often want to compute an oblique generalized inverse of X , satisfying property (iii)' of Definition 2.3.3, one may proceed as follows. Define

$$(3.2.13) \quad \dot{X} = V^{-1/2} X, \quad \text{hence } X = V^{1/2} \dot{X},$$

and

$$(3.2.14) \quad X^\dagger = \dot{X}^\dagger V^{-1/2}.$$

We verify that

$$X X^\dagger X = V^{1/2} X | \dot{X}^\dagger V^{-1/2} | V^{1/2} \dot{X} = V^{1/2} \dot{X} \dot{X}^\dagger \dot{X} = V^{1/2} \dot{X} = X$$

(where the symbols $|$ are inserted to facilitate checking the substitutions) and

$$X X^\dagger V = V^{1/2} \dot{X} | \dot{X}^\dagger V^{-1/2} | V = V^{1/2} \dot{X} \dot{X}^\dagger V^{1/2}$$

which is symmetric since $V^{1/2}$ and $\dot{X} \dot{X}^\dagger$ are symmetric.

3.3 The condition number of a matrix

The condition number of a matrix is a numerical measure of the degree to which it is ill-conditioned, i.e., “close” to having deficient rank. To develop this concept we first introduce some definitions (cf. Golub & Van Loan 1983, pp. 12–14):

DEFINITION 3.3.1. *The Hölder p -norm of an $n \times 1$ vector x is defined by*

$$(3.3.1) \quad \|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \quad \text{for } p \geq 1.$$

In the special case $p = 2$ this reduces to the Euclidean norm

$$(3.3.2) \quad \|x\|_2 \equiv \|x\| = \sqrt{x'x}.$$

Definition 3.3.1 is generalized to a matrix by

DEFINITION 3.3.2. *The Hölder p -norm of an $n \times k$ matrix X is defined by*

$$(3.3.3) \quad \|X\|_p = \sup_{b \neq 0} \frac{\|Xb\|_p}{\|b\|_p},$$

where b is of order $k \times 1$. In the special case $p = 2$ this may be written

$$(3.3.4) \quad \|X\|_2 = \sup_{b \neq 0} \left(\frac{b'X'Xb}{b'b} \right)^{1/2}.$$

Note that if $k = 1$ and $X = x$, then b is scalar and (3.3.3) reduces to (3.3.1):

$$\sup_{b \neq 0} \frac{\|xb\|_p}{\|b\|_p} = \sup_{b \neq 0} \frac{|b| \|x\|_p}{|b|} = \|x\|_p.$$

Likewise, (3.3.4) reduces to (3.3.2).

The following is a special case of the Courant-Fischer min-max theorem (cf. Courant 1922, Fischer 1905; Bellman 1960, pp. 110–115):

LEMMA 3.3.1. *Let M be a $k \times k$ symmetric nonnegative-definite matrix with eigenvalues given by $\Lambda = \text{diag}\{\lambda_i\}$ arranged in descending order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$, where $P'MP = \Lambda$, $P'P = I$. Then*

$$(3.3.5) \quad \lambda_1 = \max_b \frac{b'Mb}{b'b} \quad \text{and} \quad \lambda_k = \min_b \frac{b'Mb}{b'b}.$$

PROOF. Defining $c = P'b$ we have

$$\frac{b'Mb}{b'b} = \frac{c'P'MPc}{c'P'Pc} = \frac{c'\Lambda c}{c'c} = \frac{\sum_{i=1}^k \lambda_i c_i^2}{\sum_{i=1}^k c_i^2}.$$

Since

$$\sum_{i=1}^k \lambda_i c_i^2 \leq \lambda_1 \sum_{i=1}^k c_i^2 \quad \text{and} \quad \sum_{i=1}^k \lambda_i c_i^2 \geq \lambda_k \sum_{i=1}^k c_i^2$$

it follows that

$$(3.3.6) \quad \lambda_k \leq \frac{\sum_{i=1}^k \lambda_i c_i^2}{\sum_{i=1}^k c_i^2} \leq \lambda_1.$$

Since the equalities in (3.3.6) are attained by choice of $c = (0, \dots, 0, 1)'$ and $c = (1, 0, \dots, 0)'$ respectively, the result follows. \square

LEMMA 3.3.2. *Let X be an $n \times k$ matrix of rank m , whose positive singular values are $s_1 \geq s_2 \geq \dots \geq s_m > 0$. Then*

$$\|X\|_2 = s_1.$$

PROOF. Defining $M = X'X$ and $\lambda_i = s_i^2$, where $X = QDP'$, this follows directly from Lemma 3.3.1. \square

LEMMA 3.3.3. *Let X be an $n \times k$ matrix of rank m , whose positive singular values are $s_1 \geq s_2 \geq \dots \geq s_m > 0$. Then*

$$\|X^\dagger\|_2 = 1/s_m.$$

PROOF. From the singular-value decomposition $X = QDP'$ of X we have that of $X^\dagger = PD^\dagger Q'$, hence the singular values of X^\dagger are $1/s_1 \leq 1/s_2 \leq \dots \leq 1/s_m$. The result then follows from Lemma 3.3.1. \square

DEFINITION 3.3.3. *The condition number of an $n \times k$ matrix X of rank m is, in terms of the Hölder matrix 2-norm,*

$$\kappa(X) = \|X\|_2 \|X^\dagger\|_2.$$

From this definition and Lemmas 3.3.2 and 3.3.3 it follows immediately that $\kappa(X) = s_1/s_m$, i.e., the condition number of X is the ratio between its largest and its smallest positive singular value. Cf., e.g., Noble (1976, pp. 279, 295), Golub & Van Loan (1983, p. 140).

3.4 The Eckart-Young theorem

Eckart and Young (1936) furnished an algorithm for obtaining a best approximation of an $n \times k$ matrix X by an $n \times k$ matrix of rank less than that of X . In order to furnish a precise meaning to "best approximation" it is necessary to define a concept of distance between two $n \times k$ matrices.

DEFINITION 3.4.1. *The Frobenius norm of an $n \times k$ matrix X is defined as*

$$\|X\|_F \equiv \|X\| = \sqrt{\text{tr}(X'X)}.$$

It is equal to the square root of the sum of squares of all the elements of X . If X^1 and X^2 are two $n \times k$ matrices, we define the Frobenius distance between X^1 and X^2 as the Frobenius norm of their difference, i.e., $\|X^1 - X^2\|$.

DEFINITION 3.4.2. *Let X be any $n \times k$ matrix, and let Q and P respectively be $n \times n$ and $k \times k$ orthogonal matrices. Then a norm $\|X\|$ of X is said to be orthogonally invariant if it has the property $\|Q'XP\| = \|X\|$.*

The following is a simple extension of Lemma 1.1.1.

LEMMA 3.4.1. *The Frobenius norm is orthogonally invariant.*

PROOF. From Lemma 1.1.1,

$$\begin{aligned}\|Q'XP\| &= \sqrt{\text{tr}(P'X'QQ'XP)} \\ &= \sqrt{\text{tr}(P'X'XP)} \\ &= \sqrt{\text{tr}(X'XPP')} \\ &= \sqrt{\text{tr}(X'X)} \\ &= \|X\|. \quad \square\end{aligned}$$

It is readily seen that the Hölder matrix 2-norm is also orthogonally invariant.

