

A BOUND FOR CONDITION NUMBERS OF MATRICES*

MICHAEL I. GIL[†]

Abstract. Let A be a diagonalizable matrix; so there is an invertible matrix T and a normal matrix \hat{D} , such that $T^{-1}AT = \hat{D}$. A sharp bound for the constant $\kappa_T = \|T\|\|T^{-1}\|$ is suggested. Some applications of the obtained bound are also discussed.

Key words. Matrix, Similarity, Condition number.

AMS subject classifications. 15A45, 15A42.

1. Introduction and preliminaries. Let \mathbb{C}^n be the n -dimensional complex Euclidean space with a scalar product (\cdot, \cdot) , the Euclidean norm $\|\cdot\| = \sqrt{(\cdot, \cdot)}$ and the identity matrix I . For an $n \times n$ matrix A , $\sigma(A)$ denotes the spectrum of A , $\|A\|$ is the spectral norm; A^* is the adjoint to A ; $\|A\|_F = (\text{Trace } A^*A)^{1/2}$ is the Frobenius norm; λ_k ($k = 1, \dots, n$) are the eigenvalues of A . Everywhere below it is assumed that

$$\lambda_j \neq \lambda_m, \quad \text{whenever } j \neq m. \quad (1.1)$$

So, A is a diagonalizable matrix: There is an invertible matrix T and a normal matrix \hat{D} , such that

$$T^{-1}AT = \hat{D}. \quad (1.2)$$

The condition number $\kappa_T := \|T\|\|T^{-1}\|$ is very important for various applications, cf. [3, 16]. That number is mainly numerically calculated.

In the present paper, we suggest a sharp bound for κ_T . Applications of the obtained bound to spectrum perturbations and matrix functions are also discussed.

The following quantity (departure from normality) plays an essential role hereafter:

$$g(A) := \left[\|A\|_F^2 - \sum_{k=1}^n |\lambda_k|^2 \right]^{1/2}.$$

*Received by the editors on November 28, 2013. Accepted for publication on February 22, 2014.
 Handling Editor: Michael Tsatsomeros.

[†]Department of Mathematics, Ben Gurion University of the Negev, PO Box 653, Beer-Sheva 84105, Israel (gilmi@bezeqint.net).

$g(A)$ enjoys the following properties:

$$g^2(A) \leq 2\|A_I\|_F^2 \quad (A_I = (A - A^*)/2i) \quad \text{and} \quad g^2(A) \leq \|A\|_F^2 - |\text{Trace } A^2|, \quad (1.3)$$

cf. [10, Section 2.1]. If A is normal, then $g(A) = 0$. Put

$$\delta := \min_{j,k=1,\dots,n; k \neq j} |\lambda_j - \lambda_k|.$$

Corollary 3.6 from [11], under condition (1.1) gives us the inequality

$$\kappa_T \leq n \sum_{k=0}^{n-1} \frac{g^k(A) 2^k}{\delta^k \sqrt{k!}}. \quad (1.4)$$

That inequality is not sharp: If A is a normal matrix, then it gives $\kappa_T \leq n$, but $\kappa_T = 1$ in this case. Inequality (1.4) has been slightly improved in [12]. In this paper, we considerably refine (1.4) and the corresponding result from [12].

Put

$$\tau(A) := \sum_{k=0}^{n-2} \frac{g^{k+1}(A)}{\sqrt{k!} \delta^{k+1}}$$

and

$$\gamma(A) := \left(1 + \frac{\tau(A)}{\sqrt{n-1}}\right)^{2(n-1)}.$$

Now we are in a position to formulate the main result of this paper.

THEOREM 1.1. *Let condition (1.1) be fulfilled. Then there is an invertible matrix T , such that (1.2) holds with*

$$\kappa_T \leq \gamma(A). \quad (1.5)$$

The proof of this theorem is presented in the next two sections. Theorem 1.1 is sharp: If A is normal, then $g(A) = 0$ and $\gamma(A) = 1$. Thus, we obtain the equality $\kappa_T = 1$. So Theorem 1.1 is obviously sharper than (1.4) at least for matrices “close” to normal ones. The proof of Theorem 1.1 is absolutely different from the proof of inequality (1.4) and the proof of the corresponding result from [12].

