

## BOUNDS FOR THE COMPONENTWISE DISTANCE TO THE NEAREST SINGULAR MATRIX

S. M. RUMP <sup>†</sup>

**Abstract.** The normwise distance of a matrix  $A$  to the nearest singular matrix is well known to be equal to  $\|A\|/\text{cond}(A)$  for norms being subordinate to a vector norm. However, there is no hope to find a similar formula or even a simple algorithm for computing the componentwise distance to the nearest singular matrix for general matrices. This is because Rohn and Poljak [7] showed that this is an  $NP$ -hard problem.

Denote the minimum Bauer-Skeel condition number achievable by column scaling by  $\kappa$ . Demmel [3] showed that  $\kappa^{-1}$  is a *lower* bound for the componentwise distance to the nearest singular matrix. In this paper we prove that  $2.4 \cdot n^{1.7} \cdot \kappa^{-1}$  is an *upper* bound. This extends and proves a conjecture by N. J. Higham and J. Demmel. We give an explicit set of examples showing that an upper bound cannot be better than  $n \cdot \kappa^{-1}$ . Asymptotically, we show that  $n^{1+\ln 2+\varepsilon} \cdot \kappa^{-1}$  is a valid upper bound.

**Key words.** componentwise distance, NP-hardness, optimal Bauer-Skeel condition number

**AMS subject classifications.** 65F35, 15A60

**0. Introduction.** Let  $A$  be an  $n$  by  $n$  matrix and denote its smallest singular value by  $\sigma_n(A)$ . It is well known that the distance to the nearest singular matrix in the 2-norm or Frobenius norm is equal to  $\sigma_n(A)$ . More general, for any consistent matrix norm  $\|\cdot\|$  being subordinate to a vector norm we have

$$(1) \quad \min \{ \|\delta A\| \mid A + \delta A \text{ singular} \} = \frac{1}{\|A^{-1}\|} = \frac{\|A\|}{\text{cond}(A)}.$$

An appropriate  $\delta A$  of rank 1 can be explicitly calculated (cf. [?], [13]). Such a perturbation does, in general, alter each component of  $A$ . In many practical applications, one may be interested in leaving specific components such as system zeros unaltered, for example, if the matrix arises from some discretisation scheme. More general, this leads to the question of the componentwise distance to the nearest singular matrix. The componentwise distance may be weighted by some nonnegative matrix  $E$ . More precisely, we define

$$(2) \quad \sigma(A, E) := \min \{ \alpha \in \mathbb{R} \mid A + \tilde{E} \text{ singular where } |\tilde{E}_{ij}| \leq \alpha \cdot E_{ij} \text{ for all } i, j \}.$$

If no such  $\alpha$  exists, we set  $\sigma(A, E) := \infty$ . For singular matrices,  $\sigma(A, E) = 0$  for every weight matrix  $E$ . Specific values of  $E$  are  $E = |A|$  for relative perturbations, or  $E = (\mathbf{1})_{nn}$  for absolute perturbations. Among others, the componentwise distance to the nearest singular matrix was discussed in [8], [11], [10], and in [3]. In [8] we also find a first approach towards an estimation of the nearness to singularity in a norm not being subordinate to a vector norm, namely  $\|A\| := \max_{i,j} |A_{ij}|$ .

We cannot expect to find a formula or even a simple algorithm for calculating  $\sigma(A, E)$ . This is because Rohn and Poljak [7] proved that computation of  $\sigma(A, E)$  is  $NP$ -complete. For an outline of their proof see also [3]. Nevertheless, we may find bounds for  $\sigma(A, E)$ , and for classes of matrices even explicit formulas.

Another view of  $\sigma(A, E)$  is the maximum value such that the interval matrix  $[A - \alpha E, A + \alpha E]$  is nonsingular for  $\alpha < \sigma(A, E)$ . The interval matrix is defined as being the set of all matrices  $\tilde{A}$  with  $A_{ij} - \alpha E_{ij} \leq \tilde{A}_{ij} \leq A_{ij} + \alpha E_{ij}$  for all  $i, j$ , or in short notation  $A - \alpha E \leq \tilde{A} \leq A + \alpha E$ . The interval matrix is called nonsingular if every matrix  $\tilde{A} \in [A - \alpha E, A + \alpha E]$  is nonsingular. In a very interesting paper [9], Rohn gave 13 necessary and sufficient criteria for  $[A - E, A + E]$  being nonsingular.