DEFINITION 3.4.3. *We denote by \mathcal{X} the set of all real $n \times k$ matrices, and by \mathcal{X}_l the subset of \mathcal{X} consisting of $n \times k$ matrices of rank $\leq l$.*

In trying to approximate an $n \times k$ matrix of rank $> l$ by one of rank l , a basic problem is that the set of $n \times k$ matrices of rank exactly l is not closed. Such a set is defined by the condition that all minors (subdeterminants) of X of order greater than l vanish, and at least one minor of order l be nonvanishing. One wishes to find in this set a matrix that is closest to X in the Frobenius norm; but since this set is obviously not closed, the existence of such a matrix is not at all obvious. The procedure that is followed, therefore, is to deal with the set \mathcal{X}_l of all $n \times k$ matrices X of rank $\leq l$; this is a closed set, and can be compactified, hence a matrix in this set exists that is closest to X (see Lemma 3.4.2 below); but on the face of it, it might have rank $< l$ (see Figure 3.4.1). It has to be shown that it has rank exactly l ; this is done in Theorem 3.4.1. The Eckart-Young theorem then provides the algorithm by which this matrix is determined, namely the replacement of all but the l largest singular values of X by zeros.

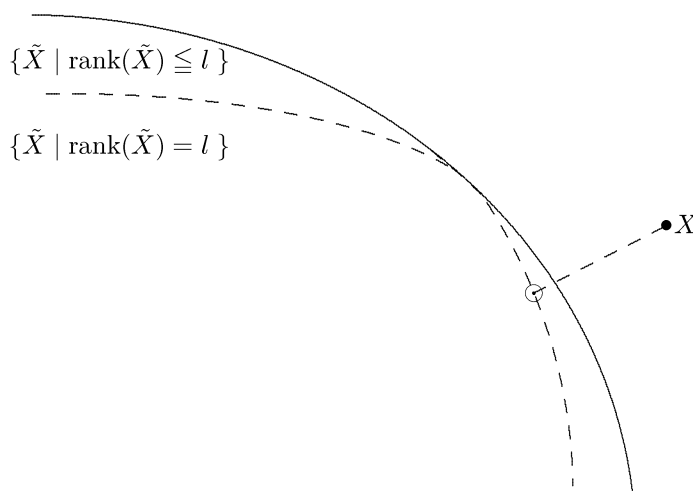


FIGURE 3.4.1

As a preliminary to the Eckart-Young theorem (Theorem 3.4.2 below), we will need Lemmas 3.4.2 and 3.4.3 and Theorem 3.4.1.

LEMMA 3.4.2. Let X be a given $n \times k$ matrix of rank $> l$, where $l < m = \min(n, k)$. Then within the class \mathcal{X}_l of $n \times k$ matrices \tilde{X} of rank $\leq l$, there exists a matrix \hat{X} closest to X in the Frobenius norm, i.e., such that

$$(3.4.1) \quad \|X - \hat{X}\| = \min_{\tilde{X} \in \mathcal{X}_l} \|X - \tilde{X}\|.$$

PROOF. The set \mathcal{X}_l of $n \times k$ matrices \tilde{X} of rank $\leq l$ is defined by the condition that all minors of order $l+1$ of such matrices are equal to zero. Since these minors are polynomials in the elements of the matrices \tilde{X} , these equations define a closed set in the nk -dimensional space of matrices X . Let B be the ball of radius $\|X\|$ in this space, centered at X . Then $B \cap \mathcal{X}_l$ is compact, and is nonempty since it contains at least the zero matrix 0 . Therefore the continuous function $f(\tilde{X}) = \|X - \tilde{X}\|$ has a minimum, \hat{X} , on $B \cap \mathcal{X}_l$, which is clearly the minimum on \mathcal{X}_l . (See Figure 3.4.2.) \square

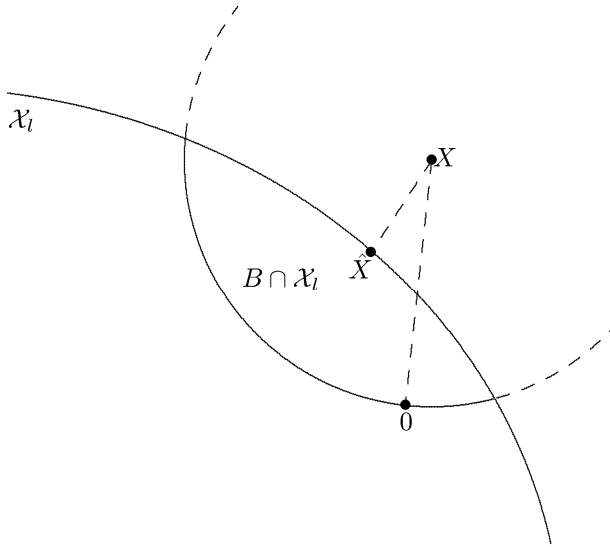


FIGURE 3.4.2

The following lemma and theorem have kindly been supplied by Joel Roberts of the School of Mathematics, University of Minnesota.

LEMMA 3.4.3 (ROBERTS). Let E_{ij} be the $n \times k$ matrix with 1 in the i, j th position, and 0s elsewhere, and let A be any $n \times k$ matrix. Then there exists a real number $\lambda \neq 0$ such that the nk matrices

$$(3.4.2) \quad \lambda E_{ij} - A \quad (i = 1, 2, \dots, n; j = 1, 2, \dots, k)$$

form a basis in nk -dimensional space.

PROOF. If the set of matrices (3.4.2) is linearly dependent, then the $nk \times nk$ matrix whose columns are the successive columns of (3.4.2) is a matrix of the form $\lambda I_{nk} - M$, where M has all its nk columns equal to the column vector of columns of

A , and λ is an eigenvalue of M . But the eigenvalues of M are 0 with a multiplicity of $nk - 1$, and just one other real number (namely, the sum of the elements of A). For any real number other than this one or 0, the set (3.4.2) is therefore linearly independent. \square

THEOREM 3.4.1 (ROBERTS). *Let X be an $n \times k$ matrix of rank $> l$, and let $\hat{X} \in \mathcal{X}$ be a matrix of rank $\leq l$ which is closest to X (in the Frobenius norm) among all matrices in the set \mathcal{X}_l of $n \times k$ matrices of rank $\leq l$. Then $\text{rank } \hat{X} = l$.*

PROOF. Suppose by way of contradiction that $\text{rank}(\hat{X}) < l$, and let λ be such as to satisfy Lemma 3.4.3. Since $\text{rank}(\lambda E_{ij}) = 1$ for $\lambda \neq 0$, we have, since multiplication of a matrix by a non-zero scalar does not affect its rank,

$$(3.4.3) \quad \text{rank}[(1-t)\hat{X} + t\lambda E_{ij}] \leq \text{rank}(\hat{X}) + \text{rank}(\lambda E_{ij}) \leq l$$

for any real number t , since the rank of the sum of two matrices is less than or equal to the sum of the ranks (because the column space of the sum of two matrices is contained in the sum of the column spaces of the two matrices, as is easily verified), and (3.4.3) holds for $t = 0$ and $t = 1$. Thus, of all points on the line $(1-t)\hat{X} + t\lambda E_{ij}$, the closest to X will be \hat{X} , since \hat{X} has been assumed to be a closest point to X of all \tilde{X} with $\tilde{X} \in \mathcal{X}_l$. Thus, the matrix $X - \hat{X}$ is perpendicular to the matrix $\hat{X} - \lambda E_{ij}$, since the shortest distance from a point to a line is along the perpendicular. But then $X - \hat{X}$ is perpendicular to all of nk -dimensional space, since $\hat{X} - \lambda E_{ij}$ is a basis for this space, by Lemma 3.4.3. Thus, $X - \hat{X} = 0$, i.e., $X = \hat{X}$. This contradiction establishes that $\text{rank}(\hat{X}) = l$. (See Figure 3.4.3.) \square

