# ON THE UNITARY COMPLETION OF A MATRIX

BY

M. MARCUS AND P. GREINER

## I. Introduction

Suppose a certain set of entries in an $n$-square array are prescribed. We consider the following question: to determine some nontrivial necessary conditions under which the rest of the entries may be constructed so that the resulting matrix is unitary. For example, a trivial necessary condition is that the sum of the squares of the absolute values of the prescribed entries in any particular row or column is at most 1. Define a *diagonal* of an $n$-square array to be a set of positions $(i, \sigma(i))$, $i = 1, \cdots, n$, where $\sigma$ is a permutation of $1, \cdots, n$. Two diagonals will overlap in $k$ places if the corresponding permutations agree on $k$ integers in $1, \cdots, n$, and we shall refer to two such diagonals as being *k-overlapping*. Our results will show for example that if $A$ is a 9-square matrix whose entries in a fixed pair of nonoverlapping diagonals are all $\frac{2}{3} + \varepsilon$, $\varepsilon > 0$, then $A$ cannot be completed to a unitary matrix. On the other hand if the absolute value of the sum of the elements in two nonoverlapping diagonals is to be no greater than 10, there exists a 9-square unitary matrix for which this value of the sum is taken on. This last statement becomes false however if 10 is replaced by $10 + \varepsilon$, $\varepsilon > 0$.

For the group of $n$-square unitary matrices we obtain the maximum and minimum over all pairs of $k$-overlapping diagonals $d_1$ and $d_2$ of the maximum absolute value of the sum of the entries in $d_1$ and $d_2$. Our main results are contained in the

THEOREM. *Consider a fixed pair of $k$-overlapping diagonals of an $n$-square array, $1 \leq k \leq n$. Let $s_k$ be the maximum taken over all $n$-square unitary matrices of the absolute value of the sum of the elements in the given pair of diagonals. Then*

(i)    $s_k \leq n$    *if*    $n = k + 2$,

(ii)    $s_k \leq n - 4 + 2 \cot \pi/8$    *if*    $n = k + 4$,

(iii)    $s_k \leq n + \alpha$    *if*    $n = k + 3\alpha$,

       $s_k \leq n + \alpha - 5 + 2 \csc \pi/10$    *if*    $n = k + 3\alpha + 5$,

       $s_k \leq n + \alpha - 7 + 2 \csc \pi/14$    *if*    $n = k + 3\alpha + 7$,

(iv)    $s_k \geq n$    *if*    $n = k + 2\alpha$,

       $s_k \geq n + 1$    *if*    $n = k + 2\alpha + 3$.

*Moreover for each of the bounds in* (i)–(iv) *there exist a pair of $k$-overlapping diagonals and a unitary matrix for which the absolute value of the sum of the elements in these diagonals is the appropriate bound.*

*In addition if $0 \leq \mu \leq s_k$, then there exists a unitary matrix such that the absolute value of the sum of the elements in the fixed pair of diagonals is $\mu$.*

We first introduce some preliminary definitions and results. Then in Section II we shall consider the case of the two $k$-overlapping diagonals. In Sections III and IV we state and prove the necessary trigonometric inequalities to complete the proof of the above result. In Section V we indicate a further application of our methods.

Let $O_n$ be the group of $n$-square unitary matrices over the complex numbers. Let $S$ be any set of pairs $(i, j)$, $1 \leq i, j \leq n$, and define the real valued function on $O_n$

$$(1.1) \qquad g_S(A) = \left| \sum_{(i,j) \epsilon S} a_{ij} \right|.$$

The general problem is to determine the maximum value of $g_S(A)$ as $A$ varies over $O_n$. Let $H_S$ be the $n$-square matrix whose $(i, j)$ element is the number of integers $t$ such that $(i, t)$ and $(j, t)$ both belong to $S$. We have

LEMMA 1. *The matrix $H_S$ is positive semidefinite, and*

$$(1.2) \qquad 0 \leq g_S(A) \leq \operatorname{tr}\left(\sqrt{H_S}\right) \qquad \qquad for \quad A \, \epsilon \, O_n$$

*where $\sqrt{\phantom{x}}$ indicates the positive semidefinite determination of the square root. Every value between the two bounds is achievable by $g_S(A)$ for an appropriate unitary $A$.*

*Proof.* Let $P_S$ be the $n$-square matrix whose $(i, j)$ element is 1 if $(j, i) \, \epsilon \, S$ and 0 otherwise. Then

$$(1.3) \qquad g_S(A) = \left| \sum_{(i,j) \epsilon S} a_{ij} \right| = \left| \operatorname{tr}\left(P_S A\right) \right|.$$