---

<sup>†</sup>Technische Informatik III, TU Hamburg-Harburg, Eißendorfer Straße 38, 21071 Hamburg, Germany

A thorough discussion of  $\sigma(A, E)$  can be found in the very interesting paper [3]. Demmel [3] proved that  $\sigma(A, E)$  is equal to the inverse of  $\min \kappa(AD, ED)$ , the minimum taken over all diagonal matrices  $D$ , where  $\kappa(A, E) := \||A^{-1}| \cdot E\|$  denotes the Bauer-Skeel condition number. For any  $p$ -norm, he proves

$$\min_D \kappa(AD, ED) = \rho(|A^{-1}| \cdot E),$$

extending a result by Bauer [1]. In other words, the minimum Bauer-Skeel condition number achievable by column scaling is equal to the inverse of  $\rho(|A^{-1}| \cdot E)$ . Demmel and N. J. Higham conjecture that  $\rho(|A^{-1}| \cdot E)$  and  $\sigma(A, E)$  are not too far apart. They conjecture for relative perturbations existence of some constant  $\gamma \in \mathbb{R}$ , possibly depending on the dimension, with

$$(3) \quad \sigma(A, |A|) \leq \frac{\gamma}{\rho(|A^{-1}| \cdot |A|)}.$$

In this paper, our main goal is to show existence of such constants  $\gamma(n)$  and to derive lower and upper bounds for  $\gamma(n)$ . First, we show that  $\sigma(A, E) \geq \sigma_n(A)$  for  $\|E\|_2 = 1$ . A corresponding result for other norms is given in §2. However, this bound can be arbitrarily weak. Following we give some new bounds for  $\sigma(A, E)$ .

In §4 a perturbation formula for determinants is stated which is the key to prove an upper bound of  $\gamma(n)$ .

In §5 we will prove  $\gamma \geq n$ . In §6, for arbitrary weight matrices  $E$  we prove

$$(4) \quad \frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \leq \frac{\gamma(n)}{\rho(|A^{-1}| \cdot E)} \quad \text{with} \quad \gamma(n) = c \cdot n^\alpha$$

for  $c = 2.4$  and  $\alpha = 1.7$ . Moreover, for  $n \rightarrow \infty$  we show that for every  $\varepsilon > 0$ ,  $\alpha$  can be replaced by  $1 + \ln 2 + \varepsilon$ . In view of  $\gamma \geq n$ , we conjecture  $\gamma = n$ .

In [3], Demmel gave reasons to be interested in the componentwise distance to the nearest singular matrix. In §2, we add a lower and upper componentwise error bound for the solution of a linear system  $Ax = b$  subject to componentwise perturbations of the matrix and the right hand side. Such upper bounds are known in the literature and are valid for nonsingular  $A$  and  $|\tilde{A} - A| \leq E$  with  $\rho(|A^{-1}| \cdot E) < 1$ . We derive a componentwise bound for the *minimum* perturbation of the solution subject to finite perturbations of  $A$  and  $b$ . (4) shows that those estimates cover perturbation matrices  $\tilde{A}$  not too far from the next singular matrix.

The paper is organized as follows. In §1 we introduce the used notation. In §2 follows a componentwise lower and upper perturbation bound for finite componentwise perturbations of a linear system. In §3, lower bounds on  $\sigma(A, E)$  are given. For orthogonal matrices we show that  $\gamma$  (see (4)) is at least of the order of  $\sqrt{n}$ .

In §4, a Sherman-Morrison-Woodbury like perturbation theorem for determinants is given. In fact, this is an *equality* for finite perturbations of a matrix. In §5 we derive upper bounds on  $\sigma(A, E)$ . For  $E$  being of rank 1, such as for absolute perturbations, we show  $\gamma(n) \leq n$ , and for relative perturbations we give a set of matrices  $A \in M_n(\mathbb{R})$  with  $\gamma(n) = n$ . For a class of matrices including  $M$ -matrices we prove  $\gamma(n) = 1$ , i.e.  $\sigma(A, E) = \rho(|A^{-1}| \cdot E)^{-1}$ .

In §6 the results are extended to obtain an explicit upper bound on  $\gamma(n)$  for general  $A$  and  $E$ , and in §7 those bounds are quantified into (4). We close with the conjecture that (4) is valid for  $\gamma(n) = n$  for all  $A, E$ . If this is true, the set of matrices given in §5 would imply that inequality (4) with  $\gamma(n) = n$  is sharp.

