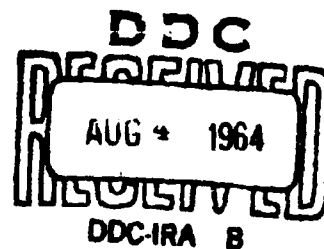


CS10
603163

ON CLOSEST NORMAL MATRICES

BY
ROBERT L. CAUSEY

TECHNICAL REPORT CS10
JUNE 30, 1964



COMPUTER SCIENCE DIVISION
School of Humanities and Sciences
STANFORD UNIVERSITY



CS10

ON CLOSEST NORMAL MATRICES

BY
ROBERT L. CAUSEY

PREPARED UNDER CONTRACT Nonr-225(37)
(NR-044-211)
OFFICE OF NAVAL RESEARCH

Reproduction in Whole or in Part is Permitted for
any Purpose of the United States Government

COMPUTER SCIENCE DIVISION
School of Humanities and Sciences
STANFORD UNIVERSITY



ACKNOWLEDGMENTS

The author wishes to express his sincere appreciation and thanks to the following persons:

Professor George E. Forsythe, Stanford University, for his encouragement and for his extremely valuable suggestions and comments in directing this dissertation.

Professor Ingram Olkin, Stanford University, for a critical appraisal of an earlier manuscript and for suggestions which led to a number of improvements, especially in Section 4.2.

Professor John G. Herriot, Stanford University, for reading the manuscript and for numerous valuable suggestions.

Dr. Patricia J. Eberlein, University of Rochester, for some valuable private communications concerning her own work on Mirsky's conjecture.

Professor J. Wallace Givens, Northwestern University and Argonne National Laboratory, for calling the author's attention (before research on this dissertation was initiated) to the useful notion of a differentiable curve of matrices and for some general comments on the problem of this dissertation.

Professor Peter Henrici, Eidgenössische Technische Hochschule, Zurich, Switzerland, Dr. J. H. Wilkinson, National Physical Laboratory, Teddington, England, and Professors Menahem M. Schiffer and Charles Loewner of Stanford University for stimulating conversations on the investigations of this thesis.

My wife, Barbara, for constant help and understanding.

Mrs. Jerri Rudnick for a most careful and competent typing of the manuscript and Mrs. Betty Jo Prine for seeing the manuscript safely through the typing and reproduction phases.

This work was supported in part by the Lockheed Missiles and Space Company Independent Research Program.

TABLE OF CONTENTS

Chapter		Page
1.	INTRODUCTION AND PRELIMINARIES	1
	1.0 Introduction	1
	1.1 Notation and Preliminary Definitions	3
	1.2 Distance Problems; Mirsky's Conjecture and Bound	7
	1.3 Differentiable Curves in Matric Space	11
	1.4 Normal Matrices	14
	1.5 Existence of ν -Minimal Matrices	19
2.	EQUIVALENCE THEOREMS CONCERNING MIRSKY'S CONJECTURE, ν -MINIMAL MATRICES, AND DISTANCE FORMULAS	22
	2.1 Mirsky's Conjecture	22
	2.2 Invariance of ν -Minimal Matrices	26
	2.3 Invariance of Distance Formulas	28
3.	IMPROVEMENTS OF MIRSKY'S BOUND	32
	3.1 A New Bound	32
	3.2 Methods for Obtaining Other Bounds	40
4.	NECESSARY CONDITIONS FOR ϵ -MINIMAL MATRICES	43
	4.1 Primary Results	43
	4.2 Use of Lagrange Multipliers in Finding Necessary Conditions	56
5.	CHARACTERIZATION OF ϵ -MINIMAL MATRICES	74
6.	ϵ -MINIMAL MATRICES OF ORDER 2	82
	6.1 Preliminaries Concerning 2 by 2 Matrices	82
	6.2 Determination of all 2 by 2 ϵ -Minimal Matrices	87
	6.3 The Maximum Problem 5.11 in the Case $n = 2$	103
7.	FURTHER NECESSARY CONDITIONS FOR ϵ -MINIMAL MATRICES	105
8.	COUNTEREXAMPLES TO MIRSKY'S CONJECTURE	113
9.	THE FIELD OF VALUES AND EIGENVALUE OF ϵ -MINIMAL MATRICES	117
10.	A GENERALIZATION OF THE JACOBI AND GOLDSTINE-HORWITZ METHODS	123

BLANK PAGE

CHAPTER 1

INTRODUCTION AND PRELIMINARIES

1.0 Introduction.

The central problem considered in this paper is the following: Given an n by n matrix A of complex elements, find those normal matrices (called ν -minimal matrices) of order n which are closest to A in the sense of a metric defined in terms of a norm ν . A related problem is that of determining the distance $d_\nu(A)$ between A and the subset \mathcal{N} of all normal matrices of order n . The historical background and precise definition of these problems are given in Section 1.2. The distance problem mentioned above was first studied by Mirsky, who offered Conjecture 1.35 as the general solution for all unitarily invariant norms.

After a proof that ν -minimal matrices always exist (Section 1.5), it is shown in Chapter 2 that the property of ν -minimality and certain other quantities are invariant under certain transformations of matrix space. An inequality of Mirsky is sharpened in Chapter 3. A number of important necessary conditions for ϵ -minimal matrices (ϵ denotes the familiar Euclidean norm) are derived in Chapters 4 and 7. In Chapter 9 it is shown that all eigenvalues of an ϵ -minimal matrix lie in the field of values $F(A)$ of A ; these eigenvalues are shown to be special extreme points of $F(A)$ in the case $n = 2$. In the final Chapter 10 an iterative computational procedure for finding $d_\epsilon(A)$ or an ϵ -minimal matrix is proposed, but its convergence is not proved.

Probably the most important results of the paper are the following:

- 1) A characterization of any ϵ -minimal matrix in terms of the Maximum Problem 5.11 (Theorem 5.13).
- 2) A determination of all ϵ -minimal matrices of order 2 (Theorem 6.24).
- 3) Mirsky's Conjecture 1.35 is shown to be true for $v = \epsilon$ and $n = 2$ (Theorem 6.80), false for $v = \epsilon$ and $n \geq 3$ (Theorem 8.5) and false for $n \geq 2$ and $v = v_p$ ($2 < p \leq \infty$), where v_p is defined by (1.15) (Theorem 8.9).

[Note: The fact that Mirsky's conjecture is false for $v = \epsilon$ and $n \geq 3$ was first proved by P. J. Eberlein.]

All results stated herein which are not specifically labeled as known or for which no reference is given are believed to be new.

1.1 Notation and Preliminary Definitions.

Let \underline{R} and \underline{C} denote respectively the real and complex number fields and let \mathcal{M}_n denote the algebra of all n by n matrices over \underline{C} , where n is a positive integer. If $X \in \mathcal{M}_n$ we denote its complex conjugate transpose by X^* . A matrix X is called normal if $X^*X = XX^*$; in particular X is called hermitian if $X^* = X$, skew-hermitian if $X^* = -X$, and unitary if $X^*X = I$, where I denotes the identity matrix of the same order as X . Let \mathcal{N}_n , \mathcal{H}_n , \mathcal{U}_n , and \mathcal{D}_n denote respectively the subsets of all normal, hermitian, unitary and diagonal matrices in \mathcal{M}_n . (We shall sometimes omit the subscript n , if the value of n need not be specified.)

For the meaning of terminology or notation not explicitly defined in this paper, the reader is referred to one of the standard textbooks on the theory of matrices (e.g., Perlis [21][†]).

A real valued function ν defined on \mathcal{M} is called a norm if the conditions

$$(1.1) \quad \nu(A) > 0 \quad \text{if } A \neq 0 ,$$

$$(1.2) \quad \nu(cA) = |c| \nu(A) ,$$

$$(1.3) \quad \nu(A + B) \leq \nu(A) + \nu(B)$$

are satisfied for all $A, B \in \mathcal{M}$ and for all $c \in \underline{C}$. A norm ν is said

†

Numbers in square brackets refer to references listed in the bibliography at the end of this paper.

to be unitarily invariant if, in addition to (1.1) - (1.3),

$$(1.4) \quad v(UA) = v(AU) = v(A)$$

holds for all $A \in \mathcal{M}$ and for all $U \in \mathcal{U}$. Moreover, a norm v is said to be multiplicative if, for arbitrary $A, B \in \mathcal{M}$,

$$(1.5) \quad v(AB) \leq v(A) v(B) .$$

Unitarily invariant norms were characterized by von Neumann [19] (see also Fan and Hoffman [9]) as follows. A real-valued function $\varphi(u) = \varphi(u_1, u_2, \dots, u_n)$ defined for all real n -vectors $u = (u_1, u_2, \dots, u_n)$ is called a symmetric gauge function if it satisfies, for arbitrary real vectors u, v and real scalars α , the following conditions

$$(1.6) \quad \varphi(u) > 0 \quad \text{if } u \neq 0 ,$$

$$(1.7) \quad \varphi(\alpha u) = |\alpha| \varphi(u) ,$$

$$(1.8) \quad \varphi(u + v) \leq \varphi(u) + \varphi(v) ,$$

$$(1.9) \quad \varphi(u_1, u_2, \dots, u_n) = \varphi(\sigma_1 u_{\pi_1}, \sigma_2 u_{\pi_2}, \dots, \sigma_n u_{\pi_n}) .$$

where σ_i can be either of the signs ± 1 ($i = 1, 2, \dots, n$) and where $(\pi_1, \pi_2, \dots, \pi_n)$ is any permutation of $(1, 2, \dots, n)$.

1.10 Definition. Let $A \in \mathcal{M}$. The nonnegative square roots of the eigenvalues of A^*A are called the singular values of A .

Remark. Since the spectrum of AA^* coincides with that of A^*A , the singular values of A^* are the same as those of A .

1.11 Theorem (von Neumann [19]). A norm v on M_n is unitarily invariant if and only if there exists a symmetric gauge function φ_v of n real variables such that

$$(1.12) \quad v(A) = \varphi_v(\alpha_1, \alpha_2, \dots, \alpha_n) \quad \text{for all } A \in M_n$$

where $\alpha_1, \alpha_2, \dots, \alpha_n$ are the singular values of A .

As examples of symmetric gauge functions we may cite

$$(1.13) \quad \varphi_p(u_1, \dots, u_n) = \left(\sum_{i=1}^n |u_i|^p \right)^{1/p} \quad \text{for } 1 \leq p < \infty.$$

As $p \rightarrow \infty$, the function in (1.13) converges to

$$(1.14) \quad \varphi_\infty(u_1, u_2, \dots, u_n) = \max_{i=1, \dots, n} (|u_i|)$$

which is also a symmetric gauge function. Thus, for $1 \leq p \leq \infty$,

$$(1.15) \quad v_p(A) = \varphi_p(\alpha_1, \dots, \alpha_n)$$

is a unitarily invariant norm, where $\alpha_1, \dots, \alpha_n$ are the singular values of A . The norm v_∞ is known as the spectral norm of A and is sometimes denoted by the symbol σ , i.e., $\sigma(A) = v_\infty(A)$ for all $A \in M_n$. For $p = 2$ $v_p(A)$ in (1.15) coincides with the familiar Euclidean norm $\epsilon(A)$ defined by

$$(1.16) \quad \epsilon^2(A) = \sum_{i,j=1}^n |a_{ij}|^2$$

where $A = (a_{ij}) \in \mathcal{M}_n$ and where by $\epsilon^2(A)$ we mean $[\epsilon(A)]^2$. It can be shown (see e.g., Faddeev and Faddeeva [7] pp. 105-111) that both σ and ϵ are multiplicative norms.

It is well known that \mathcal{M}_n is a Banach algebra with respect to the norm (1.16), that is, \mathcal{M}_n is a Banach space when considered as a linear space with norm ϵ , and the multiplication operation (transformation) $(A, B) \rightarrow AB$ is a continuous mapping from the product space $\mathcal{M}_n \times \mathcal{M}_n$ onto \mathcal{M}_n . (Cf. Hille and Phillips [13], p. 22.)

1.17 Definition. Let the eigenvalues of $M \in \mathcal{M}_n$ be denoted by $\lambda_1(M), \lambda_2(M), \dots, \lambda_n(M)$ in some order. Then $\Omega(M)$ is defined by

$$\Omega(M) = \text{diag}(\lambda_1(M), \dots, \lambda_n(M)) .$$

For any $M = (m_{ij}) \in \mathcal{M}_n$ the trace of M is defined by

$$(1.18) \quad \text{tr}(M) = \sum_{i=1}^n m_{ii} .$$

We note the following well-known properties of the trace functional:

$$(1.19) \quad \text{tr}(AB) = \text{tr}(BA) ;$$

$$(1.20) \quad \text{tr}(\alpha A + \beta B) = \alpha \text{tr}(A) + \beta \text{tr}(B) ;$$

$$(1.21) \quad \text{tr}(M^*) = \overline{\text{tr}(M)} ;$$

$$(1.22) \quad \epsilon^2(M) = \text{tr}(M^*M) = \text{tr}(MM^*) .$$

Here A, B, M are any matrices in \mathcal{M} and α, β are any complex numbers.

1.23 Definition. Let $M = (m_{ij}) \in \mathcal{M}_n$. The diagonal of M is defined by

$$(1.24) \quad \text{dg}(M) = \text{diag}(m_{11}, \dots, m_{nn}) .$$

Likewise the off-diagonal of M is defined by

$$(1.25) \quad \text{offdg}(M) = M - \text{dg}(M) .$$

1.2 Distance Problems; Mirsky's Conjecture and Bound.

Let $A \in \mathcal{M}$, let ν be any norm on \mathcal{M} , and let \mathcal{S} be any subset of \mathcal{M} . By a distance problem we mean the problem of determining the "distance"

$$(1.26) \quad \inf_{X \in \mathcal{S}} \nu(A - X)$$

between A and \mathcal{S} with respect to the norm ν . In any case in which the infimum (1.26) is attained by a matrix in \mathcal{S} we may consider the related minimum problem of finding (at least one matrix and preferably all) matrices $X_0 \in \mathcal{S}$ such that

$$\nu(A - X_0) = \min_{X \in \mathcal{S}} \nu(A - X) .$$

Several such problems have been considered and solved in the past. We now describe some of these results.

1.27 Theorem (Fan and Hoffman [9]). Let $A \in \mathcal{M}$, let \mathcal{H} denote the subset of all hermitian matrices in \mathcal{M} , and let ν denote any unitarily invariant norm on \mathcal{M} . Then

$$(1.28) \quad \min_{X \in \mathcal{H}} \nu(A - X) = \nu(A - \frac{1}{2}(A + A^*)) = \frac{1}{2} \nu(A - A^*) .$$

1.29 Theorem (Fan and Hoffman [9]). Let \mathcal{U} denote the subset of all unitary matrices in \mathcal{M} and let ν denote any unitarily invariant norm on \mathcal{M} . Let $A \in \mathcal{M}$ and suppose $A = UH$ where $H \in \mathcal{H}$ is positive semidefinite and $U \in \mathcal{U}$. Then

$$(1.30) \quad \min_{X \in \mathcal{U}} \nu(A - X) = \nu(\text{diag}(\alpha_1, \dots, \alpha_n) - I) = \nu(A - U)$$

where $\alpha_1, \dots, \alpha_n$ are the singular values of A .

The next result is apparently new, although its interpretation (see Amir-Moez and Horn [2] and [9] regarding a well-known analogy between matrices and complex numbers) and the method of proof are strictly analogous to those associated with Theorem 1.27.

1.31 Theorem. Let $A \in \mathcal{M}$, let \mathcal{S} denote the subset of all skew-hermitian matrices in \mathcal{M} , and let ν denote any unitarily invariant norm on \mathcal{M} . Then

$$(1.32) \quad \min_{X \in \mathcal{S}} \nu(A - X) = \nu(A - \frac{1}{2}(A - A^*)) = \frac{1}{2} \nu(A + A^*) .$$

Proof. Let S be any skew-hermitian matrix. We have

$$A - \frac{A-A^*}{2} = \frac{A-S}{2} + \frac{A^*-S}{2} = \frac{A-S}{2} + \frac{(A-S)^*}{2}$$

whence

$$v(A - \frac{1}{2}(A - A^*)) \leq \frac{1}{2} v(A - S) + \frac{1}{2} v((A - S)^*) .$$

By Theorem 1.11 and the remark following Definition 1.10,

$v((A - S)^*) = v(A - S)$; consequently

$$(1.33) \quad v(A - \frac{1}{2}(A - A^*)) \leq v(A - S)$$

holds for all $S \in \mathcal{S}$. This proves (1.32).

Let k denote an integer such that $1 \leq k \leq n$. In [15] Mirsky solved distance problems for the subsets $\{X ; X \in \mathcal{M}_n \text{ and } \text{rank}(X) \leq k\}$ and $\{X ; X \in \mathcal{M}_n \text{ and } \text{rank}(X) = k\}$. As in the above results, the formula for the distance can be put in the same form for all unitarily invariant norms.

In this paper we shall be primarily interested in distance and extremum problems associated with the subset \mathcal{N} of all normal matrices in \mathcal{M} . The distance problem was apparently first studied by Mirsky [15]. Let A be any fixed element of \mathcal{M} , let v be any norm on \mathcal{M} and define

$$(1.34) \quad d_v(A) = \inf_{X \in \mathcal{N}} v(A - X) .$$

Mirsky was unable to determine $d_v(A)$, even for special choices of v , but he obtained an upper bound for $d_\epsilon(A)$ (see Theorem 1.37 below) and

offered the following conjecture for the general solution when ν is unitarily invariant.

1.35 Conjecture (Mirsky). Let ν denote any unitarily invariant norm on \mathcal{M}_n . Then

$$(1.36) \quad d_{\nu}^2(A) = \frac{1}{2}(\nu^2(A) - \nu^2(\Omega(A)))$$

holds for all $A \in \mathcal{M}_n$, where $d_{\nu}^2(A) = [d_{\nu}(A)]^2$, $\nu^2(A) = [\nu(A)]^2$ and where $\Omega(A)$ is defined in Definition 1.17.

Note. By (1.9) the right side of (1.36) is independent of the order of the λ 's in $\Omega(A)$. The singular values of A are the eigenvalues of the positive semidefinite square root (denoted by $(A^*A)^{\frac{1}{2}}$) of A^*A ; consequently $\nu(A) = \nu(\Omega((A^*A)^{\frac{1}{2}}))$. A further interpretation of Mirsky's conjecture is contained in Chapter 2 where the nonnegativity of the right side of (1.36) is proved (Lemma 2.4).

1.37 Theorem (Mirsky). Let $A \in \mathcal{M}$. We have

$$(1.38) \quad d_{\epsilon}^2(A) = \inf_{X \in \mathcal{N}} \epsilon^2(A - X) \leq \frac{1}{2}(\epsilon^2(A) - |\text{tr}(A^2)|) .$$

1.39 Definition. Let $A \in \mathcal{M}_n$ and let ν be any norm on \mathcal{M}_n . A matrix $N \in \mathcal{N}_n$ such that

$$(1.40) \quad \nu(A - N) = d_{\nu}(A) = \inf_{X \in \mathcal{N}_n} \nu(A - X)$$

is called a ν -minimizing normal matrix (for A) or N is said to be ν -minimal (for A).

1.3 Differentiable Curves in Matric Space.

Let $A(t) = (a_{ij}(t))$ be a matrix function of the real variable t which is defined for $-\infty \leq a < t < b \leq \infty$. In the sequel we shall assume that each of the scalar functions $a_{ij}(t)$ is sufficiently differentiable throughout its domain of definition. We define the derivative of $A(t)$ by

$$(1.41) \quad \frac{dA}{dt} = \frac{d}{dt} A(t) = \left(\frac{da_{ij}}{dt} \right) ;$$

higher order derivatives are defined in a similar fashion. The exponential function $\exp(A)$ is defined by the power series

$$(1.42) \quad e^A = I + \sum_{k=1}^{\infty} \frac{1}{k!} A^k$$

which converges for all $A \in \mathcal{M}$.

1.43 Lemma. Let $A(t), B(t)$ be any differentiable matrices in \mathcal{M} and let C be any constant matrix in \mathcal{M} . We have

$$(1.44) \quad \frac{d}{dt} [A(t) B(t)] = \frac{dA(t)}{dt} B(t) + A(t) \frac{dB(t)}{dt} ,$$

$$(1.45) \quad \frac{d}{dt} e^{tC} = e^{tC} \cdot C = C \cdot e^{tC} ,$$

$$(1.46) \quad \frac{d}{dt} \operatorname{tr}[A(t)] = \operatorname{tr}\left(\frac{dA(t)}{dt}\right) ,$$

$$(1.47) \quad \frac{d}{dt} \operatorname{dg}(A(t)) = \operatorname{dg}\left(\frac{dA(t)}{dt}\right) ,$$

$$(1.48) \quad \frac{d}{dt} A^*(t) = \left[\frac{dA(t)}{dt} \right]^*$$

where $A^*(t) = [A(t)]^*$ and where dg is defined by (1.24).

The proof of Lemma 1.43 will be omitted since each equation is either very elementary or well known.

1.49 Definition. Let \mathcal{A} be any subset of \mathcal{M} . A matrix function $A(t)$ whose range is in \mathcal{A} and which is differentiable in some interval $a < t < b$ ($a < b$) is called a differentiable curve in \mathcal{A} .

For any differentiable curve $A(t)$, the tangent vector to $A(t)$ at any $t_0 \in (a, b)$ is defined to be the matrix $[dA(t)/dt]_{t=t_0}$.

Let $U(t)$ be an arbitrary differentiable curve in the subset \mathcal{U} of all unitary matrices in \mathcal{M} . Then $U^*(t) U(t) = I$ and $U(t) U^*(t) = I$ are identities in t . Differentiating these identities and using the rule (1.44) we obtain

$$(1.50) \quad \frac{dU^*(t)}{dt} U(t) + U^*(t) \frac{dU(t)}{dt} = 0 ,$$

$$(1.51) \quad \frac{dU(t)}{dt} U^*(t) + U(t) \frac{dU^*(t)}{dt} = 0 .$$

From either of the last two equations we obtain

$$(1.52) \quad \frac{dU(t)}{dt} = - U(t) \frac{dU^*(t)}{dt} U(t) .$$

1.53 Lemma. Let $U(t)$ be any differentiable curve in \mathcal{U} and let $S_1(t), S_2(t)$ be given by the equations

$$(1.54) \quad S_1(t) = - \frac{dU^*(t)}{dt} U(t) \quad ,$$

$$(1.55) \quad S_2(t) = - U(t) \frac{dU^*(t)}{dt} \quad .$$

Then $S_1(t)$ and $S_2(t)$ are skew-hermitian for every $t \in (a,b)$ and the tangent vector to $U(t)$ at t_0 is given by

$$(1.56) \quad \left. \frac{dU(t)}{dt} \right|_{t=t_0} = U(t_0) S_1(t_0) = S_2(t_0) U(t_0) \quad .$$

Proof. Using (1.48) and the definitions of $S_1(t)$ and $S_2(t)$ we see from (1.50) and (1.51) that $S_1(t) + S_1^*(t) = 0$, $S_2^*(t) + S_2(t) = 0$; i.e., S_1 and S_2 are skew-hermitian for each value of t . The expression (1.56) follows immediately from (1.52).

Let S be any skew-hermitian matrix. It is easy to see that $\exp(tS)$ is unitary for all finite values of the real variable t . In fact, from the definition of the exponential function, $[\exp(tS)]^* = \exp(-tS)$ and, since tS commutes with $-tS$, we have

$$(e^{tS})^* e^{tS} = e^{-tS} e^{tS} = e^{tS-tS} = e^0 = I \quad .$$

Thus, letting $U(t) = \exp(tS)$, we see from (1.45) that

$$\frac{dU(t)}{dt} = U(t) S = S U(t) \quad ;$$

consequently any skew-hermitian matrix S can occur in place of $S_1(t_0)$ and $S_2(t_0)$ in (1.56) (for some differentiable curve $U(t)$ in \mathcal{U}) and for any value of t_0 .

1.4 Normal Matrices.

In the next theorem we list several known characterizations of normal matrices, already defined in Section 1.1.

1.57 Theorem. Let the eigenvalues of $N \in \mathcal{M}_n$ be denoted by $\lambda_1, \lambda_2, \dots, \lambda_n$. Then N is normal if and only if any one of the following propositions is true:

- (a) $N = H_1 + iH_2$ where H_1 and H_2 are hermitian and $H_1H_2 = H_2H_1$.
- (b) N has a complete orthonormal set of eigenvectors.
- (c) (Toeplitz [23]) N is unitarily similar to a diagonal matrix:

$$N = U^*DU \quad (U \in \mathcal{U}, D \in \mathcal{D}) .$$

- (d) (Wintner and Murnaghan [26] and Williamson [25], see also Halmos [11], pp. 169-170) There is a positive semidefinite hermitian matrix H and a unitary matrix U such that

$$(1.58) \quad N = UH = HU .$$

- (e) (Parker [20], p. 522) There exists a unitary matrix U such that $U(N + N^*)U^* \in \mathcal{D}$ and $U(N - N^*)U^* \in \mathcal{D}$.

(f) (Parker [20], Theorem 1) The eigenvalues of NN^* are

$$|\lambda_1|^2, |\lambda_2|^2, \dots, |\lambda_n|^2.$$

(g) (Parker [20], Theorem 2) The eigenvalues of $N + N^*$ are

$$\lambda_1 + \bar{\lambda}_1, \lambda_2 + \bar{\lambda}_2, \dots, \lambda_n + \bar{\lambda}_n.$$

1.59 Theorem (Toeplitz [23] and Parker [20]). A triangular matrix in \mathcal{M} is normal if and only if it is diagonal.

We shall be interested later in utilizing differentiable curves in \mathcal{N} . One way of constructing such curves is to use Theorem 1.57 (c) and differentiable curves $U(t)$ and $D(t)$ in \mathcal{U} and \mathcal{D} respectively. Then $N(t) = U^*(t) D(t) U(t)$ is a differentiable curve in \mathcal{N} . Furthermore we can construct differentiable curves in \mathcal{N} which pass through a given normal matrix $N_0 = U_0^* D_0 U_0$ ($U_0 \in \mathcal{U}$, $D_0 \in \mathcal{D}$) for some value of t (say $t = 0$) by merely requiring that $U(0) = U_0$ and $D(0) = D_0$. For curves $U(t)$ in \mathcal{U} we shall use the formula

$$(1.60) \quad U(t) = U_0 e^{itH}$$

where H is hermitian. There is no loss of generality in restricting ourselves to the formula (1.60) since we shall be concerned with evaluations of the derivative of $U(t)$ at $t = 0$. Note that, as H runs through \mathcal{H} , iH runs through the set of all skew-hermitian matrices, so, by Lemma 1.53, all possible tangent vectors to a differentiable curve $U(t)$ in \mathcal{U} at $t = 0$ can occur for curves of the type given by (1.60)

We shall have a need later on for the following result concerning differentiable curves in \mathcal{D} .

1.61 Lemma. If $D(t)$ is any differentiable curve in \mathcal{A} , then $\left. \frac{dD(t)}{dt} \right|_{t=0}$ is in \mathcal{A} . Furthermore every $A \in \mathcal{A}$ can occur as $\left. \frac{dD(t)}{dt} \right|_{t=0}$ for some differentiable curve $D(t)$ in \mathcal{A} .

Proof. Let $D(t) = \text{diag}(\alpha_1(t), \dots, \alpha_n(t))$. Each of the scalar functions $\alpha_i(t)$ has a scalar derivative so $dD(t)/dt$ is in \mathcal{A} for every value of t . Let z_i ($i = 1, \dots, n$) be any complex numbers. Then the derivative of $D(t) = t \text{diag}(z_1, z_2, \dots, z_n)$ equals $\text{diag}(z_1, \dots, z_n)$, an arbitrary matrix in \mathcal{A} , for all t .