This equality (1.3) follows immediately upon noting that the $(t, t)$ element of $P_S A$ is the sum $\sum_s a_{st}$ for $(s, t) \, \epsilon \, S$. We next apply a result of von Neumann [2, 3] to conclude that

$$(1.4) \qquad g_S(A) = \left| \operatorname{tr}\left(P_S A\right) \right| \leq \operatorname{tr}\left(\sqrt{P_S P_S'}\right) = \operatorname{tr}\left(\sqrt{H_S}\right)$$

where $A$ is any unitary matrix. For the $(i, j)$ element of $P_S P_S'$ is precisely the number of columns in which rows $i$ and $j$ of $P_S$ have a 1 in common, and thus $P_S P_S' = H_S$. We include a short proof of (1.4) for completeness. By the polar factorization theorem it is clear that we may assume

$$g_S(A) = \left| \operatorname{tr}\left(K_S A\right) \right|$$

where $K_S = \sqrt{H_S}$. Letting $x_1, \cdots, x_n$ be an orthonormal basis and setting $y_i = A x_i$, $i = 1, \cdots, n$, we have

$$g_S(A) = \left| \operatorname{tr}\left(K_S A\right) \right| = \left| \sum_{i=1}^{n} \left(K_S A x_i, x_i\right) \right|$$

$$= \left| \sum_{i=1}^{n} \left(y_i, K_S x_i\right) \right| \leq \sum_{i=1}^{n} \left| \left(y_i, K_S x_i\right) \right|$$

$$\leq \sum_{i=1}^{n} \left(K_S x_i, x_i\right)^{1/2} \left(K_S y_i, y_i\right)^{1/2}$$

$$\leq \left\{ \sum_{i=1}^{n} \left(K_S x_i, x_i\right) \right\}^{1/2} \left\{ \sum_{i=1}^{n} \left(K_S y_i, y_i\right) \right\}^{1/2}$$

$$= \{\operatorname{tr}\left(K_S\right)\}^{1/2} \{\operatorname{tr}\left(K_S\right)\}^{1/2} = \operatorname{tr}\left(\sqrt{H_S}\right).$$

To see that $g_s(O_n)$ is precisely the closed interval $[0, \text{tr }(\sqrt{H_s})]$ we note the following facts. First, $g_s(A)$ is clearly continuous with respect to the usual distance function in $O_n$, $d(U, V) = (\sum_{i,j=1}^n |u_{ij} - v_{ij}|^2)^{1/2}$. Second, if $U$ and $V$ are in $O_n$, there exists a continuous one-parameter family $A(t) \, \epsilon \, O_n$, $0 \leq t \leq 1$, such that $A(0) = U$ and $A(1) = V$; for let $S \, \epsilon \, O_n$ be chosen so that $S^{-1}(U^{-1}V)S = \text{diag }(e^{i\theta_1}, \cdots, e^{i\theta_n})$, and set

$$A(t) = US \text{ diag }(e^{i\theta_1 t}, \cdots, e^{i\theta_n t})S^{-1}.$$

Finally, $g_s(A) = 0$ for an appropriate $A \, \epsilon \, O_n$; for $g_s(A) = |\text{ tr }(K_s A)|$, and by selecting $x_1, \cdots, x_n$ to be an orthonormal set of eigenvectors of $K_s$ and choosing $A \, \epsilon \, O_n$ such that $Ax_i = x_{i+1} \pmod n$, we have

$$g_s(A) = \left| \sum_{i=1}^n (K_s Ax_i, x_i) \right| = \left| \sum_{i=1}^n \lambda_i(x_{i+1}, x_i) \right| = 0.$$

## II. Sums down pairs of diagonals

Let $S$ be a set of pairs $(i,j)$, $1 \leq i, j \leq n$, determined by two $k$-overlapping diagonals, $0 \leq k \leq n$. In this case $P_s$ has precisely one entry 1 in each of exactly $k$ rows and columns and precisely two entries 1 in each of the remaining $n - k$ rows and columns. Then

$$(2.1) \qquad H_s = P_s P_s' \simeq (I_k \dotplus (P + Q))(I_k \dotplus (P' + Q')),$$

where in (2.1) $A \simeq B$ means $RAR' = B$ for some permutation matrix $R$, $\dotplus$ indicates direct sum, $P$ and $Q$ are both $(n - k)$-square permutation matrices, and $I_k$ is the $k$-square identity matrix. Hence from (2.1) we see that