**1. Notation.** In the following we list some notation from matrix theory, cf. for example [6], [5].  $V_n(\mathbb{R})$  denotes the set of vectors with  $n$  real components,  $M_{m,n}(\mathbb{R})$  the set of real  $m$  by  $n$  matrices, and  $M_n(\mathbb{R}) = M_{n,n}(\mathbb{R})$ . The components of a matrix  $A \in M_n(\mathbb{R})$  are referred by  $A_{ij}$  or  $A_{i,j}$ . For short notation, components of  $A^{-1}$  are referred by  $A_{ij}^{-1}$ .  $(\mathbf{1})$  denotes a vector with all components equal to 1,  $(\mathbf{1})_{nn} \in M_n(\mathbb{R})$  the matrix with all columns equal to  $(\mathbf{1})$ .

$Q_{kn}$  denotes the set of strictly increasing sequences of  $k$  integers chosen from  $\{1, \dots, n\}$ . For  $\omega \in Q_{kn}$ , we denote  $\omega = (\omega_1, \dots, \omega_k)$ . For  $C \in M_n(\mathbb{R})$ ,  $\omega \in Q_{kn}$ ,  $C[\omega] \in M_k(\mathbb{R})$  denotes the  $k$  by  $k$  submatrix of  $C$  lying in rows and columns  $\omega$ . A sequence  $\zeta = (i_1, \dots, i_k)$ ,  $k \geq 1$  of mutually different integers  $i_\nu \in \{1, \dots, n\}$  is

called a *cycle*. We identify the cycles  $(i_1, \dots, i_k)$  and  $(i_p, \dots, i_k, i_1, \dots, i_{p-1})$ , where  $1 \leq p \leq k$ . It is  $|\zeta| := k$ . A *full cycle*  $\zeta$  on  $\{1, \dots, n\}$  is a cycle  $\zeta$  with  $|\zeta| = n$ .

For  $C \in M_n(\mathbb{R})$  and a cycle  $\zeta = (i_1, \dots, i_k)$  on  $\{1, \dots, n\}$ , we put

$$\Pi_\zeta(C) := C_{i_1 i_2} \cdot \dots \cdot C_{i_{k-1} i_k} \cdot C_{i_k i_1},$$

the *cycle product* for  $\zeta$ . Note the last factor in the product. Therefore,  $|\Pi_\zeta(C)|^{1/|\zeta|}$  is the geometric mean of the elements of the cycle product. Each diagonal element  $C_{ii}$  is a cycle product, namely of the cycle  $(i)$ . (Here our definition differs from Engel/Schneider [4]).

With one exception, throughout the paper, *absolute value* and *comparison* is used *componentwise*. For example, for  $A, B \in M_n(\mathbb{R})$ ,

$$|A| \leq B \quad \text{means} \quad |A_{ij}| \leq B_{ij} \quad \text{for} \quad 1 \leq i, j \leq n.$$

The exception are cycles  $\zeta = (i_1, \dots, i_k)$ , where  $|\zeta| = k$ . The *singular values* of a matrix  $A \in M_n(\mathbb{R})$  are denoted in decreasing order with increasing indices, i.e.  $\sigma_1(A) \geq \dots \geq \sigma_n(A) \geq 0$ .

For  $A, E \in M_n(\mathbb{R})$ ,  $E \geq 0$ ,  $\sigma(A, E)$  denotes the *componentwise distance*, weighted by  $E$ , to the nearest *singular matrix* (cf. (0.2)).

For finite  $\sigma(A, E)$ , the set of all matrices  $\tilde{A} \in M_n(\mathbb{R})$  with  $|\tilde{A} - A| \leq \sigma(A, E) \cdot E$  is compact. For every nonsingular  $\tilde{A}$  there is a neighbourhood of  $\tilde{A}$  consisting only of nonsingular matrices. Therefore

$$\sigma(A, E) < \infty \quad \Rightarrow \quad \exists \delta A \in M_n(\mathbb{R}) : |\delta A| = \sigma(A, E) \cdot E \quad \text{and} \quad A + \delta A \quad \text{singular},$$

showing that we are allowed to use a minimum in the definition (0.2) of  $\sigma(A, E)$ .  $\rho$  denotes the spectral radius, whereas  $\rho_0$  denotes the *real spectral radius*:

$$B \in M_n(\mathbb{R}) : \quad \rho_0(B) := \max \{ |\lambda| \mid \lambda \in \mathbb{R} \text{ is an eigenvalue of } B \}.$$