1.62 Theorem. Let $X, Y \in \mathcal{N}$. Then $N = X + Y \in \mathcal{N}$ if and only if

$$(1.63) \quad XY^* - Y^*X + YX^* - X^*Y = 0.$$

Proof. Since X and Y are normal and have

$$\begin{aligned} NN^* - N^*N &= (X + Y)(X^* + Y^*) - (X^* + Y^*)(X + Y) \\ &= XX^* + XY^* + YX^* + YY^* - (X^*X + X^*Y + Y^*X + Y^*Y) \\ &= XY^* - Y^*X + YX^* - X^*Y \end{aligned}$$

whence N is normal if and only if (1.63) holds.

1.64 Corollary. Let z be any fixed complex number. Then $N \in \mathcal{N}$ if and only if $N + zI \in \mathcal{N}$.

Proof. Clearly $zI \in \mathcal{N}$ for all $z \in \mathbb{C}$. Setting $X = N$, $Y = zI$ we find that (1.63) is satisfied for all z :

$$XY^* - Y^*X + YX^* - X^*Y = \bar{z}N - \bar{z}N + zN^* - zN^* = 0.$$

Thus, by Theorem 1.62, $N \in \mathcal{N}$ implies $N + zI \in \mathcal{N}$. The converse implication is clear since $(N + zI) - zI = N$.

1.65 Theorem. Let $X \in \mathcal{N}$. Then $X + tY \in \mathcal{N}$ for all values of t in a real interval of positive length if and only if $Y \in \mathcal{N}$ and (1.63) holds.

Proof. Let $N = X + tY$. A short computation, using the fact that t is real and $X \in \mathcal{N}$, yields

$$(1.66) \quad NN^* - N^*N = t^2(YY^* - Y^*Y) + t(XY^* - Y^*X + YX^* - X^*Y) \quad .$$

A polynomial of second degree can have at most two zeros unless all of its coefficients vanish. Thus the assumption that $N \in \mathcal{N}$ for more than two values of t implies immediately (1.63) and $YY^* - Y^*Y = 0$ i.e., $Y \in \mathcal{N}$. Conversely $Y \in \mathcal{N}$ and (1.63) imply via (1.66) that $X + tY \in \mathcal{N}$ for all real values of t .

1.67 Theorem. Let $P_k(X) = z_0I + \sum_{i=1}^k z_i X^i$ denote a polynomial of degree k (≥ 1) in a matrix $X \in \mathcal{N}$ with arbitrary complex coefficients z_i . Then $N \in \mathcal{N}$ implies $P_k(N) \in \mathcal{N}$.

Proof. If $N \in \mathcal{N}$ then by Theorem 1.57 (c) it has a decomposition $N = U^*DU$ ($U \in \mathcal{U}$, $D \in \mathcal{D}$) and clearly $N^i = U^*D^iU$ for all positive integers i . Thus $P_k(N) = U^*P_k(D)U$, $P_k(D) \in \mathcal{D}$ so by Theorem 1.57 (c) $P_k(N)$ is normal.

1.68 Corollary. Suppose $N \in \mathcal{N}$. Then $N^k \in \mathcal{N}$ for $k = 2, 3, 4, \dots$. Furthermore if N^{-1} exists then $N^{-k} \in \mathcal{N}$ for $k = 1, 2, 3, \dots$.

Proof. That $N^k \in \mathcal{N}$ for $k \geq 2$ is obvious from Theorem 1.67. If N^{-1} exists then, by the Cayley-Hamilton Theorem, N^{-1} is a polynomial in N whence N^{-1} is normal. Applying Theorem 1.67 again we see that

$N^{-k} = (N^{-1})^k$ is normal for $k \geq 2$.

1.69 Theorem. Let $A \in \mathcal{M}$ and let α and β be complex numbers. If A is normal then $\alpha A + \beta A^* \in \mathcal{N}$ for all α, β ; if A is not normal, $\alpha A + \beta A^* \in \mathcal{N}$ if and only if $|\alpha| = |\beta|$

Proof. Setting $N = \alpha A + \beta A^*$ we find that

$$(1.70) \quad NN^* - N^*N = (|\alpha|^2 - |\beta|^2)(AA^* - A^*A) .$$

The conclusions of the theorem follow immediately from the relation (1.70).

1.71 Corollary. If $A \in \mathcal{M}$ and $\alpha, \beta \in \underline{\mathbb{C}}$ with $|\alpha| = |\beta|$, then $P_k(\alpha A + \beta A^*) \in \mathcal{N}$ where P_k is an arbitrary polynomial of degree k with coefficients in $\underline{\mathbb{C}}$.

Proof. This is an obvious consequence of Theorems 1.67 and 1.69.

1.72 Theorem. Suppose $\alpha \in \underline{\mathbb{C}}$ ($\alpha \neq 0$), $z \in \underline{\mathbb{C}}$, and $U \in \mathcal{U}$ are fixed. Then each of the transformations

$$(1.73) \quad T_\alpha(N) = \alpha N ,$$

$$(1.74) \quad T_z(N) = N + zI ,$$

$$(1.75) \quad T_U(N) = U^*NU$$

defines a one-to-one mapping of \mathcal{N} onto itself.

Proof. Let N be any normal matrix. Then from Theorem 1.67, Corollary 1.64, and Theorem 1.57 (c) we see that $\alpha^{-1}N$, $N - zI$, and UNU^* respectively are in \mathcal{N} . Thus $N = T_\alpha(\alpha^{-1}N) = T_z(N - zI) = T_U(UNU^*)$ which proves that each of the transformations (1.73), (1.74), and (1.75)

is onto. Letting N_1, N_2 denote any pair of normal matrices, one sees easily that any one of the equations $T_\alpha(N_1) = T_\alpha(N_2)$, $T_z(N_1) = T_z(N_2)$, $T_U(N_1) = T_U(N_2)$ implies $N_1 = N_2$, whence each of the transformations is one-to-one.

1.5 Existence of ν -Minimal Matrices.

The definition of a ν -minimal (or ν -minimizing normal) matrix has already been given in Section 1.2 (Definition 1.39). Let ν denote any norm on \mathcal{M} . Suppose first that $A \in \mathcal{N}$. Then there is a unique ν -minimal N_0 , namely $N_0 = A$; for if $A \in \mathcal{N}$ then $d_\nu(A) = 0$ and the infimum in (1.34) is assumed if and only if $X = A$. Our main purpose in the present section is to show that, for any $A \in \mathcal{M}$ and any norm ν , there exists a ν -minimal matrix.

Note. Any two normed linear spaces (over \mathbb{C}) of the same finite dimension are topologically isomorphic (see e.g., [13], p. 13). This implies that the norm topologies induced in \mathcal{M}_n by any two norms are the same; consequently there is only one norm topology for \mathcal{M}_n and we refer to it as the norm topology of \mathcal{M}_n .

1.76 Lemma. The set \mathcal{N}_n of all normal matrices (in \mathcal{M}_n) is closed in the norm topology of \mathcal{M}_n .

Proof. Matrix multiplication is continuous in the norm topology since it is continuous with respect to the ϵ -norm topology (cf. Section 1.1). Let N be any matrix in the closure $\overline{\mathcal{N}}$ of \mathcal{N} . Then there is a sequence $\{N_i\}$ ($N_i \in \mathcal{N}$ for $i = 1, 2, 3, \dots$) such that $N_i \rightarrow N$. Since each N_i is normal we have

$$(1.77) \quad N_1 N_1^* = N_1^* N_1 \quad (1 = 1, 2, \dots) .$$

By virtue of the continuity of multiplication we may pass to the limit in (1.77) and obtain $NN^* = N^*N$, i.e., N is normal. This implies $\overline{\mathcal{N}} \subset \mathcal{N}$ which proves that \mathcal{N} is closed.

1.78 Theorem. Let v be any norm on \mathcal{M}_n and let A be any fixed matrix in \mathcal{M}_n . Then there is a v -minimal matrix N_0 (for A).

Proof. We assume $d_v(A) > 0$ (see (1.34)) since otherwise the theorem is trivial. By the definition of $d_v(A)$, there is a sequence $\{N_1\}$ of normal matrices such that $v(A - N_1) \rightarrow d_v(A)$. The subset $\mathcal{Q} = \{X ; X \in \mathcal{M}_n \text{ and } d_v(A) \leq v(A - X) \leq 2d_v(A)\}$ is closed and bounded, hence compact, in \mathcal{M}_n . Clearly there is an index k_0 so that $N_1 \in \mathcal{Q}$ for $1 \geq k_0$. Thus, since \mathcal{Q} is compact (and therefore countably compact), there is a subsequence $\{N_{1_k}\}$ ($1_k \geq k_0$ for $k = 1, 2, 3, \dots$) which converges to a matrix $N_0 \in \mathcal{Q}$. This N_0 is a point of closure of \mathcal{N} so, by Lemma 1.76, it is normal. Finally

$$(1.79) \quad d_v(A) \leq v(A - N_0) \leq v(A - N_{1_k}) + v(N_{1_k} - N_0) ;$$

and, since $v(A - N_{1_k}) \rightarrow d_v(A)$, we have

$$v(A - N_{1_k}) = d_v(A) + \delta_k \geq d_v(A)$$

where $\delta_k \rightarrow 0$ as $k \rightarrow \infty$. Thus, given $\epsilon > 0$, there is an index k_1 such that the right side of (1.79) is less than $d_v(A) + \epsilon$ for $k \geq k_1$. Since

ϵ is arbitrary we have $d_v(A) \leq v(A - N_0) \leq d_v(A)$ which completes the proof of Theorem 1.78.

Theorem 1.78 shows that the distance problem of finding $d_v(A)$ is actually a minimum problem for all v and for all A . This suggests the possibility of finding $d_v(A)$ by determining a v -minimal matrix. We shall investigate this aspect of the distance problem in subsequent sections of this paper.

Since \mathcal{N}_n is not a convex set (the sum of two normal matrices is not necessarily normal), we naturally expect that there might exist matrices $A \in \mathcal{N}_n$ for which there is no unique v -minimal matrix. We shall show, at least for $v = \epsilon$ and $n = 2$, that this is the case.

CHAPTER 2

EQUIVALENCE THEOREMS CONCERNING MIRSKY'S CONJECTURE, ν -MINIMAL MATRICES, AND DISTANCE FORMULAS.

Now that the existence of ν -minimal matrices has been established, it is of interest to determine what transformations of \mathcal{M}_n leave the property of minimality invariant. In this section we shall give three results of this type and we shall also prove three closely parallel results concerning distance formulas resembling Mirsky's formula (1.36). We begin with an examination of the meaning of Mirsky's conjecture.

2.1 Mirsky's Conjecture.

2.1 Definition. A norm ν on \mathcal{M}_n such that

$$(2.2) \quad \nu(A) = \nu(\Omega(A)) \quad \text{for all } A \in \mathcal{M}_n$$

and

$$(2.3) \quad \nu(A) > \nu(\Omega(A)) \quad \text{for all } A \in \mathcal{M}_n, A \notin \mathcal{N}_n,$$

where $\Omega(A)$ is defined in Definition 1.17, is said to have property S.

2.4 Lemma. The Euclidean norm ϵ has property S. Furthermore, for any unitarily invariant norm ν , we have

$$(2.5) \quad \nu(\Omega((A^*A)^{\frac{1}{2}})) = \nu(A) \geq \nu(\Omega(A)) \quad \text{for all } A \in \mathcal{M}_n$$

with equality holding in (2.5) for all $A \in \mathcal{N}_n$.

Proof. By a well-known theorem of Schur [22] (see e.g., [16], p. 307) every $A \in \mathcal{M}_n$ is unitarily similar to a triangular matrix:

$$(2.6) \quad VAV^* = \Omega(A) + M \quad (V \in \mathcal{U}_n)$$

where in M only elements above the principal diagonal may be different from zero. Thus, by (1.16), we have

$$\epsilon^2(A) = \epsilon^2(VAV^*) = \epsilon^2(\Omega(A)) + \epsilon^2(M)$$

so

$$(2.7) \quad \epsilon^2(A) - \epsilon^2(\Omega(A)) = \epsilon^2(M)$$

By Theorem 1.59 $M = 0$ if and only if $A \in \mathcal{N}_n$; hence from (2.7) we see that ϵ has property S.

In order to prove (2.5) we shall need the following two results.

2.8 Theorem (Fan [8], Theorem 4). Let $a_1 \geq a_2 \geq \dots \geq a_n \geq 0$, $b_1 \geq b_2 \geq \dots \geq b_n \geq 0$. Then

$$(2.9) \quad \varphi(a_1, a_2, \dots, a_n) \leq \varphi(b_1, b_2, \dots, b_n)$$

holds for all symmetric gauge functions φ of n real variables if and only if

$$(2.10) \quad \sum_{i=1}^k a_i \leq \sum_{i=1}^k b_i \quad (1 \leq k \leq n) \quad .$$

2.11 Theorem (Weyl [24], p. 409). Let the eigenvalues λ_1 and singular values α_1 of $A \in \mathcal{M}_n$ be arranged so that $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$, $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$. Then, for any real exponent $s > 0$,

$$(2.12) \quad \sum_{i=1}^k |\lambda_i|^s \leq \sum_{i=1}^k \alpha_i^s \quad (1 \leq k \leq n) .$$

Since $\Omega(A)$ is diagonal (and therefore normal) its singular values are $|\lambda_1|, \dots, |\lambda_n|$ (see Theorem 1.57 (f)). The inequality (2.5) now follows immediately by combining Theorem 1.11, Theorem 2.8 (here we set $a_1 = |\lambda_1|, b_1 = \alpha_1$), and Theorem 2.11 (use $s = 1$). If A is normal it is unitarily similar to $\Omega(A)$ whence equality holds in (2.5) for all $A \in \mathcal{N}_n$.

The question of which unitarily invariant norms $\nu \neq \epsilon$ have property S (i.e., which ones satisfy (2.3)) is apparently open. We shall presently show that the spectral norm σ does not have property S.

2.13 Definition. Let ν be any norm on \mathcal{M}_n which has property S. Then we define

$$(2.14) \quad p_{\nu,n}(A) = \frac{d_{\nu}^2(A)}{\nu^2(A) - \nu^2(\Omega(A))} \quad \text{for } A \notin \mathcal{N}_n .$$

For any ν with property S we see from the definition of $d_{\nu}(A)$ and (2.3) that $p_{\nu,n}(A)$ is a well-defined positive quantity so that $d_{\nu}(A)$ can be expressed in the form

$$(2.15) \quad d_{\nu}^2(A) = p_{\nu,n}(A)(\nu^2(A) - \nu^2(\Omega(A))) \quad \text{for all } A \in \mathcal{M}_n, A \notin \mathcal{N}_n .$$

If we assigned to $p_{v,n}(A)$ some convenient finite value for $A \in \mathcal{N}_n$, then (2.15) would be valid for all $A \in \mathcal{M}_n$. Note that (2.15) has the same general form as Mirsky's formula (1.36). We now prove the following characterization of Mirsky's conjecture for a particular norm v .

2.16 Theorem. Let v be a norm on \mathcal{M}_n . Mirsky's conjectured formula (1.36) for $d_v(A)$ holds for all $A \in \mathcal{M}_n$ if and only if v has property S and

$$(2.17) \quad p_{v,n}(A) = \frac{1}{2}$$

for all $A \in \mathcal{M}_n$, $A \notin \mathcal{N}_n$ and for $n = 2, 3, \dots$.

Proof. Since \mathcal{N}_n is closed as a subset of \mathcal{M}_n , we see from (1.34) that $d_v(A)$ is zero when $A \in \mathcal{N}_n$ and strictly positive when $A \notin \mathcal{N}_n$. Thus, if (1.36) holds for all $A \in \mathcal{M}_n$, v has property S and (2.17) holds. The converse is obvious.

Let the n by n matrix A ($n \geq 3$) be given by

$$(2.18) \quad A = \begin{pmatrix} B & 0 \\ 0 & 0 \end{pmatrix} \quad \text{where } B = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}.$$

The nonzero singular values of A are $2, \sqrt{2}$ and the nonzero singular values of $\Omega(A)$ are $2, 1$. From (1.14) and (1.15) we have $\sigma(A) = \sigma(\Omega(A)) = 2$ while, by Theorem 1.59, A is not normal. Thus (2.3) does not hold for $v = \sigma$; consequently σ fails to have property S for $n \geq 3$. We have proved

2.19 Theorem. The spectral norm σ does not have property S for $n \geq 3$. Mirsky's Conjecture 1.35 is false for $v = \sigma$ and for $n \geq 3$.

In Chapter 8 where we discuss some other counterexamples we shall prove that Mirsky's conjecture also fails for $v = \sigma$ and $n = 2$. (See Theorem 8.9).

2.2 Invariance of v -Minimal Matrices.

The content of the next three theorems is, roughly speaking, the following: v -minimality is generally invariant under the transformations (1.73) - (1.75). The accompanying corollaries show that uniqueness of a v -minimal matrix is also invariant under the same transformations.

2.20 Theorem. Let v be any norm on \mathcal{M}_n and suppose $\alpha \in \underline{\mathbb{C}}$, $\alpha \neq 0$. Then N_0 is v -minimal for $A \in \mathcal{M}_n$ if and only if αN_0 is v -minimal for αA . Furthermore

$$(2.21) \quad d_v(\alpha A) = |\alpha| d_v(A) \quad .$$

Proof. By Theorem 1.72 αN runs through \mathcal{N} as N runs through \mathcal{N} . Consequently $v(A - N)$ assumes its minimum for $N = N_0$ if and only if $|\alpha| \cdot v(A - N) = v(\alpha A - \alpha N)$ assumes its minimum for $N = N_0$. If αN_0 is v -minimal for αA then

$$d_v(\alpha A) = v(\alpha A - \alpha N_0) = |\alpha| v(A - N_0) = |\alpha| d_v(A) \quad .$$

2.22 Corollary. Let v be any norm on \mathcal{M}_n and suppose $\alpha \in \underline{\mathbb{C}}$, $\alpha \neq 0$. $A \in \mathcal{M}_n$ has a unique v -minimal matrix if and only if αA has a unique v -minimal matrix.

Proof. If N_0 is the unique v -minimal matrix for A and if N_1, N_2 are v -minimal for αA , then by Theorem 2.20 $N_0 = \alpha^{-1}N_1 = \alpha^{-1}N_2$ whence $N_1 = N_2$. The converse is proved in a similar manner.

2.23 Theorem. Let v denote any norm on \mathcal{M} . If N_0 is v -minimal for $A \in \mathcal{M}$ then $N_0 + zI$ is v -minimal for $A + zI$ for all $z \in \underline{\mathbb{C}}$. Conversely, if $N_0 + zI$ is v -minimal for $A + zI$ for one value of $z \in \underline{\mathbb{C}}$ then N_0 is v -minimal for A . Furthermore, for all $z \in \underline{\mathbb{C}}$, we have

$$(2.24) \quad d_v(A + zI) = d_v(A) \quad .$$

Proof. Obviously

$$(2.25) \quad v(A - N) = v[(A + zI) - (N + zI)]$$

holds for any norm v , for all $A, N \in \mathcal{M}$, and for all $z \in \underline{\mathbb{C}}$. By Theorem 1.72 $T_z(N) = N + zI$ runs through \mathcal{M} in a one-to-one manner as N runs through \mathcal{M} . Thus, as N runs through \mathcal{M} , the left and right sides of (2.25) assume their minima simultaneously. This proves the first two statements of Theorem 2.23. The relation (2.24) follows immediately from (2.25) if one assumes that N is v -minimal for A .

2.26 Corollary. Let v be any norm on \mathcal{M} and suppose $z \in \underline{\mathbb{C}}$. Then $A \in \mathcal{M}$ has a unique v -minimal matrix if and only if $A + zI$ has a unique v -minimal matrix.

Proof. The proof is strictly analogous to that for Corollary 2.22.

2.27 Theorem. Let v denote any unitarily invariant norm on \mathcal{M} and let $U \in \mathcal{U}$ be fixed. Then N_0 is v -minimal for $A \in \mathcal{M}$ if and only if

$U^* N_0 U$ is ν -minimal for $U^* A U$. Furthermore

$$(2.28) \quad d_\nu(U^* A U) = d_\nu(A) \quad .$$

Proof. Since ν is unitarily invariant

$$(2.29) \quad \nu(A - N) = \nu(U^* A U - U^* N U)$$

holds for all $A, N \in \mathcal{M}$ and for all $U \in \mathcal{U}$. By Theorem 1.72 $T_U(N) = U^* N U$ runs through \mathcal{N} in a one-to-one manner as N runs through \mathcal{N} . Therefore, as N runs through \mathcal{N} , the left and right sides of (2.29) assume their minima simultaneously. This proves the first assertion in Theorem 2.27. Equation (2.28) follows from (2.29) if it is assumed that N is ν -minimal for A .

2.30 Corollary. Let ν be any unitarily invariant norm on \mathcal{M} and let $U \in \mathcal{U}$. Then $A \in \mathcal{M}$ has a unique ν -minimal matrix if and only if $U^* A U$ has a unique ν -minimal matrix.

The proof is analogous to that for Corollary 2.22 and is therefore omitted.

2.3 Invariance of Distance Formulas.

2.31 Definition. A transformation T whose domain space is \mathcal{M}_n and whose range space is contained in \mathcal{M}_n is said to be discriminating if

$$T(M) \notin \mathcal{N}_n \quad \text{for } M \notin \mathcal{N}_n$$

and

$$T(M) \in \mathcal{M}_n \quad \text{for } M \in \mathcal{M}_n .$$

We have already seen that, for any norm ν with property S, $d_\nu(A)$ can be represented by the formula

$$(2.32) \quad d_\nu^2(A) = p_{\nu,n}(A)(\nu^2(A) - \nu^2(\Omega(A))), \quad (A \in \mathcal{M}_n, A \notin \mathcal{N}_n)$$

where $p_{\nu,n}(A)$ is defined by (2.14). The question naturally arises as to what happens to $p_{\nu,n}(A)$ when A is subjected to various discriminating transformations. The next three theorems provide some answers to this question for the discriminating transformations (1.73) - (1.75).

2.33 Theorem. Let ν be any norm on \mathcal{M}_n which has property S. Let $\alpha \in \mathbb{C}$, $\alpha \neq 0$ and let $A \in \mathcal{M}_n$, $A \notin \mathcal{N}_n$. Then

$$(2.34) \quad p_{\nu,n}(\alpha A) = p_{\nu,n}(A) .$$

Proof. If $\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of A then $\alpha\lambda_1, \alpha\lambda_2, \dots, \alpha\lambda_n$ are the eigenvalues of αA ; consequently (cf. Definition 1.17) $\Omega(\alpha A) = \alpha\Omega(A)$. Using that fact, (2.32), and (2.21) we obtain

$$\begin{aligned} d_\nu^2(\alpha A) &= |\alpha|^2 d_\nu^2(A) \\ &= p_{\nu,n}(A)(\nu^2(\alpha A) - \nu^2(\alpha\Omega(A))) \\ &= p_{\nu,n}(A)(\nu^2(\alpha A) - \nu^2(\Omega(\alpha A))) . \end{aligned}$$

The first factor in the last line must be $p_{v,n}(\alpha A)$ which proves (2.34).

2.35 Theorem. Let ϵ denote the Euclidean norm (1.16) on \mathcal{M}_n . Let A be any nonnormal matrix in \mathcal{M}_n and let $z \in \mathbb{C}$. Then

$$(2.36) \quad p_{\epsilon,n}(A + zI) = p_{\epsilon,n}(A) \quad .$$

Proof. Let λ_i ($i = 1, 2, \dots, n$) denote the eigenvalues of $A = (a_{ij})$. From (1.16) we have

$$\begin{aligned} \epsilon^2(A + zI) &= \sum_{i \neq j} |a_{ij}|^2 + \sum_{i=1}^n |a_{ii} + z|^2 \\ &= \epsilon^2(A) + \bar{z} \sum_{i=1}^n a_{ii} + z \sum_{i=1}^n \bar{a}_{ii} + n|z|^2 \quad ; \end{aligned}$$

and, since $\sum_{i=1}^n \lambda_i = \sum_{i=1}^n a_{ii}$, we obtain

$$(2.37) \quad \epsilon^2(A + zI) = \epsilon^2(A) + \bar{z} \sum_{i=1}^n \lambda_i + z \sum_{i=1}^n \bar{\lambda}_i + n|z|^2 \quad .$$

Furthermore, since $(\lambda_i + z)$ ($i = 1, \dots, n$) are the eigenvalues of $A + zI$,

$$\begin{aligned} (2.38) \quad \epsilon^2(\Omega(A + zI)) &= \sum_{i=1}^n (\lambda_i + z)(\bar{\lambda}_i + \bar{z}) \\ &= \epsilon^2(\Omega(A)) + \bar{z} \sum_{i=1}^n \lambda_i + z \sum_{i=1}^n \bar{\lambda}_i + n|z|^2 \quad . \end{aligned}$$

Combining (2.37) and (2.38) we find that

$$(2.39) \quad \epsilon^2(A) - \epsilon^2(\Omega(A)) = \epsilon^2(A + zI) - \epsilon^2(\Omega(A + zI)) \quad .$$

Applying (2.24) in the case $\nu = \epsilon$ and using (2.39) we obtain

$$\begin{aligned} d_{\epsilon}^2(A + zI) &= d_{\epsilon}^2(A) = p_{\epsilon, n}(A)(\epsilon^2(A) - \epsilon^2(\Omega(A))) \\ &= p_{\epsilon, n}(A)(\epsilon^2(A + zI) - \epsilon^2(\Omega(A + zI))) \quad . \end{aligned}$$

The relation (2.36) now follows from the last equation.

2.40 Theorem. Let ν be any unitarily invariant norm on \mathcal{M}_n which has property S. Let A be any nonnormal matrix in \mathcal{M}_n and let $U \in \mathcal{U}_n$. Then

$$(2.41) \quad p_{\nu, n}(U^*AU) = p_{\nu, n}(A) \quad .$$

Proof. The eigenvalues of a matrix are invariant under a unitary similarity transformation so $\Omega(U^*AU) = \Omega(A)$. Consequently, since ν is unitarily invariant, we have

$$(2.42) \quad \nu^2(A) - \nu^2(\Omega(A)) = \nu^2(U^*AU) - \nu^2(\Omega(U^*AU)) \quad .$$

From (2.28) we have

$$d_{\nu}^2(U^*AU) = d_{\nu}^2(A) = p_{\nu, n}(A)(\nu^2(A) - \nu^2(\Omega(A)))$$

and, combining this with (2.42) we obtain (2.41).

CHAPTER 3

IMPROVEMENTS OF MIRSKY'S BOUND

3.1 A New Bound.

In this section we shall obtain an upper bound for $d_e^2(A)$ which is sharper than the bound (1.38) obtained by Mirsky. Before doing this we prove a lemma which sheds some light on Mirsky's result (Theorem 1.37). This lemma furnishes at least a partial answer to the question: given $A \in \mathcal{M}_n$, what normal matrices lie at the distance from A which is given by Mirsky's bound?

