

THE SUM OF THE ELEMENTS OF THE POWERS OF A MATRIX

MARVIN MARCUS AND MORRIS NEWMAN

1. Introduction and results. In the first two sections of this paper A will be assumed to be an irreducible nonnegative n -square matrix; $A \geq 0$. Let $s_k = s_k(A)$ denote the sum of the entries in the matrix A^k , where k is a positive integer. The problem considered in the first section is the convergence of the ratio s_k/s_{k-1} as $k \rightarrow \infty$. In §3 we obtain an inequality relating the s_k for various k in the case A is a Hermitian matrix and in §4 we discuss convexity properties of s_2/s_1 .

Let λ_1 be the dominant positive characteristic root of A which can be taken as 1 for the purposes of our subsequent arguments. If h is the number of characteristic roots of A of modulus 1, then they are the roots of $\lambda^h - 1 = 0$ and are all simple [3]. Let $\varepsilon = e^{2\pi i/h}$ so that $1, \varepsilon, \varepsilon^2, \dots, \varepsilon^{h-1}$ are the roots of modulus 1. Choose permutation matrices P and Q so that

$$(1) \quad PAP^T = \begin{bmatrix} 0 & A_1 & & & 0 \\ & 0 & A_2 & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & A_{h-1} \\ A_h & & & & 0 \end{bmatrix}$$

and

$$(2) \quad QA^TQ^T = \begin{bmatrix} 0 & B_1 & & & 0 \\ & 0 & B_2 & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & B_{h-1} \\ B_h & & & & 0 \end{bmatrix}$$

where the zero blocks down the main diagonal in both (1) and (2) are square. We shall assume henceforth that A is in this *Frobenius normal form*. In other words we assume A is already in the form given on the right in (1). Let u_1, \dots, u_h and v_1, \dots, v_h be the characteristic vectors of A and A^T corresponding to $1, \varepsilon, \dots, \varepsilon^{h-1}$ respectively. We write for the *maximal characteristic vector*

$$(3) \quad u_1 = z_1 + \dots + z_h,$$

Received September 27, 1961. This work was supported in part by the Office of Naval Research. The authors wish to express their thanks to K. Goldberg for many valuable suggestions.

where $\dot{+}$ indicates direct sum and the partitioning in (3) is conformal to the partitioning into main diagonal blocks in (1). That is, if the square blocks down the main diagonal in (1) are of sizes n_1, \dots, n_h then z_i has coordinates numbered $n_{i-1} + 1, \dots, n_i$ ($n_0 = 0$) all positive, the rest zero. Thus no two of the z_i have positive coordinates in the same position and the direct sum notation is appropriate.

Since Qv_1 is the maximal characteristic vector of QA^xQ^x we can partition Qv_1 into a direct sum exactly as was done with u_1 : $Qv_1 = m_1 \dot{+} \dots \dot{+} m_h$. Then $v_1 = Q^x m_1 \dot{+} \dots \dot{+} Q^x m_h$ and we set $Q^x m_t = w_t$, $t = 1, \dots, h$, to obtain

$$(4) \quad v_1 = w_1 \dot{+} \dots \dot{+} w_h .$$

Let $r = r(A)$ and $c = c(A)$ be the n -tuples of row and column sums of A respectively. Partition r and c conformally with v_1 and u_1 respectively as in (4) and (3):

$$\begin{aligned} r &= r_1 \dot{+} \dots \dot{+} r_h \\ c &= c_1 \dot{+} \dots \dot{+} c_h . \end{aligned}$$

The notations (z_i, c_i) and (w_i, r_i) will be used for the ordinary euclidean inner product.

Our main result is in terms of these inner products.

THEOREM 1. $\lim_{k \rightarrow \infty} s_k / s_{k-1}$ exists and is equal to the dominant characteristic root if and only if the numbers

$$(5) \quad \sum_{i=1}^h (z_i, c_i)(w_{\alpha-i+1}, r_{\alpha-i+1})$$

are all equal for $\alpha = 1, \dots, h$.

We remark that the indices in (5) are to be reduced modulo h .

In case A is symmetric then the roots of modulus 1 can be only 1 or -1 . Thus $h = 1$ or 2. In case $h = 1$ (A is primitive) then Theorem 1 automatically holds since there is only one item (5). In case $h = 2$ we have.