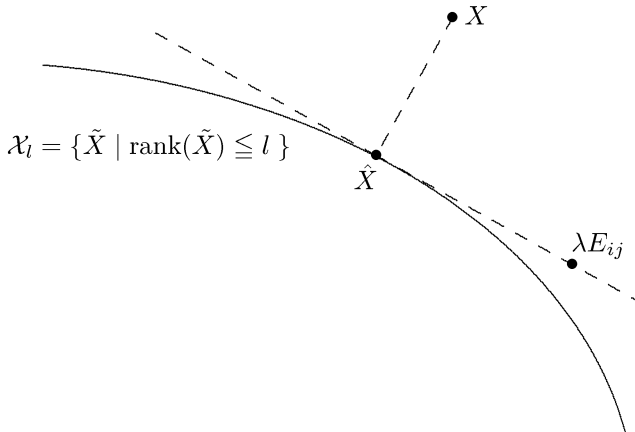


FIGURE 3.4.3

The Eckart-Young theorem states that any $n \times k$ matrix X of rank $p \leq m = \min(n, k)$ can be best approximated in terms of the Frobenius norm by an $n \times k$ matrix of rank $l < p$, and that when the singular values of X are arranged in descending order, this approximating matrix can be obtained by striking out the singular values $s_{l+1}, s_{l+2}, \dots, s_p$ of X . The theorem goes back to Eckart & Young

(1936, 1939), who tacitly assumed the first part to be true, and for the second part provided only a heuristic proof. This was followed by a more complete proof by Householder & Young (1938) and a much more detailed but still somewhat incomplete proof by Stewart (1973, pp. 322–3) (cf. Chipman 1997). A number of incorrect proofs have also appeared.¹ It was subsequently pointed out by Stewart & Sun (1990, pp. 208–210) that the theorem had already been proved (for integral operators and the Hilbert-Schmidt norm) by Schmidt (1907), and subsequently by Mirsky (1960) for unitarily invariant norms. Thus they refer to it as the “Schmidt-Mirsky theorem”. Mirsky’s proof (1960, Theorem 3) relied on the apparatus of symmetric gauge functions introduced by von Neumann (1937). For further discussion see Meyer (1993).

The following proof is based on the method followed in Stewart (1973).

THEOREM 3.4.2 (ECKART-YOUNG, SCHMIDT, MIRSKY). *Let X be a given $n \times k$ matrix of rank $p > l$, where $p \leq m = \min(n, k)$. Then a matrix \tilde{X} that minimizes $\|X - \tilde{X}\|$ over the set $\mathcal{X}_l = \{\tilde{X} \mid \text{rank}(\tilde{X}) \leq l\}$ is given by*

$$(3.4.4) \quad \hat{X} = QD_{(l)}P',$$

where

$$(3.4.5) \quad X = QDP'$$

is a singular-value decomposition of X , and the $n \times k$ matrix $D_{(l)}$ is obtained from D by replacing all but a set of its l largest diagonal elements by 0s. Further, \hat{X} is a minimizer if and only if it is obtained in this way.

PROOF. By Lemma 3.4.2 such a matrix \tilde{X} exists, and by Theorem 3.4.1 it has rank l .

Let a singular-value decomposition of \hat{X} be denoted

$$(3.4.6) \quad \hat{X} = \hat{Q}\hat{D}\hat{P}',$$

where D is an $n \times k$ diagonal matrix of the form

$$(3.4.7) \quad \hat{D} = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix},$$

where in turn S is an $l \times l$ diagonal matrix $\text{diag}(s_1, s_2, \dots, s_l)$ with $s_1 \geq s_2 \geq \dots \geq s_l > 0$. The main task of the proof is to show that $\hat{D} = D_{(l)}$, where the latter is obtained from D in the manner described in the statement of the theorem.

Define

$$(3.4.8) \quad \bar{D} = \hat{Q}'X\hat{P}$$

and partition it conformably with \hat{D} , as

$$(3.4.9) \quad \bar{D} = \begin{bmatrix} \bar{D}_{11} & \bar{D}_{12} \\ \bar{D}_{21} & \bar{D}_{22} \end{bmatrix},$$

¹Cf. Golub & Kahan (1965, p. 220); Rao (1965, p. 56; 1973, p. 70); Ben-Israel & Greville (1974, pp. 246–9; 1980, pp. 246–9). These are discussed in Chipman (1997).

\bar{D}_{11} being of order $l \times l$ and \bar{D}_{22} of order $(n-l) \times (k-l)$. \bar{D} has rank p . Owing to the orthogonal invariance of the Frobenius norm, it follows from (3.4.6) and (3.4.8) that (3.4.1) is equivalent to

$$(3.4.10) \quad \|\bar{D} - \hat{D}\| = \min_{\text{rank}(\tilde{D}) \leq l} \|\bar{D} - \tilde{D}\|.$$

We now show in Steps 1–3 below that the matrix \bar{D} of (3.4.9) must be of the form

$$(3.4.11) \quad \bar{D} = \begin{bmatrix} S & 0 \\ 0 & \bar{D}_{22} \end{bmatrix}.$$

In Step 4 we will show that \bar{D} is orthogonally equivalent to a diagonal matrix, D .

Step 1. First we show that $\bar{D}_{12} = 0$. Suppose not. Then the matrix

$$\tilde{D} = \begin{bmatrix} S & \bar{D}_{12} \\ 0 & 0 \end{bmatrix}$$

has the same rank as S , which is l , and

$$\|\bar{D} - \tilde{D}\| = \left\| \begin{bmatrix} \bar{D}_{11} - S & 0 \\ \bar{D}_{21} & \bar{D}_{22} \end{bmatrix} \right\| < \left\| \begin{bmatrix} \bar{D}_{11} - S & \bar{D}_{12} \\ \bar{D}_{21} & \bar{D}_{22} \end{bmatrix} \right\| = \|\bar{D} - \hat{D}\|.$$

Now define $\tilde{X} = \tilde{Q}\tilde{D}\tilde{P}'$; this matrix also has rank l . Then by the orthogonal invariance of the Frobenius norm we have

$$(3.4.12) \quad \|X - \tilde{X}\| = \|D - \tilde{D}\| < \|D - \hat{D}\| = \|X - \hat{X}\|.$$

Therefore \tilde{X} , which has rank l , is closer to X than \hat{X} . But this contradicts the hypothesis that \hat{X} is a closest matrix to X among all $n \times k$ matrices \tilde{X} of rank l . Therefore $\bar{D}_{12} = 0$.

Step 2. That $\bar{D}_{21} = 0$ is proved in similar fashion.

Step 3a. Now we show that $\text{rank}(\bar{D}_{11}) = l$. Suppose not; then since $\text{rank}(\bar{D}) = p > l$, and \bar{D}_{12} and \bar{D}_{21} have been shown to be zero, we can find a partition

$$\bar{D}_{22} = \begin{bmatrix} \bar{D}_{22,11} & \bar{D}_{22,12} \\ \bar{D}_{22,21} & \bar{D}_{22,22} \end{bmatrix}$$

of the $(n-l) \times (n-l)$ matrix \bar{D}_{22} (if necessary by temporarily interchanging its rows and columns) such that $\bar{D}_{22,11}$ is of order $(p-l) \times (p-l)$ and rank $p-l$, so that the $n \times k$ matrix

$$\tilde{D} = \begin{bmatrix} \bar{D}_{11} & 0 & 0 \\ 0 & \bar{D}_{22,11} & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

has rank l . Accordingly,

$$\begin{aligned} \|\bar{D} - \tilde{D}\| &= \left\| \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \bar{D}_{22,12} \\ 0 & \bar{D}_{22,21} & \bar{D}_{22,22} \end{bmatrix} \right\| \\ &< \left\| \begin{bmatrix} \bar{D}_{11} - S & 0 & 0 \\ 0 & \bar{D}_{22,11} & \bar{D}_{22,12} \\ 0 & \bar{D}_{22,21} & \bar{D}_{22,22} \end{bmatrix} \right\| = \|\bar{D} - \hat{D}\|. \end{aligned}$$

As before, we define $\tilde{X} = \hat{Q}\tilde{D}\hat{P}'$; this matrix also has rank l . Then with this new \tilde{X} , (3.4.12) holds as before and we arrive at a contradiction. Therefore $\text{rank}(\bar{D}_{11}) = l$.