3.1 Definition. For any $A \in \mathcal{M}_n$ we define

$$(3.2) \quad \eta(A) = \begin{cases} \{\operatorname{tr}(A^2)/|\operatorname{tr}(A^2)|\} & \text{if } \operatorname{tr}(A^2) \neq 0 \\ \{z ; z \in \mathbb{C} \text{ and } |z| = 1\} & \text{if } \operatorname{tr}(A^2) = 0 ; \end{cases}$$

$$(3.3) \quad M_0(A) = \left\{ \frac{1}{2} (A + \eta A^*) ; \eta \in \eta(A) \right\} .$$

Note. $\eta(A)$ is a set of complex numbers of unit modulus. If $\operatorname{tr}(A^2) \neq 0$ then $\eta(A)$ contains a single well-defined number but, if $\operatorname{tr}(A^2) = 0$, $\eta(A)$ consists of all complex numbers on the boundary of the unit disk. An analogous remark applies to the set $M_0(A)$.

3.4 Lemma. Let A be any matrix in \mathcal{M}_n . Then every matrix in the set $M_0(A)$ is normal. Furthermore

$$(3.5) \quad \epsilon^2(A - X) = \frac{1}{2} (\epsilon^2(A) - |\text{tr}(A^2)|) \text{ for all } X \in M_0(A) .$$

Proof. The fact that $M_0(A) \subset \mathcal{M}_n$ follows immediately from Theorem 1.69. Assuming that $|\eta| = 1$ and using the properties (1.20) and (1.22) of the trace, we obtain

$$(3.6) \quad \begin{aligned} \epsilon^2(A - \frac{1}{2}(A + \eta A^*)) &= \frac{1}{4} \text{tr}[(A^* - \bar{\eta}A)(A - \eta A^*)] \\ &= \frac{1}{4} [\text{tr}(A^*A) - \eta \text{tr}(A^{*2}) - \bar{\eta} \text{tr}(A^2) + \text{tr}(AA^*)] . \end{aligned}$$

Using (1.22) again and (1.21) we see that the last equation becomes

$$(3.7) \quad \epsilon^2(A - \frac{1}{2}(A + \eta A^*)) = \frac{1}{2} \epsilon^2(A) - \frac{1}{4} (\eta \overline{\text{tr}(A^2)} + \bar{\eta} \text{tr}(A^2)) .$$

If $\text{tr}(A^2) \neq 0$ and $\eta \in \eta(A)$ then (3.2) implies $\eta \overline{\text{tr}(A^2)} = \bar{\eta} \text{tr}(A^2) = |\text{tr}(A^2)|$, whence (3.7) implies (3.5). If $\text{tr}(A^2) = 0$ then the last term on the right side of (3.7) vanishes for all $\eta \in \eta(A)$ and the right sides of both (3.5) and (3.7) reduce to $(1/2) \epsilon^2(A)$. This proves Lemma 3.4.

Remark. The author believes that the set $M_0(A)$ contains all matrices of the form $\alpha A + \beta A^*$ (with $|\alpha| = |\beta|$) which satisfy

$$\epsilon^2[A - (\alpha A + \beta A^*)] = \frac{1}{2} (\epsilon^2(A) - |\text{tr}(A^2)|) .$$

However, no attempt will be made here to prove that assertion.

3.8 Definition. For any $A \in \mathcal{M}_n$ we define

$$(3.9) \quad \zeta(A) = \begin{cases} ((\operatorname{tr}(A^2) - (1/n)(\operatorname{tr}(A))^2) / |\operatorname{tr}(A^2) - (1/n)(\operatorname{tr}(A))^2|) & \text{if } \operatorname{tr}(A^2) \neq \frac{1}{n} (\operatorname{tr}(A))^2 \\ [z ; z \in \mathbb{C} \text{ and } |z| = 1] & \text{if } \operatorname{tr}(A^2) = \frac{1}{n} (\operatorname{tr}(A))^2 ; \end{cases}$$

$$(3.10) \quad M(A) = \left\{ \frac{1}{2} (A + \zeta A^*) + \frac{1}{2n} \operatorname{tr}(A - \zeta A^*) I ; \zeta \in \zeta(A) \right\} .$$

Note. The situation with $\zeta(A)$ and $M(A)$ is exactly the same as that for $\eta(A)$ and $M_\eta(A)$ (see the note following Definition 3.1). If $\operatorname{tr}(A^2) \neq (1/n)(\operatorname{tr}(A))^2$ we use the symbol $N(A)$ to denote the (single) matrix in $M(A)$. In the ambiguous case $\operatorname{tr}(A^2) = (1/n)(\operatorname{tr}(A))^2$ of (3.9) we shall use the notation $N_\theta(A)$ (or $N_\zeta(A)$) to denote the particular matrix in $M(A)$ which corresponds to the element $\zeta = \exp(i\theta)$ of $\zeta(A)$ (here θ is real).

3.11 Theorem. Let $A \in \mathcal{M}_n$ and let ϵ denote the Euclidean norm (1.16) on \mathcal{M}_n . Then every matrix in the set $M(A)$ is normal. Furthermore

$$(3.12) \quad \epsilon^2(A - X) = \frac{1}{2} \left(\epsilon^2(A) - \frac{|\operatorname{tr}(A)|^2}{n} - \left| \operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n} \right| \right)$$

for all $X \in M(A)$.

Proof. The fact that $M(A) \subset \mathcal{M}_n$ follows immediately from Corollary 1.71. Let

$$B = A - \frac{\operatorname{tr}(A)}{n} I .$$

Then

$$B^2 = A^2 - \frac{2 \operatorname{tr}(A)}{n} A + \frac{(\operatorname{tr}(A))^2}{n^2} I$$

whence

$$(3.13) \quad \operatorname{tr}(B^2) = \operatorname{tr}(A^2) - \frac{1}{n} (\operatorname{tr}(A))^2 .$$

A comparison of the definitions of the sets (3.2) and (3.9) reveals that

$$(3.13) \quad \eta(B) = \zeta(A) .$$

Furthermore

$$\begin{aligned} A - \left[\frac{1}{2} (A + \zeta A^*) + \frac{1}{2n} \operatorname{tr}(A - \zeta A^*) I \right] &= \frac{1}{2} \left(A - \frac{\operatorname{tr}(A)}{n} I \right) - \frac{1}{2} \left(\zeta A^* - \frac{\operatorname{tr}(A^*)}{n} I \right) \\ &= \frac{1}{2} (B - \zeta B^*) = B - \frac{1}{2} (B + \zeta B^*) ; \end{aligned}$$

consequently, from (3.5) and (3.14), we obtain

$$(3.15) \quad \epsilon^2 \left[A - \left(\frac{1}{2} (A + \zeta A^*) + \frac{1}{2n} \operatorname{tr}(A - \zeta A^*) I \right) \right] = \frac{1}{2} (\epsilon^2(B) - |\operatorname{tr}(B^2)|)$$

for all $\zeta \in \zeta(A)$.

Now

$$\begin{aligned} \epsilon^2(B) &= \operatorname{tr} \left[\left(A^* - \frac{\operatorname{tr}(A^*)}{n} I \right) \left(A - \frac{\operatorname{tr}(A)}{n} I \right) \right] \\ (3.16) \quad &= \operatorname{tr} \left[A^* A - \frac{\operatorname{tr}(A)}{n} A^* - \frac{\operatorname{tr}(A^*)}{n} A + \frac{|\operatorname{tr}(A)|^2}{n^2} I \right] \\ &= \epsilon^2(A) - \frac{|\operatorname{tr}(A)|^2}{n} . \end{aligned}$$

Combining (3.13) - (3.16) we obtain (3.12) and this proves Theorem 3.11.

An obvious consequence of Theorem 3.11 is

3.17 Theorem. We have, for all $A \in \mathcal{M}_n$,

$$(3.18) \quad d_{\epsilon}^2(A) \leq \frac{1}{2} \left(\epsilon^2(A) - \frac{|\operatorname{tr}(A)|^2}{n} - |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n}| \right) .$$

We shall show that the bound (3.18) is sharper than Mirsky's bound

(1.38). We first prove

3.19 Lemma. If $A \in \mathcal{M}_n$ then

$$(3.20) \quad \epsilon^2(A) - \frac{|\operatorname{tr}(A)|^2}{n} - |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n}| \leq \epsilon^2(A) - |\operatorname{tr}(A^2)| .$$

If $\operatorname{tr}(A) = 0$ we have equality in (3.20) and, if $\operatorname{tr}(A) \neq 0$, equality holds in (3.20) if and only if

$$(3.21) \quad \frac{n \operatorname{tr}(A^2)}{(\operatorname{tr}(A))^2} \geq 1 .$$

Note. The inequality (3.21) is to be interpreted as follows: it is satisfied if and only if the left hand side of (3.21) is both real and greater than or equal to unity.

Proof. The inequality (3.20) and the conditions of equality in it for the case $\operatorname{tr}(A) \neq 0$ are an immediate consequence of the triangle inequality and its conditions of equality (see e.g., [1] pp. 8-9).

While Lemma 3.19 provides, in a certain sense, a complete answer to the questions of equality and inequality of the upper bounds in (1.38) and (3.18), there remain the more interesting questions of what relationship

these bounds have to the actual distance d_ϵ and to the conjectured distance (1.36) in the case $v = \epsilon$. We shall discuss the former question in connection with some counterexamples in Chapter 8. The next three results delineate a partial answer to the latter question.

3.22 Lemma. Let $\lambda_1, \lambda_2, \dots, \lambda_n$ denote complex numbers. Then

$|\sum_{k=1}^n \lambda_k^2| \leq \sum_{k=1}^n |\lambda_k|^2$ with equality holding if and only if all nonzero λ 's lie on a single straight line through the origin in the complex plane.

Proof. By the triangle inequality (see [1], p. 9) the desired inequality holds with equality holding if and only if the ratio of the squares of any pair of nonzero λ 's is positive:

$$(3.23) \quad \frac{\lambda_k^2}{\lambda_j^2} = \left(\frac{\lambda_k}{\lambda_j} \right)^2 > 0, \quad (\lambda_k \neq 0, \lambda_j \neq 0).$$

Obviously (3.23) holds if and only if λ_k/λ_j is real so that, if $\lambda_k = r_k \exp(i\theta_k)$ ($k = 1, 2, \dots, n$) where $r_k \geq 0$ and θ_k is the principal value of $\arg(\lambda_k)$, $\exp[i(\theta_k - \theta_j)] = \pm 1$ whence either $\theta_k - \theta_j \equiv 0 \pmod{2\pi}$ or $\theta_k - \theta_j \equiv \pi \pmod{2\pi}$ i.e., either $\theta_k = \theta_j$ or $\theta_k = \theta_j \pm \pi$.

3.24 Corollary. Let $A \in \mathcal{M}_n$. We have

$$(3.25) \quad |\operatorname{tr}(A^2)| \leq \epsilon^2(\Omega(A))$$

with equality holding if and only if all nonzero eigenvalues of A lie on a single straight line through the origin in the complex plane.

Proof. Let $\lambda_1, \dots, \lambda_n$ denote the eigenvalues of A . Then by Lemma 3.22

$$(3.26) \quad |\operatorname{tr}(A^2)| = \left| \sum_{k=1}^n \lambda_k^2 \right| \leq \sum_{k=1}^n |\lambda_k|^2 = \epsilon^2(\Omega(A))$$

with equality holding as stated in Corollary 3.24.

3.27 Lemma. Let $A \in \mathcal{M}_n$ be given. If, for this particular matrix A , Mirsky's bound in (1.38) is equal to the conjectured distance (1.36) (with $v = \epsilon$), i.e., if

$$(3.28) \quad |\operatorname{tr}(A^2)| = \epsilon^2(\Omega(A)) \quad ,$$

then the bound in (3.18) is also equal to the conjectured distance, i.e.,

$$(3.29) \quad \frac{|\operatorname{tr}(A)|^2}{n} + |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n}| = \epsilon^2(\Omega(A)) \quad ,$$

and

$$(3.30) \quad \epsilon^2(A - X) = \frac{1}{2} (\epsilon^2(A) - \epsilon^2(\Omega(A))) \quad \text{for all } X \in M(A)$$

where $M(A)$ is defined in Definition 3.8. .

Proof. By Corollary 3.24 (3.28) holds if and only if the eigenvalues of A can be written in the form $\lambda_k = r_k \exp[i(\theta + p_k)]$ ($k = 1, 2, \dots, n$) where $r_k \geq 0$, $\exp[ip_k] = \pm 1$ ($k = 1, 2, \dots, n$) and θ is a real constant. Thus

$$(3.31) \quad \operatorname{tr}(A^2) = \sum_{k=1}^n \lambda_k^2 = e^{12\theta} \sum_{k=1}^n r_k^2 e^{12p_k} = e^{12\theta} \sum_{k=1}^n r_k^2 \quad ,$$

$$(3.32) \quad (\operatorname{tr}(A))^2 = \left(\sum_{k=1}^n \lambda_k \right)^2 = \left(e^{i\theta} \sum_{k=1}^n r_k e^{ip_k} \right)^2 = e^{i2\theta} \left(\sum_{k=1}^n r_k e^{ip_k} \right)^2 .$$

Clearly

$$(3.33) \quad \left(\sum_{k=1}^n r_k e^{ip_k} \right)^2 \leq \left(\sum_{k=1}^n r_k \right)^2 \leq n \sum_{k=1}^n r_k^2 .$$

By Lemma 3.19 and (3.28), (3.29) holds if $\operatorname{tr}(A) = 0$. If $\operatorname{tr}(A) \neq 0$, then at least one r_k is positive and, from (3.31), (3.32), and (3.33), we have

$$\frac{n \operatorname{tr}(A^2)}{(\operatorname{tr}(A))^2} = \frac{n \sum r_k^2}{\left(\sum r_k \exp(ip_k) \right)^2} \geq \frac{n \sum r_k^2}{n \sum r_k^2} = 1 .$$

Therefore, if $\operatorname{tr}(A) \neq 0$, (3.21) holds so that equality holds in (3.20). This proves (3.29). Equation (3.30) then follows immediately from (3.12). This completes the proof of Lemma 3.27.

Remark. The equality (3.29) can hold also in many cases in which (3.28) is not satisfied; furthermore the bound in (3.18) can lie between the bound in (1.38) and the conjectured distance (1.36) (with $\nu = \epsilon$) as the following examples show. Suppose A is of order 4 and has eigenvalues $-1 + i$, $-1 + i$, $-1 + i$, and $2 + 6i$. Then straightforward computations show that

$$\frac{|\operatorname{tr}(A)|^2}{n} + |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n}| = \epsilon^2(\Omega(A)) = 46$$

while

$$36.71 < |\operatorname{tr}(A^2)| < 36.72 \quad .$$

Again, if A is of order 4 and has eigenvalues 1, 1, 1, and -1, then

$$\epsilon^2(\Omega(A)) = 4, \quad \frac{|\operatorname{tr}(A)|^2}{n} + |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{n}| = 2, \quad |\operatorname{tr}(A^2)| = 0 \quad .$$

We shall show later (Lemma 6.5) that (3.29) and (3.30) hold for all $A \in \mathcal{M}_2$.

3.2 Methods for Obtaining Other Bounds.

Another upper bound for $d_\epsilon(A)$ will be obtained in Chapter 5 (Theorem 5.24) in connection with a maximum problem which is closely related to ϵ -minimal matrices.

We shall now describe some problems whose solutions, if they were known, would give rise to upper bounds for $d_\epsilon(A)$.

Let $A \in \mathcal{M}_n$. For each natural number $k \geq 1$ we define the following subsets of \mathcal{M}_n :

$$(3.34) \quad \mathcal{L}_k(A) = \{z_0 I + \sum_{\alpha=1}^k z_\alpha (A + \zeta A^*)^\alpha; \zeta \in \underline{\mathbb{C}}, |\zeta| = 1, z_\alpha \in \underline{\mathbb{C}} \\ (\alpha = 0, 1, \dots, k)\} \quad ,$$

$$(3.35) \quad \mathcal{L}_k(A; \zeta) = \{z_0 I + \sum_{\alpha=1}^k z_\alpha (A + \zeta A^*)^\alpha; z_\alpha \in \underline{\mathbb{C}} \quad (\alpha = 0, 1, \dots, k)\}$$

where in (3.35) we require $\xi \in \underline{\mathbb{C}}$ and $|\xi| = 1$. From Theorems 1.67 and 1.69 we obtain

$$(3.36) \quad \mathcal{L}_k(A) \subset \mathcal{M}_n \quad \text{for } A \in \mathcal{M}_n,$$

$$(3.37) \quad \mathcal{L}_k(A; \xi) \subset \mathcal{L}_k(A) \quad \text{for all } |\xi| = 1.$$

Also, obviously,

$$(3.38) \quad \mathcal{L}_{k-1}(A) \subset \mathcal{L}_k(A), \mathcal{L}_{k-1}(A; \xi) \subset \mathcal{L}_k(A; \xi) \quad (k \geq 2).$$

From the Cayley-Hamilton Theorem one may deduce

$$(3.39) \quad \mathcal{L}_k(A) \subset \mathcal{L}_{n-1}(A), \mathcal{L}_k(A; \xi) \subset \mathcal{L}_{n-1}(A; \xi) \quad \text{where } A \in \mathcal{M}_n, \\ k \geq n.$$

3.40 Problem. Let $A \in \mathcal{M}_n$. Find

$$(3.41) \quad \min_{X \in \mathcal{L}_k(A)} \epsilon^2(A - X) \quad (1 \leq k \leq n - 1)$$

and find all matrices in $\mathcal{L}_k(A)$ for which the minimum in (3.41) is assumed.

A simpler problem is the following

3.42 Problem. Let $A \in \mathcal{M}_n$ and let ξ be a fixed complex number satisfying $|\xi| = 1$. Find

$$(3.43) \quad \min_{X \in \mathcal{L}_k(A; \xi)} \epsilon^2(A - X) \quad (1 \leq k \leq n - 1)$$

and find all matrices in $\mathcal{L}_k(A; \xi)$ for which the minimum in (3.43) is assumed.

By virtue of (3.36) and (3.37) a solution to either of the above problems would yield an upper bound for $d_\epsilon(A)$.

Remark. The author believes that the set $M(A)$ defined in Definition 3.8 provides a complete solution to Problem 3.40 for $k = 1$. No attempt to prove that assertion will be made here.

CHAPTER 4

NECESSARY CONDITIONS FOR ϵ -MINIMAL MATRICES

4.1 Primary Results.

In this section we shall derive a number of conditions which a matrix N_0 must satisfy if it is ϵ -minimal. One method which we employ involves the notion of a differentiable curve $N(t)$ in \mathcal{N} . If N_0 is v -minimal for A then, by elementary calculus, we know that

$$(4.1) \quad \left. \frac{d}{dt} v(A - N(t)) \right|_{t=0} = 0 \quad \text{or} \quad \left. \frac{d}{dt} v^2(A - N(t)) \right|_{t=0} = 0$$

must hold for every differentiable curve $N(t)$ such that $N(0) = N_0$, provided the derivatives indicated in (4.1) exist in an interval containing $t = 0$ in its interior. This brings up the (possibly difficult) question of how to differentiate an arbitrary (unitarily invariant) norm or some other function whose minima coincide with those of $v(A - N(t))$ (such as $v^2(A - N(t))$). In the case where v is the Euclidean norm, the derivatives of all orders can easily be computed. Indeed, from (1.22), we have

$$\begin{aligned} (4.2) \quad \frac{d}{dt} \epsilon^2(A - N(t)) &= \frac{d}{dt} \operatorname{tr}[(A - N(t))^* (A - N(t))] \\ &= \frac{d}{dt} \operatorname{tr}[A^* A + N^*(t)N(t) - A^* N(t) - N^*(t)A] \end{aligned}$$

Using (1.20) and (1.46), we have

$$(4.3) \quad \frac{d}{dt} \epsilon^2(A - N(t)) = \frac{d}{dt} \operatorname{tr}(N^*(t)N(t)) = \operatorname{tr}\left(A^* \frac{dN(t)}{dt} + \frac{dN^*(t)}{dt} A\right).$$

In the proof of the next theorem we shall need the following three elementary lemmas.

4.4 Lemma. Let $N(t) = U^*(t) D_0 U(t)$ where D_0 is any fixed element of \mathcal{O}_n and $U(t)$ is any differentiable curve in \mathcal{U}_n . Then

$$(4.5) \quad \operatorname{tr}(N^*(t)N(t)) = \epsilon^2(D_0)$$

Proof. This follows immediately from (1.22) and the unitary invariance of ϵ .

4.6 Lemma. (von Neumann [19], p. 290) Let $A \in \mathcal{M}_n$. Then $\operatorname{tr}(AH) = 0$ for all $H \in \mathcal{N}_n$ if and only if $A = 0$.

4.7 Lemma. Let $A = (a_{ij}) \in \mathcal{M}_n$. Then $\operatorname{Re} \operatorname{tr}(AA) = 0$ for all $A \in \mathcal{O}_n$ if and only if

$$(4.8) \quad a_{ii} = 0 \quad (i = 1, 2, \dots, n).$$

Proof. Letting $A = \operatorname{diag}(\alpha_1, \alpha_2, \dots, \alpha_n)$, the sufficiency of (4.8) is obvious from

$$(4.9) \quad \operatorname{Re} \operatorname{tr}(AA) = \operatorname{Re} \sum_{i=1}^n \alpha_i a_{ii} = \sum_{i=1}^n \operatorname{Re}(\alpha_i a_{ii}).$$

For the necessity we set $A = \operatorname{diag}(\bar{a}_{11}, \bar{a}_{22}, \dots, \bar{a}_{nn})$, obtaining $\operatorname{Re} \operatorname{tr}(AA) = \sum_{i=1}^n |a_{ii}|^2 = 0$ which implies (4.8).

4.10 Theorem. Let $A \in \mathcal{M}$ and define, for all $X \in \mathcal{M}$, the operator $L_A(X)$ by

$$(4.11) \quad L_A(X) = XA^* - A^*X + X^*A - AX^* .$$

If $N_0 = U_0^* D_0 U_0$, where $U_0 \in \mathcal{U}$ and $D_0 \in \mathcal{D}$, is ϵ -minimal for A then

$$(4.12) \quad L_A(N_0) = 0$$

and

$$(4.13) \quad D_0 = dg(U_0 A U_0^*)$$

where the function dg is defined by (1.24).

Proof. We first utilize differentiable curves of the form $N(t) = U^*(t) D_0 U(t)$, where $U(t)$ is given by (1.60) and H is any matrix in \mathcal{H} . For every $H \in \mathcal{H}$ we have $N(0) = N_0$ and

$$(4.14) \quad \frac{dU(t)}{dt} = iU_0 H e^{itH}, \quad \left. \frac{dU(t)}{dt} \right|_{t=0} = iU_0 H ;$$

$$\frac{dN(t)}{dt} = \frac{dU^*(t)}{dt} D_0 U(t) + U^*(t) D_0 \frac{dU(t)}{dt} ,$$

$$(4.15) \quad \left. \frac{dN(t)}{dt} \right|_{t=0} = i(N_0 H - H N_0) .$$

Using (4.3), Lemma 4.4, (1.48), and (4.15) we obtain

$$\begin{aligned}
(4.16) \quad \frac{d}{dt} \epsilon^2(A - N(t)) \Big|_{t=0} &= -\operatorname{tr}[A^* i(N_0 H - H N_0) - i(H N_0^* - N_0^* H) A] \\
&= i \operatorname{tr}(A^* H N_0 - A^* N_0 H + H N_0^* A - N_0^* H A) \\
&= i \operatorname{tr}[(N_0 A^* - A^* N_0 + N_0^* A - A N_0^*) H]
\end{aligned}$$

where, in the last step, we used (1.19) twice. According to (4.1) the last expression in (4.16) must vanish for all $H \in \mathcal{H}$. Therefore, by Lemma 4.6, we obtain

$$(4.17) \quad N_0 A^* - A^* N_0 + N_0^* A - A N_0^* = 0 ,$$

i.e., (4.12) holds.

Next let $N(t) = U_0^* D(t) U_0$ where $D(t)$ is a differentiable curve in \mathcal{D} satisfying $D(0) = D_0$. In using (4.1) with this type of curve $N(t)$ we are concerned only with the derivative of $D(t)$ (evaluated at $t = 0$) and not $D(t)$ itself. By Lemma 1.61 we obtain complete generality by setting

$$(4.18) \quad \frac{dD(t)}{dt} \Big|_{t=0} = \Lambda$$

where Λ denotes an arbitrary diagonal matrix. We define N_1 by $N_1 = U_0^* \Lambda U_0$ and note that

$$(4.19) \quad \frac{dN(t)}{dt} \Big|_{t=0} = U_0^* \frac{dD(t)}{dt} \Big|_{t=0} U_0 = U_0^* \Lambda U_0 = N_1 .$$

Thus, by (1.44),

$$(4.20) \quad \left. \frac{d}{dt} [N^*(t)N(t)] \right|_{t=0} = N_1^* N_0 + N_0^* N_1$$

and, using (4.3), (4.19), and (4.20), we obtain

$$\begin{aligned} \left. \frac{d}{dt} \epsilon^2(A - N(t)) \right|_{t=0} &= \text{tr}(N_1^* N_0 + N_0^* N_1 - A^* N_1 - N_1^* A) \\ &= \text{tr}(\Lambda^* D_0 + D_0^* \Lambda - U_0^* A^* U_0^* \Lambda - \Lambda^* U_0 A U_0^*) \\ (4.21) \quad &= \text{tr}(\Lambda^* (D_0 - U_0 A U_0^*) + [\Lambda^* (D_0 - U_0 A U_0^*)]^*) \\ &= 2 \text{Re tr}[\Lambda^* (D_0 - U_0 A U_0^*)] . \end{aligned}$$

It follows from (4.1) that (4.21) must vanish for all $\Lambda \in \mathcal{A}$; consequently, by Lemma 4.7, all diagonal elements of the matrix $D_0 - U_0 A U_0^*$ vanish, i.e., (4.13) holds.

Remark. It is instructive to investigate the result of applying (4.1) in the more general situation in which $N(t) = U^*(t)D(t)U(t)$, where $U(t)$ is again given by (1.60) and $D(t)$ is any differentiable curve in \mathcal{A} such that $D(0) = D_0$. The computations can be summarized as follows:

$$(4.22) \quad \frac{dN(t)}{dt} = U^*(t) \frac{dD(t)}{dt} U(t) + \frac{dU^*(t)}{dt} D(t)U(t) + U^*(t)D(t) \frac{dU(t)}{dt} ,$$

$$(4.23) \quad \left. \frac{dN(t)}{dt} \right|_{t=0} = N_1 + i(N_0 H - H^* N_0) ,$$

where N_1 is given by (4.19), and

$$(4.24) \quad \left. \frac{d}{dt} \epsilon^2(A - N(t)) \right|_{t=0} = 2 \operatorname{Re} \operatorname{tr}[A^*(D_0 - U_0 A U_0^*)] \\ + i \operatorname{tr}[(N_0 A^* - A^* N_0 + N_0^* A - A N_0^*)H] .$$

Of course the right side of (4.24) must vanish for all $A \in \mathcal{A}$ and for all $H \in \mathcal{H}$; but, from (4.16) and (4.21), we see that this is implied by the separate arguments used in the proof of Theorem 4.10. Therefore nothing new can be obtained by using the more general type of differentiable curves in \mathcal{N} .