THEOREM 2. If

$$A = \begin{bmatrix} O_p & C \\ C^x & O_q \end{bmatrix}$$

is irreducible and has maximal characteristic vector

$$u_1 = z_1 \dot{+} z_2 = (a_1, \dots, a_p) \dot{+} (b_1, \dots, b_q)$$

$$(9) \quad s_k = \text{tr}(JA^k) = \text{tr}\left(JS\left(\sum_{t=1}^h \varepsilon^{(t-1)k} E_{tt}\right)S^{-1}\right) + \delta_k .$$

Since $|\lambda_t| < 1$ for $t \geq h + 1$ we conclude that $\lim_{k \rightarrow \infty} \delta_k = 0$. If x is an n -vector it will be convenient to denote by $\sigma(x)$ the sum of the coordinates of x . Moreover, S_t and S^t will designate the t th row and column of the n -square matrix S respectively. Let $\mu_t = (S^{-1}JS)_{tt}$, $t = 1, \dots, h$ and note that

$$(10) \quad \begin{aligned} s_k - \delta_k &= \sum_{t=1}^h \text{tr}(E_{tt}S^{-1}JS)\varepsilon^{(t-1)k} \\ &= \sum_{t=1}^h (S^{-1}JS)_{tt}\varepsilon^{(t-1)k} \\ &= \sum_{t=1}^h \mu_t \varepsilon^{(t-1)k} , \end{aligned}$$

and

$$(11) \quad \begin{aligned} \mu_t &= \sum_{\alpha, \beta=1}^n (S^{-1})_{t\alpha} J_{\alpha\beta} S_{\beta t} \\ &= \sum_{\alpha, \beta} (S^{-1})_{t\alpha} S_{\beta t} \\ &= \sigma(S_t^{-1})\sigma(S^t) , \end{aligned} \quad t = 1, \dots, h .$$

From (8) we have

$$AS^t = \varepsilon^{t-1}S^t , \quad t = 1, \dots, h$$

and since the dimension of the null space of $A - \varepsilon^{t-1}I$ is ≤ 1 for $t = 1, \dots, h$ we conclude that

$$(12) \quad S^t = c_t u_t , \quad t = 1, \dots, h$$

for appropriate nonzero scalars c_t . Similarly

$$(13) \quad S_t^{-1} = d_t v_t , \quad t = 1, \dots, h .$$

From (11), (12) and (13) we have, for e the n -tuple all of whose coordinates is 1,

$$(14) \quad \begin{aligned} \mu_t &= \sigma(S_t^{-1})\sigma(S_t) = d_t c_t \sigma(v_t)\sigma(u_t) \\ &= d_t c_t (v_t, e)(u_t, e) \\ &= d_t c_t (A^T(v_t/\varepsilon^{t-1}), e)(A(u_t/\varepsilon^{t-1}), e) \\ &= d_t c_t \varepsilon^{2(1-t)}(v_t, Ae)(u_t, A^T e) \\ &= d_t c_t \varepsilon^{2(1-t)}(v_t, r)(u_t, c) . \end{aligned}$$

The vectors u_t and v_t , $t = 2, \dots, h$, have explicit representations in terms of u_1, v_1 and ε as follows:

$$\begin{aligned}
 u_t &= z_1 + \varepsilon^{t-1}z_2 + \dots + \varepsilon^{(t-1)(h-1)}z_h \\
 v_t &= w_1 + \varepsilon^{(t-1)}w_2 + \dots + \varepsilon^{(t-1)(h-1)}w_h \quad t = 2, \dots, h.
 \end{aligned}$$

Let $k_t = d_t c_t \varepsilon^{2(1-t)} \neq 0$, $\zeta_i = (z_i, c_i)$, $\eta_i = (w_i, r_i)$, $i = \dots, h$ and we compute from (14) and the fact that $\varepsilon^{(t-1)h} = 1$ that

$$\begin{aligned}
 (15) \quad \mu_t &= k_t(v_t, r)(u_t, c) \\
 &= k_t \sum_{i=1}^h (w_i, r_i) \varepsilon^{(t-1)(i-1)} \sum_{i=1}^h (z_i, c_i) \varepsilon^{(t-1)(i-1)} \\
 &= k_t \sum_{\alpha=1}^h \left(\sum_{i=1}^h \zeta_i \eta_{\alpha-i+1} \right) \varepsilon^{(t-1)(\alpha-1)},
 \end{aligned}$$

where the subscripts are always reduced modulo h .

LEMMA 1. $\mu_2 = \dots = \mu_h = 0$ if and only if the sums $\sum_{i=1}^h \zeta_i \eta_{\alpha-i+1}$ are all equal for $\alpha = 1, \dots, h$.