Step 3b. Next we show that $\bar{D}_{11} = S$. Suppose not. Then define

$$\tilde{D} = \begin{bmatrix} \bar{D}_{11} & 0 \\ 0 & 0 \end{bmatrix}.$$

As just shown, this matrix has rank l , and

$$\|\bar{D} - \tilde{D}\| = \left\| \begin{bmatrix} 0 & 0 \\ 0 & \bar{D}_{22} \end{bmatrix} \right\| < \left\| \begin{bmatrix} \bar{D}_{11} - S & 0 \\ 0 & \bar{D}_{22} \end{bmatrix} \right\| = \|\bar{D} - \hat{D}\|,$$

leading to a contradiction, as before. Therefore $\bar{D}_{11} = S$.

From Steps 1–3 it follows that \bar{D} must be of the form (3.4.11).

Step 4. Now let

$$(3.4.13) \quad \bar{D}_{22} = Q_{22}RP'_{22}$$

be a singular-value decomposition of the $(n-l) \times (k-l)$ matrix \bar{D}_{22} , where Q_{22} and P_{22} are, respectively, $(n-l) \times (n-l)$ and $(k-l) \times (k-l)$ orthogonal matrices, and R is an $(n-l) \times (k-l)$ diagonal matrix of singular values of \bar{D}_{22} . Define further the partitions $\hat{P} = [\hat{P}_1, \hat{P}_2]$ and $\hat{Q} = [\hat{Q}_1, \hat{Q}_2]$ of \hat{P} and \hat{Q} into their first l and last $k-l$ and $n-l$ columns, respectively. Finally, define the rectangular $n \times k$ diagonal matrix

$$(3.4.14) \quad D = \begin{bmatrix} S & 0 \\ 0 & R \end{bmatrix}$$

and the $k \times k$ and $n \times n$ matrices

$$(3.4.15) \quad P = [\hat{P}_1, \hat{P}_2] \begin{bmatrix} I_l & 0 \\ 0 & P_{22} \end{bmatrix}, \quad Q = [\hat{Q}_1, \hat{Q}_2] \begin{bmatrix} I_l & 0 \\ 0 & Q_{22} \end{bmatrix},$$

which are readily verified to be orthogonal. Then we verify from (3.4.15), (3.4.14), (3.4.13), (3.4.10), and (3.4.8) that

$$(3.4.16) \quad QDP' = \hat{Q}\bar{D}\hat{P}' = X;$$

thus, QDP' is a singular-value decomposition of X , in accordance with (3.4.5); and from the orthogonal invariance of the Frobenius norm, D is orthogonally equivalent to \bar{D} (and to X):

$$(3.4.17) \quad \|D\| = \|\bar{D}\| = \|X\|.$$

On the other hand, it is clear from (3.4.15), (3.4.7), and (3.4.6) that

$$(3.4.18) \quad Q\hat{D}P' = \hat{Q}\hat{D}\hat{P}' = \hat{X},$$

so that $Q\hat{D}P'$ is a singular-value decomposition of \hat{X} . It remains to show that $\hat{D} = D_{(r)}$, establishing (3.4.4).

From (3.4.16), (3.4.14), and (3.4.13) we have

$$(3.4.19) \quad \|X\|^2 = \|S\|^2 + \|R\|^2 = \|S\|^2 + \|\bar{D}_{22}\|^2,$$

so that X has the diagonal elements of S as l of its singular values, and the sum of squares of its remaining $m - l$ singular values is equal to $\|\bar{D}_{22}\|^2$. From (3.4.16), (3.4.18), (3.4.14), (3.4.7), and (3.4.13) we have

$$(3.4.20) \quad \|X - \hat{X}\| = \|D - \hat{D}\| = \|\bar{D}_{22}\| = \|R\|.$$

Since by hypothesis, (3.4.20) is a minimum (satisfying (3.4.1)), this can only be the case if, in (3.4.19), the diagonal elements of S are the l largest singular values of X , and those of R are the $m - l$ smallest (with possible ties). It follows that, if the singular values of X are ordered as $s_1 \geq s_2 \geq \dots \geq s_l \geq s_{l+1} \geq \dots \geq s_m$, S must contain s_1, s_2, \dots, s_l , and R must contain s_{l+1}, \dots, s_m . (If $s_l = s_{l+1}$, \hat{X} is not unique.) Applying this requirement to (3.4.7) and (3.4.14) we have $\hat{D} = D_{(r)}$ and the main part of the theorem is proved.

We finally come to the last statement of the theorem. Let $X = \check{Q}\check{D}\check{P}'$ be any other singular-value decomposition of X , and let $\check{D}_{(r)}$ be obtained from \check{D} by replacing all but a set of its r largest singular values by 0s. Define $\check{X} = \check{Q}\check{D}_{(r)}\check{P}'$. Then by the orthogonal invariance of the Frobenius norm we have

$$\|X - \check{X}\| = \|\check{D} - \check{D}_{(r)}\| = \|D - D_{(r)}\| = \|X - \hat{X}\|. \quad \square$$

The following elementary proof of Theorem 3.4.2 is based on that of Neudecker in Chipman (1997, pp. 80–81) and Magnus & Neudecker (1999, pp. 359–361). First, a simple lemma will be used.

LEMMA 3.4.4. *Let A, B be $n \times n$ matrices, where B is symmetric. Then $\text{tr}(AB) = \text{tr}(A'B)$.*

PROOF. From the symmetry of B and Lemma 1.1.1,

$$\text{tr}(AB) = \text{tr}(AB)' = \text{tr}(B'A') = \text{tr}(BA') = \text{tr}(A'B). \quad \square$$

ALTERNATIVE PROOF OF THEOREM 3.4.2 (NEUDECKER).² Let $\hat{X} \in \mathcal{X}_l$ be closest to the given $n \times k$ matrix X in the Frobenius norm. (The existence of \hat{X} is assured by Lemma 3.4.2.) Its rows may be expressed without loss of generality as linear combinations of l $1 \times k$ orthonormal vectors, i.e.,

$$\hat{X} = AB', \quad B'B = I_l$$

where A is $n \times l$ and B is $k \times l$. By Theorem 3.4.1, \hat{X} has rank l ; therefore A must have rank l . We wish to find A and B that solve the problem

$$\text{Minimize}_{A,B} \text{tr}[(X - AB)'(X - AB)] \quad \text{subject to} \quad B'B = I,$$

²This proof, like the previous one, assumes the truth of Lemma 3.4.2 and Theorem 3.4.1.

or equivalently,

$$\text{Maximize}_{A,B} \psi \equiv 2 \operatorname{tr}(BA'X) - \operatorname{tr}(A'A) \quad \text{subject to} \quad B'B = I.$$

Setting up the Lagrangean expression

$$\varphi = 2 \operatorname{tr}(BA'X) - \operatorname{tr}(A'A) - \operatorname{tr}[L(B'B - I)],$$

we note that since $B'B$ is symmetric, from Lemma 3.4.4 we may without loss of generality assume the Lagrangean multiplier matrix L to be symmetric (or replaced by $\frac{1}{2}(L + L')$). Using this symmetry we obtain for variations in A and B

$$\begin{aligned} d\varphi &= 2 \operatorname{tr}(B'X')dA + 2 \operatorname{tr}(A'X)dB - 2 \operatorname{tr}(A')dA - \operatorname{tr}(LB')dB - \operatorname{tr}(L'B')dB \\ &= 2 \operatorname{tr}(B'X' - A')dA + 2 \operatorname{tr}(A'X - LB')dB. \end{aligned}$$

Setting $d\varphi = 0$ for arbitrary dA and dB yields, with the given constraint,

$$\begin{aligned} \text{(i)} \quad & XB = A \\ \text{(ii)} \quad & A'X = LB' \\ \text{(iii)} \quad & B'B = I. \end{aligned}$$

From these three equations we obtain

$$\text{(iv)} \quad A' \underline{A} \stackrel{\text{(i)}}{=} \underline{A'XB} \stackrel{\text{(ii)}}{=} \underline{LB'B} \stackrel{\text{(iii)}}{=} L.$$

Since $\operatorname{rank}(A) = l$, $L = A'A$ is positive-definite. From equations (i) and (ii) and the symmetry of L we obtain

$$\text{(v)} \quad X' \underline{XB} \stackrel{\text{(i)}}{=} \underline{X'A'} \stackrel{\text{(ii)}}{=} BL' = BL.$$

From (ii) and (iv) it follows that

$$\psi = 2 \operatorname{tr}(BA'X) - \operatorname{tr}(A'A) = 2 \operatorname{tr}(BLB') - \operatorname{tr}(L) = \operatorname{tr}(L),$$

which is to be a maximum.