4.25 Theorem. Let $A \in \mathcal{M}$. If N_0 is ϵ -minimal for A then

$$(4.26) \quad \operatorname{tr}(A) = \operatorname{tr}(N_0) ,$$

$$(4.27) \quad \operatorname{tr}[N_0^*(A - N_0)] = \operatorname{tr}[(A - N_0)^* N_0] = 0 .$$

Remark. The reader will note from the following proof that both (4.26) and (4.27) are consequences of (4.13), as are also the following equalities which are equivalent with (4.27):

$$(4.28) \quad \epsilon^2(N_0) = \operatorname{tr}(A^* N_0) = \operatorname{tr}(N_0^* A) ,$$

$$(4.29) \quad d_\epsilon^2(A) = \epsilon^2(A - N_0) = \epsilon^2(A) - \epsilon^2(N_0) .$$

Proof. Since $\operatorname{tr}(D_0) = \operatorname{tr}(N_0)$, (4.26) is an immediate consequence of (4.13). Writing $N_0 = U_0^* D_0 U_0$ as before, we have

$$\begin{aligned}
(4.30) \quad \operatorname{tr}[N_0^*(A - N_0)] &= \operatorname{tr}[U_0 N_0^* U_0^* U_0 (A - N_0) U_0^*] \\
&= \operatorname{tr}[D_0^*(U_0 A U_0^* - D_0)] .
\end{aligned}$$

By (4.13) all diagonal elements of the matrix $U_0 A U_0^* - D_0$ are zero, consequently the same holds for the matrix $D_0^*(U_0 A U_0^* - D_0)$ hence $\operatorname{tr}[D_0^*(U_0 A U_0^* - D_0)] = 0$. Thus, from (4.30), we have

$$(4.31) \quad \operatorname{tr}[N_0^*(A - N_0)] = 0$$

The last equality implies that $\operatorname{tr}[N_0^*(A - N_0)]$ is real; consequently the rest of (4.27) follows from (1.21).

Remark. Let $A, B \in \mathcal{M}_n$. The set \mathcal{M}_n can be considered as a complex Hilbert space with respect to the inner product

$$(4.32) \quad (A, B) = \operatorname{tr}(AB^*) .$$

Comparing (4.32) with (4.27) we see that, if N_0 is ϵ -minimal for A ,

$$(4.33) \quad (A - N_0, N_0) = 0 ;$$

that is, N_0 is orthogonal to $A - N_0$. This geometrical reformulation of the necessary condition (4.27) is not too surprising. Indeed, if N_0 is ϵ -minimal, then it must also be ϵ -minimal among all members of the one-dimensional linear subspace (of \mathcal{M}_n):

$$\mathcal{L} = \{X ; X = zN_0 \text{ where } z \in \underline{\mathbb{C}}\}$$

which is contained in \mathcal{N}_n . \mathcal{L} is precisely the "straight line" through N_0 and the null matrix. Since no point on \mathcal{L} can be closer to A than N_0 , we see that N_0 must be the point of intersection of \mathcal{L} and the line through A which is perpendicular to \mathcal{L} . This gives immediately the result (4.33) so we have proved (4.27) a second time without using differentiable curves and calculus. One can also prove (4.27) directly (without using (4.13)) by using (4.1) and the particular differentiable curves $N(t) = (1 - t)N_0$ and $N(t) = (1 - it)N_0$. The latter proof, which will not be worked out in detail here, appears to be essentially an analytic reformulation of the geometric considerations just mentioned.

We shall prove later (Theorem 5.13) that (4.13) and another condition on U_0 constitute necessary and sufficient conditions for ϵ -minimality. It appears that (4.13) is a much more stringent condition than (4.12); however, the latter is much easier to check than (4.13) and this fact enhances its value. We now describe some consequences of the necessary condition (4.12).

4.34 Theorem. Let $A \in \mathcal{M}$, $N_0 \in \mathcal{M}$. Then the following statements are equivalent.

- (a) $L_A(N_0) = 0$;
- (b) $N_0 A^* + N_0^* A \in \mathcal{H}$;
- (c) $(N_0 - \alpha A)^* + (N_0 - \beta A)^* A \in \mathcal{H}$ for any $\alpha, \beta \in \underline{\mathbb{R}}$;
- (d) $N_0(A - \alpha N_0)^* + N_0^*(A - \beta N_0) \in \mathcal{H}$ for any $\alpha, \beta \in \underline{\mathbb{R}}$;
- (e) $N_0(A - zN_0)^* + N_0^*(A - wN_0) = [N_0(A - wN_0)^* + N_0^*(A - zN_0)]^*$
holds for any $z, w \in \underline{\mathbb{C}}$ provided $N_0 \in \mathcal{N}$.

Proof. Statement (a) is equivalent to (4.17) and we shall show that (b), (c), (d), and (e) are equivalent to (4.17). First note that (4.17) can be rewritten in the form

$$N_O A^* + N_O^* A = A N_O^* + A^* N_O .$$

The right side of the last equation obviously equals $[N_O A^* + N_O^* A]^*$ which shows that (b) is equivalent to (4.17). For any real numbers α and β we note that

$$N_O A^* - \alpha N_O N_O^* - A^* N_O + \beta N_O^* N_O + N_O^* A - \beta N_O^* N_O - A N_O + \alpha N_O N_O^* = 0 ,$$

$$N_O A^* - \alpha A A^* - A^* N_O + \beta A^* A + N_O^* A - \beta A^* A - A N_O + \alpha A A^* = 0$$

are equivalent to (4.17). Clearly the last two equations are equivalent respectively to

$$N_O (A - \alpha N_O)^* - (A - \beta N_O)^* N_O + N_O^* (A - \beta N_O) - (A - \alpha N_O) N_O^* = 0 ,$$

$$(N_O - \alpha A) A^* - A^* (N_O - \beta A) + (N_O - \beta A)^* A - A (N_O - \alpha A)^* = 0$$

which are equivalent respectively to statements (c) and (d). If we assume that N_O is normal and let z and w be any complex numbers, we obtain from (4.17)

$$N_O A^* - \bar{z} N_O N_O^* - A^* N_O + \bar{z} N_O^* N_O + N_O^* A - w N_O^* N_O - A N_O^* + w N_O N_O^* = 0$$

or

$$N_O (A - z N_O)^* - (A - z N_O)^* N_O + N_O^* (A - w N_O) - (A - w N_O) N_O^* = 0$$

which is equivalent to (e). This proves Theorem 4.34.

In terms of the operator (4.11) we observe the following rather obvious consequences of Theorem 4.34. If $A, N_0 \in \mathcal{M}$, then the following statements are equivalent to (4.12):

$$(4.35) \quad L_A(N_0 - \alpha A) = 0 \quad \text{for any } \alpha \in \underline{\mathbb{R}},$$

$$(4.36) \quad L_{A-\alpha N_0}(N_0) = 0 \quad \text{for any } \alpha \in \underline{\mathbb{R}},$$

$$(4.37) \quad L_{A-zN_0}(N_0) = 0 \quad \text{for any } z \in \underline{\mathbb{C}} \text{ provided } N_0 \in \mathcal{M}.$$

4.38 Lemma. Let $A \in \mathcal{M}$ and suppose the operator $L_A(X)$ is defined by (4.11). Then, for all $X, Y \in \mathcal{M}$, we have

$$(4.39) \quad L_A(\alpha X) = \alpha L_A(X) \quad \text{for all } \alpha \in \underline{\mathbb{R}},$$

$$(4.40) \quad L_A(X + Y) = L_A(X) + L_A(Y).$$

Proof. For any $\alpha \in \underline{\mathbb{C}}$ we have

$$(4.41) \quad L_A(\alpha X) = \alpha(XA^* - A^*X) + \bar{\alpha}(X^*A - AX^*)$$

whence (4.39) holds if α is real. The verification of (4.40) is equally simple and need not be given here.

Remark. Equations (4.39) and (4.40) show that, given any $A \in \mathcal{M}$,

$$\{X ; X \in \mathcal{M} \text{ and } L_A(X) = 0\}$$

is a real linear subspace of \mathcal{M} , i.e., a linear subspace of \mathcal{M} over the field $\underline{\mathbb{R}}$ of real scalars.

Page Intentionally Left Blank

BLANK PAGE

The result (4.39) cannot in general be extended to all complex scalars α as the following simple example shows. Suppose $A \notin \mathcal{M}$. Then from (4.11) and (4.41) we find that $L_A(A) = 0$ but $L_A(iA) = 2i(AA^* - A^*A) \neq 0$. Of course, in certain special situations, (4.39) can hold for all $\alpha \in \mathbb{C}$. Examples of the latter can be deduced from the following easily verified results which are valid for any $A \in \mathcal{M}_n$:

$$(4.42) \quad L_A(zI) = 0 \quad \text{for all } z \in \mathbb{C} ;$$

$$(4.43) \quad L_A(\alpha A) = 0 \quad \text{for all } \alpha \in \mathbb{R} ;$$

$$(4.44) \quad L_A(z(A^*)^k) = 0 \quad \text{for all } z \in \mathbb{C} \text{ and for } k = 1, 2, \dots ;$$

$$(4.45) \quad L_A(zA^2 + \bar{z}(AA^* + A^*A)) = 0 \quad \text{for all } z \in \mathbb{C} ;$$

$$(4.46) \quad L_A(\alpha A + zA^*) = 0 \quad \text{for all } \alpha \in \mathbb{R}, \quad z \in \mathbb{C} ;$$

$$(4.47) \quad L_A[z(A + \frac{\bar{z}}{z}A^*)^2] = 0 \quad \text{for all } z \in \mathbb{C} \text{ with } z \neq 0 .$$

Note. The content of the next lemma is that every matrix in the set $M(A)$ defined in Definition 3.8 satisfies (for any value of n) all of the necessary conditions given in Theorems 4.10 and 4.25 with the possible exception of (4.13).

Remark. While condition (4.13) is apparently the most crucial test of a candidate for ϵ -minimality of all the necessary conditions derived so far, it cannot be checked unless the normal matrix in question is first diagonalized by a unitary transformation.

4.48 Lemma. Let $A \in \mathcal{M}_n$ and let $M(A)$ be defined as in Definition 3.8. Then the necessary conditions (4.12), (4.26), and (4.27) are satisfied when N_0 is replaced by any matrix in $M(A)$.

Proof. Let

$$(4.49) \quad P_{\zeta}(A) = \frac{1}{2} (A + \zeta A^*) + \frac{1}{2n} \operatorname{tr}(A - \zeta A^*) I.$$

The fact that $P_{\zeta}(A)$ satisfies (4.12) for all $\zeta \in \underline{\mathbb{C}}$ follows immediately from (4.40), (4.42), and (4.46). This implies that every matrix in $M(A)$ satisfies (4.12). Furthermore

$$\operatorname{tr}[P_{\zeta}(A)] = \frac{1}{2} \operatorname{tr}(A) + \frac{1}{2} \zeta \operatorname{tr}(A^*) + \frac{1}{2} \operatorname{tr}(A - \zeta A^*) = \operatorname{tr}(A)$$

for all $\zeta \in \underline{\mathbb{C}}$ hence every matrix in $M(A)$ satisfies (4.26). Let

$$(4.50) \quad z = \frac{1}{n} \operatorname{tr}(A - \zeta A^*)$$

Then $\bar{z} = (1/n) \operatorname{tr}(A^* - \bar{\zeta} A)$, $P_{\zeta}(A) = (1/2)(A + \zeta A^* + zI)$, and $A - P_{\zeta}(A) = (1/2)(A - \zeta A^* - zI)$. Consequently, for $|\zeta| = 1$,

$$\begin{aligned} P_{\zeta}^*(A)(A - P_{\zeta}(A)) &= \frac{1}{4} (A^* + \bar{\zeta} A + \bar{z} I)(A - \zeta A^* - zI) \\ &= \frac{1}{4} (A^* A - \zeta A^{*2} + \bar{\zeta} A^2 - A A^* - z(A^* + \bar{\zeta} A) + \bar{z}(A - \zeta A^*) - z\bar{z} I) \end{aligned}$$

so that

$$(4.51) \quad 4 \operatorname{tr}[P_{\xi}^*(A)(A - P_{\xi}(A))] = \bar{\xi} \operatorname{tr}(A^2) - \xi \operatorname{tr}(A^{*2}) - z \operatorname{tr}(A^* + \bar{\xi}A) \\ + \bar{z} \operatorname{tr}(A - \xi A^*) - nz\bar{z}, \quad (|\xi| = 1).$$

Using (4.50) in the last term of equation (4.51) we see that the last two terms on the right side of (4.51) cancel. Using (4.50) again we find, for $|\xi| = 1$,

$$(4.52) \quad z \operatorname{tr}(A^* + \bar{\xi}A) = \frac{1}{n} [\operatorname{tr}(A) - \xi \operatorname{tr}(A^*)][\operatorname{tr}(A^*) + \bar{\xi} \operatorname{tr}(A)] \\ = \frac{1}{n} [|\operatorname{tr}(A)|^2 + \bar{\xi}(\operatorname{tr}(A))^2 - \xi(\operatorname{tr}(A^*))^2 - |\operatorname{tr}(A)|^2] \\ = \frac{1}{n} [\bar{\xi}(\operatorname{tr}(A))^2 - \xi(\operatorname{tr}(A^*))^2].$$

Combining (4.51) and (4.52) we obtain

$$(4.53) \quad 4 \operatorname{tr}[P_{\xi}^*(A)(A - P_{\xi}(A))] \\ = \bar{\xi}(\operatorname{tr}(A^2) - \frac{1}{n}(\operatorname{tr}(A))^2) - \xi(\operatorname{tr}(A^2) - \frac{1}{n}(\operatorname{tr}(A))^2).$$

If $\operatorname{tr}(A^2) = (1/n)(\operatorname{tr}(A))^2$ then the right side of (4.53) vanishes for all $\xi \in \xi(A)$. On the other hand, if $\operatorname{tr}(A^2) \neq (1/n)(\operatorname{tr}(A))^2$ then by (3.9) the right side of (4.53) becomes

$$|\operatorname{tr}(A^2) - \frac{1}{n}(\operatorname{tr}(A))^2| - |\operatorname{tr}(A^2) - \frac{1}{n}(\operatorname{tr}(A))^2| = 0.$$

Thus every matrix in $M(A)$ satisfies (4.27) and this proves Lemma 4.48.

In the light of Lemma 4.48, it is natural to ask whether the matrices in the set $M(A)$ can ever be ϵ -minimal. We shall answer this question in the affirmative in Section 6.2 below.

4.2 Use of Lagrange Multipliers in Finding Necessary Conditions.

As was observed in Section 1.5, the problem of finding a v -minimal matrix can be considered as a problem of minimizing a real valued function of several complex variables. The minimum must be taken over normal matrices only, so certain constraints on the variables are inevitable. This point of view suggests that the method of Lagrange multipliers might be useful in deriving necessary conditions for a v -minimal matrix. Since the variables and constraints are complex, it first appears that the initial computational labor of deriving such conditions would be prohibitive. However, by using conjugate complex coordinates (see e.g., [18] pp. 16-21) and the complex (matrix) differential calculus, the algebraic manipulations can be made almost inconsequential. In the next two paragraphs we outline briefly the mathematical basis of the technique we shall use in deriving necessary conditions for ϵ -minimal matrices. The equations which will be derived appear to be very difficult to solve in general, primarily because they involve a large number of auxiliary unknown variables (Lagrange multipliers). As a consequence these equations are not of immediate value in finding ϵ -minimal matrices. Nevertheless, the equations themselves and the technique of deriving

them seem interesting. Furthermore as byproducts we obtain alternative derivations of the primary necessary conditions (4.12), (4.13) and also additional necessary conditions for ϵ -minimal matrices involving theoretically interesting interpretations of certain Lagrange multipliers.

Let x and y denote a pair of real variables (coordinates). Following Nehari [18] we define $z = x + iy$ and $\bar{z} = x - iy$ to be the corresponding pair of conjugate complex variables (coordinates). We define formally a pair of "partial differential operators" by the expressions

$$(4.54) \quad \frac{\partial}{\partial z} = \frac{1}{2} \left(\frac{\partial}{\partial x} - i \frac{\partial}{\partial y} \right) ,$$

$$(4.55) \quad \frac{\partial}{\partial \bar{z}} = \frac{1}{2} \left(\frac{\partial}{\partial x} + i \frac{\partial}{\partial y} \right) .$$

As shown by Nehari, these operations can be carried out, at least for a wide class of functions of x and y (or equivalently of z and \bar{z}), by treating z and \bar{z} as independent variables and differentiating by the usual rules of calculus. As a simple illustration of the method, suppose $g(x,y) = f(z,\bar{z})$ is a real valued function for which the stationary points are desired. The standard method of elementary calculus involves solving the simultaneous pair of equations

$$(4.56) \quad \frac{\partial g}{\partial x} = 0 , \quad \frac{\partial g}{\partial y} = 0 ;$$

whereas in the conjugate variables method one merely works with the single (complex) equation

$$(4.57) \quad \frac{\partial f}{\partial \bar{z}} = 0$$

which is equivalent to the pair of equations (4.56). It is obvious how one may extend this procedure to the problem of finding stationary points of real valued functions of $2k$ real variables where $k = 2, 3, \dots$.

An alternative technique, which is often easier to apply in matrix problems, is based on the use of differentials. Nehari showed that the operators (4.54) and (4.55) behave like genuine partial derivatives in the sense that the differential of a real function $g(x, y) = f(z, \bar{z})$ is given by

$$(4.58) \quad dg = df = \frac{\partial f}{\partial z} dz + \frac{\partial f}{\partial \bar{z}} d\bar{z}$$

where $dz = dx + i dy$ and $d\bar{z} = dx - i dy$. In elementary calculus the equivalent formulation of (4.56) in terms of differentials is the statement that $dg = (\partial g / \partial x) dx + (\partial g / \partial y) dy$ vanishes for all values of the independent differentials dx and dy . This does not correspond exactly to the statement that (4.58) vanishes for all dz and $d\bar{z}$ considered as independent differentials. For, from that statement, one would obtain both $\partial f / \partial z = 0$ and $\partial f / \partial \bar{z} = 0$ which would imply (4.56) twice since $\partial f / \partial z = 0$ if and only if $\partial f / \partial \bar{z} = 0$. It is sufficient to compute the complex "half differential"

$$(dg)_{\bar{z}} = (df)_{\bar{z}} = \frac{\partial f}{\partial \bar{z}} d\bar{z}$$

and say that it must vanish for all $d\bar{z}$. In the case of functions of more than two real variables, this rule for deriving necessary conditions can be generalized as follows.

4.59 Theorem. Let $z = (z_1, z_2, \dots, z_k)$ denote a vector with complex components and suppose $\phi(z, \bar{z}) \equiv \phi(z_1, z_2, \dots, z_k, \bar{z}_1, \bar{z}_2, \dots, \bar{z}_k)$ is a real valued differentiable function of the $2k$ real variables x_α, y_α ($\alpha = 1, \dots, k$) where $z_\alpha = x_\alpha + i y_\alpha$ ($\alpha = 1, \dots, k$). Let the half differential $(d\phi)_{\bar{z}}$ be defined by

$$(4.60) \quad (d\phi)_{\bar{z}}(z, \bar{z}) \equiv \frac{\partial \phi(z, \bar{z})}{\partial \bar{z}_1} d\bar{z}_1 + \frac{\partial \phi(z, \bar{z})}{\partial \bar{z}_2} d\bar{z}_2 + \dots + \frac{\partial \phi(z, \bar{z})}{\partial \bar{z}_k} d\bar{z}_k.$$

If ϕ has an extreme value at the point $z_0 = (z_1^{(0)}, \dots, z_k^{(0)})$, then $(d\phi)_{\bar{z}}(z_0, \bar{z}_0)$ must vanish for all values of the independent differentials $d\bar{z}_1, d\bar{z}_2, \dots, d\bar{z}_k$.

We refrain from giving a proof of Theorem 4.59.

Let $A = (a_{ij}) \in \mathcal{M}_n$ and $N = (n_{ij}) \in \mathcal{M}_n$. We first consider the straightforward problem of minimizing $\epsilon^2(N - A)$ subject to the constraint

$$(4.61) \quad N^*N - NN^* = 0.$$

Since the matrix $N^*N - NN^*$ is hermitian, the matrix equation (4.61) amounts to: (1) at most $(1/2)n(n-1)$ independent complex constraints on the n_{ij} (these correspond to all elements of (4.61) which lie on one side of the diagonal), and (2) at most n independent real constraints (these correspond to the n (real) diagonal elements of (4.61)). Thus

(4.61) involves at most n^2 real constraints on the n^2 complex quantities n_{ij} . Clearly then, we need at most n^2 real Lagrange multipliers. Nevertheless, it is convenient in deriving the necessary conditions to introduce at the beginning $2n^2$ real Lagrange multipliers, one for each of the real constraints

$$(4.62) \quad \operatorname{Re}[(N^*N - NN^*)_{ij}] = 0, \quad \operatorname{Im}[(N^*N - NN^*)_{ij}] = 0$$

where $1 \leq i \leq n$, $1 \leq j \leq n$. The introduction of superfluous Lagrange multipliers causes no difficulty since the superfluous ones either drop out or combine automatically with others to form the maximal number of independent multipliers. We let λ_{ij} , μ_{ij} be real Lagrange multipliers which correspond (in that order) to the constraints (4.62). We now define the scalar quantities

$$(4.63) \quad \tau_{ij} = \lambda_{ij} + i \mu_{ij} \quad (1 \leq i, j \leq n)$$

and call the τ_{ij} complex Lagrange multipliers. For later use we also define the matrix

$$(4.64) \quad T = (\tau_{ij})$$

In the classical way of applying the method of Lagrange multipliers one would consider the function

$$(4.65) \quad \epsilon^2(N - A) + \sum_{i,j=1}^n (\lambda_{ij} \operatorname{Re}[(N^*N - NN^*)_{ij}] + \mu_{ij} \operatorname{Im}[(N^*N - NN^*)_{ij}])$$

and differentiate it partially with respect to the real variables $\text{Re}(n_{rc}), \text{Im}(n_{rc})$ ($1 \leq r, c \leq n$). As was suggested above, this procedure is unnecessarily complicated and can be greatly simplified and shortened by using complex analysis.

Note that

$$\lambda_{1j} \text{Re}[(N^*N - NN^*)_{1j}] + \mu_{1j} \text{Im}[(N^*N - NN^*)_{1j}] = \text{Re}[\bar{\tau}_{1j}(N^*N - NN^*)_{1j}]$$

so that the sum on the right side of (4.65) can be written in the form

$$\begin{aligned} \text{Re}\left\{\sum_{j=1}^n \bar{\tau}_{1j}(N^*N - NN^*)_{1j}\right\} &= \text{Re} \text{tr}[(N^*N - NN^*)T^*] \\ (4.66) \qquad &= \frac{1}{2} (\text{tr}[(N^*N - NN^*)T^*] + \overline{\text{tr}[(N^*N - NN^*)T^*]}) \\ &= \frac{1}{2} \text{tr}[(N^*N - NN^*)(T + T^*)] \end{aligned}$$

where T is given by (4.64). In view of the form of the constraint equations (4.61) the factor $(1/2)$ in (4.66) can be omitted upon substitution of (4.66) into (4.65). Thus, writing $\epsilon^2(N - A) = \text{tr}[(N - A)^*(N - A)]$, we consider the problem of minimizing the function

$$(4.67) \quad \Phi(N, N^*) = \text{tr}[(N^* - A^*)(N - A)] + \text{tr}[(N^*N - NN^*)H]$$

where H is given by

$$(4.68) \quad H = T + T^*.$$

Equation (4.67) shows that, as far as the Lagrange multipliers (4.63) are concerned, it is the matrix H which has a definite meaning (i.e., which has to be determined using (4.61)) and not T . Clearly the Lagrange multipliers μ_{ii} ($i = 1, \dots, n$) which correspond to nonexistent constraints in (4.61) have dropped out and the others have combined to form the n^2 "actual Lagrange multipliers":

$$\begin{aligned} 2 \lambda_{ii} & \quad (i = 1, \dots, n) , \\ \lambda_{ij} + \lambda_{ji} , \quad \mu_{ij} - \mu_{ji} & \quad (1 \leq i < j \leq n) . \end{aligned}$$

Differentiating (4.67) we obtain

$$\begin{aligned} (d\phi)_{N^*}^* &= \text{tr}[dN^*(N - A) + (dN^*N - NdN^*)H] \\ (4.69) \quad &= \text{tr}[(N - A + NH - HN)dN^*] . \end{aligned}$$

By Theorem 4.59 ϕ has a minimum at the point (matrix) N only if $(d\phi)_{N^*}^*$ vanishes for all matrices dN^* . It is easy to show from von Neumann's Lemma 4.6 that (4.69) vanishes for all $dN^* \in \mathcal{M}_n$ if and only if

$$N - A + NH - HN = 0 .$$

If we rewrite the last equation and the constraint (4.61) and note that (4.68) implies $H \in \mathcal{H}_n$, we obtain the following set of three simultaneous matrix equations for the unknown matrices N and H .

$$(4.70) \quad N + NH - HN = A ,$$

$$(4.71) \quad H = H^* ,$$

$$(4.72) \quad N^* N - N N^* = 0 .$$

4.73 Theorem. Let $A \in \mathcal{M}_n$. If N_0 is ϵ -minimal for A then there exists a matrix $H_0 \in \mathcal{M}_n$ such that (4.70), (4.71), and (4.72) are satisfied when N and H are replaced respectively by N_0 and H_0 .

Proof. If N_0 is ϵ -minimal it minimizes the functional $\epsilon^2(N - A)$ and satisfies (4.72). Consequently N_0 is a stationary point of $\epsilon^2(N - A)$ which implies the existence of complex Lagrange multipliers τ_{ij} such that N_0 is also a stationary point of (4.67). By the above derivation of equation (4.70) there exists a hermitian matrix H_0 such that $N_0 + N_0 H_0 - H_0 N_0 = A$. This proves Theorem 4.73.

Remark. It should be noted that the hermitian matrix H of Lagrange multipliers need not necessarily be of the form of the most general hermitian matrix. For example in the case $n = 2$ the diagonal elements of the left side of (4.72) differ only in sign (cf. Lemma 6.1 below) so one additional real Lagrange multiplier can be eliminated, allowing H to assume the form

$$H = \begin{pmatrix} \lambda_1 & \lambda + i\mu \\ \lambda - i\mu & -\lambda_1 \end{pmatrix}$$

where λ_1 , λ , and μ are real. It should also be noted, however, that H may be assumed to be a general hermitian matrix. There is no harm in repeating constraints in the method of Lagrange multipliers.

Since at least one solution of the necessary condition described in Theorem 4.73 is ϵ -minimal, it is natural to investigate the problem of solving the system (4.70) - (4.72). In theory one could carry out the following procedure. Assuming that H is hermitian, solve (4.70) for N as a function of A and H , i.e., find

$$(4.74) \quad N = F(A, H) \quad .$$

Then substitute (4.74) into (4.72) and solve the latter equation for the hermitian matrix H . Once H was determined it would be eliminated by substitution into (4.74) which would yield a normal matrix.

Actually, it is possible to solve the system (4.70) - (4.72) formally, although not explicitly, by a method which runs somewhat along the lines of the procedure outlined in the previous paragraph. Since the matrix H of Lagrange multipliers is hermitian it can be diagonalized by a unitary transformation:

$$(4.75) \quad UH U^* = D \quad (U \in \mathcal{U}_n, D \in \mathcal{D}_n)$$

where D is real. By means of (4.75) we replace the unknown hermitian matrix H by two other matrices, namely an unknown unitary matrix U and an unknown diagonal matrix D . Now let

$$(4.76) \quad UN U^* = M = (m_{ij}) \quad , \quad UAU^* = B = (b_{ij})$$

where U is the same matrix which occurs in (4.75). By premultiplication by U and postmultiplication by U^* we see that the equations (4.70) and (4.72) are equivalent respectively to

$$(4.77) \quad M + MD - DM = B$$

and

$$(4.78) \quad M^* M - M M^* = 0 \quad .$$

Note that N is normal if and only if M is normal.