If  $B$  has no real eigenvalues, we set  $\rho_0(B) := 0$ .  $I$  denotes the identity matrix of proper dimension, especially  $I_k \in M_k(\mathbb{R})$  denotes the  $k$  by  $k$  identity matrix. A *signature matrix*  $S$  is a diagonal matrix with diagonal entries  $+1$  or  $-1$ , i.e.  $|S| = I$ .

We frequently use standard results from matrix and Perron-Frobenius theory such as

$$(5) \quad A \in M_{nk}(\mathbb{R}), B \in M_{kn}(\mathbb{R}) \quad \Rightarrow \quad \text{The set of nonzero eigenvalues of } AB \text{ and } BA \text{ are identical,}$$

cf. Theorem 1.3.20 in [5], and

$$(6) \quad A \in M_n(\mathbb{R}) \text{ and } A \geq 0, x \in V_n(\mathbb{R}) \text{ with } x > 0 \quad \Rightarrow \quad \min_i \frac{(Ax)_i}{x_i} \leq \rho(A) \leq \max_i \frac{(Ax)_i}{x_i}.$$

The latter can be found in [2].

**2. Finite perturbations for a linear system.** Calculating bounds on  $\sigma(A, E)$  can be motivated, for example, by looking at linear systems with finite perturbations of the input data. For a linear system  $Ax = b$  consider the perturbed system  $\tilde{A}\tilde{x} = \tilde{b}$  with  $\delta A := \tilde{A} - A$ ,  $\delta b := \tilde{b} - b$ ,  $\delta x := \tilde{x} - x$ . Then for nonsingular  $A$ ,

$$(7) \quad A \cdot (I + A^{-1} \cdot \delta A) \cdot (\tilde{x} - x) = \tilde{A} \cdot (\tilde{x} - x) = \tilde{b} - \tilde{A}x = \delta b - \delta A \cdot x.$$

If  $\rho(A^{-1} \cdot \delta A) < 1$ , then  $I + A^{-1} \cdot \delta A$  and  $\tilde{A} = A \cdot (I + A^{-1} \cdot \delta A)$  are nonsingular, and (7) implies

$$(8) \quad \delta x = (I + A^{-1} \cdot \delta A)^{-1} \cdot A^{-1} \cdot (\delta b - \delta A \cdot x).$$

If  $\rho(|A^{-1}| \cdot \Delta A) < 1$  then  $I - |A^{-1}| \cdot \Delta A$  is an  $M$ -matrix. If the perturbations  $\delta A$ ,  $\delta b$  are componentwise bounded by  $|\delta A| \leq \Delta A$ ,  $|\delta b| \leq \Delta b$  then (2.2) implies

$$(9) \quad |\delta x| \leq (I - |A^{-1}| \cdot \Delta A)^{-1} \cdot |A^{-1}| \cdot (\Delta b + \Delta A \cdot |x|).$$

For given weight matrix  $\Delta A$ , consider the set of matrices with componentwise distance from  $A$  weighted by  $\Delta A$  not greater than  $\sigma$ :

$$\tilde{A} \in U_\sigma(A, \Delta A) \Leftrightarrow |\tilde{A} - A| \leq \sigma \cdot \Delta A.$$

For  $\sigma \leq \rho(|A^{-1}| \cdot \Delta A)$ , Perron-Frobenius-Theory yields

$$\rho(I - A^{-1} \cdot \tilde{A}) = \rho(A^{-1} \cdot (A - \tilde{A})) \leq \rho(|A^{-1}| \cdot \Delta A) < 1,$$

and therefore regularity of all  $\tilde{A} \in U_\sigma(A, \Delta A)$ . The bound (2.3) requires  $|A^{-1}| \cdot \Delta A$  to be convergent, whereas (8) is valid for  $\rho(A^{-1} \cdot \delta A) < 1$ . Therefore we may ask, how far a matrix  $\tilde{A}$  with  $\rho(|A^{-1}| \cdot |\tilde{A} - A|) \geq 1$  can be from the nearest singular matrix. An answer to this question shows how strong the assumption  $\rho(|A^{-1}| \cdot \Delta A) < 1$  is.