Note from (i) that

$$\hat{X} = AB' = XBB',$$

which states (since BB' is idempotent and symmetric) that \hat{X} is a perpendicular projection of the rows of X onto an l -dimensional subspace.

Let T be an orthogonal matrix such that

$$L = T\Lambda T'$$

where Λ is diagonal, and define

$$\tilde{A} = AT, \quad \tilde{B} = BT.$$

Then

$$\tilde{A}'\tilde{A} = T'A'AT = T'LT = \Lambda$$

and

$$\tilde{B}'\tilde{B} = T'B'BT = T'T = I.$$

Equations (i) to (iii) above then become

$$\begin{aligned} \text{(i')} \quad & X\tilde{B} = \tilde{A} \\ \text{(ii')} \quad & \tilde{A}'X = T'LB' = T'LT\tilde{B}' = \Lambda\tilde{B}' \\ \text{(iii')} \quad & \tilde{B}'\tilde{B} = I. \end{aligned}$$

From these equations it follows that

$$X'X\tilde{B} = \tilde{B}\Lambda \quad \text{and} \quad \tilde{B}'\tilde{B} = I.$$

Thus, Λ , whose trace is to be maximized (being equal to the trace of L), is a diagonal matrix of l eigenvalues of $X'X$, and \tilde{B} is the matrix whose l columns constitute an associated orthonormal set of l eigenvectors of $X'X$. $\text{tr}(\Lambda)$ is maximized when these l eigenvalues are a set of l largest eigenvalues of $X'X$.

Now let X have the singular-value decomposition (3.4.4); its singular values s_i (the diagonal elements of D) are the positive square roots of the eigenvalues λ_i of Λ (from $X'X = PD'DP'$), so since \tilde{X} has been assumed to be the closest to X , its singular values must be the l largest singular values of D . \square

3.5 Reduced-rank estimation

It was proposed by Marquardt (1970)—with respect to the regression model

$$(3.5.1) \quad y = X\beta + \varepsilon; \quad E\{\varepsilon\} = 0; \quad \text{Var}\{\varepsilon\} = \sigma^2\Omega,$$

in the special case $\Omega = I$ —that when the observation matrix X is ill-conditioned (as defined by its condition number—see Definition 1.3.3), one can obtain an estimator of β in this regression model which has lower scalar mean-square error than that of the least-squares estimator, by finding the best approximation to X by a matrix $X_{(l)}$ of rank $l < k$ and then replacing the least-squares estimator $\tilde{\beta} = X^\dagger y$ by the estimator $\hat{\beta}_{(l)} = X_{(l)}^\dagger y$. The theory behind this procedure—and a generalization—will be developed in this section.

LEMMA 3.5.1. *Let the $n \times k$ matrix X , of rank $k < n$, have singular-value decomposition $X = Q_1SP'$ (as in (3.2.8)), where the singular values of X are arranged in descending order $s_1 \geq s_2 \geq \dots \geq s_k > 0$, and let its best approximation by an $n \times k$ matrix of rank l , in terms of the Frobenius norm of Definition 3.4.1, be given by $X_{(l)} = Q_1S_{(l)}P'$, where $S_{(l)}$ is obtained from S by replacing s_i by zero for $i = l + 1, \dots, k$. Then*

$$(3.5.2) \quad X_{(l)} = XP_1P_1' = X[I - P_2P_2'],$$

where $P = [P_1, P_2]$ is a partition of the $k \times k$ orthogonal matrix P into its first l and last $r = k - l$ columns. Further, the Moore-Penrose generalized inverse of $X_{(l)}$ is given by

$$(3.5.3) \quad X_{(l)}^\dagger = P_1P_1'X^\dagger = [I - P_2P_2']X^\dagger.$$

PROOF. We denote the partitioned matrices

$$(3.5.4) \quad S = \begin{bmatrix} S_1 & 0 \\ 0 & S_2 \end{bmatrix} \quad \text{and} \quad S_{(l)} = \begin{bmatrix} S_1 & 0 \\ 0 & 0 \end{bmatrix},$$

where S_1 and S_2 are diagonal matrices of orders $l \times l$ and $r \times r$ respectively. Then

$$X_{(l)} = Q_1 \begin{bmatrix} S_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} P_1' \\ P_2' \end{bmatrix} = Q_1 \begin{bmatrix} S_1 P_1' \\ 0 \end{bmatrix}$$

while

$$X P_1 P_1' = Q_1 \begin{bmatrix} S_1 & 0 \\ 0 & S_2 \end{bmatrix} \begin{bmatrix} P_1' \\ P_2' \end{bmatrix} P_1 P_1' = Q_1 \begin{bmatrix} S_1 & 0 \\ 0 & S_2 \end{bmatrix} \begin{bmatrix} P_1' \\ 0 \end{bmatrix} = Q_1 \begin{bmatrix} S_1 P_1' \\ 0 \end{bmatrix},$$

and these are the same, establishing (3.5.2). Likewise,

$$X_{(l)}^\dagger = P S_{(l)}^\dagger Q_1' = [P_1 \quad P_2] \begin{bmatrix} S_1^{-1} & 0 \\ 0 & 0 \end{bmatrix} Q_1' = [P_1 S_1^{-1} \quad 0] Q_1',$$

while

$$\begin{aligned} P_1 P_1' X^\dagger &= P_1 P_1' [P_1 \quad P_2] \begin{bmatrix} S_1^{-1} & 0 \\ 0 & S_2^{-1} \end{bmatrix} Q_1' \\ &= [P_1 \quad 0] \begin{bmatrix} S_1^{-1} & 0 \\ 0 & S_2^{-1} \end{bmatrix} Q_1' = [P_1 S_1^{-1} \quad 0] Q_1', \end{aligned}$$

and these are the same, establishing (3.5.3). \square

In order to pursue Marquardt's result in the general case $\text{Var}\{\varepsilon\} = \sigma^2 \Omega$, and in terms of the matrix concept of mean-square error, it is clear that we will want to employ an estimator that is an oblique generalized inverse of a reduced-rank matrix with respect to Ω . This may be accomplished by premultiplying the variables y , X , and ε in (3.5.1) by $\Omega^{-1/2}$, where $\Omega^{1/2}$ is a symmetric positive-definite square root of Ω , and defining

$$(3.5.6) \quad \dot{y} = \Omega^{-1/2} y, \quad \dot{X} = \Omega^{-1/2} X, \quad \dot{\varepsilon} = \Omega^{-1/2} \varepsilon.$$