4.79 Definition. Let the matrix (4.75) be given by $D = \text{diag}(d_1, d_2, \dots, d_n)$,

let $\Delta_{ij} = d_j - d_i$ ($1 \leq i, j \leq n$), and let $\Delta = \{\Delta_{ij}; 1 \leq i < j \leq n\}$.

Then $M(A, U, \Delta)$ is defined to be the matrix

$$(4.80) \quad M(A, U, \Delta) = \left(\frac{b_{ij}}{1 + \Delta_{ij}} \right)$$

where $B = (b_{ij})$ is defined by (4.76).

Using the notation of (4.76) and the preceding definition we see that (4.77) (which can be rewritten in the form $M(I + D) - DM = B$) amounts to a set of scalar equations:

$$m_{ij}(1 + d_j) - d_i m_{ij} = b_{ij}$$

or

$$(4.81) \quad m_{ij}(1 + d_j - d_i) = m_{ij}(1 + \Delta_{ij}) = b_{ij} \quad (1 \leq i, j \leq n) \quad .$$

If we choose the d_i 's in such a way that $1 + d_j - d_i = 1 + \Delta_{ij} \neq 0$ for all pairs (i, j) then we can solve (4.81) for the m_{ij} . From (4.80) we obtain $M = M(A, U, \Delta)$, a known function of A , U , and Δ . The set Δ is determined by D so the matrix (4.90) is actually a function of A , U , and D . Thus, if we choose U and D so as to make the matrix (4.80) normal then

BLANK PAGE

$$(4.88) \quad T_{U_0, \Delta_0}(A) \in \mathcal{N}_n$$

and

$$(4.89) \quad \epsilon^2(U_0 A U_0^* - T_{U_0, \Delta_0}(A)) = \text{Min } \epsilon^2(U A U^* - T_{U, \Delta}(A))$$

where the minimum is taken over all U, D such that $T_{U, \Delta}(A) \in \mathcal{N}_n$ and where Δ and Δ_0 are determined respectively by D and D_0 as in Definition 4.79.

In order to examine this problem we note that the transformation (4.86) is the composition $T_{\Delta} \cdot T_{U^*}$ where the transformation on the right is to be carried out first, where T_U is defined by (1.75), and where

$$(4.90) \quad T_{\Delta}(M) = \left(\frac{m_{ij}}{1 + \Delta_{ij}} \right), \quad (M = (m_{ij}) \in \mathcal{M}_n, \Delta_{ij} \neq -1).$$

Thus $T_{U, \Delta}(A) = T_{\Delta}(T_{U^*}(A))$. If A is not normal, then neither is $B = T_{U^*}(A)$; hence the second transformation T_{Δ} plays an indispensable role in achieving the normality of $T_{U, \Delta}(A)$. On the other hand, it can be shown by examples that $T_{U, \Delta}(A) \in \mathcal{N}_n$ can hold for all Δ 's which arise from real diagonal matrices D . Thus we have gained some insight into the meaning and role of the matrix H of Lagrange multipliers:

The possibility that

$$(4.91) \quad T_{U, \Delta}(A) \in \mathcal{N}_n$$

for some Δ depends on the choice of U ; the actual realization of (4.91) depends on the choice of Δ and not on U . If we translate back into the coordinate system of the set of equations (4.70) - (4.72), we can summarize those observations in the following manner.

4.92 Theorem. Let $A \in \mathcal{M}_n$, $A \notin \mathcal{N}_n$ and let $H \in \mathcal{H}_n$. Then the possibility that (4.70) has a normal solution N (for some choice of the eigenvalues of H) depends on the eigenvectors of H . The choice of the eigenvalues (and not the eigenvectors) of H determines whether or not (4.70) has a normal solution N .

Proof. We need only observe from (4.75) that the eigenvectors of H are the columns of U^* and that Δ is determined by the eigenvalues d_1 of H .

Instead of Problem 4.87 we may consider the simpler

4.93 Problem. Let $A \in \mathcal{M}_n$, $A \notin \mathcal{N}_n$ and define

$$\mathcal{S}(A) = \{U ; U \in \mathcal{U}_n \text{ and } T_{U,\Delta}(A) \in \mathcal{N}_n \text{ for some admissible } \Delta\}.$$

Find a $U \in \mathcal{S}(A)$. Then, for this U , find a $D_0 \in \mathcal{D}_n$ such that

$$T_{U,\Delta_0}(A) \in \mathcal{N}_n$$

and

$$\epsilon^2(UAU^* - T_{U,\Delta_0}(A)) = \text{Min } \epsilon^2(UAU^* - T_{U,\Delta}(A))$$

where the minimum is taken over all D such that $T_{U,\Delta}(A) \in \mathcal{N}_n$ and where Δ and Δ_0 are determined respectively by D and D_0 as in Definition 4.79.

The explicit determination of solutions to either Problem 4.87 or Problem 4.93 appears to be difficult to carry out for $n \geq 3$. The author has found solutions to Problem 4.93 for all $A \in \mathcal{M}_2$ and, on the basis of the results in Chapter 6 below, these have all been ϵ -minimal. No further investigation of Problems 4.87 or 4.93 will be made here except for the following simple observation.

4.94 Lemma. Let $U \in \mathcal{U}_n$ and let $\Delta = [\Delta_{ij}]$ ($1 \leq i, j \leq n$) where $\Delta_{ij} \neq -1$ for all i and j . Then the transformation $T_{U,\Delta}(A)$ defined by (4.86) is a linear transformation on \mathcal{M}_n .

Proof. The transformations T_U of (1.75) and T_Δ of (4.90) are both obviously linear. Thus $T_{U,\Delta}$, being a composition of linear transformations, is linear too.

Note. We shall find another application for equation (4.70) in Chapter 7 below.

We consider next another approach to the use of Lagrange multipliers. By Theorem 1.57 (c) $N = U^*DU$ runs through \mathcal{N}_n as U and D run independently through \mathcal{U}_n and \mathcal{D}_n respectively. Thus we can find an ϵ -minimal matrix for A by minimizing

$$(4.95) \quad \epsilon^2(U^*DU - A) = \text{tr}[(U^*D^*U - A^*)(U^*DU - A)]$$

subject to the constraints

$$(4.96) \quad U^*U = I ,$$

$$(4.97) \quad D \in \mathcal{D}_n .$$

As in the case of the constraint (4.61) we introduce complex Lagrange multipliers τ_{ij} corresponding to the (i,j) elements $(U^*U - I)_{ij}$ of the constraint equation (4.96) and write $T = (\tau_{ij})$ as before. If we merely keep in mind that D is diagonal, we need not introduce any Lagrange multipliers for the constraint (4.97). We have

$$2 \operatorname{Re} \operatorname{tr}(U^*UT^*) = \operatorname{tr}[U^*U(T + T^*)] .$$

Here, as in the previous constrained minimization problem, it is not T which has to be determined but the hermitian matrix $H = T + T^*$. Therefore the function we wish to minimize in the method of Lagrange multipliers is

$$(4.98) \quad \Phi(U, U^*, D, D^*) = \operatorname{tr}[(U^*D^*U - A^*)(U^*DU - A) + U^*UH] .$$

Differentiating this we obtain

$$\begin{aligned} (d\Phi)_{U^*, D^*} &= \operatorname{tr}[dU^*D^*U(U^*DU - A) + (U^*D^*U - A)dU^*DU + dU^*UH] \\ &\quad + \operatorname{tr}[U^*dD^*U(U^*DU - A)] \\ &= \operatorname{tr}([D^*U(U^*DU - A) + DU(U^*D^*U - A^*) + UH]dU^*) \\ &\quad + \operatorname{tr}[dD^*U(U^*DU - A)U^*] . \end{aligned}$$

By Theorem 4.59 Φ has a minimum at the point U^*DU only if $(d\Phi)_{U^*, D^*}$ vanishes for all $dU^* \in \mathcal{M}_n$ and $dD^* \in \mathcal{D}_n$. Setting $dU^* = 0$, varying dD^* in \mathcal{D}_n , and using Lemma 4.7 we obtain (4.99) below.

Setting $dD^* = 0$ and varying dU^* we obtain (4.100) below. If we write down all the relevant constraints with these equations we obtain the following complete set of five simultaneous matrix equations for the unknown matrices U , D , and H .

$$(4.99) \quad dg[U(U^*DU - A)U^*] = 0 \quad ,$$

$$(4.100) \quad D^*U(U^*DU - A) + DU(U^*DU - A)^* + UH = 0 \quad ,$$

$$(4.101) \quad U^*U = I \quad ,$$

$$(4.102) \quad H = H^* \quad ,$$

$$(4.103) \quad \text{offdg}(D) = 0 \quad .$$

4.104 Theorem. Let $A \in \mathcal{M}_n$. If $N_0 = U_0^* D_0 U_0$ where $U_0 \in \mathcal{U}_n$ and $D_0 \in \mathcal{H}_n$ is ϵ -minimal for A , then there exists a matrix $H_0 \in \mathcal{M}_n$ such that (4.99) through (4.103) are satisfied when U, D , and H are replaced by U_0, D_0 , and H_0 respectively.

Proof. The proof of Theorem 4.104 follows along the same lines as the proof of Theorem 4.73.

Remark. We can use Theorem 4.104 to derive again the necessary conditions (4.12) and (4.13). Let $N_0 = U_0^* D_0 U_0$ be ϵ -minimal as in Theorem 4.104 and let H_0 be a hermitian matrix of Lagrange multipliers which corresponds to it. Then using the conclusion of Theorem 4.104 we can immediately combine (4.99) and (4.101) to obtain $dg(D_0 - U_0 A U_0^*) = 0$ which is the same as (4.13). Premultiplying (4.100) by U_0^* and again

using (4.101) we find, upon setting $N_O = U_O^* D_O U_O$, that

$$N_O^* N_O + N_O N_O^* + H = N_O^* A + N_O A^* .$$

The last equation shows clearly that $N_O^* A + N_O A^*$ is hermitian so, by Theorem 4.34, we have again derived the necessary condition (4.12).

Still another approach to the use of Lagrange multipliers is suggested by the representation of normal matrices given by Theorem 1.57 (a). Here we set

$$(4.105) \quad N = H + iK .$$

By Theorem 1.57, (4.105) runs through \mathcal{N} as H and K run through the set of all pairs of commuting hermitian matrices. Thus we try to minimize $\epsilon^2(H + iK - A)$ subject to the constraints

$$(4.106) \quad H = H^* , \quad K = K^* , \quad \text{and} \quad HK = KH .$$

We introduce matrices Z, Π , and T of complex Lagrange multipliers which correspond respectively to the matrix constraints (4.106). Without going through the computational details we simply give here the results of the differentiations, which again turn out to be matrix equations:

$$(4.107) \quad H + iK - A + TK^* - K^*T + (Z - Z^*) = 0 ,$$

$$(4.108) \quad -i(H + iK - A) + H^*T - TH^* + (\Pi - \Pi^*) = 0 .$$

Here, by analogy with what happened in the previous derivations, it is not the matrices Z and Π which have to be determined but the

skew-hermitian matrices $S = Z - Z^*$ and $P = \Pi - \Pi^*$ instead. Note, however, that no combinations of elements of T occurs so that every element of T is an "actual Lagrange multiplier". This was to be expected, since the equation $HK = KH$ involves no duplications.

If we rewrite (4.107) and (4.108), replacing $Z - Z^*$ and $\Pi - \Pi^*$ by S and P respectively, and write down all the relevant auxiliary equations we obtain the following set of simultaneous matrix equations for the five unknown matrices $H, K, T, S,$ and P .

$$(4.109) \quad H + iK - A + TK^* - K^*T + S = 0 ,$$

$$(4.110) \quad -iH + K + iA + H^*T - TH^* + P = 0 ,$$

$$(4.111) \quad H = H^* ,$$

$$(4.112) \quad K = K^* ,$$

$$(4.113) \quad HK = KH ,$$

$$(4.114) \quad S + S^* = 0 ,$$

$$(4.115) \quad P + P^* = 0 .$$

As in the previous constrained minimization problems (cf. Theorems 4.73 and 4.104) we can prove

4.116 Theorem. Let $A \in \mathcal{M}_n$. If $N_0 = H_0 + iK_0$ where $H_0, K_0 \in \mathcal{H}_n$ and $H_0K_0 = K_0H_0$ is ϵ -minimal for A then there exist matrices $T_0, S_0, P_0 \in \mathcal{M}_n$ such that (4.109) through (4.115) are satisfied when $H, K, T, S,$ and P are replaced by $H_0, K_0, T_0, S_0,$ and P_0 respectively.

CHAPTER 5

CHARACTERIZATION OF ϵ -MINIMAL MATRICES

We begin the present chapter by proving the following simple

5.1 Lemma. Let v denote any norm on \mathcal{M} . We have

$$\begin{aligned}
 d_v(A) &= \inf_{\substack{U \in \mathcal{U} \\ D \in \mathcal{D}}} v(A - U^* D U) = \inf_{U \in \mathcal{U}} \left(\inf_{D \in \mathcal{D}} v(A - U^* D U) \right) \\
 &= \inf_{D \in \mathcal{D}} \left(\inf_{U \in \mathcal{U}} v(A - U^* D U) \right) .
 \end{aligned}
 \tag{5.2}$$

Proof. Let

$$\begin{aligned}
 d' &= \inf_{U \in \mathcal{U}} \left(\inf_{D \in \mathcal{D}} v(A - U^* D U) \right) , \\
 \alpha(U) &= \inf_{D \in \mathcal{D}} v(A - U^* D U) .
 \end{aligned}$$

We want to show that $d' = d_v(A)$. Clearly $d' \leq \alpha(U)$ and

$$d_v(A) \leq \alpha(U)
 \tag{5.3}$$

hold for all $U \in \mathcal{U}$. Taking the infimum over \mathcal{U} on the right side of (5.3) we obtain $d_v(A) \leq d'$. Suppose $d_v(A) < d'$; then there is a pair of matrices $U_0 \in \mathcal{U}$, $D_0 \in \mathcal{D}$ such that $v(A - U_0^* D_0 U_0) < d'$. But this implies $\alpha(U_0) < d'$ which contradicts the inequality $d' \leq \alpha(U_0)$. Thus $d_v(A) = d'$ which proves the second equality in (5.2). A similar argument will establish the last equality in (5.2).

By very much the same arguments used in proving Theorem 1.78, it can be shown that any of the infima over \mathcal{D} in (5.2) are assumed for some $D \in \mathcal{D}$. Furthermore, since v is continuous and \mathcal{U} is compact in the norm topology of \mathcal{M} , all infima over \mathcal{U} are likewise assumed. Thus, in (5.2), we can replace "inf" by "Min" in every instance. We shall use this fact in proving

5.4 Theorem. Let $A \in \mathcal{M}$ and let ϵ denote the Euclidean norm (1.16) on \mathcal{M} . Then

$$(5.5) \quad d_{\epsilon}^2(A) = \epsilon^2(A) - \max_{U \in \mathcal{U}} \epsilon^2(\text{dg}(UAU^*)) .$$

Proof. Since ϵ is unitarily invariant we have

$$(5.6) \quad \epsilon^2(A - U^*DU) = \epsilon^2(UAU^* - D) .$$

From Definition 1.23 and (1.16) we can easily deduce the following equalities which hold for any $M \in \mathcal{M}$:

$$(5.7) \quad \min_{D \in \mathcal{D}} \epsilon^2(M - D) = \epsilon^2(\text{offdg}(M)) ,$$

$$(5.8) \quad \epsilon^2(M) = \epsilon^2(\text{offdg}(M) + \text{dg}(M)) = \epsilon^2(\text{offdg}(M)) + \epsilon^2(\text{dg}(M))$$

where, in (5.7), the minimum is assumed if and only if $D = \text{dg}(M)$. Thus, for every $U \in \mathcal{U}$, we have

$$(5.9) \quad \min_{D \in \mathcal{D}} \epsilon^2(UAU^* - D) = \epsilon^2(UAU^* - \text{dg}(UAU^*))$$

whence, using (5.2),

$$(5.10) \quad d_{\epsilon}^2(A) = \min_{U \in \mathcal{U}} \epsilon^2(UAU^* - dg(UAU^*)) .$$

Noting that $\epsilon^2(UAU^*) = \epsilon^2(A)$ and combining (5.8) and (5.10) we obtain

$$\begin{aligned} d_{\epsilon}^2(A) &= \min_{U \in \mathcal{U}} [\epsilon^2(UAU^*) - \epsilon^2(dg(UAU^*))] \\ &= \min_{U \in \mathcal{U}} [\epsilon^2(A) - \epsilon^2(dg(UAU^*))] \\ &= \epsilon^2(A) - \max_{U \in \mathcal{U}} \epsilon^2(dg(UAU^*)) \end{aligned}$$

which proves (5.5).

5.11 Maximum Problem. Let $A \in \mathcal{M}_n$. Find a $U_0 \in \mathcal{U}_n$ such that

$$(5.12) \quad \epsilon^2(dg(U_0AU_0^*)) = \max_{U \in \mathcal{U}_n} \epsilon^2(dg(UAU^*)) .$$

5.13 Theorem (Characterization Theorem). Let $A \in \mathcal{M}_n$ and let $N_0 = U_0^* D_0 U_0$ where $U_0 \in \mathcal{U}_n$ and $D_0 \in \mathcal{A}_n$. Then N_0 is ϵ -minimal for A if and only if U_0 satisfies (5.12) (i.e., U_0 solves the Maximum Problem 5.11) and

$$(5.14) \quad D_0 = dg(U_0AU_0^*) .$$

Proof. Let U_0 be any unitary matrix which satisfies (5.12) and let D_0 be given by (5.14). Then, from (5.5), we have

$$\begin{aligned}
d_{\epsilon}^2(A) &= \epsilon^2(A) - \epsilon^2(\text{dg}(U_0 A U_0^*)) \\
&= \epsilon^2(U_0 A U_0^*) - \epsilon^2(\text{dg}(U_0 A U_0^*)) \\
&= \epsilon^2(\text{offdg}(U_0 A U_0^*)) \\
&= \epsilon^2(U_0 A U_0^* - \text{dg}(U_0 A U_0^*)) = \epsilon^2(U_0 A U_0^* - D_0) \\
&= \epsilon^2(A - U_0^* D_0 U_0) = \epsilon^2(A - N_0)
\end{aligned}$$

so that N_0 is ϵ -minimal for A . This proves the sufficiency of the conditions (5.12) and (5.14). For their necessity we first note from Theorem 4.10 that the necessity of (5.14) has already been established. Thus, if N_0 is ϵ -minimal for A , we have

$$\begin{aligned}
d_{\epsilon}^2(A) &= \epsilon^2(A - N_0) = \epsilon^2(U_0 A U_0^* - D_0) \\
&= \epsilon^2(U_0 A U_0^* - \text{dg}(U_0 A U_0^*)) \\
&= \epsilon^2(\text{offdg}(U_0 A U_0^*)) \\
&= \epsilon^2(U_0 A U_0^*) - \epsilon^2(\text{dg}(U_0 A U_0^*)) \\
&= \epsilon^2(A) - \epsilon^2(\text{dg}(U_0 A U_0^*)) .
\end{aligned}$$

Comparing the last equation with (5.5) we see that U_0 satisfies (5.12). This completes the proof of Theorem 5.13.

An immediate corollary of Theorem 5.13 is the following result which indicates the importance of the Maximum Problem 5.11 in so far as ϵ -minimal matrices are concerned.

5.15 Theorem. Let U_0 be a solution of the Maximum Problem 5.11. Then

$$(5.16) \quad U_0^* \operatorname{dg}(U_0 A U_0^*) U_0$$

is ϵ -minimal for A .

Let u and v denote complex column n -vectors. The (Euclidean) inner product of u and v is defined by

$$(u, v) = v^* u$$

where v^* denotes the conjugate transpose of v . We call u and v orthogonal if $(u, v) = 0$. A set of vectors u_1, u_2, \dots, u_k ($k \leq n$) is called orthonormal if

$$(u_i, u_j) = \delta_{ij} = \begin{cases} 1 & , \quad i = j \\ 0 & , \quad i \neq j \end{cases} .$$

5.17 Maximum Problem. Let $A \in \mathcal{M}_n$. Find an orthonormal set of column n -vectors u_1, u_2, \dots, u_n such that

$$\sum_{i=1}^n |(Au_i, u_i)|^2 = \operatorname{Max} \sum_{i=1}^n |(Av_i, v_i)|^2$$

where the maximum is taken over all orthonormal sets v_1, v_2, \dots, v_n of column n -vectors.

Suppose $U \in \mathcal{M}_n$ and let v_1, \dots, v_n denote the columns of the matrix U^* . Clearly U is unitary if and only if v_1, \dots, v_n is orthonormal. Since $v^* A u = (Au, v)$ for all vectors u, v , we see from

(1.24) that

$$\text{dg}(\text{UAU}^*) = \text{diag}((Av_1, v_1), \dots, (Av_n, v_n))$$

whence, by (1.16),

$$(5.19) \quad \epsilon^2(\text{dg}(\text{UAU}^*)) = \sum_{i=1}^n |(Av_i, v_i)|^2.$$

From (5.12), (5.18), and (5.19) one can immediately deduce

5.20 Theorem. The Maximum Problems 5.11 and 5.17 are equivalent in the sense that a solution of one yields immediately a solution of the other.

From the decomposition $N_0 = U_0^* D_0 U_0$ one finds that the columns of U_0^* are eigenvectors of N_0 . Therefore, by Theorem 5.13, we have

5.21 Theorem. Let $A \in \mathcal{M}_n$ and let $N_0 = U_0^* D_0 U_0$ where $U_0 \in \mathcal{U}_n$ and $D_0 \in \mathcal{D}_n$. Then N_0 is ϵ -minimal for A if and only if it has an orthonormal set of eigenvectors which solves the Maximum Problem 5.17.

The next result exhibits the relationship between the maximum (5.12) and distance formulas of the type (2.32).

5.22 Lemma. Let $A \in \mathcal{M}_n$, $A \notin \mathcal{M}_n$ and let $p_{\epsilon, n}(A)$ be defined by Definition 2.13 in the case $v = \epsilon$. Then

$$(5.23) \quad \max_{U \in \mathcal{U}_n} \epsilon^2(\text{dg}(\text{UAU}^*)) = (1 - p_{\epsilon, n}(A)) \epsilon^2(A) + p_{\epsilon, n}(A) \epsilon^2(\Omega(A)).$$

Proof. Since ϵ has property S (see Definition 2.1), (2.32) holds for $v = \epsilon$. That equation and (5.5) yield immediately (5.23).

Remark. Theorem 5.4 opens up a new avenue for obtaining estimates of $d_\epsilon(A)$. Thus, a lower bound for the maximum (5.12) will yield an

upper bound for $d_\epsilon(A)$ and an upper bound for (5.12) will furnish a lower bound for $d_\epsilon(A)$. The next theorem is a result of this type.

5.24 Theorem. Let $A \in \mathcal{M}_n$ and let $A = WH$ where H is positive semidefinite hermitian and $W \in \mathcal{U}_n$. Furthermore let y_1, y_2, \dots, y_n be an orthonormal set of eigenvectors of H , i.e., $Hy_1 = \alpha_1 y_1$ ($1 \leq i \leq n$) where α_1 are the singular values of A and $(y_i, y_j) = \delta_{ij}$ ($1 \leq i, j \leq n$). Then

$$(5.25) \quad d_\epsilon^2(A) \leq \sum_{i=1}^n \alpha_i^2 (1 - |(Wy_1, y_1)|^2) \quad .$$

Proof. We have

$$(Ay_1, y_1) = (WHy_1, y_1) = \alpha_1 (Wy_1, y_1) \quad \text{for } i = 1, 2, \dots, n \quad .$$

Since y_1, \dots, y_n is an orthonormal set, we see from (5.18) and (5.19) that

$$(5.26) \quad \sum_{i=1}^n |(Ay_1, y_1)|^2 = \sum_{i=1}^n \alpha_i^2 |(Wy_1, y_1)|^2$$

is a lower bound for the maximum (5.12). As was observed in Section 1.1, $\epsilon(A)$ coincides with (1.15) for $p = 2$; consequently

$$(5.27) \quad \epsilon^2(A) = \sum_{i=1}^n \alpha_i^2 \quad .$$

Therefore, from (5.5), (5.26), and (5.27) we obtain

$$d_{\epsilon}^2(A) \leq \sum_{i=1}^n \alpha_i^2 - \sum_{i=1}^n \alpha_i^2 |(wy_i, y_i)|^2$$

which is the same as (5.25).

The next result provides a slight simplification of the problem of finding a solution to the Maximum Problem 5.11.

5.28 Lemma. Let \mathcal{U}_n^+ denote the set of all unitary matrices of order n which have nonnegative diagonal elements. Then, for any $A \in \mathcal{M}_n$, we have

$$(5.29) \quad \max_{U \in \mathcal{U}_n} \epsilon^2(\text{dg}(UAU^*)) = \max_{U \in \mathcal{U}_n^+} \epsilon^2(\text{dg}(UAU^*)) .$$

Proof. Let U be any unitary matrix of order n and let $u_{kk} = r_k \exp(i\theta_k)$ be its k -th diagonal element where $r_k \geq 0$ and θ_k is real. By factoring out $\exp(i\theta_k)$ from the k -th row of U ($k = 1, 2, \dots, n$) we can write $U = \Lambda V$ where $V \in \mathcal{U}_n^+$ and $\Lambda = \text{diag}(\exp(i\theta_1), \dots, \exp(i\theta_n))$ is a diagonal unitary matrix, i.e., an element of $\mathcal{N}_n \cap \mathcal{U}_n$. Noting that $\text{dg}(\Lambda M \Lambda^*) = \text{dg}(M)$ holds for all $M \in \mathcal{M}_n$ and for all $\Lambda \in \mathcal{N}_n \cap \mathcal{U}_n$ we find that

$$\epsilon^2(\text{dg}(UAU^*)) = \epsilon^2(\text{dg}(\Lambda V \Lambda^*)) = \epsilon^2(\text{dg}(VAV^*)) .$$

The last equation shows that, in computing the maximum (5.12), it suffices to consider only unitary matrices in \mathcal{U}_n^+ . The proof of Lemma 5.28 is now complete.

CHAPTER 6

ϵ -MINIMAL MATRICES OF ORDER 2

6.1 Preliminaries Concerning 2 by 2 Matrices.

6.1 Lemma. Let $N \in \mathcal{M}_2$ be given by

$$(6.2) \quad N = \begin{pmatrix} n_1 & n_2 \\ n_3 & n_4 \end{pmatrix}.$$

Then N is normal if and only if

$$(6.3) \quad |n_2|^2 = |n_3|^2$$

and

$$(6.4) \quad (n_1 - n_4)\bar{n}_3 = n_2(\bar{n}_1 - \bar{n}_4) \quad .$$

Proof. We have

$$NN^* = \begin{pmatrix} |n_1|^2 + |n_2|^2 & n_1\bar{n}_3 + n_2\bar{n}_4 \\ \bar{n}_1n_3 + \bar{n}_2n_4 & |n_3|^2 + |n_4|^2 \end{pmatrix},$$

$$N^*N = \begin{pmatrix} |n_1|^2 + |n_3|^2 & \bar{n}_1n_2 + \bar{n}_3n_4 \\ n_1\bar{n}_2 + n_3\bar{n}_4 & |n_2|^2 + |n_4|^2 \end{pmatrix}$$

whence $NN^* - N^*N = 0$ if and only if (6.3) and (6.4) hold.