Theorem 1.1 supplements the interesting recent investigations of the similarity of matrices, cf. [4, 5, 9, 13] and references therein.

2. Auxiliary results. Let matrix A have in \mathbb{C}^n a chain of invariant projections P_k ($k = 1, \dots, m$; $m \leq n$):

$$0 \subset P_1\mathbb{C}^n \subset P_2\mathbb{C}^n \subset \dots \subset P_m\mathbb{C}^n = \mathbb{C}^n \quad (2.1)$$

and

$$P_k A P_k = A P_k \quad (k = 1, \dots, m). \quad (2.2)$$

Put $\Delta P_k = P_k - P_{k-1}$ ($P_0 = 0$), $A_k = \Delta P_k A \Delta P_k$,

$$Q_k = I - P_k, \quad B_k = Q_k A Q_k \quad \text{and} \quad C_k = \Delta P_k A Q_k.$$

It is assumed that the spectra $\sigma(A_k)$ of A_k in $\Delta P_k\mathbb{C}^n$ satisfies the condition

$$\sigma(A_k) \cap \sigma(A_j) = \emptyset \quad (j \neq k). \quad (2.3)$$

LEMMA 2.1. *One has*

$$\sigma(A) = \cup_{k=1}^m \sigma(A_k).$$

Proof. Put

$$S = \sum_{k=1}^m A_k \quad \text{and} \quad W = A - S.$$

Due to (2.2), we have $W P_k = P_{k-1} W P_k$. Hence,

$$\begin{aligned} W^m &= W^m P_m = W^{m-1} P_{m-1} W P_m = W^{m-2} P_{m-2} W P_{m-1} W P_m \\ &= W^{m-2} P_{m-2} W^2 = W^{m-3} P_{m-3} W^3 = \dots = P_0 W^m = 0. \end{aligned}$$

So, W is nilpotent. Similarly, taking into account that

$$(S - \lambda I)^{-1} W P_k = P_{k-1} (S - \lambda I)^{-1} W P_k,$$

we prove that $((S - \lambda I)^{-1} W)^m = 0$ ($\lambda \notin \sigma(S)$). Thus,

$$\begin{aligned} (A - \lambda I)^{-1} &= (S + W - \lambda I)^{-1} = (I + (S - \lambda I)^{-1} W)^{-1} (S - \lambda I)^{-1} \\ &= \sum_{k=0}^{m-1} (-1)^k ((S - \lambda I)^{-1} W)^k (S - \lambda I)^{-1}. \end{aligned}$$

Hence, it easily follows that $\sigma(S) = \sigma(A)$. This proves the lemma. \square

Since B_j is a block triangular matrix, according to the previous lemma, we have

$$\sigma(B_j) = \cup_{k=j+1}^m \sigma(A_k) \quad (j = 0, \dots, m-1).$$

So, due to (2.3),

$$\sigma(B_j) \cap \sigma(A_j) = \emptyset.$$

Under this condition, the equation

$$A_j X_j - X_j B_j = -C_j \quad (j = 1, \dots, m-1). \quad (2.4)$$

has a unique solution, e.g., [1, Section VII.2] or [2].

LEMMA 2.2. *Let condition (2.3) hold and X_j be a solution to (2.4). Then*

$$\begin{aligned} (I - X_{m-1})(I - X_{m-2}) \cdots (I - X_1) A (I + X_1)(I + X_2) \cdots (I + X_{m-1}) \\ = A_1 + A_2 + \cdots + A_m. \end{aligned} \quad (2.5)$$