We see that

$$\text{Var}\{\dot{\varepsilon}\} = \sigma^2 I.$$

Defining

$$(3.5.7) \quad X^\ddagger = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1},$$

we verify that

$$(3.5.8) \quad X^\ddagger = \dot{X}^\dagger \Omega^{-1/2}.$$

This transformation of variables amounts to replacing the Frobenius norm $\|X\|$ of Definition 3.4.1 by the Ω -norm defined by

$$(3.5.9) \quad \|X\|_{\Omega} = \sqrt{\text{tr}(X'\Omega^{-1}X)}.$$

In terms of this norm the “best approximation” of the matrix $X = \Omega^{1/2}\dot{X}$ by an $n \times k$ matrix of rank l is then

$$(3.5.10) \quad X_{[l]} = \Omega^{1/2}\dot{X}_{(l)} = \Omega^{1/2}\dot{X}\dot{P}_1\dot{P}'_1 = X\dot{P}_1\dot{P}'_1 = X(I - \dot{P}_2\dot{P}'_2),$$

where

$$(3.5.11) \quad \dot{X} = \dot{Q}_1\dot{S}\dot{P}'$$

is a singular-value decomposition of \dot{X} and $\dot{P} = [\dot{P}_1, \dot{P}_2]$ is a partition of the orthogonal matrix \dot{P} into its first l and last $r = k - l$ columns. Then, defining

$$(3.5.12) \quad X_{[l]}^{\ddagger} = \dot{X}_{(l)}^{\ddagger}\Omega^{-1/2} = \dot{P}_1\dot{P}'_1\dot{X}^{\ddagger}\Omega^{-1/2} = \dot{P}_1\dot{P}'_1X^{\ddagger} = (I - \dot{P}_2\dot{P}'_2)X^{\ddagger},$$

we verify that $X_{[l]}^{\ddagger}$ is the oblique generalized inverse of $X_{[l]}$ with respect to $U = I$ and $V = \Omega$.

The following generalizes Marquardt’s theorem (1970, pp. 591–612):

THEOREM 3.5.1. *Let $X_{[l]}$ be the best approximation of X by a matrix of rank $l < k$ in terms of the Ω -norm (3.5.9), and let $X_{[l]}^{\ddagger}$ be defined by (3.5.12). Then:*

(a) *a necessary and sufficient condition for the reduced-rank estimator*

$$(3.5.13) \quad \hat{\beta}_{[l]} = X_{[l]}^{\ddagger}y = (I - \dot{P}_2\dot{P}'_2)X^{\ddagger}y$$

to have matrix mean-square error no greater than that of the generalized least-squares (Gauss-Markov) estimator

$$(3.5.14) \quad \tilde{\beta} = X^{\ddagger}y$$

is that

$$(3.5.15) \quad \sigma^{-2}\dot{P}'_2\beta\beta'\dot{P}_2 \preceq \dot{S}_2^{-2}.$$

(b) *A sufficient condition for (3.5.15) to hold is that*

$$(3.5.16) \quad \lambda \equiv \frac{\beta'\dot{P}_2\dot{S}_2^2\dot{P}'_2\beta}{\sigma^2} \leq 1.$$

(c) *If $r = k - l > 1$, condition (3.5.16) is sufficient for the inequality (3.5.15) to be strict, and if $r = 1$ then conditions (3.5.15) and (3.5.16) are equivalent, hence both necessary and sufficient.*

PROOF. (a) From the singular-value decomposition of X we readily compute the mean-square error of the generalized least-squares estimator,

$$(3.5.17) \quad \text{Risk}\{\tilde{\beta}\} = \sigma^2\dot{P}\dot{S}^{-2}\dot{P}' = \sigma^2(\dot{P}_1\dot{S}_1^{-2}\dot{P}'_1 + \dot{P}_2\dot{S}_2^{-2}\dot{P}'_2).$$

To compute that of the reduced-rank estimator we observe that

$$\begin{aligned}\hat{\beta}_{[l]} - \beta &= X_{[l]}^\dagger \varepsilon - (I - \dot{P}_1 \dot{P}_1') \beta \\ &= \dot{P}_1 \dot{P}_1' X_{[l]}^\dagger \varepsilon - \dot{P}_2 \dot{P}_2' \beta.\end{aligned}$$

Consequently, from the singular-value decomposition (3.5.11) we have

$$\begin{aligned}\text{Risk}\{\hat{\beta}_{[l]}\} &= \sigma^2 \dot{P}_1 \dot{P}_1' (\dot{X}' \dot{X})^{-1} \dot{P}_1 \dot{P}_1' + \dot{P}_2 \dot{P}_2' \beta \beta' \dot{P}_2 \dot{P}_2' \\ (3.5.18) \quad &= \sigma^2 \dot{P}_1 \dot{P}_1' [\dot{P}_1, \dot{P}_2] \begin{bmatrix} \dot{S}_1^{-2} & 0 \\ 0 & \dot{S}_2^{-2} \end{bmatrix} \begin{bmatrix} \dot{P}_1 \\ \dot{P}_2 \end{bmatrix} \dot{P}_1 \dot{P}_1' + \dot{P}_2 \dot{P}_2' \beta \beta' \dot{P}_2 \dot{P}_2' \\ &= \sigma^2 \dot{P}_1 \dot{S}_1^{-2} \dot{P}_1' + \dot{P}_2 \dot{P}_2' \beta \beta' \dot{P}_2 \dot{P}_2' .\end{aligned}$$

The difference between the two mean-square errors (3.5.17) and (3.5.18) is therefore

$$(3.5.19) \quad \text{Risk}\{\tilde{\beta}\} - \text{Risk}\{\hat{\beta}_{[l]}\} = \sigma^2 \dot{P}_2 [\dot{S}_2^{-2} - \sigma^{-2} \dot{P}_2' \beta \beta' \dot{P}_2] \dot{P}_2' .$$

This is nonnegative-definite if and only if (3.5.15) holds.

(b) The sufficient condition (3.5.16) is obtained by making use of the generalized matrix Cauchy-Schwarz inequality $AVA' \succcurlyeq AX(X'V^{-1}X)^{-1}X'A'$ (Lemma 2.4.1) with I_r substituted for “A”, \dot{S}_2^{-2} for “V”, and $\dot{P}_2' \beta$ for “X”, yielding the first inequality in

$$(3.5.20) \quad \dot{S}_2^{-2} \succcurlyeq \dot{P}_2' \beta (\beta' \dot{P}_2 \dot{S}_2^2 \dot{P}_2 \beta)^{-1} \beta' \dot{P}_2 = \lambda^{-1} \sigma^{-2} \dot{P}_2' \beta \beta' \dot{P}_2 \succcurlyeq \sigma^{-2} \dot{P}_2' \beta \beta' \dot{P}_2 .$$

The equality in (3.5.20) follows from the definition of λ in (3.5.16), and the second inequality in (3.5.20) follows from $\lambda^{-1} \geq 1$. The entire inequality in (3.5.20) is simply (3.5.15).

(c) Since the matrix \dot{S}_2^{-2} on the right side of the inequality (3.5.15) has rank r , whereas the matrix $\beta \beta'$ on the left has rank 1, it is impossible for equality to hold in (3.5.15) unless $r = 1$. Thus, the inequality is necessarily strict if $r > 1$. If $r = 1$, since \dot{P}_2 is $k \times r$ and \dot{S}_2 is $r \times r$, both sides of (3.5.15) are scalars, hence (3.5.15) reduces to (3.5.16), which is then necessary and sufficient. \square

REMARK 3.5.1. *It will be shown in the next chapter (Theorem 4.2.3 and formula (4.2.30)) that the reduced-rank estimator (3.5.13) is the minimum-variance conditionally unbiased estimator of β subject to the homogeneous linear restriction $\dot{P}_2' \beta = 0$.*

The duality referred to in Remark 3.5.1 was apparently first noticed in an unpublished paper by Johnson & Wallace (1969).

COROLLARY 3.5.1. *A sufficient condition for $\text{Risk}\{\hat{\beta}_{[l]}\} \preccurlyeq \text{Risk}\{\tilde{\beta}\}$ is*

$$(3.5.21) \quad \frac{\beta' \dot{P}_2 \dot{P}_2' \beta}{\sigma^2} \leq \frac{1}{\dot{s}_{l+1}^2} .$$

A sufficient condition for (3.5.21) is in turn

$$(3.5.22) \quad \frac{\beta' \beta}{\sigma^2} \leq \frac{1}{\dot{s}_{l+1}^2} ,$$

where the \dot{s}_i^2 are the eigenvalues of $X'\Omega^{-1}X$, in descending order.