6.5 Lemma. Let $A \in \mathcal{M}_2$. Let $\Omega(A)$ and $M(A)$ be defined as in Definitions 1.17 and 3.8 respectively. Then

$$(6.6) \quad \epsilon^2(A - X) = \frac{1}{2} (\epsilon^2(A) - \epsilon^2(\Omega(A))) \quad \text{for all } X \in M(A) .$$

Proof. Let the eigenvalues of A be denoted by λ_1, λ_2 . Then the eigenvalues of A^2 are λ_1^2 and λ_2^2 so we have $\text{tr}(A) = \lambda_1 + \lambda_2$ and $\text{tr}(A^2) = \lambda_1^2 + \lambda_2^2$. Thus

$$\begin{aligned} \frac{1}{2} |\text{tr}(A)|^2 &= \frac{1}{2} (\lambda_1 + \lambda_2)(\bar{\lambda}_1 + \bar{\lambda}_2) \\ &= \frac{1}{2} (|\lambda_1|^2 + |\lambda_2|^2 + \lambda_1 \bar{\lambda}_2 + \bar{\lambda}_1 \lambda_2) \end{aligned}$$

or

$$(6.7) \quad \frac{1}{2} |\text{tr}(A)|^2 = \frac{1}{2} \epsilon^2(\Omega(A)) + \text{Re}(\lambda_1 \bar{\lambda}_2) .$$

Furthermore

$$\begin{aligned} \text{tr}(A^2) - \frac{1}{2} (\text{tr}(A))^2 &= \lambda_1^2 + \lambda_2^2 - \frac{1}{2} (\lambda_1 + \lambda_2)^2 \\ &= \lambda_1^2 + \lambda_2^2 - \frac{1}{2} (\lambda_1^2 + \lambda_2^2 + 2\lambda_1 \lambda_2) \end{aligned}$$

or

$$(6.8) \quad \text{tr}(A^2) - \frac{1}{2} (\text{tr}(A))^2 = \frac{1}{2} (\lambda_1 - \lambda_2)^2 .$$

From the last equation we obtain

$$\begin{aligned}
|\operatorname{tr}(A^2) - \frac{1}{2} (\operatorname{tr}(A))^2| &= \frac{1}{2} |\lambda_1 - \lambda_2|^2 \\
&= \frac{1}{2} (\lambda_1 - \lambda_2)(\bar{\lambda}_1 - \bar{\lambda}_2) \\
&= \frac{1}{2} (|\lambda_1|^2 + |\lambda_2|^2 - \lambda_1 \bar{\lambda}_2 - \bar{\lambda}_1 \lambda_2)
\end{aligned}$$

or

$$(6.9) \quad |\operatorname{tr}(A^2) - \frac{1}{2} (\operatorname{tr}(A))^2| = \frac{1}{2} \epsilon^2(\Omega(A)) - \operatorname{Re}(\lambda_1 \bar{\lambda}_2) .$$

Combining (6.7) and (6.9) we find that

$$(6.10) \quad \frac{|\operatorname{tr}(A)|^2}{2} + |\operatorname{tr}(A^2) - \frac{(\operatorname{tr}(A))^2}{2}| = \epsilon^2(\Omega(A)) .$$

The equation (6.6) now follows immediately from (6.10) and Theorem 3.11.

The next lemma shows how the simplification provided by Lemma

5.28 works out in the case $n = 2$.

6.11 Lemma. Let $A \in \mathcal{M}_2$. Let $\mu \in \underline{\mathbb{C}}$ and let

$$(6.12) \quad W(\mu) = \frac{1}{\sqrt{1 + \mu\bar{\mu}}} \begin{pmatrix} 1 & -\bar{\mu} \\ \mu & 1 \end{pmatrix} .$$

Then $W(\mu) \in \mathcal{U}_2$ for all $\mu \in \underline{\mathbb{C}}$ and

$$(6.13) \quad \max_{U \in \mathcal{U}_2} \epsilon^2(\operatorname{dg}(UAU^*)) = \max_{\mu \in \underline{\mathbb{C}}} \epsilon^2(\operatorname{dg}(W(\mu)AW^*(\mu))) .$$

Proof. It was shown in [3] (see also [17] for another parametrization of 2 by 2 unitary matrices) that every $U \in \mathcal{U}_2$ can be obtained from the formula

$$(6.14) \quad U = \begin{pmatrix} e^{i\alpha} \cos \theta & -e^{i\beta} \sin \theta \\ e^{i\gamma} \sin \theta & e^{i\delta} \cos \theta \end{pmatrix}$$

where $\theta \in \underline{\mathbb{R}}$ and where $\alpha, \beta, \gamma, \delta$ are real numbers satisfying

$$(6.15) \quad \alpha - \beta - \gamma + \delta \equiv 0 \pmod{2\pi}.$$

By factoring out $\exp(i\alpha)$ and $\exp(i\delta)$ respectively from the first and second rows of (6.14) and using (6.15) we see that (6.14) can be written in the form

$$(6.16) \quad U = \begin{pmatrix} e^{i\alpha} & 0 \\ 0 & e^{i\delta} \end{pmatrix} \begin{pmatrix} \cos \theta & -e^{-i(\gamma-\delta)} \sin \theta \\ e^{i(\gamma-\delta)} \sin \theta & \cos \theta \end{pmatrix} = \Lambda W(\mu)$$

where $\Lambda = \text{diag}(\exp(i\alpha), \exp(i\delta))$ and where $W(\mu)$ is given by (6.12) with $\mu = \exp[i(\gamma-\delta)]\tan \theta$. Since the real quantities $\theta, \alpha, \beta, \gamma, \delta$ in (6.14) are arbitrary except for the constraint (6.15) which was used in obtaining (and which is automatically satisfied by) the product $\Lambda W(\mu)$ in (6.16), one sees easily that (6.16) furnishes a parametrization of all $U \in \mathcal{U}_2$ in the following manner. If we consider θ, α, δ , and $\psi = \gamma - \delta$ as independent real variables, or, alternatively, if we consider α and δ as independent real variables and $\mu = \exp(i\psi)\tan \theta$ as an independent complex variable, then, as μ runs through $\underline{\mathbb{C}}$ and as α, δ run through $\underline{\mathbb{R}}$, $U = \Lambda W(\mu)$ runs through \mathcal{U}_2 . As was shown in the proof of Lemma 5.28, the factor Λ in (6.16) can be disregarded

in computing $\epsilon^2(\text{dg}(\text{UAU}^*))$. Therefore only the complex parameter μ matters and we have (6.13).

6.17 Lemma. Let $A \in \mathcal{M}_2$, $A \notin \mathcal{N}_2$. Then A is unitarily similar either to a matrix of the form

$$(6.18) \quad A_1 = \begin{pmatrix} a_1 & a_2 \\ 0 & a_1 \end{pmatrix} \quad \text{with } a_2 \neq 0$$

or to a matrix of the form

$$(6.19) \quad A_2 = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_1 \end{pmatrix} \quad \text{with } a_2 \neq 0, a_3 \neq 0.$$

Proof. By Schur's theorem [22] (or see [16] p. 307) every $A \in \mathcal{M}_2$ is unitarily similar to a triangular matrix:

$$(6.20) \quad VAV^* = B = \begin{pmatrix} b_1 & b_2 \\ 0 & b_4 \end{pmatrix} \quad (V \in \mathcal{U}_2);$$

and, if $A \notin \mathcal{N}_2$, we have $b_2 \neq 0$. If $b_1 = b_4$, i.e., if the two eigenvalues of A are equal, then (6.20) is already of the form (6.18). It remains only to show that, if $b_1 \neq b_4$, A is unitarily similar to a matrix of the form (6.19). To prove the latter statement we consider

$$(6.21) \quad W(\mu) B W^*(\mu) = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix}$$

where $W(\mu)$ is given by (6.12) and where $\mu \in \mathbb{C}$. We have

$$a_1 - a_4 = (1 + \mu\bar{\mu})^{-1} [(b_1 - b_4)(1 - \mu\bar{\mu}) - 2\mu b_2]$$

hence $a_1 = a_4$ if and only if $2\mu b_2 = (b_1 - b_4)(1 - \mu\bar{\mu})$ or

$$2|\mu|e^{i \arg \mu} |b_2|e^{i \arg b_2} = |b_1 - b_4|e^{i \arg(b_1 - b_4)}(1 - |\mu|^2) .$$

The last equation will be satisfied if

$$(6.22) \quad \arg \mu \equiv \arg(b_1 - b_4) - \arg b_2 \pmod{2\pi} \quad \text{and}$$

$$(6.23) \quad |\mu| = \frac{-|b_2| + \sqrt{|b_2|^2 + |b_1 - b_4|^2}}{|b_1 - b_4|} .$$

Thus, if μ is chosen to satisfy (6.22) and (6.23), $a_1 = a_4$ holds in (6.21). Furthermore, for this choice of μ , neither a_2 nor a_3 in (6.21) can vanish since otherwise the eigenvalues of A , namely b_1 and b_4 , would have to be equal. This completes the proof of Lemma 6.17.

6.2 Determination of all 2 by 2 ϵ -Minimal Matrices.

6.24 Theorem. Let $A \in \mathcal{M}_2$ and let $M(A)$ be defined as in Definition 3.8. Then the set $M(A)$ represents the totality of all ϵ -minimal matrices for A .

Proof. Suppose first that $A \in \mathcal{M}_2$. Then by Lemma 2.4 $\epsilon^2(A) = \epsilon^2(\Omega(A))$ whence Lemma 6.5 implies

$$(6.25) \quad \epsilon^2(A - X) = 0 \quad \text{for all } X \in M(A) \text{ where } A \in \mathcal{M}_2 .$$

The meaning of the last equation is that every matrix in the set (3.10) is equal to A :

$$(6.26) \quad M(A) = \{A\} \quad (A \in \mathcal{N}_2)$$

that is, if $A \in \mathcal{N}_2$, then the set (3.10) contains exactly one matrix, namely A itself! This proves Theorem 6.24 in the case $A \in \mathcal{N}_2$.

We assume henceforth in this proof that $A \notin \mathcal{N}_2$. The principal tools we shall employ in the case $A \notin \mathcal{N}_2$ are Theorem 5.13 and Lemma 6.11. In order to keep the computations manageable we shall employ a change of coordinates defined by a unitary matrix which transforms A into one of the forms (6.18) or (6.19). This procedure is justified in the next paragraph. The general outline of our proof is as follows. We shall first show that every matrix in the set $M(A)$ is ϵ -minimal. Following this we show that there are no other ϵ -minimal matrices.

From (3.9) and the invariance of the trace under unitary similarity we have

$$(6.27) \quad \zeta(A) = \zeta(UAU^*) \quad (A \in \mathcal{M}_n, U \in \mathcal{U}_n).$$

In like manner we obtain the identity

$$(6.28) \quad P_\zeta(UAU^*) = UP_\zeta(A)U^* \quad (A \in \mathcal{M}_n, U \in \mathcal{U}_n)$$

for all $\zeta \in \underline{\mathbb{C}}$ where P_ζ is defined by (4.49). In view of (3.10), (6.28), and Theorem 2.27 it will suffice to prove the conclusion of Theorem 6.24 for some unitary transform of A . As mentioned in the

preceding paragraph we shall use the forms (6.18) and (6.19) for this purpose.

We first observe from (6.8) that the case in which A is unitarily similar to (6.18) corresponds to the ambiguous case of (3.9). Similarly the case in which A is similar to (6.19) corresponds to the non-ambiguous case of (3.9) where $\zeta(A)$ contains exactly one number. As stated in the note following Definition 3.8, we use the notation $N_\theta(A_1)$ to denote one of the matrices in $M(A_1)$ and $N(A_2)$ to denote the (single) matrix in $M(A_2)$. Straightforward calculations based on Definition 3.8 yield

$$(6.29) \quad N_\theta(A_1) = \begin{pmatrix} a_1 & \frac{1}{2} a_2 \\ \frac{1}{2} \bar{a}_2 e^{i\theta} & a_1 \end{pmatrix} \quad (\theta \in \underline{\mathbb{R}}) ,$$

$$(6.30) \quad N(A_2) = \begin{pmatrix} a_1 & \frac{1}{2} (a_2 + \zeta \bar{a}_3) \\ \frac{1}{2} (a_3 + \zeta \bar{a}_2) & a_1 \end{pmatrix} \quad \text{where } \zeta = \frac{a_2 a_3}{|a_2 a_3|} .$$

In order to prove that $N_\theta(A_1)$ is ϵ -minimal for A_1 for all $\theta \in \underline{\mathbb{R}}$ and that $N(A_2)$ is ϵ -minimal for A_2 , we shall diagonalize (6.29) and (6.30) by unitary matrices of the form (6.12) and verify the sufficient conditions of Theorem 5.13.

Let

$$M = \begin{pmatrix} m_1 & m_2 \\ m_3 & m_4 \end{pmatrix}$$

and let $W(\mu)$ be given by (6.12); then

$$(6.31) \quad W(\mu)MW^*(\mu) = \frac{1}{1 + \mu\bar{\mu}} \begin{pmatrix} m_1 - \mu m_2 - \bar{\mu} m_3 + \mu\bar{\mu} m_4 & \bar{\mu} m_1 + m_2 - \bar{\mu}^2 m_3 - \bar{\mu} m_4 \\ \mu m_1 - \mu^2 m_2 + m_3 - \mu m_4 & \mu\bar{\mu} m_1 + \mu m_2 + \bar{\mu} m_3 + m_4 \end{pmatrix}.$$

In order to diagonalize $N_\theta(A_1)$ and $N(A_2)$ it will suffice, by Theorem 1.59, to triangularize them. Thus, from (6.29) and (6.31) we see that $W(\mu_{1\theta})N_\theta(A_1)W^*(\mu_{1\theta})$ is diagonal if and only if

$$(6.32) \quad \mu_{1\theta}^2 = \frac{\bar{a}_2 e^{i\theta}}{a_2}.$$

Similarly, $W(\mu_2)N(A_2)W^*(\mu_2)$ is diagonal if and only if

$$(6.33) \quad \mu_2^2 = \frac{a_3 + \bar{a}_2}{a_2 + \bar{a}_3} = \frac{a_3}{a_2} \cdot \frac{|a_2 a_3| + |a_2|^2}{|a_2 a_3| + |a_3|^2} = \frac{a_3 |a_2|}{a_2 |a_3|}.$$

For notational simplicity we shall usually omit the subscript θ from $\mu_{1\theta}$ and write only μ_1 . However, it should always be understood that μ_1 depends on the parameter θ . Clearly, each of equations (6.32) and (6.33) has two solutions and we are at liberty to choose either solution in diagonalizing $N_\theta(A_1)$ or $N(A_2)$. In what follows it will not matter which solution is chosen, so we use the symbols μ_1 and μ_2 to denote any solutions of (6.32) and (6.33). We note from (6.32) that

$$(6.34) \quad |\mu_{1\theta}| = 1,$$

$$(6.35) \quad \mu_{1\theta}^2 a_2^2 = |a_2|^2 e^{i\theta}.$$

From (6.33) we have

$$(6.36) \quad \mu_2^2 a_2 \bar{a}_3 = \bar{\mu}_2^2 \bar{a}_2 a_3 = |a_2 a_3| > 0 ;$$

and, from (6.30), (6.33), and Lemma 6.1 (equation (6.3)), we obtain

$$(6.37) \quad |\mu_2| = 1 .$$

By straightforward calculations we find that

$$(6.38) \quad W(\mu_1) N_\theta(A_1) W^*(\mu_1) = \frac{1}{2} \begin{pmatrix} 2a_1 - \frac{1}{2} (\mu_1 a_2 + \bar{\mu}_1 \bar{a}_2 e^{1\theta}) & 0 \\ 0 & 2a_1 + \frac{1}{2} (\mu_1 a_2 + \bar{\mu}_1 \bar{a}_2 e^{1\theta}) \end{pmatrix} ,$$

$$(6.39) \quad W(\mu_2) N(A_2) W^*(\mu_2) = \frac{1}{2} \begin{pmatrix} 2a_1 - \frac{1}{2} [\mu_2 (a_2 + \zeta \bar{a}_3) + \bar{\mu}_2 (a_3 + \zeta \bar{a}_2)] & 0 \\ 0 & 2a_1 + \frac{1}{2} [\mu_2 (a_2 + \zeta \bar{a}_3) + \bar{\mu}_2 (a_3 + \zeta \bar{a}_2)] \end{pmatrix} ,$$

$$(6.40) \quad W(\mu_1) A_1 W^*(\mu_1) = \frac{1}{2} \begin{pmatrix} 2a_1 - \mu_1 a_2 & a_2 \\ -\mu_1^2 a_2 & 2a_1 + \mu_1 a_2 \end{pmatrix} ,$$

$$(6.41) \quad W(\mu_2) A_2 W^*(\mu_2) = \frac{1}{2} \begin{pmatrix} 2a_1 - (\mu_2 a_2 + \bar{\mu}_2 a_3) & a_2 - \bar{\mu}_2^2 a_3 \\ a_3 - \mu_2^2 a_2 & 2a_1 + (\mu_2 a_2 + \bar{\mu}_2 a_3) \end{pmatrix} .$$

From (6.32) and (6.34) we obtain $\mu_1 a_2 = \bar{\mu}_1 \bar{a}_2 \exp(i\theta)$ whence, using (6.38) and (6.40),

$$(6.42) \quad W(\mu_1) N_\theta(A_1) W^*(\mu_1) = \frac{1}{2} \text{diag}(2a_1 - \mu_1 a_2, 2a_1 + \mu_1 a_2) = \text{dg}(W(\mu_1) A_1 W^*(\mu_1)) .$$

Using (6.37) we can write

$$\mu_2(a_2 + \zeta \bar{a}_3) + \bar{\mu}_2(a_3 + \zeta \bar{a}_2) = \bar{\mu}_2[\mu_2^2 a_2 + \mu_2^2 \zeta \bar{a}_3 + a_3 + \zeta \bar{a}_2]$$

and, by (6.30) and (6.36), we have

$$\mu_2^2 \zeta \bar{a}_3 = \mu_2^2 a_2 \bar{a}_3 \frac{1}{|a_2 a_3|} a_3 = a_3 ,$$

$$\zeta \bar{a}_2 = a_2 \frac{\bar{a}_2 a_3}{|a_2 a_3|} = a_2 \frac{1}{\bar{a}_2} = \mu_2^2 a_2$$

so that

$$(6.43) \quad \mu_2(a_2 + \zeta \bar{a}_3) + \bar{\mu}_2(a_3 + \zeta \bar{a}_2) = 2\bar{\mu}_2(\mu_2^2 a_2 + a_3) = 2(\mu_2 a_2 + \bar{\mu}_2 a_3) .$$

Combining (6.39), (6.41) and (6.43) we obtain

$$(6.44) \quad \begin{aligned} W(\mu_2) N(A_2) W^*(\mu_2) &= \frac{1}{2} \text{diag}[2a_1 - (\mu_2 a_2 + \bar{\mu}_2 a_3), 2a_1 + (\mu_2 a_2 + \bar{\mu}_2 a_3)] \\ &= \text{dg}(W(\mu_2) A_2 W^*(\mu_2)) . \end{aligned}$$

Equations (6.42) and (6.44) merely express the fact that the decompositions of $N_\theta(A_1)$ and $N(A_2)$ implied in (6.38) and (6.39) satisfy the condition (5.14) of Theorem 5.13.

In order to show that every matrix in the set $M(A)$ is ϵ -minimal for A , it remains only to prove that $W(\mu_1)$ solves the Maximum Problem 5.11 for A_1 ($i = 1, 2$). By Lemma 6.11 it will suffice to prove that

$$(6.45) \quad \sigma_1(\lambda) \leq \epsilon^2 [\text{dg}(W(\mu_1)A_1W^*(\mu_1))] \quad (i = 1, 2)$$

hold for all $\lambda \in \underline{\mathbb{C}}$ where

$$(6.46) \quad \sigma_1(\lambda) = \epsilon^2 [\text{dg}(W(\lambda)W(\mu_1)A_1W^*(\mu_1)W^*(\lambda))] \quad (i = 1, 2) .$$

Straightforward calculations using (6.12), (6.32) through (6.37), (6.40), and (6.41) yield

$$(6.47) \quad \epsilon^2 [\text{dg}(W(\mu_1)A_1W^*(\mu_1))] = 2 |a_1|^2 + \frac{1}{2} |a_2|^2 ,$$

$$(6.48) \quad \epsilon^2 [\text{dg}(W(\mu_2)A_2W^*(\mu_2))] = 2 |a_1|^2 + \frac{1}{2} (|a_2| + |a_3|)^2 ,$$

$$(6.49) \quad \sigma_1(\lambda) = \frac{1}{\alpha^2} (|\alpha a_1 + \beta_1 a_2|^2 + |\alpha a_1 - \beta_1 a_2|^2) ,$$

$$(6.50) \quad \sigma_2(\lambda) = \frac{1}{\alpha^2} (|\alpha a_1 + \beta_2 a_2 + \bar{\beta}_2 a_3|^2 + |\alpha a_1 - \beta_2 a_2 - \bar{\beta}_2 a_3|^2)$$

where

$$(6.51) \quad \alpha = 2(1 + \lambda\bar{\lambda})$$

and

$$(6.52) \quad \beta_1 = \mu_1^2 \bar{\lambda} + \mu_1 \lambda \bar{\lambda} - \mu_1 - \lambda \quad (i = 1, 2) .$$

After further calculations we find that

$$(6.53) \quad \sigma_1(\lambda) = 2|a_1|^2 + \frac{2|\beta_1|^2}{\alpha^2} |a_2|^2 ,$$

$$(6.54) \quad \sigma_2(\lambda) = 2|a_1|^2 + \frac{2}{\alpha^2} |\beta_2 a_2 + \bar{\beta}_2 a_3|^2 ,$$

$$(6.55) \quad \beta_i = -\mu_i(\mu_i + \lambda)(\bar{\mu}_i - \bar{\lambda}) \quad (i = 1, 2) ,$$

$$(6.56) \quad |\beta_i| = |\mu_i^2 - \lambda^2| \quad (i = 1, 2) .$$

Rewriting (6.45) using (6.47) and (6.48) we obtain

$$(6.57) \quad \sigma_1(\lambda) \leq 2|a_1|^2 + \frac{1}{2} |a_2|^2 ,$$

$$(6.58) \quad \sigma_2(\lambda) \leq 2|a_1|^2 + \frac{1}{2} (|a_2| + |a_3|)^2 .$$

We wish to show that (6.57) and (6.58) hold for all $\lambda \in \underline{\mathbb{C}}$. One sees from (6.53) that (6.57) holds if and only if $2|\beta_1|^2/\alpha^2 \leq 1/2$ that is

$$(6.59) \quad \frac{|\beta_1|^2}{(1 + \lambda \bar{\lambda})^2} \leq 1 \quad \text{for all } \lambda \in \underline{\mathbb{C}} .$$

Furthermore equality holds in (6.57) if and only if equality holds in (6.59). From (6.34) and (6.37) we have $|\mu_i| = 1$ ($i = 1, 2$);

consequently

$$|\mu_1^2 - \lambda^2| \leq |\mu_1|^2 + |\lambda|^2 = 1 + \lambda\bar{\lambda} \quad (1 = 1, 2)$$

so, from (6.56) and the last inequality, we see that

$$(6.60) \quad \frac{|\beta_1|^2}{(1 + \lambda\bar{\lambda})^2} = \frac{|\mu_1^2 - \lambda^2|^2}{(1 + \lambda\bar{\lambda})^2} \leq \frac{(1 + \lambda\bar{\lambda})^2}{(1 + \lambda\bar{\lambda})^2} = 1 \quad (1 = 1, 2)$$

hold for all $\lambda \in \mathbb{C}$. Equality holds in (6.60) if and only if $-\frac{\lambda^2}{\mu_1^2}\lambda^2$ is real and nonnegative, i.e., $-\frac{\lambda^2}{\mu_1^2}\lambda^2 \geq 0$. By the triangle inequality

$$(6.61) \quad |\beta_2 a_2 + \bar{\beta}_2 a_3|^2 \leq |\beta_2|^2 (|a_2| + |a_3|)^2$$

with equality holding if and only if $\beta_2^2 a_2 \bar{a}_3 \geq 0$. From (6.36) and (6.55) we have

$$(6.62) \quad \beta_2^2 a_2 \bar{a}_3 = \mu_2^2 a_2 \bar{a}_3 (\mu_2 + \lambda)^2 (\bar{\mu}_2 - \bar{\lambda})^2 = |a_2 a_3| (\mu_2 + \lambda)^2 (\bar{\mu}_2 - \bar{\lambda})^2.$$

From (6.51), (6.54), and (6.61) we obtain

$$\sigma_2(\lambda) \leq 2|a_1|^2 + \frac{|\beta_2|^2}{(1 + \lambda\bar{\lambda})^2} \cdot \frac{1}{2} (|a_2| + |a_3|)^2$$

consequently, by (6.60) and (6.62), (6.58) is valid for all $\lambda \in \mathbb{C}$ with equality if and only if

$$(6.63) \quad (\mu_2 + \lambda)^2 (\bar{\mu}_2 - \bar{\lambda})^2 \geq 0$$

and

$$(6.64) \quad -\frac{2}{\mu_2} \lambda^2 \geq 0 .$$

Furthermore, from (6.60) and (6.59), we see that (6.57) holds for all $\lambda \in \underline{\mathbb{C}}$ with equality if and only if

$$(6.65) \quad -\frac{2}{\mu_1} \lambda^2 \geq 0 .$$

This completes the proof of the fact that every matrix in the set $M(A)$ is ϵ -minimal for A .

In order to determine all ϵ -minimal matrices for A it will suffice, by Theorem 2.27 and Lemma 6.17, to determine all ϵ -minimal matrices for $B_i = W(\mu_i)A_iW^*(\mu_i)$ ($i = 1, 2$). The following lemma will provide help in making that determination.

6.66 Lemma. Let $B \in \mathcal{M}_2$ and let $N_0 = U_0^* D_0 U_0$, where $U_0 \in \mathcal{U}_2$ and $D_0 \in \mathcal{A}_2$, be ϵ -minimal for B . Then there exists a $\lambda_0 \in \underline{\mathbb{C}}$ such that

$$N_0 = W^*(\lambda_0) \operatorname{dg}(W(\lambda_0)BW^*(\lambda_0))W(\lambda_0) .$$

Proof. By the parametric representation of \mathcal{U}_2 developed in the proof of Lemma 6.11 there exist $\Lambda_0 \in \mathcal{A}_2 \cap \mathcal{U}_2$ and $\lambda_0 \in \underline{\mathbb{C}}$ such that $U_0 = \Lambda_0 W(\lambda_0)$. Thus, recalling that $\operatorname{dg}(\Lambda_0 M \Lambda_0^*) = \operatorname{dg}(M)$ for any $M \in \mathcal{M}_2$, we have

$$\operatorname{dg}(U_0 B U_0^*) = \operatorname{dg}(W(\lambda_0) B W^*(\lambda_0)) ;$$

consequently by Theorem 5.13 and the fact that diagonal matrices commute

$$\begin{aligned} N_0 &= U_0^* \operatorname{dg}(W(\lambda_0) B W^*(\lambda_0)) U_0 \\ &= W^*(\lambda_0) \Lambda_0^* \operatorname{dg}(W(\lambda_0) B W^*(\lambda_0)) \Lambda_0 W(\lambda_0) \\ &= W^*(\lambda_0) \operatorname{dg}(W(\lambda_0) B W^*(\lambda_0)) W(\lambda_0) \end{aligned}$$

as desired.