Proof. Set $f_\alpha = \sum_{i=1}^h \zeta_i \eta_{\alpha-i+1}$ and from (15) the conditions $\mu_t = 0$, $t = 2, \dots, h$ are equivalent to the system of linear equations

$$(16) \quad \sum_{\alpha=1}^h f_\alpha \varepsilon^{(\alpha-1)(t-1)} = 0 \quad t = 2, \dots, h.$$

Since $\sum_{\alpha=1}^h \varepsilon^{(\alpha-1)(t-1)} = 0$ ($1 \leq t-1 < h$), each of the equations (16) has the solution $f_1 = \dots = f_h$. On the other hand, the $(h-1)$ -square submatrix of coefficients in (16) obtained by deleting the first column in the coefficient matrix has as its determinant the Vandermonde of $\varepsilon, \varepsilon^2, \dots, \varepsilon^{h-2}$ to within a nonzero constant multiple. Thus the system (16) has rank $h-1$ and $f_1 = \dots = f_h$ is the *only* solution of (16). The proof of Theorem 1 will then be complete if we establish

LEMMA 2. $\lim_{k \rightarrow \infty} s_k/s_{k-1}$ exists if and only if $\mu_2 = \dots = \mu_h = 0$. If it exists it has value 1.

Proof. From (10) we have

$$s_k/s_{k-1} = \left(\sum_{t=1}^h \mu_t \varepsilon^{(t-1)k} + \delta_k \right) / \left(\sum_{t=1}^h \mu_t \varepsilon^{(t-1)(k-1)} + \delta_{k-1} \right),$$

and since $\lim_{k \rightarrow \infty} \delta_k = 0$, $\lim_{k \rightarrow \infty} s_k/s_{k-1}$ exists if and only if

$$g_k = m_k/m_{k-1} = \frac{\sum_{t=1}^h \mu_t \varepsilon^{(t-1)k}}{\sum_{t=1}^h \mu_t \varepsilon^{(t-1)(k-1)}}$$

approaches a limit. Note that m_k is periodic of period h . Also $\mu_1 = c_1 d_1 \sigma(u_1) \sigma(v_1) \neq 0$ follows from (3) and (4) and so the condition is clearly sufficient. Since g_k takes on only a finite number of values it follows

that if g_k approaches a limit l then $g_k = l$ for all k . Moreover if g_k exists,

$$\begin{aligned} g_k &= m_k/m_{k-1} = m_k/m_{k-h-1} \\ &= \prod_{\alpha=1}^{h+1} (m_{k-\alpha+1}/m_{k-\alpha}) \\ &= g_k^{h+1}, \end{aligned}$$

and thus $l = 1$. But then $m_k = m_{k-1}$ and we conclude that

$$(17) \quad \sum_{t=2}^h \mu_t (1 - \varepsilon^{(1-t)}) \varepsilon^{(t-1)k} = 0.$$

Letting $k = 0, \dots, h - 2$ successively in (17) and noting that $\prod_{0 \leq i < j \leq h-2} (\varepsilon^i - \varepsilon^j) \neq 0$ we conclude that $\mu_t (1 - \varepsilon^{1-t}) = 0, t = 2, \dots, h$, and hence that $\mu_2 = \dots = \mu_h = 0$.

To proceed to the proof of Theorem 2 note that the maximal characteristic vectors of A and $A^x = A$ are given by

$$\begin{aligned} u_1 &= v_1 = z_1 + z_2 = (a_1, \dots, a_p) + (b_1, \dots, b_q). \\ c_1 &= (\sigma(C_1), \dots, \sigma(C_p)) \\ c_2 &= (\sigma(C^1), \dots, \sigma(C^q)), \text{ and} \\ r_1 &= c_1, r_2 = c_2. \end{aligned}$$

The condition that the items (5) be equal for $\alpha = 1, 2$ becomes, in succession,

$$\begin{aligned} (z_1, c_1)(z_2, c_2) + (z_2, c_2)(z_1, c_1) &= (z_1, c_1)^2 + (z_2, c_2)^2, \\ (z_1, c_1) &= (z_2, c_2), \\ \sum_{i=1}^p a_i \sigma(c_i) &= \sum_{i=1}^q b_i \sigma(C^i), \\ (18) \quad \sigma\left(\sum_{i=1}^p a_i C_i\right) &= \sigma\left(\sum_{i=1}^q b_i C^i\right). \end{aligned}$$

Now $Cb = a, C^x a = b$ and hence $a = \sum_{i=1}^q b_i C^i, b = \sum_{i=1}^p a_i C_i$. We then have from (18) that $\sigma(a) = \sigma(b)$ is equivalent to (5) in the case A symmetric and $h = 2$.