PROOF. Applying the matrix Cauchy-Schwarz inequality (Lemma 2.4.1) for $A = V = I$, we have $I \succcurlyeq X(X'X)^{-1}X'$, hence replacing “ X ” by $\dot{P}'_2\beta$ we obtain

$$(3.5.23) \quad I_r \succcurlyeq \dot{P}'_2\beta(\beta'\dot{P}_2\dot{P}'_2\beta)^{-1}\beta'\dot{P}_2, \quad \text{or} \quad (\beta'\dot{P}_2\dot{P}'_2\beta)I_r \succcurlyeq \dot{P}'_2\beta\beta'\dot{P}_2,$$

Assuming (3.5.21) to hold it follows that

$$\dot{S}_2^{-2} \succcurlyeq \dot{s}_{l+1}^{-2}I_r \succcurlyeq \sigma^{-2}(\beta'\dot{P}_2\dot{P}'_2\beta)I_r \succcurlyeq \sigma^{-2}\dot{P}'_2\beta\beta'\dot{P}_2,$$

i.e., (3.5.15) holds. Since $I - \dot{P}_2\dot{P}'_2 = \dot{P}_1\dot{P}'_1 \succcurlyeq 0$, $\beta'\dot{P}_2\dot{P}'_2\beta \leq \beta'\beta$, hence (3.5.21) follows from (3.5.22). \square

Thus, if one is able to impose an *a priori* upper bound on $\beta'\beta/\sigma^2$, and is one has data on the singular values \dot{s}_i of $\dot{X} = \Omega^{-1/2}X$, one may select the appropriate dimensionality l by choosing the largest singular value that satisfies (3.5.22)—or still more accurately, (3.5.21).

We may note that condition (3.5.21) is stronger than, since it implies, (3.5.16); for $\dot{S}_2^2 \preccurlyeq \dot{s}_{l+1}^2 I_r$ from the definition of \dot{S}_2^2 , hence (3.5.21) implies

$$\lambda = \frac{\beta'\dot{P}_2\dot{S}_2^2\dot{P}'_2\beta}{\sigma^2} \leq \dot{s}_{l+1}^2 \frac{\beta'\dot{P}_2\dot{P}'_2\beta}{\sigma^2} \leq 1.$$

While the still stronger condition (3.5.22) has the advantage of not requiring computation of \dot{P}_2 , condition (3.5.16) has the advantage—as will become clear in the next chapter—that it can be tested from the data (provided of course ε is normally distributed), since λ is the noncentrality parameter of the noncentral F -distribution used for testing the null hypothesis $\lambda = 1$ ($\dot{P}'_2\beta = 0$) against the alternative hypothesis $\lambda > 1$ ($\dot{P}'_2\beta \neq 0$); see Remark 3.5.1 above and Theorem 4.3.1 below.

Let us now consider Marquardt's criterion. Marquardt used a scalar criterion of mean-square error, equal to the trace of the matrix measure used here. If an estimator $\hat{\beta}_{[l]}$ has lower matrix-mean-square error than $\tilde{\beta}$, this means that *each component* of $\hat{\beta}_{[l]}$ has lower mean-square error than the corresponding component of $\tilde{\beta}$. With the scalar definition, this need not be the case; some components could have lower, others higher, mean-square error; only an average of them has lower mean-square error. It is to be expected, then, that the matrix definition requires more stringent conditions. Taking the trace of both sides of (3.5.15) we obtain

$$(3.5.24) \quad \frac{\beta'\dot{P}_2\dot{P}'_2\beta}{\sigma^2} \leq \sum_{i=l+1}^k \frac{1}{\dot{s}_i^2}.$$

This is clearly a much weaker condition than (3.5.21); Marquardt (1970, p. 601) actually specified the somewhat more stringent condition

$$(3.5.25) \quad \frac{\beta'\beta}{\sigma^2} \leq \sum_{i=l+1}^k \frac{1}{\dot{s}_i^2},$$

which implies (3.5.24). Clearly, (3.5.21) and (3.5.22) imply (3.5.24) and (3.5.25) respectively. However, Marquardt's condition has a curious interpretation: if the inequality is satisfied for some $l < k$, then it is satisfied for any $l' < l$; in particular, if it is satisfied for $l = k - 1$, then it is also satisfied for $l = 1$. Thus it provides no guide for choosing the appropriate rank, l . On the other hand, (3.5.22), or better still (3.5.21), provides just such a guide.

3.6 Exercises

1. An investigator wishes to estimate the parameters β_1 and β_2 in the model

$$y_t = x_{t1}\beta_1 + x_{t2}\beta_2 + \varepsilon_t, \quad E\{\varepsilon_t\} = 0, \quad E\{\varepsilon_t\varepsilon_{t'}\} = \delta_{tt'}\sigma^2$$

$$\text{where } t, t' = 1, 2, \dots, n \quad \text{and} \quad \delta_{tt'} = \begin{cases} 1 & \text{for } t = t', \\ 0 & \text{for } t \neq t', \end{cases}$$

where the column vectors $x^j = (x_{1j}, x_{2j}, \dots, x_{nj})'$ have been normalized to have length 1 for $j = 1, 2$. This investigator notices that the correlation $r = \sum_{t=1}^n x_{t1}x_{t2}$ between the two independent variables is rather close to 1, and therefore decides to estimate the two parameters by finding the best approximation of $X = [x^1, x^2]$ by an $n \times 2$ matrix $X_{(1)}$ of rank 1 (the distance between two matrices X and X^* being defined by the Frobenius norm $\|Z\|$ of their difference $Z = X - X^*$), and then estimating the 2×1 vector β by $\hat{\beta}_{(1)} = X_{(1)}^\dagger y$ (the ‘‘Marquardt estimator’’), where \dagger denotes the Moore-Penrose generalized inverse.

- Show that the best approximation of X by a matrix of rank 1 is the $n \times 2$ matrix each of whose columns is the average of the two.
- Obtain the formula for the Marquardt estimator.
- Show that the Marquardt estimator estimates each parameter β_j by the simple average of the least-squares estimators $\tilde{\beta}_1$ and $\tilde{\beta}_2$.
- Find the expressions for the (scalar) mean-square errors

$$E\{(\hat{\beta}_{(1)} - \beta)'(\hat{\beta}_{(1)} - \beta)\} \quad \text{and} \quad E\{(\tilde{\beta} - \beta)'(\tilde{\beta} - \beta)\}$$

of $\hat{\beta}_{(1)}$ and $\tilde{\beta}$ respectively.

- Show that $\hat{\beta}_{(1)}$ has lower mean-square error than $\tilde{\beta}$ if and only if

$$\frac{(\beta_1 - \beta_2)^2}{2\sigma^2} < \frac{1}{1 - r}.$$

2. What is the condition number of the matrix X of the previous question?

3.7 References

- BELLMAN, RICHARD. *Introduction to Matrix Analysis*. New York: McGraw-Hill Book Company, Inc., 1960. xx, 328 pp.
- BEN ISRAEL, ADI, and GREVILLE, THOMAS N. E. *Generalized Inverses: Theory and Applications*. New York: John Wiley & Sons, Inc., 1974. xi, [3], 395 pp. Reprint edition with corrections, Huntington, New York: Robert E. Krieger Publishing Company, 1980.
- CHIPMAN, JOHN S. ‘‘On Least Squares with Insufficient Observations.’’ *Journal of the American Statistical Association*, 59 (December 1964), 1078–1111. Corrigendum: 60 (December 1965), 1249.
- CHIPMAN, JOHN S. ‘‘Estimation and Aggregation in Econometrics: An Application of the Theory of Generalized Inverses.’’ In Nashed (1976), 549–769.