Since $W(\mu_1)$ solves the Maximum Problem 5.11 for A_1 ($i = 1, 2$) we see from Lemma 6.66 that, if we determine all $\lambda_0 \in \underline{\mathbb{C}}$ such that

$$\epsilon^2(\operatorname{dg}(B_1)) = \epsilon^2(\operatorname{dg}(W(\lambda_0) B_1 W^*(\lambda_0))) ,$$

then we will have determined all ϵ -minimal matrices for B_1 ($i = 1, 2$).

By the definition (6.46) of $\sigma_1(\lambda)$ this amounts to determining all cases of equality in (6.57) and (6.58). We shall do this by using the next two lemmas.

6.67 Lemma. The inequalities

$$(6.68) \quad -\frac{2}{\mu_k} \lambda^2 \geq 0 \quad (k = 1, 2)$$

hold if and only if

$$(6.69) \quad \lambda = \rho i \mu_k \quad (k = 1, 2)$$

where ρ is real.

Proof. Let $\lambda = r \exp(i\varphi)$ where $r \geq 0$ and φ is real. Since $|\mu_k| = 1$ ($k = 1, 2$) we can set $\mu_k = \exp(i\varphi_k)$. Thus

$$(6.70) \quad -\frac{2}{\mu_k} \lambda^2 = -r^2 e^{i 2(\varphi - \varphi_k)} .$$

Obviously (6.68) holds if $r = 0$. If $r > 0$ then, from (6.70), we see that (6.68) holds if and only if $\exp[i 2(\varphi - \varphi_k)] = -1$ i.e., if and only if $\exp(i\varphi) = \pm i \mu_k$ or $\lambda = \pm r i \mu_k$. Setting $\rho = \pm r$ we see that (6.68) holds if and only if λ has the form (6.69) where ρ is any real number.

6.71 Lemma. The inequality (6.63) holds if and only if

$$(6.72) \quad \lambda = \sigma \mu_2$$

where σ is real.

Proof. Straightforward calculations using (6.37) yield

$$(6.73) \quad (\mu_2 + \lambda)^2 (\bar{\mu}_2 - \bar{\lambda})^2 = [1 - 4|\lambda|^2 + |\lambda|^4 + 2 \operatorname{Re}(\lambda^2 \bar{\mu}_2^2)] + [2(|\lambda|^2 - 1)(\mu_2 \bar{\lambda} - \bar{\mu}_2 \lambda)] .$$

The first term in square brackets on the right of (6.73) is real while the second such term is pure imaginary. Therefore in order that (6.63) hold it is necessary and sufficient that

$$(6.74) \quad (|\lambda|^2 - 1)(\mu_2 \bar{\lambda} - \bar{\mu}_2 \lambda) = 0$$

and

$$(6.75) \quad f(\lambda, \mu_2) = 1 - 4|\lambda|^2 + |\lambda|^4 + 2 \operatorname{Re}(\lambda^2 \bar{\mu}_2^2) \geq 0 .$$

Obviously (6.74) holds if and only if either $|\lambda| = 1$ or $\mu_2 \bar{\lambda} = \bar{\mu}_2 \lambda$.

If $|\lambda| = 1$ then

$$f(\lambda, \mu_2) = 2(\operatorname{Re}(\lambda^2 \bar{\mu}_2^2) - 1) \leq 0$$

whence (6.75) holds if and only if equality holds in the last inequality, i.e.,

$$(6.76) \quad \lambda = \pm \mu_2.$$

If $|\lambda| \neq 1$ then (6.74) holds if and only if $\mu_2 \bar{\lambda} = \bar{\mu}_2 \lambda$, i.e., $\lambda = \mu_2^2 \bar{\lambda}$. Let $\lambda = r \exp(i\varphi)$ where $r \geq 0$ and φ is real. We need consider only the case $r > 0$ since (6.63) holds if $\lambda = 0$. Thus $r \exp(i\varphi) = \mu_2^2 r \exp(-i\varphi)$ or $\exp(i2\varphi) = \mu_2^2$ or $\exp(i\varphi) = \pm \mu_2$. Consequently

$$(6.77) \quad \lambda = \pm r \mu_2.$$

Substituting this into $f(\lambda, \mu_2)$ defined in (6.75) we have

$$f(\lambda, \mu_2) = 1 - 4r^2 + r^4 + 2r^2 = r^4 - 2r^2 + 1 = (r^2 - 1)^2 \geq 0$$

for all r . Setting $\sigma = \pm r$ and combining (6.76) and (6.77) we see that (6.74) and (6.75) hold simultaneously if and only if λ has the form (6.72) where σ is any real number. This completes the proof of Lemma 6.71.

Consider first the case in which $M(A)$ contains exactly one matrix, i.e., the case in which A is unitarily similar to A_2 . Here we wish to determine all cases of equality in (6.58). We have previously

observed that equality holds in (6.58) if and only if (6.63) and (6.64) hold simultaneously. Therefore by Lemmas 6.67 and 6.71 we must have $\lambda = \sigma\mu_2$ and $\lambda = \rho i\mu_2$ where ρ and σ are real. But the last two equations represent straight lines in the complex plane which intersect only for $\rho = \sigma = 0$, i.e., only for $\lambda = 0$. Thus equality holds in (6.58) only for $\lambda = 0$. Since $W(0) = I$ we see from Theorem 5.13 and Lemma 6.66 that

$$N(A_2) = W^*(\mu_2) \operatorname{dg}(W(\mu_2)A_2W^*(\mu_2))W(\mu_2)$$

is the only ϵ -minimal matrix for A_2 . By Corollary 2.30 A has exactly one ϵ -minimal matrix, namely $N(A)$.

Now consider the case in which A is unitarily similar to A_1 . Here $M(A)$ contains infinitely many matrices and we wish to determine all cases of equality in (6.57). We observed previously that equality in (6.57) if and only if (6.65) holds. By Lemma 6.67 equality holds in (6.57) if and only if $\lambda = \rho i\mu_{1\theta}$ where ρ is real and $\mu_{1\theta}$ is defined by (6.32). Using (6.34) and (6.12) we obtain

$$\begin{aligned} W(\lambda)W(\mu_1) &= \frac{1}{\sqrt{2}\sqrt{1+\rho^2}} \begin{pmatrix} 1 & \rho i\bar{\mu}_1 \\ \rho i\mu_1 & 1 \end{pmatrix} \begin{pmatrix} 1 & -\bar{\mu}_1 \\ \mu_1 & 1 \end{pmatrix} \\ &= \frac{1}{\sqrt{2}\sqrt{1+\rho^2}} \begin{pmatrix} 1+\rho i & \bar{\mu}_1(-1+\rho i) \\ \mu_1(1+\rho i) & 1-\rho i \end{pmatrix} \\ &= \frac{1}{\sqrt{1+\rho^2}} \begin{pmatrix} 1+\rho i & 0 \\ 0 & 1-\rho i \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -\bar{\mu}_1 \frac{(1+\rho i)}{1-\rho i} \\ \mu_1 \frac{1+\rho i}{1-\rho i} & 1 \end{pmatrix} \\ &= \Lambda_{\rho} W(\mu_1 e^{i\theta} \rho^{1/2}) \end{aligned}$$

where

$$\Lambda_{\rho} = \begin{pmatrix} \frac{1 + \rho 1}{\sqrt{1 + \rho^2}} & 0 \\ 0 & \frac{1 - \rho 1}{\sqrt{1 + \rho^2}} \end{pmatrix} \in \mathcal{N}_2 \cap \mathcal{U}_2$$

and where

$$e^{1\theta_{\rho}/2} = \frac{1 + \rho 1}{1 - \rho 1}.$$

From (6.32) we have

$$\mu_{1\theta} e^{1\theta_{\rho}/2} = \sqrt{a_2/a_1} e^{1 \frac{\theta+\theta_{\rho}}{2}} = \mu_{1, \theta+\theta_{\rho}}.$$

Therefore by Lemma 6.66 every ϵ -minimal matrix for A_1 is of the form

$$\begin{aligned} & W^*(\mu_1) W^*(\lambda) \operatorname{dg}(W(\lambda) W(\mu_1) A_1 W^*(\mu_1) W^*(\lambda)) W(\lambda) W(\mu_1) \\ &= W^*(\mu_{1, \theta+\theta_{\rho}}) \Lambda_{\rho}^* \operatorname{dg}(\Lambda_{\rho} W(\mu_{1, \theta+\theta_{\rho}}) A_1 W^*(\mu_{1, \theta+\theta_{\rho}}) \Lambda_{\rho}^*) \Lambda_{\rho} W(\mu_{1, \theta+\theta_{\rho}}) \\ &= W^*(\mu_{1, \theta+\theta_{\rho}}) \operatorname{dg}(W(\mu_{1, \theta+\theta_{\rho}}) A_1 W^*(\mu_{1, \theta+\theta_{\rho}})) W(\mu_{1, \theta+\theta_{\rho}}) \\ &= N_{\theta+\theta_{\rho}}(A_1) \end{aligned}$$

which is just another one of the matrices in the set $M(A_1)$ no matter what the value of ρ is. Thus $M(A_1)$ contains all the ϵ -minimal matrices for A_1 . This completes the proof of Theorem 6.24.

Some immediate consequences of Theorem 6.24 are the following.

6.78 Corollary. If $A \in \mathcal{M}_2$ then the set (3.10) contains exactly one matrix, namely A .

6.79 Theorem. Let λ_1 and λ_2 denote the eigenvalues of $A \in \mathcal{M}_2$. If $A \notin \mathcal{M}_2$ then there is a unique ϵ -minimal matrix for A if and only if $\lambda_1 \neq \lambda_2$; if $\lambda_1 = \lambda_2$ there are infinitely many ϵ -minimal matrices for A .

Proof. By Theorem 6.24 A has a unique ϵ -minimal matrix if and only if $\text{tr}(A^2) - (1/2)(\text{tr}(A))^2 \neq 0$ (cf. Definition 3.8). By (6.8) this happens if and only if $\lambda_1 \neq \lambda_2$. Similarly, if $\lambda_1 = \lambda_2$, $M(A)$ contains an infinite number of matrices.

6.80 Theorem. Mirsky's Conjecture 1.35 is true for $v = \epsilon$ and $n = 2$.

Proof. This follows immediately from Theorem 6.24, Lemma 6.5, and the definition (1.34) of $d_\epsilon(A)$.

6.81 Theorem. The set $M(A)$ of (3.10) provides a complete solution to Problem 3.40 for $n = 2$ and $k = 1$.

Proof. Clearly $M(A) \subset \mathcal{L}_1(A) \subset \mathcal{M}_2$ so by Theorem 6.24 we have the desired conclusion.

For any subset \mathcal{S} of \mathcal{M}_n we denote the set of all real matrices in \mathcal{S} by $\mathcal{S}/\underline{\mathbb{R}}$ (read: \mathcal{S} restricted to $\underline{\mathbb{R}}$). Theorem 6.24 provides a complete solution to the distance problem (and associated minimum problem) of finding (1.26) where $v = \epsilon$, $A \in \mathcal{M}_2/\underline{\mathbb{R}}$ and where \mathcal{S} is replaced by $\mathcal{M}_2/\underline{\mathbb{R}}$. For, if $A \in \mathcal{M}_2/\underline{\mathbb{R}}$ has eigenvalues λ_1, λ_2 , then the set $M(A)$ contains all matrices in \mathcal{M}_2 which are ϵ -minimal for A . If $\lambda_1 \neq \lambda_2$ (cf. (6.8)) then by (3.9) $\xi(A)$ contains exactly one number and that is real; consequently by (3.10) $M(A)$ contains exactly one matrix and that is real. If $\lambda_1 = \lambda_2$ then $\xi(A)$ contains all

complex numbers of unit modulus and by (3.10) $N_\zeta(A)$ is real if and only if either $\zeta = +1$ or $\zeta = -1$. We have proved

6.82 Theorem. Let $A \in \mathcal{M}_2/\mathbb{R}$, $A \notin \mathcal{M}_2$ and let the eigenvalues of A be denoted by λ_1, λ_2 . If $\lambda_1 \neq \lambda_2$ there is a unique real ϵ -minimal matrix for A given by (3.10). If $\lambda_1 = \lambda_2$ there are exactly two real ϵ -minimal matrices for A , namely

$$\frac{1}{2} (A \pm A^*) + \frac{1}{4} \operatorname{tr}(A \mp A^*) I .$$

6.3 The Maximum Problem 5.11 in the Case $n = 2$.

6.83 Theorem. Let $A \in \mathcal{M}_2$. The identity matrix I solves the Maximum Problem 5.11 for A , i.e., A satisfies

$$(6.84) \quad \epsilon^2(\operatorname{dg}(A)) = \max_{U \in \mathcal{U}_2} \epsilon^2(\operatorname{dg}(UAU^*)) ,$$

if and only if

$$(6.85) \quad A + \zeta A^* \in \mathcal{O}_2$$

for some ζ in the set $\zeta(A)$ defined by (3.9).

Proof. If (6.85) holds for $\zeta \in \zeta(A)$ then by (3.10) $N_\zeta(A) \in \mathcal{O}_2$ whence I diagonalizes $N_\zeta(A)$. By Theorem 6.24 $N_\zeta(A)$ is ϵ -minimal so by Theorem 5.13 I solves the Maximum Problem 5.11. Suppose now that (6.84) holds. Then by Theorem 5.13 $\operatorname{dg}(A)$ is ϵ -minimal for A whence $\operatorname{dg}(A) \in M(A)$ by Theorem 6.24. Thus there is a $\zeta \in \zeta(A)$ such that

$N_{\zeta}(A) = dg(A)$, i.e., $N_{\zeta}(A) \in \mathcal{A}_2$, and this implies (6.85).

6.86 Theorem. Let $A \in \mathcal{M}_2$. U_0 solves the Maximum Problem 5.11 for A if and only if $U_0 \in \mathcal{U}_2$ and

$$(6.87) \quad U_0(A + \zeta A^*)U_0^* \in \mathcal{A}_2$$

for some ζ in the set $\zeta(A)$ defined by (3.9).

Proof. Every $V \in \mathcal{U}_2$ can be written uniquely as $V = UU_0$ and V runs through \mathcal{U}_2 if and only if U runs through \mathcal{U}_2 . Thus

$$(6.88) \quad \epsilon^2(dg(U_0AU_0^*)) = \max_{U \in \mathcal{U}_2} \epsilon^2(dg(U(U_0AU_0^*)U^*))$$

whence $U_0AU_0^*$ satisfies (6.88) if and only if U_0 solves the Maximum Problem 5.11. By virtue of (6.27) Theorem 6.86 follows immediately from Theorem 6.83.

Remark. The significance of Theorem 6.86 lies in the fact that it characterizes any solution of the Maximum Problem 5.11 (for $n = 2$) in terms of an algebraic condition which is very easy to check.

Note. A more precise determination of the values of ζ for which (6.85) and (6.87) hold in the ambiguous case of (3.9) will be made in Chapter 7 (Theorems 7.24 and 7.26).

CHAPTER 7

FURTHER NECESSARY CONDITIONS FOR ϵ -MINIMAL MATRICES

According to Theorem 5.13 all ϵ -minimal matrices are determined by solutions to the Maximum Problem 5.11. Therefore a necessary condition on a unitary matrix solving this maximum problem will, indirectly, be a necessary condition on an ϵ -minimal matrix. By working through the Maximum Problem 5.11 we shall be able to derive some additional necessary conditions in the present chapter.

The identity matrix I solves the Maximum Problem 5.11 for $B \in \mathcal{M}_n$ if and only if B satisfies the condition

$$(7.1) \quad \epsilon^2(\text{dg}(B)) = \max_{U \in \mathcal{U}_n} \epsilon^2(\text{dg}(UBU^*)) .$$

If B satisfies (7.1) then, by Theorem 5.15, $\text{dg}(B)$ is ϵ -minimal for B and, by Theorem 4.10, $\text{dg}(B)$ satisfies the necessary condition (4.12), i.e.,

$$(7.2) \quad \text{dg}(B)B^* - B^* \text{dg}(B) + \text{dg}(B^*)B - B \text{dg}(B^*) = 0 .$$

Letting $B = (b_{ij})$ we can express the last equation in terms of the elements of B as follows

$$(7.3) \quad \bar{b}_{ji}(b_{ii} - b_{jj}) + b_{ij}(\bar{b}_{ii} - \bar{b}_{jj}) = 0 \quad (1 \leq i, j \leq n) .$$

In like manner we find from Theorem 4.73 that there is a hermitian matrix H such that

$$dg(B) + dg(B)H - H dg(B) = B = dg(B) + offdg(B)$$

or

$$offdg(B) = dg(B)H - H dg(B) .$$

Writing $H = (h_{ij})$ we can express the last equation in terms of the elements of B and H as follows

$$b_{ij} = h_{ij}(b_{ii} - b_{jj}) \quad (i \neq j) .$$

Since H is hermitian we have $h_{ji} = \bar{h}_{ij}$ so $b_{ji} = \bar{h}_{ij}(b_{jj} - b_{ii}) = -\bar{h}_{ij}(b_{ii} - b_{jj})$. Therefore we find that every 2 by 2 principal submatrix B_{ij} of B is of the form

$$(7.4) \quad B_{ij} = \begin{pmatrix} b_{ii} & b_{ij} \\ b_{ji} & b_{jj} \end{pmatrix} = \begin{pmatrix} b_{ii} & h_{ij}(b_{ii} - b_{jj}) \\ -\bar{h}_{ij}(b_{ii} - b_{jj}) & b_{jj} \end{pmatrix} \quad (1 \leq i < j \leq n).$$

Consider now unitary matrices $U = (u_{rc}) \in \mathcal{U}_n$ of the following special type. For any pair (i, j) of row and column indices satisfying $1 \leq i < j \leq n$ we let the 2 by 2 principal submatrix

$$U_{ij} = \begin{pmatrix} u_{ii} & u_{ij} \\ u_{ji} & u_{jj} \end{pmatrix}$$

be unrestricted (except for the requirement $U \in \mathcal{U}_n$) and specify that all other elements of U satisfy

$$u_{rc} = \delta_{rc} = \begin{cases} 1, & r = c \\ 0, & r \neq c \end{cases}.$$

One may easily verify that U is unitary if and only if U_{ij} is unitary. If B satisfies (7.1) then certainly $\epsilon^2(\text{dg}(B)) \geq \epsilon^2(\text{dg}(UBU^*))$ holds for all U of the special type just described. The effect of the transformation UBU^* on the submatrix B_{ij} is that of replacing B_{ij} by $U_{ij} B_{ij} U_{ij}^*$; therefore, since U_{ij} can be any matrix in \mathcal{U}_2 , we see that B_{ij} itself satisfies (7.1) with n replaced by 2:

$$\epsilon^2(\text{dg}(B_{ij})) = \max_{U \in \mathcal{U}_2} \epsilon^2(\text{dg}(UB_{ij}U^*)) \quad (1 \leq i < j \leq n).$$

We now find from Theorem 6.83 that

$$(7.5) \quad B_{ij} + \xi_{ij} B_{ij}^* \in \mathcal{N}_2 \quad (1 \leq i < j \leq n)$$

holds for some $\xi_{ij} \in \zeta(B_{ij})$. Expressed in terms of elements, the condition (7.5) states that

$$(7.6) \quad b_{ij} + \xi_{ij} \bar{b}_{ji} = b_{ji} + \xi_{ij} \bar{b}_{ij} = 0 \quad (1 \leq i < j \leq n).$$

Since $|\xi_{ij}| = 1$ the last equation implies $|b_{ij}| = |b_{ji}|$.

We can obtain some information about the complex Lagrange multipliers h_{ij} in (7.4) as follows. From (7.3) and (7.4) we obtain

$$-\bar{h}_{ij}(b_{ii} - b_{jj})^2 + h_{ij} |b_{ii} - b_{jj}|^2 = 0$$

so, unless $b_{11} = b_{jj}$ (in which case the value of h_{1j} in (7.4) does not matter at all) or $h_{1j} = 0$, we have

$$\frac{h_{1j}^2}{|h_{1j}|^2} = \frac{(b_{11} - b_{jj})^2}{|b_{11} - b_{jj}|^2}$$

or

$$\frac{h_{1j}}{|h_{1j}|} = \pm \frac{b_{11} - b_{jj}}{|b_{11} - b_{jj}|} .$$

Straightforward computations using (7.4) show that

$$(7.7) \quad \text{tr}(B_{1j}^2) - \frac{1}{2} (\text{tr}(B_{1j}))^2 = \left(\frac{1}{2} - 2|h_{1j}|^2\right)(b_{11} - b_{jj})^2 .$$

If (7.7) does not vanish we see from (3.9) that ζ_{1j} in (7.5) must have the value

$$(7.8) \quad \zeta_{1j} = s_{1j} \frac{b_{11} - b_{jj}}{\overline{b_{11}} - \overline{b_{jj}}}$$

where

$$s_{1j} = \text{sgn}\left(\frac{1}{2} - 2|h_{1j}|^2\right)$$

and where for $a \in \mathbb{R}$ $\text{sgn}(a)$ denotes the sign of a :

$$\text{sgn}(a) = \begin{cases} +1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \end{cases} .$$

Using (7.4) and (7.8) we obtain

$$b_{ij} + \zeta_{ij} \bar{b}_{ji} = h_{ij}(1 - s_{ij})(b_{ii} - b_{jj}) ,$$

$$b_{ji} + \zeta_{ij} \bar{b}_{ij} = \bar{h}_{ij}(s_{ij} - 1)(b_{ii} - b_{jj})$$

whence (7.6) holds if and only if either $h_{ij} = 0$ or $s_{ij} = 1$. Clearly $s_{ij} = 1$ holds if and only if $(1/2) - 2|h_{ij}|^2 > 0$, i.e.,

$$(7.9) \quad |h_{ij}| < \frac{1}{2} .$$

The other possibility is that $h_{ij} = 0$ in which case (7.9) still holds.

In the ambiguous case of (3.9) where (7.7) vanishes we obtain

$$(7.10) \quad b_{ij} + \zeta_{ij} \bar{b}_{ji} = h_{ij}[(b_{ii} - b_{jj}) - \zeta_{ij} \overline{(b_{ii} - b_{jj})}] ,$$

$$(7.11) \quad b_{ji} + \zeta_{ij} \bar{b}_{ij} = -\bar{h}_{ij}[(b_{ii} - b_{jj}) - \zeta_{ij} \overline{(b_{ii} - b_{jj})}] .$$

If we disregard the uninteresting case $b_{ii} = b_{jj}$ then obviously (7.7) vanishes if and only if $|h_{ij}| = (1/2)$. Thus, from (7.10) and (7.11), (7.6) holds if and only if

$$\zeta_{ij} = \frac{b_{ii} - b_{jj}}{\bar{b}_{ii} - \bar{b}_{jj}} .$$

We now summarize the above results in the following

7.12 Theorem. Let $B = (b_{ij}) \in \mathcal{M}_n$ satisfy the condition (7.1).

For each pair of indices (i, j) satisfying $1 \leq i < j \leq n$ the following statements are true.

$$(7.13) \quad b_{1j} \overline{(b_{11} - b_{jj})} = \bar{b}_{j1} (b_{11} - b_{jj}) .$$

$$(7.14) \quad |b_{1j}| = |b_{j1}| .$$

There exists a complex number (= complex Lagrange multiplier) h_{1j} such that

$$(7.15) \quad B_{1j} = \begin{pmatrix} b_{11} & b_{1j} \\ b_{j1} & b_{jj} \end{pmatrix} = \begin{pmatrix} b_{11} & h_{1j}(b_{11}-b_{jj}) \\ -\bar{h}_{1j}(b_{11}-b_{jj}) & b_{jj} \end{pmatrix} .$$

If $b_{11} = b_{jj}$ then $b_{1j} = b_{j1} = 0$. If $b_{11} \neq b_{jj}$ then

$$(7.16) \quad |h_{1j}| \leq \frac{1}{2} .$$

More precisely, if

$$(7.17) \quad \text{tr}(B_{1j}^2) - \frac{1}{2} (\text{tr}(B_{1j}))^2 \neq 0 ,$$

then

$$(7.18) \quad |h_{1j}| < \frac{1}{2}$$

and, if

$$(7.19) \quad \text{tr}(B_{1j}^2) - \frac{1}{2} (\text{tr}(B_{1j}))^2 = 0$$

and $b_{11} \neq b_{jj}$, then

$$(7.20) \quad |h_{ij}| = \frac{1}{2} .$$

If $h_{ij} \neq 0$ and $b_{ii} \neq b_{jj}$ then

$$(7.21) \quad \frac{h_{ij}}{|h_{ij}|} = \pm \frac{b_{ii} - b_{jj}}{|b_{ii} - b_{jj}|} .$$

Furthermore, if $b_{ii} \neq b_{jj}$, there is a uniquely determined complex number ξ_{ij} satisfying $|\xi_{ij}| = 1$ such that

$$(7.22) \quad b_{ij} + \xi_{ij} \bar{b}_{ji} = b_{ji} + \xi_{ij} \bar{b}_{ij} = 0 .$$

If (7.17) holds $\xi_{ij} \in \xi(B_{ij})$ where the set $\xi(B_{ij})$ is defined by (3.9). If (7.19) holds then

$$(7.23) \quad \xi_{ij} = \frac{b_{ii} - b_{jj}}{\bar{b}_{ii} - \bar{b}_{jj}} .$$

Specializing Theorem 7.12 to the case $n = 2$ we can improve Theorems 6.83 and 6.86 as follows.

7.24 Theorem. Let $A = (a_{ij}) \in \mathcal{M}_2$ satisfy (6.84). If $a_{11} = a_{22}$ then $A = a_{11}I$ and (6.85) holds for all $\xi \in \xi(A)$. If $a_{11} \neq a_{22}$ and and if $\text{tr}(A^2) - (1/2)(\text{tr}(A))^2 = 0$ then (6.85) holds only for

$$(7.25) \quad \xi = \frac{a_{11} - a_{22}}{\bar{a}_{11} - \bar{a}_{22}} .$$

7.26 Theorem. Let $A \in \mathcal{M}_2$, let U_0 solve the Maximum Problem 5.11 for A , and let $B = (b_{ij}) = U_0 A U_0^*$. If $b_{11} = b_{22}$ then $B = b_{11}I$ and (6.87) holds for all $\zeta \in \zeta(A)$. If $b_{11} \neq b_{22}$ and if $\text{tr}(A^2) - (1/2)(\text{tr}(A))^2 = 0$ then (6.87) holds only for

$$(7.27) \quad \zeta = \frac{b_{11} - b_{22}}{\bar{b}_{11} - \bar{b}_{22}}.$$

7.28 Theorem. Let $A \in \mathcal{M}_n$, let U_0 solve the Maximum Problem 5.11 for A , and let $B = (b_{ij}) = U_0 A U_0^*$. Then all the conclusions of Theorem 7.12 hold.

Proof. If U_0 solves the Maximum Problem 5.11 for A then, by an argument similar to that used in the proof of Theorem 6.86, B satisfies (7.1).

CHAPTER 8

COUNTEREXAMPLES TO MIRSKY'S CONJECTURE

In the present chapter we shall present some selected examples which shed some light on the distance problem of finding $d_v(A)$ and which also show that Mirsky's Conjecture 1.35 is incorrect in a number of instances.

Consider first the following class of matrices of order n . We define $A = (a_{kl}) \in \mathcal{M}_n$ as follows: $a_{kl} = 0$ ($1 \leq k, l \leq n$) with the following exceptions

$$(8.1) \quad a_{k,k+1} = e^{i\theta_k} \quad (k = 1, 2, \dots, n-1)$$

where the θ_k are arbitrary real numbers. Now define $N_\varphi = (b_{kl}) \in \mathcal{M}_n$ by

$$(8.2) \quad b_{nl} = e^{i\varphi}, \quad b_{kl} = a_{kl} \text{ otherwise}$$

where φ is any real number. Clearly the rows and columns of N_φ form orthonormal sets so that N_φ is unitary (for all $\theta_1, \dots, \theta_{n-1}, \varphi \in \mathbb{R}$) and therefore normal. (Actually N_φ is the product of a diagonal unitary matrix with a permutation matrix.)