$$(2.2) \qquad H_s \simeq I_k \dotplus (2I_{n-k} + PQ' + (PQ')').$$

Thus from (2.2) we conclude that the eigenvalues of $\sqrt{H_s}$ are 1 with multiplicity $k$ and $(2 + \lambda_i + \lambda_i^{-1})^{1/2}$, $i = 1, \cdots, n - k$, where the $\lambda_i$ are the eigenvalues of the permutation matrix $PQ'$. Let $\sigma$ and $\gamma$ be the permutations on $n - k$ symbols corresponding to $P$ and $Q$ respectively. Then $\sigma\gamma^{-1}$ holds no symbol fixed; otherwise the two diagonals would be at least $(k + 1)$-overlapping. Thus $\sigma\gamma^{-1}$ has a decomposition into the product of disjoint cycles each of which has length at least 2. Let $m_1, \cdots, m_p$ be the cycle lengths in this decomposition,

$$\sum_{j=1}^p m_j = n - k, \qquad m_i \geq 2.$$

Then it is well known that the characteristic polynomial of $PQ'$ is $\sum_{i=1}^p (x^{m_i} - 1)$. Hence $PQ' + (PQ')'$ has as eigenvalues the numbers

$$(2.3) \qquad \begin{aligned} e^{2\pi i k/m_t} + e^{-2\pi i k/m_t} &= 2\cos(2\pi k/m_t), \\ k = 0, \cdots, m_t - 1, \quad &t = 1, \cdots, p. \end{aligned}$$

Since $(2 + 2\cos 2\theta)^{1/2} = 2 |\cos \theta|$ and

$$\frac{1}{2} + \sum_{t=1}^n \cos t\theta = \frac{\sin(n + \frac{1}{2})\theta}{2\sin(\theta/2)},$$

we conclude from (2.3) that

(2.4) $$\text{tr } (\sqrt{H_s}) = k + 2 \sum_{t=1}^{p} f_{m_t}$$

where

(2.5) $$f_\alpha = \sum_{k=0}^{\alpha-1} | \cos (\pi k/\alpha) | = \begin{cases} \cot (\pi/2\alpha) \text{ for } \alpha \text{ even} \\ \csc (\pi/2\alpha) \text{ for } \alpha \text{ odd.} \end{cases}$$

From (2.4) we see that the proof of the theorem depends upon finding the maximum and minimum values of

(2.6) $$\psi(\gamma) = \sum_{t=1}^{p} f_{m_t}$$

where $\gamma$ varies over all partitions of the form

(2.7)
$$\gamma : m_1 + \cdots + m_p = n - k,$$
$$m_i \geq 2 \text{ for each } i = 1, \cdots, p.$$

In the next section we state and prove the inequalities necessary to evaluate the extreme values of $\psi(\gamma)$.

## III. Some inequalities

The necessary inequalities are contained in the following lemmas.

LEMMA 2.  *If $\alpha$ and $\beta$ are even, then*

(3.1) $$f_\alpha + f_\beta \leq f_{\alpha+\beta} .$$

LEMMA 3.  *If $\alpha$ and $\beta$ are odd, then*

(3.2) $$f_\alpha + f_\beta \geq f_{\alpha+\beta} .$$

LEMMA 4.  *If $\alpha \geq 13$, then*

(3.3) $$f_3 + f_{\alpha-3} \geq f_\alpha .$$

LEMMA 5.  *If $\alpha$ is odd, then*

(3.4) $$f_2 + f_\alpha \leq f_{\alpha+2} .$$

We prove Lemma 2 first.   Let $g(x) = \cot (\pi/x)$, and note that

$$g''(x) = \frac{2\pi}{x^3} \csc^2 \frac{\pi}{x} \left( \frac{\pi}{x} \cot \frac{\pi}{x} - 1 \right),$$

and since $(\pi/x) \cot (\pi/x) < 1$ for $x > 2$, we conclude that $g(x)$ is a concave function for $x \geq 2$.   Thus

$$2g(\alpha + \beta) \geq g(2\alpha) + g(2\beta),$$

$$2 \cot \frac{\pi}{\alpha + \beta} \geq \cot \frac{\pi}{2\alpha} + \cot \frac{\pi}{2\beta},$$

and finally

$$\cot \frac{\pi}{2\alpha + 2\beta} \geq 2 \cot \frac{\pi}{\alpha + \beta} \geq \cot \frac{\pi}{2\alpha} + \cot \frac{\pi}{2\beta},$$

$$f_{\alpha+\beta} \geq f_\alpha + f_\beta .$$

This completes the proof of Lemma 2.