Proof. Since $X_j = \Delta P_j X_j Q_j$, we have $X_j A_j = B_j X_j = X_j C_j = C_j X_j = 0$. Due to (2.2), $Q_j A P_j = 0$. Thus, $A = A_1 + B_1 + C_1$, and consequently,

$$(I - X_1)A(I + X_1) = (I - X_1)(A_1 + B_1 + C_1)(I + X_1) =$$

$$A_1 + B_1 + C_1 - X_1 B_1 + A_1 X_1 = A_1 + B_1.$$

Furthermore, $B_1 = A_2 + B_2 + C_2$. Hence,

$$(Q_1 - X_2)B_1(Q_1 + X_2) = (Q_1 - X_1)(A_2 + B_2 + C_2)(Q_1 + X_1) =$$

$$A_2 + B_2 + C_2 - X_2 B_2 + A_2 X_2 = A_2 + B_2.$$

Therefore,

$$(I - X_2)(A_1 + B_1)(I + X_2) = (P_1 + Q_1 - X_2)(A_1 + B_1)(P_1 + Q_1 + X_2) =$$

$$A_1 + (Q_1 - X_2)(A_1 + B_1)(Q_1 + X_2) = A_1 + A_2 + B_2.$$

Consequently,

$$(I - X_2)(A_1 + B_1)(I + X_2) = (I - X_2)(I - X_1)A(I + X_1)(I + X_2) = A_1 + A_2 + B_2.$$

Continuing this process and taking into account that $B_{m-1} = A_m$, we obtain the required result. \square

Take

$$T = (I + X_1)(I + X_2) \cdots (I + X_{m-1}). \quad (2.6)$$

It is simple to see that the inverse to $I + X_j$ is the matrix $I - X_j$. Thus,

$$T^{-1} = (I - X_{m-1})(I - X_{m-2}) \cdots (I - X_1) \quad (2.7)$$

and (2.5) can be written as

$$T^{-1}AT = \text{diag}(A_{kk})_{k=1}^m. \quad (2.8)$$

By the inequalities between the arithmetic and geometric means, we get

$$\|T\| \leq \prod_{k=1}^{m-1} (1 + \|X_k\|) \leq \left(1 + \frac{1}{m-1} \sum_{k=1}^{m-1} \|X_k\|\right)^{m-1} \quad (2.9)$$

and

$$\|T^{-1}\| \leq \left(1 + \frac{1}{m-1} \sum_{k=1}^{m-1} \|X_k\|\right)^{m-1}. \quad (2.10)$$

3. Proof of Theorem 1.1. Let $\{e_k\}$ be the Schur basis (the orthogonal normal basis of the triangular representation) of matrix A :

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & a_{22} & a_{23} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & a_{nn} \end{pmatrix}$$

with $a_{jj} = \lambda_j$ in that basis. Besides,

$$\sum_{k=2}^n \sum_{i=1}^{k-1} |a_{ik}|^2 = g^2(A).$$

Take $P_j = \sum_{k=1}^j (\cdot, e_k)e_k$. Then one can apply Lemma 2.2 with $m = n$, $\Delta P_k = (\cdot, e_k)e_k$,

$$Q_j = \sum_{k=j+1}^n (\cdot, e_k)e_k, A_k = \Delta P_k A \Delta P_k = \lambda_k \Delta P_k, \quad \text{diag}(A_{kk})_{k=1}^n = \text{diag}(\lambda_k)_{k=1}^n,$$

$$B_j = Q_j A Q_j = \begin{pmatrix} a_{j+1,j+1} & a_{j+1,j+2} & \cdots & a_{j+1,n} \\ 0 & a_{j+2,j+2} & \cdots & a_{j+2,n} \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{pmatrix}$$

and

$$C_j = \Delta P_j A Q_j = \begin{pmatrix} a_{j,j+1} & a_{j,j+2} & \cdots & a_{j,n} \end{pmatrix}.$$

Besides,

$$A = \begin{pmatrix} \lambda_1 & C_1 \\ 0 & B_1 \end{pmatrix}, B_1 = \begin{pmatrix} \lambda_2 & C_2 \\ 0 & B_2 \end{pmatrix}, \dots, B_j = \begin{pmatrix} \lambda_{j+1} & C_{j+1} \\ 0 & B_{j+1} \end{pmatrix} \quad (j < n).$$