The convergence of s_k/s_{k-1} in Theorem 3 is clear since $h = 1$. If A is positive semi-definite and $\alpha \geq 0$ let A^α be the unique positive semi-definite determination. Then if p and q are nonnegative,

$$\begin{aligned} (19) \quad s_{(p+q)/2}^2 &= (A^{p+q/2}e, e)^2 \\ &= (A^{p/2}e, A^{q/2}e)^2 \\ &\leq (A^p e, e)(A^q e, e) \\ &= s_p s_q \end{aligned}$$

with equality if and only if $A^{p/2}e$ and $A^{q/2}e$ are linearly dependent. Set $p = k - 1$ and $q = k + 1$ to finish the argument.

The Corollary follows from Theorem 3.

3. The Hermitian case. In this section A is assumed to be an n -square Hermitian matrix with characteristic roots $\lambda_1, \dots, \lambda_n$. We have

THEOREM 4. *Let p, q, m, t be nonnegative integers and assume that $t = \min(p, q, m, t)$ is even and $p + q = m + t$ is even. Then*

$$(20) \quad s_p s_q \leq s_m s_t .$$

Proof. Let $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ and choose a unitary matrix U such that $U^*AU = D$. Then

$$s(A) = (1/n)\text{tr}(JAJ) = (1/n)\text{tr}(JUDU^*J) .$$

It is not difficult to compute that each main diagonal element of $JUDU^*J$ is $\sum_{i=1}^n \lambda_i |\sigma(U^i)|^2$. Let $w_i = |\sigma(U^i)|^2$ and then

$$(21) \quad s(A) = \sum_{i=1}^n \lambda_i w_i .$$

Replacing A by A^p in (21) yields

$$s_p = \sum_{i=1}^n \lambda_i^p w_i$$

and (20) is equivalent to

$$(22) \quad \sum_{i=1}^n \lambda_i^m w_i \sum_{i=1}^n \lambda_i^t w_i - \sum_{i=1}^n \lambda_i^p w_i \sum_{i=1}^n \lambda_i^q w_i \geq 0 .$$

The left side of (22) becomes, after symmetrizing the sums,

$$\begin{aligned} & (1/2) \sum_{i,j} w_i w_j (\lambda_i \lambda_j)^t [\lambda_j^{m-t} + \lambda_i^{m-t} - \lambda_i^{p-t} \lambda_j^{q-t} - \lambda_i^{q-t} \lambda_j^{p-t}] \\ &= (1/2) \sum_{i,j} w_i w_j (\lambda_i \lambda_j)^t (\lambda_j^{p-t} - \lambda_i^{p-t})(\lambda_j^{q-t} - \lambda_i^{q-t}) \\ &= (1/2) \sum_{i,j} w_i w_j (\lambda_i \lambda_j)^t \lambda_i^{p+q-2t} ((\lambda_j/\lambda_i)^{p-t} - 1)((\lambda_j/\lambda_i)^{q-t} - 1) . \end{aligned}$$

Since t and $p + q - 2t$ are even and moreover $p - t \geq 0, q - t \geq 0$ it follows that $(\lambda_i \lambda_j)^t \lambda_i^{p+q-2t} ((\lambda_j/\lambda_i)^{p-t} - 1)((\lambda_j/\lambda_i)^{q-t} - 1) \geq 0$ for all i, j .

Setting $t = 0$ in Theorem 4 yields the

COROLLARY. *If p and q are nonnegative integers and $p + q$ is even then*

$$(23) \quad s_p s_q \leq n s_{p+q} .$$

we remark that formula (21) doesn't yield $s_0 = \sum_{i=1}^n w_i = n$ unless A is nonsingular. But the singular case follows from the nonsingular one by the standard continuity argument.

In case A has nonnegative entries then a specialization of an inequality in [1] implies that $n^2 s_3 \geq s_1$. We conjecture that $n s_3 \geq s_1 s_2$ in analogy with (23).