- CHIPMAN, JOHN S. “‘Proofs’ and Proofs of the Eckart-Young Theorem,” with an Appendix by Heinz Neudecker. In Jerome A. Goldstein, Neil E. Gretskey, and J. J. Uhl, Jr. (eds.), *Stochastic Processes and Functional Analysis. In Celebration of M. M. Rao’s 65th Birthday* (New York: Marcel Dekker, Inc., 1997), 71–83.
- CHIPMAN, JOHN S. “Linear Restrictions, Rank Reduction, and Biased Estimation in Linear Regression.” *Linear Analysis and Its Applications*, 289 (1999), 55–74.
- COURANT, R. “Zur Theorie der kleinen Schwingungen.” *Zeitschrift für angewandte Mathematik und Mechanik*, 2 (1922), 278–285.
- ECKART, CARL, and YOUNG, GALE. “The Approximation of One Matrix by Another of Lower Rank.” *Psychometrika*, 1 (September 1936), 211–218.
- ECKART, CARL, and YOUNG, GALE. “A Principal Axis Transformation for Non-Hermitian Matrices.” *Bulletin of the American Mathematical Society*, 45 (February 1939), 118–121.
- FISCHER, ERNST. “Über quadratische Formen mit reellen Koeffizienten.” *Monatshrift für Mathematik und Physik*, 16 (1905), 234–249.
- FRISCH, RAGNAR. *Statistical Confluence Analysis by Means of Complete Regression Systems*. Oslo: Universitetets Økonomisk Institutt, Publikasjon nr. 5, 1934. 192 pp.
- GOLUB, GENE H., and KAHAN, W. “Calculating the Singular Values and Pseudo-Inverse of a Matrix.” *SIAM Journal of Numerical Analysis* [B], 2 (1965), 205–224.
- GOLUB, GENE H., and REINSCH, C. “Singular Value Decomposition and Least Squares Solutions.” *Numerische Mathematik*, 14 (1970), 403–420. Reprinted in Wilkinson & Reinsch (1971), 134–151.
- GOLUB, GENE H., and STYAN, GEORGE P. H. “Numerical Computations for Univariate Linear Models.” *Journal of Statistical Computation and Simulation*, 2 (1973), 253–274.
- GOLUB, GENE H., and VAN LOAN, CHARLES F. *Matrix Computations*. Baltimore: The Johns Hopkins University Press, 1983. xvi, [2], 476, [1] pp.
- GUNST, RICHARD F., and MASON, ROBERT L. “Generalized Mean Squared Error Properties of Regression Estimators.” *Communications in Statistics*, A5 (1976), 1501–1508.
- GUNST, RICHARD F., and MASON, ROBERT L. “Biased Estimation in Regression: An Evaluation using Mean Squared Error.” *Journal of the American Statistical Association*, 72 (September 1977), 616–628.
- GUNST, RICHARD F., WEBSTER, JOHN T., and MASON, ROBERT L. “A Comparison of Least Squares and Latent Root Regression Estimators.” *Technometrics*, 18 (1976), 75–83.
- HAWKINS, DOUGLAS M. “On the Investigation of Alternative Regressions by Principal Component Analysis.” *Applied Statistics*, 22 (1973), 275–286.
- HOUSEHOLDER, A. S., and YOUNG, GALE. “Matrix Approximations and Latent

- Roots." *American Mathematical Monthly*, 45 (March 1938), 165–171.
- JOHNSON, THOMAS, and WALLACE, T. D. "Principal Components and Multicollinearity." Workshop Discussion Paper, Department of Economics, North Carolina State University, Raleigh, NC, June 1969.
- LOTT, WILLIAM F. "The Optimal Set of Principal Component Restrictions on a Least-Squares Regression." *Communications in Statistics*, 2 (1973), 449–463.
- MAGNUS, JAN R., and NEUDECKER, HEINZ. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Chichester and New York: John Wiley & Sons, 1991, 1994. Revised edition, c. 1999. xviii, 395, [8] pp.
- MANSFIELD, EDWARD R., WEBSTER, JOHN T., and GUNST, RICHARD F. "An Analytic Variable Selection Technique for Principal Component Regression." *Applied Statistics*, 26 (1977), 34–40.
- MARQUARDT, DONALD W. "Generalized Inverses, Ridge Regression, Biased Linear Estimation, and Nonlinear Estimation." *Technometrics*, 12 (August 1970), 591–612.
- MASON, ROBERT L., GUNST, RICHARD F., and WEBSTER, JOHN T. "Regression Analysis and Problems of Multicollinearity." *Communications in Statistics*, 4 (1975), 277–292.
- MASSY, WILLIAM F. "Principal Components Regression in Exploratory Statistical Research." *Journal of the American Statistical Association*, 60 (March 1965), 234–256.
- MEYER, RENATE. *Matrix-Approximation in der multivariaten Statistik*. Aachen: Verlag der Augustinus Buchhandlung, 1993. [6], 131, [2] pp.
- MIRSKY, L. "Symmetric Gauge Functions and Unitarily Invariant Norms." *Quarterly Journal of Mathematics*, Oxford Second Series, 11 (March 1960), 50–59.
- NASHED, M. ZUHAIR (ed.). *Generalized Inverses and Applications*. New York: Academic Press, 1976. xiv, 1054 pp.
- NEUMANN, J[OH]N VON. "Some Matrix-Inequalities and Metrization of Metric-Space." *Tomsk University Review*, 1 (1937), 286–300.
- NOBLE, B. "Methods for Computing the Moore-Penrose Generalized Inverse, and Related Matters." In Nashed (1976), 245–301.
- PENROSE, R. "A Generalized Inverse for Matrices." *Proceedings of the Cambridge Philosophical Society*, 51 (1955), 406–413.
- RAO, C. RADHAKRISHNA. *Linear Statistical Inference and Its Applications*. New York: John Wiley & Sons, 1965. xviii, [2], 522 pp. 2nd edition, 1973. xx, [2], 625 pp.
- RAO, C. RADHAKRISHNA. "Separation Theorems for Singular Values of Matrices and Their Applications in Multivariate Analysis." *Journal of Multivariate Analysis*, 9 (1979), 362–377.
- RAO, C. RADHAKRISHNA. "Matrix Approximations and Reduction of Dimen-

- sionality in Multivariate Statistical Analysis." In Paruchuri R. Krishnaiah, ed., *Multivariate Analysis – V. Proceedings of the Fifth International Symposium on Multivariate Analysis* (Amsterdam: North-Holland Publishing Company, 1980), 3–22.
- SCHMIDT, ERHARD. "Zur Theorie der linearen und nichtlinearen Integralgleichungen. I. Teil: Entwicklung willkürlicher Funktionen nach Systemen vorgeschriebener." *Mathematische Annalen*, 63 (1907), 433–476.
- SILVEY, S. D. "Multicollinearity and Imprecise Estimation." *Journal of the Royal Statistical Society [B]*, 31 (1969), 539–552.
- SONDERMANN, DIETER. "Best Approximate Solutions to Matrix Equations under Rank Restrictions." Report No. 23/80, Institute for Advanced Studies, The Hebrew University of Jerusalem, Mount Scopus, Israel, August 1980.
- STEWART, G. W. *Introduction to Matrix Computations*. New York: Academic Press, 1973. xi, [3], 441 pp.
- STEWART, G. W., and SUN, JI-GUANG. *Matrix Perturbation Theory*. San Diego: Academic Press, Inc., 1990. xvi, 365 pp.
- WEBSTER, JOHN T., GUNST, RICHARD F., and MASON, ROBERT L. "Latent Root Regression Analysis." *Technometrics*, 16 (1974), 513–522.
- WILKINSON, J. H., and REINSCH, C. H. (eds.). *Linear Algebra*. New York: Springer-Verlag, 1971. VIII, [2], 439, [3] pp.