So, B_j is an upper-triangular $(n-j) \times (n-j)$ matrix. Equation (2.4) takes the form $\lambda_j X_j - X_j B_j = -C_j$. Since $X_j = X_j Q_j$, we can write $X_j(\lambda_j Q_j - B_j) = C_j$. Therefore,

$$X_j = C_j(\lambda_j Q_j - B_j)^{-1}. \quad (3.1)$$

The inverse matrix is understood in the sense of subspace $Q_j \mathbb{C}^n$. Hence,

$$\|X_j\| \leq \|C_j\| \|(\lambda_j Q_j - B_j)^{-1}\|.$$

Besides,

$$\|C_j\|^2 = \sum_{k=j+1}^n |a_{jk}|^2,$$

and due to [10, Corollary 2.2.2], we have

$$\|(\lambda_j Q_j - B_j)^{-1}\| \leq \sum_{k=0}^{n-j-1} \frac{g^k(B_j)}{\sqrt{k!} \delta^{k+1}} \quad (j = 1, 2, \dots, n-1).$$

But

$$g^2(B_j) = g^2(Q_j A Q_j) = \sum_{k=j+2}^n \sum_{i=j+1}^{k-1} |a_{ik}|^2 \leq g^2(A).$$

So, with the notation

$$\tau_1(A) := \sum_{k=0}^{n-2} \frac{g^k(A)}{\sqrt{k!} \delta^{k+1}},$$

we have

$$\|(\lambda_j Q_j - B_j)^{-1}\| \leq \tau_1(A) \quad \text{and} \quad \|X_j\| \leq \|C_j\| \tau_1(A).$$

Take T as is in (2.6) with X_k defined by (3.1). Besides (2.9) and (2.10), imply

$$\|T\| \leq \left(1 + \frac{1}{n-1} \sum_{k=1}^{n-1} \|X_k\|\right)^{n-1} \leq \left(1 + \frac{\tau_1(A)}{n-1} \sum_{k=1}^{n-1} \|C_k\|\right)^{n-1}$$

and

$$\|T^{-1}\| \leq \left(1 + \frac{\tau_1(A)}{n-1} \sum_{j=1}^{n-1} \|C_j\|\right)^{n-1}.$$

But, by the Schwarz inequality,

$$\left(\sum_{j=1}^{n-1} \|C_j\|\right)^2 \leq (n-1) \sum_{j=1}^{n-1} \|C_j\|^2 = (n-1) \sum_{j=1}^{n-1} \sum_{k=j+1}^n |a_{jk}|^2 = (n-1)g^2(A).$$

Thus,

$$\|T\|^2 \leq \left(1 + \frac{\tau_1(A)}{\sqrt{n-1}}g(A)\right)^{2(n-1)} = \left(1 + \frac{\tau(A)}{\sqrt{n-1}}\right)^{2(n-1)} = \gamma(A)$$

and $\|T^{-1}\|^2 \leq \gamma(A)$. Now (2.8) proves the theorem. \square

4. Applications of Theorem 1.1. Theorem 1.1 immediately implies the following.

COROLLARY 4.1. *Let condition (1.1) hold and $f(z)$ be a scalar function defined on the spectrum of A . Then $\|f(A)\| \leq \gamma(A) \max_k |f(\lambda_k)|$.*

Let A and \tilde{A} be complex $n \times n$ matrices whose eigenvalues λ_k and $\tilde{\lambda}_k$, respectively, are taken with their algebraic multiplicities. Recall that

$$sv_A(\tilde{A}) := \max_k \min_j |\tilde{\lambda}_k - \lambda_j|$$

is the spectral variation of \tilde{A} with respect to A .

COROLLARY 4.2. *Let condition (1.1) hold. Then $sv_A(\tilde{A}) \leq \gamma(A)\|A - \tilde{A}\|$.*

Indeed, the matrix $\hat{D} = TAT^{-1}$ is normal. Put $B = T\tilde{A}T^{-1}$. Thanks to the well-known Corollary 3.4 [16], $sv_{\hat{D}}(B) \leq \|\hat{D} - B\|$. Now the required result is due to Theorem 1.1.