**3. Lower bounds on  $\sigma(A, E)$ .** A simple and well-known lower bound on  $\sigma(A, E)$  is

$$(10) \quad \frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \quad \text{for all nonsingular } A \in M_n(\mathbb{R}), \quad 0 \leq E \in M_n(\mathbb{R}).$$

This can be seen using Perron-Frobenius Theory and

$$\begin{aligned} \rho(|A^{-1}| \cdot E) < 1 &\Rightarrow \rho(A^{-1} \cdot \delta A) < 1 \quad \text{for all } |\delta A| \leq E \\ &\Rightarrow A + \delta A = A \cdot (I + A^{-1} \cdot \delta A) \quad \text{is nonsingular.} \end{aligned}$$

Another lower bound is (cf. [12], Theorem 1.8, p. 75)

$$(11) \quad \frac{\sigma_n(A)}{\sigma_1(E)} \leq \sigma(A, E).$$

This can be generalized in the following way.

**THEOREM 3.1.** *Let  $\|\cdot\|$  be a matrix norm subordinate to an absolute vector norm  $\|\cdot\|$ . Then for nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$ ,*

$$(12) \quad \frac{1}{\|A^{-1}\| \cdot \|E\|} \leq \sigma(A, E).$$

(12) is especially valid for all  $p$ -norms. For absolute norms such as 1-norm and  $\infty$ -norm,

$$(13) \quad \frac{1}{\|A^{-1}\| \cdot \|E\|} \leq \frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E),$$

whereas for the 2-norm

$$(14) \quad \frac{1}{\|A^{-1}\|_2 \cdot \|E\|_2} \leq \frac{\sqrt{n}}{\rho(|A^{-1}| \cdot E)}.$$

*Proof.* To prove (12), let  $\delta A \in M_n(\mathbb{R})$  with  $|\delta A| \leq \sigma \cdot E$  for  $\sigma < (\|A^{-1}\| \cdot \|E\|)^{-1}$ . The vector norm is absolute implying  $\|x\| = \||x|\|$  and  $|x| \leq |y| \Rightarrow \|x\| \leq \|y\|$  for  $x, y \in V_n(\mathbb{R})$  (cf. [13], Theorem II.1.2). Let  $x \in V_n(\mathbb{R})$  with  $\|x\| = 1$  and  $\|\delta A\| = \|\delta A \cdot x\|$ . Then

$$(15) \quad \begin{aligned} \|\delta A\| &= \|\delta A \cdot x\| = \||\delta A \cdot x|\| \leq \|\sigma \cdot E \cdot |x|\| \leq \|\sigma \cdot E\| \cdot \| |x| \| \\ &= \sigma \cdot \|E\| < \|A^{-1}\|^{-1}. \end{aligned}$$

For every  $0 \neq y \in V_n(\mathbb{R})$  holds  $\|y\| \leq \|A^{-1}\| \cdot \|Ay\|$ , and (15) yields

$$\|\delta A \cdot y\| \leq \|\delta A\| \cdot \|y\| < \|A^{-1}\|^{-1} \cdot \|y\| \leq \|Ay\|, \quad \text{and therefore } (A + \delta A) \cdot y \neq 0.$$

Therefore,  $A + \delta A$  is nonsingular for  $|\delta A| \leq \sigma \cdot E$ , and  $\sigma < (\|A^{-1}\| \cdot \|E\|)^{-1}$ , proving (12). For absolute matrix norms,

$$\rho(|A^{-1}| \cdot E) \leq \| |A^{-1}| \cdot E \| \leq \|A^{-1}\| \cdot \|E\|$$

proving (13). For the 2-norm holds

$$\rho(|A^{-1}| \cdot E) \leq \| |A^{-1}| \|_2 \cdot \|E\|_2 \leq \| |A^{-1}| \|_F \cdot \|E\|_2 = \|A^{-1}\|_F \cdot \|E\|_2 \leq \sqrt{n} \cdot \|A^{-1}\|_2 \cdot \|E\|_2,$$

proving (14) and the theorem.  $\square$

(13) shows that for absolute matrix norms such as the 1-norm or  $\infty$ -norm, the bound (12) cannot be better than (10). The 2-norm is not absolute, and (14) shows that the lower bound (11) for  $\sigma(A, E)$  may be better up to a factor  $\sqrt{n}$  than (10). In fact, we can identify a class of matrices for which this improvement is approximately achieved.