To prove Lemma 3 let $g(x) = \csc(\pi/x)$, and note that

$$g''(x) = \frac{\pi}{x^3}\csc\frac{\pi}{x}\left(\frac{\pi}{x}\csc^2\frac{\pi}{x} + \frac{\pi}{x}\cot^2\frac{\pi}{x} - 2\cot\frac{\pi}{x}\right).$$

Let $\theta = \pi/x > 0$, and observe that

$$\theta\left(\csc^2\theta + \cot^2\theta\right) - 2\cot\theta > 0$$

if and only if

$$h(\theta) = \theta\left(1 + \cos^2\theta\right) - \sin 2\theta > 0.$$

Now $h(0) = 0$, and $h'(\theta) > 0$ if and only if $\tan\theta > \frac{2}{3}\theta$. This shows that $g(x) = \csc(\pi/x)$ is convex for $x \geq 2$, and hence

$$f_{\alpha+\beta} = \cot\frac{\pi}{2\alpha+2\beta} \leq 2\csc\frac{\pi}{\alpha+\beta} \leq \csc\frac{\pi}{2\alpha} + \csc\frac{\pi}{2\beta} = f_\alpha + f_\beta.$$

A somewhat less direct argument is necessary for Lemma 4. First note that for $\alpha$ even Lemma 4 follows from Lemma 3; thus our problem is to show that for $\alpha$ odd

(3.5) $$2 + \cot\frac{\pi}{2(\alpha-3)} \geq \csc\frac{\pi}{2\alpha}, \qquad \alpha \geq 13.$$

We write a sequence of inequalities each of which implies its predecessor;

$$2 + \frac{2(\alpha-3)}{\pi} \geq \frac{1}{\sin(\pi/2\alpha)\cos(\pi/2(\alpha-3))},$$

$$\sin\frac{\pi(2\alpha-3)}{2\alpha(\alpha-3)} - \sin\frac{3\pi}{2\alpha(\alpha-3)} \geq \frac{\pi}{\pi-3+\alpha},$$

and from the series expansion for the sine function,

$$\frac{\pi(2\alpha-3)}{2\alpha(\alpha-3)} - \frac{\pi^3(2\alpha-3)^3}{48\alpha^3(\alpha-3)^3} \geq \frac{\pi}{\pi-3+\alpha} + \frac{3\pi}{2\alpha(\alpha-3)},$$

(3.6) $$\frac{\pi-3}{\alpha(\pi-3+\alpha)} \geq \frac{\pi^2}{6(\alpha-3)^3}.$$

Now (3.6) holds if

$$C(\alpha) = 6(\alpha-3)^3(\pi-3) - \pi^2\alpha(\pi-3+\alpha) \geq 0.$$

Now $C(23) > 0$ may be checked, and it may also be directly verified that the largest root of $C'(\alpha)$ is less than 23. Hence $C(\alpha) > 0$ for $\alpha \geq 23$, and we check separately the values $\alpha = 13, 15, 17, 19, 21$ to complete the proof. We remark that to check these values requires a table containing hundredths of a degree. We omit the similar proof of Lemma 5.

## IV. The proof of the theorem

Now (i) is clear, and (ii) follows from Lemma 2 since the only partition of 4 is $4 = 2 + 2$. Next, every integer greater than or equal to 4 is either even and greater than or equal to 6, or odd and at least 9 except for the integers 3, 4, 5, 7.

Therefore since $f_9 \leq 3f_3$ and $f_{11} \leq 2f_3 + f_5$, we see by repeated applications of (3.1), (3.2), (3.3), and (3.4) that $\psi(\gamma)$ is dominated by the value of $\psi$ on

a partition of $n - k$ which involves only the integers 3, 4, 5, 7 with appropriate multiplicities.  Checking separately that

$$f_3 + f_4 \leqq f_7, \qquad\qquad f_5 + f_5 \leqq f_3 + f_7,$$
$$f_4 + f_5 \leqq 3f_3, \qquad\qquad f_5 + f_7 \leqq 4f_3,$$
$$f_4 + f_7 \leqq 2f_3 + f_5, \qquad f_7 + f_7 \leqq 3f_3 + f_5,$$

we conclude that the value of $\psi$ on any partition of $n - k$ is dominated by its value on a partition consisting of all 3's, or all 3's and a 5, or all 3's and 7. This representation is of course unique since $n - k \equiv 0, 5,$ or 7 (mod 3). This completes the proof of (iii).  The proof of (iv) proceeds in an analogous way.  The last statement in the theorem is precisely the content of Lemma 1.