Furthermore, let us suppose that λ_k are real and

$$\lambda_1 > \lambda_2 > \cdots > \lambda_n. \quad (4.1)$$

Then (1.2) holds with a Hermitian matrix \hat{D} . Again put $B = T\tilde{A}T^{-1}$. Then, due to Theorem 1.1, we have

$$\|\hat{D} - B\|_F = \|TAT^{-1} - T\tilde{A}T^{-1}\|_F \leq \|T\|\|A - \tilde{A}\|_F\|T^{-1}\| \leq \gamma(A)\|A - \tilde{A}\|_F. \quad (4.2)$$

The eigenvalues of B coincide with the eigenvalues $\tilde{\lambda}_k$ of \tilde{A} . Denote $\mu_k = \operatorname{Re} \tilde{\lambda}_k$ and assume that $\tilde{\lambda}_k$ are ordered in such a way that

$$\mu_1 \geq \mu_2 \geq \cdots \geq \mu_n. \quad (4.3)$$

Due to the Kahan theorem [16, Theorem IV.5.2, p. 213, inequality (5.4)], we can write

$$\left[\sum_{k=1}^n |\tilde{\lambda}_k - \lambda_k|^2 \right]^{1/2} \leq \sqrt{2}\|\hat{D} - B\|_F.$$

Hence, taking into account (4.2), we arrive at our next result.

COROLLARY 4.3. *Let the inequalities (4.1) and (4.3) hold. Then*

$$\left[\sum_{k=1}^n |\tilde{\lambda}_k - \lambda_k|^2 \right]^{1/2} \leq \sqrt{2}\gamma(A)\|A - \tilde{A}\|_F.$$

In addition, note that Theorem 4.5.4 in [16, p. 215] and Theorem 1.1 yield the following corollary.

COROLLARY 4.4. *Let A and \tilde{A} be diagonalizable $n \times n$ matrices having purely real eigenvalues:*

$$\lambda_1 < \lambda_2 < \cdots < \lambda_n \quad \text{and} \quad \tilde{\lambda}_1 < \tilde{\lambda}_2 < \cdots < \tilde{\lambda}_n, \quad \text{respectively}.$$

Then

$$|\tilde{\lambda}_j - \lambda_j| \leq \gamma(A)\gamma(\tilde{A})\|A - \tilde{A}\| \quad (j = 1, \dots, n).$$

To consider an additional application of Theorem 1.1, put

$$md(A, \tilde{A}) := \min_{\pi} \left[\sum_{k=1}^n |\tilde{\lambda}_k - \lambda_k|^2 \right]^{1/2},$$

where π ranges over all permutations of the integers $1, 2, \dots, n$, cf. [16]. Let us use Theorem 4.5.5 [16, p. 216]. That theorem together with Theorem 1.1 implies our next result.

COROLLARY 4.5. *Let the conditions (1.1) and*

$$\tilde{\lambda}_j \neq \tilde{\lambda}_m, \quad \text{whenever } j \neq m \quad (4.4)$$

be fulfilled. Then

$$md(A, \tilde{A}) \leq \gamma(A)\gamma(\tilde{A})\|A - \tilde{A}\|_F.$$

About the interesting recent publications on spectrum perturbations see for instance [6, 8].

Finally, note that Corollaries 2.2 and 2.3 from [11] and the above proved Theorem 1.1 yield the following corollary.

COROLLARY 4.6. *Let conditions (1.1) and (4.4) be fulfilled, and f be a function defined on $\sigma(A) \cup \sigma(\tilde{A})$. Then the inequalities*

$$\|f(A) - f(\tilde{A})\|_F \leq \gamma(A)\gamma(\tilde{A}) \max_{j,k} \left| \frac{f(\lambda_k) - f(\tilde{\lambda}_j)}{\lambda_k - \tilde{\lambda}_j} \right| \|A - \tilde{A}\|_F,$$

and

$$\|f(A) - f(\tilde{A})\|_F \leq \gamma(A)\gamma(\tilde{A}) \max_{j,k} |f(\lambda_k) - f(\tilde{\lambda}_j)|$$

are valid.

The recent interesting results devoted to matrix-valued functions can be found in [7, 14].