4. Some remarks on s_2/s_1 .

Let $f(t) = f(t_1, \dots, t_n) = \sum_{1 \leq i < j \leq n} t_i t_j / \sum_{i=1}^n t_i$ and note that

$$\begin{aligned} f(a+b) - f(a) - f(b) \\ = \sum_{i=1}^n \left(a_i \sum_{j=1}^n b_j - b_i \sum_{j=1}^n a_j \right)^2 / 2 \sum_{i=1}^n (a_i + b_i) \sum_{i=1}^n a_i \sum_{i=1}^n b_i . \end{aligned}$$

It follows that if $\sum_{i=1}^n a_i > 0$ and $\sum_{i=1}^n b_i > 0$ then $f(a+b) \geq f(a) + f(b)$ with equality if and only if the sets (a_1, \dots, a_n) and (b_1, \dots, b_n) are proportional. Define the functions

$$\begin{aligned} g(t) &= \frac{\sum_{i=1}^n t_i^2}{\sum_{i=1}^n t_i} , \\ h(t) &= \left(\sum_{i=1}^n t_i^2 + \sum_{1 \leq i < j \leq n} t_i t_j \right) / \sum_{i=1}^n t_i , \end{aligned}$$

and observe that

$$\begin{aligned} g(t) &= \sum_{i=1}^n t_i - 2f(t) , \\ h(t) &= \sum_{i=1}^n t_i - f(t) . \end{aligned}$$

Then if $\sum_{i=1}^n a_i > 0$, $\sum_{i=1}^n b_i > 0$,

$$\begin{aligned} (24) \quad g(a+b) &= \sum_{i=1}^n a_i + \sum_{i=1}^n b_i - 2f(a+b) \\ &\leq \sum_{i=1}^n a_i - 2f(a) + \sum_{i=1}^n b_i - 2f(b) \\ &= g(a) + g(b) , \end{aligned}$$

and similarly $h(a+b) \leq h(a) + h(b)$.

Equality holds in the preceding two inequalities if and only if the sets a and b are proportional. From the inequality (24) we can then prove

THEOREM 5. *If A and B are symmetric n -square matrices satisfying $\sum_{i,j=1}^n a_{ij} > 0$, $\sum_{i,j=1}^n b_{ij} > 0$ then*

$$(25) \quad s_2(A+B)/s_1(A+B) \leq s_2(A)/s_1(A) + s_2(B)/s_1(B)$$

with equality if and only if $r(A)$ and $r(B)$ are proportional.

Proof. From the formula

$$s_2(A)/s_1(A) = \frac{\sum_{i=1}^n \sigma(A_i)^2}{\sum_{i=1}^n \sigma(A_i)} = g(r(A))$$

we compute that

$$\begin{aligned} s_2(A+B)/s_1(A+B) &= g(r(A+B)) = g(r(A) + r(B)) \\ &\leq g(r(A)) + g(r(B)) \\ &= s_2(A)/s_1(A) + s_2(B)/s_1(B). \end{aligned}$$

A similar result can be formulated for the function h . It might be conjectured that a convexity result like (25) is true for the functions $s_r(A)/s_{r-1}(A)$, $r > 2$. This is not the case: take

$$a_1 = \frac{1}{3}, a_2 = \dots = a_{n-1} = 0, a_n = 1, b_1 = \dots = b_{n-1} = 0, \\ b_n = 1 \text{ and observe that}$$

$$\frac{\sum_{i=1}^n (a_i + b_i)^r}{\sum_{i=1}^n (a_i + b_i)^{r-1}} = (3^{-r} + 2^r)/(3^{-(r-1)} + 2^{r-1})$$

whereas

$$\frac{\sum_{i=1}^n a_i^r}{\sum_{i=1}^n a_i^{r-1}} + \frac{\sum_{i=1}^n b_i^r}{\sum_{i=1}^n b_i^{r-1}} = (3^{-r} + 1)(3^{-(r-1)} + 1) + 1$$

and it is simple to check that

$$(3^{-r} + 2^r)/(3^{-(r-1)} + 2^{r-1}) > (3^{-r} + 1)/(3^{-(r-1)} + 1) + 1$$

for $r \geq 3$.

The referee suggests that the arguments of the paper could be rephrased in terms of the vector e , where $e^x = (1, 1, \dots, 1)$. Thus $s_k = e^x A^k e$, $J = e e^x$, $\sigma(x) = e^x x$, etc. He also notes that e could be replaced by any other positive vector with conclusions similar to those obtained in the paper. We have not thought it advisable to pursue the matter further.

REFERENCES

1. F. V. Atkinson, G. A. Watterson, P. A. P. Moran, *A matrix inequality*, Quart J. Math., Oxford (2), **11**, (1960), 137-140.
2. G. H. Hardy, J. E. Littlewood, G. Pólya, *Inequalities* (Cambridge, Second Edition, 1952).
3. H. Wielandt, *Unzerlegbare, nicht negative Matrizen*, Math Zeit. **52** (1950), 642-648.

THE UNIVERSITY OF CALIFORNIA, SANTA BARBARA AND
THE UNITED STATES NATIONAL BUREAU OF STANDARDS