## V. Another application

In problem 4845 of the advanced problem section of the American Mathematical Monthly [1] the following question is posed: Find the maximum of $g_S(A)$ for $A \in O_n$ where $S$ is the set of $(i, j)$ satisfying $i \geq j$. We answer a generalization of this question in which we assume $i \geq j + p$, $p$ a fixed non-negative integer.  In this case $P_S$ becomes the $n$-square matrix whose $i$th row is $(0, \cdots, 0, 1, 1, \cdots, 1)$ where the first 1 appears in the $i + p$ position; if $i + p > n$, then the $i$th row of $P_S$ is the zero vector.  Then it is easy to check that

$$H_S = \begin{bmatrix} 0\cdots0\cdots\cdots\cdots\cdots0 \\ \cdot \quad \cdot \qquad\qquad\qquad \cdot \\ \cdot \quad \cdot \qquad\qquad\qquad \cdot \\ \cdot \quad \cdot \qquad\qquad\qquad \cdot \\ 0\cdots0\cdots\cdots\cdots\cdots0 \\ \cdot \quad \cdot \quad 1\cdots\cdots\cdots1 \\ \cdot \quad \cdot \quad \cdot \quad 2\cdots\cdots2 \\ \cdot \quad \cdot \quad \cdot \quad \cdot \quad 3\cdots3 \\ \cdot \quad \cdot \quad \cdot \quad \cdot \qquad \cdot \\ \cdot \quad \cdot \quad \cdot \quad \cdot \qquad \cdot \\ \cdot \quad \cdot \quad \cdot \quad \cdot \qquad \cdot \\ 0\cdots0 \quad 1 \quad 2 \quad 3\cdots n-p \end{bmatrix}.$$

Now we observe that if

$$C_S = \begin{bmatrix} 0\cdots\cdots0\cdots\cdots\cdots\cdots0 \\ \cdot \qquad\quad \cdot \qquad\qquad\qquad \cdot \\ \cdot \qquad\quad \cdot \qquad\qquad\qquad \cdot \\ \cdot \qquad\quad \cdot \qquad\qquad\qquad 0 \\ 0\cdots\cdots0\cdots\cdots\cdots0 \quad 1 \\ \cdot \qquad\quad \cdot \qquad\qquad 1 \quad 1 \\ \cdot \qquad\quad \cdot \qquad\qquad\quad \cdot \\ \cdot \qquad\quad \cdot \qquad \cdot \qquad \cdot \\ \cdot \qquad\quad 0 \quad 1 \qquad\qquad \cdot \\ 0\cdots0 \quad 1 \quad 1\cdots\cdots\cdots1 \end{bmatrix},$$

then $C_S^2 = H_S$.

Thus by Lemma 1 we know that the maximum of the function $g_s(A)$ for $A \epsilon O_n$ is $\sum_{j=1}^{n-p} | \lambda_j |$ where $\lambda_j$, $j = 1, \cdots n - p$, are the $n - p$ nonzero eigenvalues of $C_s$. In a private communication Professor A. C. Aitken proved the following result:

$$\lambda_j = \frac{(-1)^{j-1}}{2} \csc \frac{(2j - 1)\pi}{4(n - p) + 2}, \qquad j = 1, \cdots, n - p.$$

These values were obtained by the very elegant observation that the inverse of the lower right nonsingular $(n - p)$-square block of $C_s$ is a differencing matrix whose eigenvectors can be readily computed.

We then can conclude finally that

$$\max_{A \epsilon O_n} \left| \sum_{i \geq j+p} a_{ij} \right| = \frac{1}{2} \sum_{j=1}^{n-p} \left| \csc \frac{(2j - 1)\pi}{4(n - p) + 2} \right|.$$

REFERENCES

1. F. KOEHLER, *Problem 4845*, Amer. Math. Monthly, vol. 66 (1959), p. 426.
2. M. MARCUS AND B. N. MOYLS, *On the maximum principle of Ky Fan*, Canadian J. Math., vol. 9 (1957), pp. 313–320.
3. J. VON NEUMANN, *Some matrix-inequalities and metrization of matric-space*, Tomsk. Gos. Univ. Uč. Zap., vol. 1 (1937), pp. 286–299.

THE UNIVERSITY OF BRITISH COLUMBIA
VANCOUVER, CANADA