5. Example. To illustrate Theorem 1.1, consider the simple matrix

$$A = \begin{pmatrix} 5 & 0 & -\frac{1}{3} \\ 0 & 7 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

It is simple to check that $\lambda_1 = 5, \lambda_2 = 7, \lambda_3 = 3$. So, $\delta = 2$. In addition, due to (1.3), we have $g(A) \leq \|A - A^*\|_F / \sqrt{2} = \frac{1}{3}$. Thus,

$$\tau(A) = \sum_{k=0}^1 \frac{g^{k+1}(A)}{\sqrt{k!} \delta^{k+1}} \leq \frac{1}{6} + \frac{1}{36} \approx 0.1944$$

and

$$\gamma(A) \leq \left(1 + \frac{0.1944}{\sqrt{2}}\right)^4 \approx 1.6741.$$

On the other hand, it is not hard to check that the matrix

$$T = \begin{pmatrix} 1 & 0 & \frac{1}{6} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \text{ has the inverse one } T^{-1} = \begin{pmatrix} 1 & 0 & -\frac{1}{6} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

and satisfies relation (1.2) with

$$\hat{D} = \begin{pmatrix} 5 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

The direct calculations gives us $\kappa_T \approx 1.1815$.

Acknowledgment. I am very grateful to the referee of the present paper for his/her really deep and helpful remarks.

REFERENCES

- [1] R. Bhatia. *Matrix Analysis*. Springer, New York, 1997.
- [2] R. Bhatia and P. Rosenthal. How and why to solve the matrix equation $AX - XB = Y$. *Bull. London Math. Soc.*, 29:1–21, 1997.
- [3] L. Collatz. *Functional Analysis and Numerical Mathematics*. Academic Press, New York, 1966.
- [4] J.A. Dias da Silva and C.R. Johnson. Cospectrality and similarity for a pair of matrices under multiplicative and additive composition with diagonal matrices. *Linear Algebra Appl.*, 326(1-3):15–25, 2001.
- [5] D. Djokovic. Universal zero patterns for simultaneous similarity of several matrices. *Oper. Matrices*, 1(1):113–119, 2007.
- [6] K. Du. Note on structured indefinite perturbations to Hermitian matrices. *J. Comput. Appl. Math.*, 202(2):258–265, 2007.
- [7] B. Fritzsche, B. Kirstein, and A. Lasarow. Orthogonal rational matrix-valued functions on the unit circle: Recurrence relations and a Favard-type theorem. *Math. Nachr.*, 279(5-6):513–542, 2006.
- [8] P. Forrester and E. Rains. Jacobians and rank 1 perturbations relating to unitary Hessenberg matrices. *Int. Math. Res. Not.*, 2006(5):Article ID 48306 (36 pages), 2006.
- [9] A. George and K. Ikramov. Unitary similarity of matrices with quadratic minimal polynomials. *Linear Algebra Appl.*, 349(1-3):11–16, 2002.
- [10] M.I. Gil'. *Operator Functions and Localization of Spectra*. Lecture Notes in Mathematics, Vol. 1830, Springer-Verlag, Berlin, 2003.
- [11] M.I. Gil'. Perturbations of functions of diagonalizable matrices. *Electron. J. Linear Algebra*, 20:303–313, 2010.
- [12] M.I. Gil'. Matrix equations with diagonalizable coefficients. *Gulf Journal of Mathematics*, 1:98–104, 2013.
- [13] T. Jiang, X. Cheng, and L. Chen. An algebraic relation between consimilarity and similarity of complex matrices and its applications. *J. Phys. A, Math. Gen.*, 39(29):9215–9222, 2006.
- [14] A. Lasarow. Dual Szegő pairs of sequences of rational matrix-valued functions. *Int. J. Math. Math. Sci.*, 2006(5):1–37, 2006.
- [15] C. Rajian and T. Chelvam. On similarity invariants of *EP* matrices. *East Asian Math. J.*, 23(2):207–212, 2007.
- [16] G.W. Stewart and J.-G. Sun. *Matrix Perturbation Theory*. Academic Press, New York, 1990.