For all $\theta_1, \dots, \theta_{n-1}$ we have $\Omega(A) = 0$ and $\epsilon^2(A) = n-1$ so

$$(8.3) \quad \epsilon^2(A) - \epsilon^2(\Omega(A)) = n-1.$$

An elementary calculation shows that $f(c) = \epsilon^2(A - cN_\varphi)$, where c is real, assumes its absolute minimum if and only if $c = (n-1)/n$.

Furthermore, by (8.3)

$$(8.4) \quad \epsilon^2(A - \frac{n-1}{n} N_\varphi) = \frac{n-1}{n} = \frac{1}{n} (\epsilon^2(A) - \epsilon^2(N(A)))$$

holds for all $\theta_1, \dots, \theta_{n-1}, \varphi \in \mathbb{R}$. This proves (cf. Theorem 2.16).

8.5 Theorem. Mirsky's Conjecture 1.35 is false for $\nu = \epsilon$ and $n \geq 3$.

There exist matrices A in \mathcal{M}_n such that

$$(8.6) \quad p_{\epsilon,n}(A) \leq \frac{1}{n} \quad (n \geq 3)$$

where $p_{\epsilon,n}(A)$ is defined (for $\nu = \epsilon$) by (2.14).

Note. An example of order 3 similar to the pair $A, (n-1)/n N_\varphi$ in the case $n = 3$ is due to Eberlein [6]. She also obtains counterexamples to Mirsky's conjecture for $n \geq 4$ by bordering her 3 by 3 example with zeros.

Consider next the class of 2 by 2 matrices of the form

$$(8.7) \quad A = \begin{pmatrix} \lambda & m \\ 0 & 0 \end{pmatrix} \quad \text{where } \lambda, m \in \mathbb{C}, \quad m \neq 0.$$

Let N be any ϵ -minimal matrix for A , that is, any one of the matrices

$$(8.8) \quad \frac{1}{2} \begin{pmatrix} 2\lambda & m \\ \xi \bar{m} & 0 \end{pmatrix} \quad (\xi \in \zeta(A))$$

where $\zeta(A)$ is defined by (3.9) (cf. Theorem 6.24). One finds easily

$$AA^* = \text{diag}(|\lambda|^2 + |m|^2, 0), \quad (A - N)(A - N)^* = \frac{1}{4} \text{diag}(|m|^2, |m|^2)$$

whence, by Definition 1.10, the singular values of A are

$$\sqrt{|\lambda|^2 + |m|^2} \quad , \quad 0$$

and the singular values of $A - N$ are

$$\frac{1}{2} |m| \quad , \quad \frac{1}{2} |m| \quad .$$

Using the defining formulas (1.13) - (1.15) for the unitarily invariant norms v_p , we find that the following hold for all p ($1 \leq p \leq \infty$):

$$v_p(A) = [(\sqrt{|\lambda|^2 + |m|^2})^p]^{1/p} = \sqrt{|\lambda|^2 + |m|^2} \quad ,$$

$$v_p(\Omega(A)) = |\lambda| \quad ,$$

$$v_p(A - N) = [2(\frac{|m|}{2})^p]^{1/p} = 2^{1/p-1} |m|$$

so

$$v_p^2(A) - v_p^2(\Omega(A)) = |m|^2 \quad ,$$

$$v_p^2(A - N) = 2^{2/p-2} |m|^2$$

or

$$v_p^2(A - N) = 2^{2/p-2} (v_p^2(A) - v_p^2(\Omega(A))) \quad .$$

Now $2^{2/p-2} < (1/2)$ if and only if $2^{2/p} < 2$ and this happens if and only if $2/p < 1$, i.e., $p > 2$. Thus the pair A, N provides a

counterexample to Mirsky's conjecture for $n = 2$ and for all $p > 2$. By bordering the matrices A, N with zeros and by carrying out the relevant computations with (1.13) - (1.15) for higher values of n , we then obtain counterexamples to the conjecture for all $n \geq 3$. We have proved

8.9 Theorem. Mirsky's Conjecture 1.35 is false for $\nu = \nu_p$ ($2 < p \leq \infty$) and $n \geq 2$ where ν_p is given by (1.15).

Remark. Since $\nu_\infty = \sigma$ we have again proved the second statement of Theorem 2.19.

CHAPTER 9

THE FIELD OF VALUES AND EIGENVALUES OF ϵ -MINIMAL MATRICES

The field of values (or numerical range) of a matrix $A \in \mathcal{M}_n$ is defined to be the following set of complex numbers.

$$(9.1) \quad F(A) = \{(Ax, x); (x, x) = 1 \text{ and } x \text{ is a complex column vector}\}.$$

The following are known facts concerning $F(A)$:

9.2 Theorem. (Toeplitz [23] and Hausdorff [12]) For any $A \in \mathcal{M}_n$

$F(A)$ is a closed, bounded, connected, convex subset of \mathbb{C} .

9.3 Theorem. (Toeplitz [23]) Let $A \in \mathcal{M}_n$. All eigenvalues of A are in $F(A)$. If $A \in \mathcal{H}_n$ then $F(A)$ coincides with the convex hull $C(A)$ of the eigenvalues of A .

9.4 Theorem. (Hausdorff [12]) The field of values is invariant under a unitary similarity transformation:

$$(9.5) \quad F(A) = F(UAU^*) \text{ where } A \in \mathcal{M}_n, U \in \mathcal{U}_n.$$

A simple consequence of Theorems 9.2 and 9.3 is the following

9.6 Theorem. Let $A \in \mathcal{M}_n$. Then the convex hull $C(A)$ of the eigenvalues of A is contained in $F(A)$:

$$(9.7) \quad C(A) \subset F(A).$$

An elementary computation using Theorem 9.4 and results of Toeplitz [23] and Donoghue [5] yields the following result whose proof is omitted.

9.8 Theorem. Let $A \in \mathcal{M}_2$ and let

$$(9.9) \quad VAV^* = \begin{pmatrix} \lambda_1 & m \\ 0 & \lambda_2 \end{pmatrix} \quad (V \in \mathcal{U}_2)$$

be a Schur triangular form for A . If A is not normal, i.e., if $m \neq 0$, and if $\lambda_1 \neq \lambda_2$, then $F(A)$ is the interior and boundary of an ellipse whose foci are λ_1 and λ_2 , whose minor axis has length $|m|$, and whose major axis has length $(|m|^2 + |\lambda_1 - \lambda_2|^2)^{1/2}$; if $\lambda_1 = \lambda_2$ then $F(A)$ is the interior and boundary of a circle with center at λ_1 and diameter $|m|$. If A is normal, i.e., if $m = 0$, and if $\lambda_1 \neq \lambda_2$, then $F(A)$ is the straight line segment connecting λ_1 and λ_2 ; if $\lambda_1 = \lambda_2$, then $F(A)$ reduces to a single point, namely λ_1 .

Our first objective in the present chapter is to prove

9.10 Theorem. Let $A \in \mathcal{M}_2$, $A \notin \mathcal{N}_2$ and let (9.9) be a Schur form for A . If $\lambda_1 \neq \lambda_2$ then the eigenvalues of the (unique) ϵ -minimal matrix for A are the end points of the major axis of the ellipse which is the boundary of $F(A)$. If $\lambda_1 = \lambda_2$ and if $N_\theta(A)$ (see the note following Definition 3.8) is one of the ϵ -minimal matrices for A , then the eigenvalues of $N_\theta(A)$ are

$$(9.11) \quad \lambda_1 \pm \frac{1}{2} |m| e^{i\theta/2} ;$$

given any diameter of the circle which is the boundary of $F(A)$, there is a $\theta \in \mathbb{R}$ such that the eigenvalues of $N_\theta(A)$ are the endpoints of that diameter.

Proof. In view of Theorem 9.4, (6.27), and (6.28) we need only prove the result for some unitary transform of A . We shall use the form (9.9) for this purpose. Consider first the case $\lambda_1 = \lambda_2$. Then by putting $a_1 = \lambda_1$, $a_2 = m$ in (6.29) and (6.40) we see that

$$N_\theta(VAV^*) = \begin{pmatrix} \lambda_1 & \frac{1}{2}m \\ \frac{1}{2}\bar{m}e^{i\theta} & \lambda_1 \end{pmatrix} ;$$

furthermore, by Theorem 5.13, the eigenvalues of N_θ are the diagonal elements of (6.40), namely

$$(9.12) \quad \lambda_1 \pm \frac{1}{2} \mu_{1\theta} m ,$$

where, by (6.35), $(\mu_{1\theta} m)^2 = |m|^2 \exp(i\theta)$ or

$$(9.13) \quad \mu_{1\theta} m = \pm |m| e^{i\theta/2} .$$

Whatever sign is chosen in (9.13) we see from (9.12) that the eigenvalues of N_θ are given by (9.11). By Theorem 9.8 the boundary of $F(A)$ is a circle with center λ_1 and radius $(1/2)|m|$. As θ increases from 0 to 2π it is clear from (9.11) that the eigenvalues of N_θ are the endpoints of a diameter of that circle which rotates through an angle of π . That is, every diameter of the circle which is the boundary of $F(A)$ is included as one of those which can occur as the line segment connecting the eigenvalues of some ϵ -minimal matrix for A .

We now consider the case $\lambda_1 \neq \lambda_2$. Straightforward computations

based on Definition 3.8 show that

$$(9.14) \quad N(VAV^*) = \begin{pmatrix} \lambda_1 & \frac{1}{2} m \\ \frac{1}{2} \bar{\zeta} \bar{m} & \lambda_2 \end{pmatrix}, \quad \text{where } \zeta = \frac{\lambda_1 - \lambda_2}{\bar{\lambda}_1 - \bar{\lambda}_2}$$

and, according to Theorems 6.24 and 6.79, this is the only ϵ -minimal matrix for (9.9). The eigenvalues μ_1, μ_2 of $N(VAV^*)$ are the roots of the quadratic equation $\mu^2 - (\lambda_1 + \lambda_2)\mu + \lambda_1\lambda_2 - (1/4)|m|^2\zeta = 0$. Using the quadratic formula and the expression for ζ in (9.14) we find

$$(9.15) \quad \begin{cases} \mu_1 = \frac{1}{2}(\lambda_1 + \lambda_2) + \frac{1}{2}\sqrt{\zeta} \sqrt{|m|^2 + |\lambda_1 - \lambda_2|^2} \\ \mu_2 = \frac{1}{2}(\lambda_1 + \lambda_2) - \frac{1}{2}\sqrt{\zeta} \sqrt{|m|^2 + |\lambda_1 - \lambda_2|^2} \end{cases}.$$

Now

$$\zeta = \frac{\lambda_1 - \lambda_2}{\bar{\lambda}_1 - \bar{\lambda}_2} = \frac{(\lambda_1 - \lambda_2)^2}{|\lambda_1 - \lambda_2|^2}$$

whence

$$(9.16) \quad \arg(\sqrt{\zeta}) = \arg(\lambda_1 - \lambda_2).$$

According to Theorem 9.8 the boundary of $F(A)$ is an ellipse with foci λ_1, λ_2 , center $(1/2)(\lambda_1 + \lambda_2)$ and major axis of length $(|m|^2 + |\lambda_1 - \lambda_2|^2)^{1/2}$. From (9.15) and (9.16) we find

$$(9.17) \quad \arg(\lambda_1 - \frac{1}{2}(\lambda_1 + \lambda_2)) = \arg(\lambda_1 - \lambda_2) = \arg(\mu_1 - \frac{1}{2}(\lambda_1 + \lambda_2)),$$

BLANK PAGE

$$(9.18) \quad \arg(\lambda_2 - \frac{1}{2}(\lambda_1 + \lambda_2)) = \arg[-(\lambda_1 - \lambda_2)] = \arg(\mu_2 - \frac{1}{2}(\lambda_1 + \lambda_2)) .$$

Relations (9.17) and (9.18) show that $\lambda_1, \lambda_2, \mu_1$, and μ_2 all lie on the same straight line through the center $(1/2)(\lambda_1 + \lambda_2)$ of the ellipse. Since λ_1 and λ_2 lie on the major axis, since $|\xi| = 1$, and since the length of the major axis is $(|m|^2 + |\lambda_1 - \lambda_2|^2)^{1/2}$, we see from (9.15) that μ_1 and μ_2 are the endpoints of the major axis. This completes the proof of Theorem 9.10.

9.19 Corollary. Let $A \in \mathcal{M}_2$, $A \notin \mathcal{N}_2$. If $F(A)$ is a circular disk, then given any diameter of that disk there is a $\theta \in \mathbb{R}$ such that $F(N_\theta(A))$ coincides with that diameter. If $F(A)$ is the interior and boundary of an ellipse (not a circle) then $F(N(A))$ coincides with the major axis of that ellipse.

Proof. This follows immediately from Theorems 9.8 and 9.10, since every matrix in $M(A)$ is normal.

9.20 Corollary. Let $A \in \mathcal{M}_2$, $A \notin \mathcal{N}_2$ and let N_0 be ϵ -minimal for A . Then the eigenvalues μ_1, μ_2 of N_0 are extreme points of $F(A)$. Furthermore $|\mu_1 - \mu_2| = \text{diam}(F(A)) = \sup_{z, w \in F(A)} |z - w|$.

9.21 Theorem. Let $A \in \mathcal{M}_n$ and let N_0 be ϵ -minimal for A . Then every eigenvalue of N_0 belongs to $F(A)$ and

$$(9.22) \quad F(N_0) \subset F(A) .$$

Proof. Let $N_0 = U_0^* D_0 U_0$ where $U_0 \in \mathcal{U}_n$ and $D_0 \in \mathcal{D}_n$. According to Theorem 5.13 the eigenvalues of N_0 are the diagonal elements of the matrix $U_0 A U_0^*$. If u_k is the k -th column of U_0^* then the k -th

diagonal element of $U_0 A U_0^*$ is given by $(A u_k, u_k)$ and, by (9.1), this is in $F(A)$. Since $F(A)$ is convex and since $F(N_0)$ is the convex hull of the eigenvalues of N_0 , we obtain immediately (9.22).

CHAPTER 10

A GENERALIZATION OF THE JACOBI AND GOLDSTINE-HORWITZ METHODS

Long ago Jacobi [14] devised a method for diagonalizing a real symmetric matrix. The method utilized coordinate-plane rotations and was essentially dependent only on an elementary technique for diagonalizing a real symmetric matrix of order 2 using an orthogonal similarity transformation. Since 1950 Jacobi's method has been extensively studied and generalized (see e.g., [3], [4]). In the present chapter we shall describe still another generalization of Jacobi's method which amounts to a computational technique for solving the Maximum Problem 5.11. This new technique also generalizes and simplifies a method devised by Goldstine and Horwitz [10].

Consider the following computational algorithm. Let $A \in \mathcal{M}_n$ and set $A_0 = A$. One calculates a sequence of matrices A_1, A_2, \dots , $A_k = (a_{rc}^{(k)})$, ... which are unitarily similar to A by the recurrence relation

$$(10.1) \quad A_{k+1} = U_k A_k U_k^* \quad (k = 0, 1, 2, \dots) .$$

The $U_k = (u_{rc}^{(k)})$ are special unitary matrices of order n . For every value of k there is specified a pair $\pi_k = (i_k, j_k) = (i, j)$ of indices (we omit the subscript k in the sequel for notational simplicity) satisfying $1 \leq i < j \leq n$ such that the 2 by 2 matrix

$$(10.2) \quad V_k = \begin{pmatrix} u_{ii}^{(k)} & u_{ij}^{(k)} \\ u_{ji}^{(k)} & u_{jj}^{(k)} \end{pmatrix},$$

which is a principal submatrix of U_k , is unitary. All other elements of U_k satisfy

$$u_{rc}^{(k)} = \delta_{rc} = \begin{cases} 1, & r = c \\ 0, & r \neq c \end{cases}.$$

The matrices U_k are completely determined by the pairs π_k and the 2 by 2 unitary matrices V_k . We shall always take V_k to be a matrix of the form (6.12).

Any set of rules for choosing the sequence $\{U_k\}$ will be called a method of Jacobi type. The following example which is defined only for $A \in \mathcal{H}_n$ and which is a straightforward extension of Jacobi's original method to hermitian matrices will be referred to as the classical Jacobi method (cf. [4]). Here one chooses π_k such that

$$(10.3) \quad |a_{ij}^{(k)}| = \max_{r \neq c} |a_{rc}^{(k)}|$$

and chooses V_k to be a matrix of the form (6.12) such that

$$(10.4) \quad a_{ij}^{(k+1)} = 0.$$

The new generalization of Jacobi's method which was announced in the first paragraph of this chapter is a method of Jacobi type with the

following set of rules for determining the U_k of (10.1). Let

$$(10.5) \quad A_{rc}^{(k)} = \begin{pmatrix} a_{rr}^{(k)} & a_{rc}^{(k)} \\ a_{cr}^{(k)} & a_{cc}^{(k)} \end{pmatrix} \quad (1 \leq r < c \leq n)$$

and define

$$(10.6) \quad \Delta_k(r, c) = \max_{U \in \mathcal{U}_2} \epsilon^2(\text{dg}(UA_{rc}^{(k)}U^*)) - \epsilon^2(\text{dg}(A_{rc}^{(k)})) .$$

Choose λ_k so that

$$(10.7) \quad \Delta_k(1, j) = \max_{1 \leq r < c \leq n} \Delta_k(r, c)$$

and choose V_k to be a matrix of the form (6.12) which solves the Maximum Problem 5.11 for $A_{1j}^{(k)}$. We show in the next paragraph how to calculate easily $\Delta_k(r, c)$ and V_k .

By Theorems 5.13 and 6.24 a unitary matrix U_0 solves the Maximum Problem 5.11 for $A \in \mathcal{M}_2$ if and only if $U_0 N_\zeta(A) U_0^*$ is diagonal and, by (3.10), that happens if and only if $U_0 (A + \zeta A^*) U_0^*$ is diagonal where $\zeta \in \zeta(A)$. Therefore all we need to do in order to solve the Maximum Problem 5.11 for any of the submatrices $A_{rc}^{(k)}$ is to find a $\mu_{rc}^{(k)} \in \mathbb{C}$ such that

$$(10.8) \quad W(\mu_{rc}^{(k)})(A_{rc}^{(k)}) + \zeta A_{rc}^{(k)} W^*(\mu_{rc}^{(k)}) \in \mathcal{D}_2$$

where ζ is any number in $\zeta(A_{rc}^{(k)})$. Thus from (10.6) we have

$$(10.9) \quad \Delta_k(r,c) = \epsilon^2(\text{dg}(W(\mu_{rc}^{(k)})A_{rc}^{(k)}W^*(\mu_{rc}^{(k)}))) - \epsilon^2(\text{dg}(A_{rc}^{(k)}))$$

and we set

$$(10.10) \quad V_k = W(\mu_{1j}^{(k)}) \quad .$$

We now show that if $A \in \mathcal{H}_n$ the set of rules (10.7) and (10.10) coincide respectively with (10.3) and (10.4). This will show that our new method is a generalization of the classical Jacobi method. If $A \in \mathcal{H}_n$ then, from (10.1), $A_k \in \mathcal{H}_n$ ($k = 0, 1, 2, \dots$) which implies $A_{rc}^{(k)} \in \mathcal{H}_2$ for $1 \leq r < c \leq n$ and for all k . Since $A_{rc}^{(k)}$ is hermitian it can be diagonalized by a unitary transformation and the value of the maximum in (10.6) is given by $\epsilon^2(A_{rc}^{(k)})$. Thus we find that

$$\Delta_k(r,c) = 2|a_{rc}^{(k)}| \quad (1 \leq r < c \leq n)$$

whence the rule (10.7) reduces to (10.3). Furthermore, by Corollary 6.78, $M(A_{rc}^{(k)}) = \{A_{rc}^{(k)}\}$ whence the V_k given by (10.10) diagonalizes $A_{1j}^{(k)}$ i.e., (10.4) holds.

Since the classical Jacobi method diagonalizes any $A \in \mathcal{H}_n$ we have proved that the method of Jacobi type (10.7), (10.10) will diagonalize any hermitian matrix. By a similar argument we could show that the method will also diagonalize any skew-hermitian matrix.

In [10] Goldstine and Horwitz devised a method of Jacobi type which was applicable to any normal matrix $A \in \mathcal{M}_n$. At each stage they sought, by a very complicated procedure, to determine a V_k of the form (6.12)

such that

$$\tau^2(A_{1j}^{(k+1)}) = \epsilon^2(\text{offdg}(V_k A_{1j}^{(k)} V_k^*))$$

was minimized. It is easy to see that this is the same as determining V_k so that $\epsilon^2(\text{dg}(V_k A_{1j}^{(k)} V_k^*))$ is maximized i.e., solving the Maximum Problem 5.11 for $A_{1j}^{(k)}$. Thus our technique provides a simple solution to the problem studied by Goldstine and Horwitz and at the same time makes its application to arbitrary $A \in \mathcal{M}_n$ meaningful in the context of solving the Maximum Problem 5.11 for A .

Of course the main question here is whether or not the method of Jacobi type (10.7), (10.10) will actually solve the Maximum Problem 5.11 for any $A \in \mathcal{M}_n$, i.e., whether or not $\epsilon^2(\text{dg}(A_k))$ converges to the maximum (5.12). A related question is whether or not the infinite product $\cdots U_k U_{k-1} \cdots U_1 U_0$ of unitary matrices converges to a unitary matrix which solves the Maximum Problem 5.11 (cf. [4]). If the answers to these questions were affirmative then we would have a constructive method of computing ϵ -minimal matrices for any $A \in \mathcal{M}_n$. If the answer to the first question were affirmative then, by Theorem 5.4, we would have a constructive method of computing $d_\epsilon(A)$. These convergence questions appear to be rather difficult and we content ourselves here by proving only the following

10.11 Lemma. Let $A \in \mathcal{M}_n$ and let the matrices U_k of (10.1) be determined for each k by the rules (10.7) and (10.10). Then

$$(10.12) \quad \Delta_k(i, j) \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

Proof. Let $\sigma_k = \epsilon^2(\text{dg}(A_k))$ for $k = 0, 1, 2, \dots$. Since each A_k is unitarily similar to A , we have

$$(10.13) \quad \sigma_k \leq \epsilon^2(A) \quad (k = 0, 1, 2, \dots)$$

The only diagonal elements of A_k affected by the transformation (10.1) are the ones in the i -th and j -th rows. Moreover, since the transformation $A_{ij}^{(k+1)} = V_k A_{ij}^{(k)} V_k^*$ increases the function $\epsilon^2(\text{dg}(A_{ij}))$ by the amount $\Delta_k(i, j)$, we have

$$(10.14) \quad \sigma_{k+1} = \sigma_k + \Delta_k(i, j) \geq \sigma_k \quad \text{for } k = 0, 1, 2, \dots$$

The relations (10.13) and (10.14) show that $\{\sigma_k\}$ is a monotonically increasing sequence of positive numbers which is bounded from above. Consequently $\sigma_k \rightarrow \sigma \leq \epsilon^2(A)$ as $k \rightarrow \infty$ and from (10.14) we obtain the desired conclusion (10.12).

BIBLIOGRAPHY

- [1] L. V. Ahlfors, Complex Analysis, McGraw-Hill, New York, 1953.
- [2] A. R. Amir-Moéz and A. Horn, "Singular values of a matrix," Amer. Math. Monthly, vol. 65, (1958), pp. 742-748.
- [3] R. L. Causey, "Computing eigenvalues of non-hermitian matrices by methods of Jacobi type," J. Soc. Indust. Appl. Math., vol. 6, (1958), pp. 172-181.
- [4] R. L. Causey and P. Henrici, "Convergence of approximate eigenvectors in Jacobi methods," Numerische Math., vol. 2, (1960), pp. 67-78.
- [5] W. F. Donoghue, Jr., "On the numerical range of a bounded operator," Mich. Math. J., vol. 4, (1957), pp. 261-263.
- [6] P. J. Eberlein, "On a conjecture of Mirsky," submitted for publication.
- [7] D. K. Faddeev and V. N. Faddeeva, Computational Methods of Linear Algebra (translated from Russian by Robert C. Williams), W. H. Freeman and Co., San Francisco, 1963.
- [8] Ky Fan, "Maximum properties and inequalities for eigenvalues of completely continuous operators," Proc. Nat. Acad. Sci. U.S.A., vol. 37, (1951), pp. 760-766.
- [9] Ky Fan and A. J. Hoffman, "Some metric inequalities in the space of matrices," Proc. Amer. Math. Soc., vol. 6, (1951), pp. 760-766.
- [10] H. H. Goldstine and L. P. Horwitz, "A procedure for the diagonalization of normal matrices," J. Assoc. Comput. Mach., vol. 6, (1959), pp. 176-195.
- [11] P. R. Halmos, Finite-Dimensional Vector Spaces, 2nd Ed., Van Nostrand, New York, 1958.

- [12] F. Hausdorff, "Der Wertvorrat einer Bilinearform," Math. Zeitschrift, vol. 3, (1919), pp. 314-316.
- [13] E. Hille and R. Phillips, Functional Analysis and Semi-groups, Amer. Math. Soc. Colloq. Publ., vol. XXXI, New York, 1957.
- [14] C. G. J. Jacobi, "Über ein leichtes Verfahren, die in der Theorie der Säkularstörungen vorkommenden Gleichungen numerisch aufzulösen," J. reine angew. Math., vol. 30, (1846), pp. 51-95.
- [15] L. Mirsky, "Symmetric gauge functions and unitarily invariant norms," Quart. J. Math. Oxford (2), vol. 11, (1960), pp. 50-59.
- [16] L. Mirsky, An Introduction to Linear Algebra, Oxford Univ. Press, London, 1955.
- [17] F. D. Murnaghan, "On a convenient system of parameters for the unitary group," Proc. Nat. Acad. Sci. U.S.A., vol. 38, (1952), pp. 127-129.
- [18] Z. Nehari, Introduction to Complex Analysis, Allyn and Bacon, Boston, 1961.
- [19] J. von Neumann, "Some matrix-inequalities and metrization of metric space," Tomsk Univ. Review, vol. 1, (1937), pp. 286-300; Collected Works, vol. IV, p. 205.
- [20] W. V. Parker, "The characteristic roots of matrices," Duke Math. J., vol. 12, (1945), pp. 519-526.
- [21] S. Perlis, Theory of Matrices, Addison-Wesley Publishing Co., Inc., Cambridge, Mass., 1952.
- [22] I. Schur, "Über die charakteristischen Wurzeln einer linearen Substitution mit einer Anwendung auf die Theorie der Integralgleichungen," Math. Annalen, vol. 66, (1904), pp. 488-510.

- [23] O. Toeplitz, "Das algebraische Analogon zu einem Satze von Fejér," Math. Zeitschrift, vol. 2, (1918), pp. 187-197.
- [24] H. Weyl, "Inequalities between the two kinds of eigenvalues of a linear transformation," Proc. Nat. Acad. Sci. U.S.A., vol. 35, (1949), pp. 408-411.
- [25] J. Williamson, "A polar representation of singular matrices," Bull. Amer. Math. Soc., vol. 41, (1935), pp. 118-123.
- [26] A. Wintner and F. D. Murnaghan, "On a polar representation of non-singular square matrices," Proc. Nat. Acad. Sci. U.S.A., vol. 17, (1931), pp. 676-678.