Let  $Q \in M_n(\mathbb{R})$  be orthogonal, and consider absolute perturbations  $E = (\mathbf{1})_{nn}$ . Then (11) yields

$$\frac{\sigma_n(Q)}{\sigma_1((\mathbf{1})_{nn})} = \frac{1}{n} \leq \sigma(Q, E).$$

On the other hand,  $E = (\mathbf{1})_{nn} \in M_n(\mathbb{R})$  and  $x = (\mathbf{1}) \in V_n(\mathbb{R})$  imply

$$\{ |Q^{-1}| \cdot E \}^T \cdot x = E \cdot |Q| \cdot x = \left( \sum_{i,j} |Q_{ij}| \right) \cdot x,$$

and (6) yields  $\rho(|Q^{-1}| \cdot E) = \sum_{i,j} |Q_{ij}|$ . If  $Q$  is an orthogonalized random matrix with components uniformly distributed in  $[-1, 1]$ , then  $|Q_{ij}| \approx n^{-1/2}$ . Thus, for the ratio between the two lower bounds (11) and (10) we obtain

$$\{ \sigma_n(Q) / \sigma_1(E) \} / \{ 1 / \rho(|Q^{-1}| \cdot E) \} \approx n^{-1} \cdot n^2 \cdot n^{-1/2} = \sqrt{n}.$$

The same heuristic holds for  $E = |Q|$ , cf. [12]. For every Hadamard matrix ( $H \in M_n(\mathbb{R})$  with  $H^T H = n \cdot I$ ) the ratio is equal to  $\sqrt{n}$ . This sheds a first light on a possible quantity  $\gamma(n)$  such that (4) holds. In §5 we will prove  $\gamma(n) \geq n$ .

EXAMPLE 3.2. *The lower bound (11) may be arbitrarily weak. Consider*

$$A = \begin{pmatrix} 2\varepsilon & -\varepsilon \\ -\varepsilon & 1 \end{pmatrix} \quad \text{and} \quad E = |A| \quad \text{for some } \varepsilon > 0.$$

*A is a diagonally dominant M-matrix. As we will see in (5.5), A being M-matrix implies equality in (10), i.e.  $\sigma(A, |A|) = \rho(|A^{-1}| \cdot |A|) = 1 + 0(\sqrt{\varepsilon})$ . On the other hand,  $\sigma_2(A) / \sigma_1(|A|) = 2\varepsilon + 0(\varepsilon^2)$  underestimates  $\sigma(A, |A|)$  arbitrarily. This corresponds to  $\sigma_2(A) = 2\varepsilon + 0(\varepsilon^2)$ . That means, the normwise distance in the 2-norm or Frobenius norm to the nearest singular matrix can be arbitrarily small compared to a componentwise distance.*

**4. A perturbation theorem for determinants.** A lower bound on  $\sigma(A, E)$  is obtained by proving *regularity* of a set of matrices. This was done in §3 by using spectral properties. To obtain an *upper* bound on  $\sigma(A, E)$ , we may construct a specific perturbation  $\delta A$  with  $|\delta A| \leq \sigma_0 \cdot E$ ,  $\sigma_0 \in \mathbb{R}$  such that  $A + \delta A$  is singular. This proves  $\sigma(A, E) \leq \sigma_0$ . Another possibility to obtain an upper bound on  $\sigma(A, E)$  is the following. If  $|\delta A| \leq \sigma_0 \cdot E$  and  $\det(A) \cdot \det(A + \delta A) \leq 0$ , then a continuity argument yields  $\sigma(A, E) \leq \sigma_0$ . Therefore we state the following explicit formula for the relative change of the determinant of a matrix subject to a rank- $k$  perturbation. It is a Sherman-Morrison-Woodbury like perturbation formula for determinants.

LEMMA 4.1. *Let  $A \in M_n(\mathbb{R})$  and  $U, V \in M_{n,k}(\mathbb{R})$  be given. Then for nonsingular A,*

$$(16) \quad \det(A + UV^T) = \det(A) \cdot \det(I_k + V^T A^{-1} U),$$

*where  $I_k$  denotes the  $k$  by  $k$  identity matrix.*

*Proof.* It is

$$\det(A + UV^T) = \det(A) \cdot \det(I_n + A^{-1}UV^T).$$

Denoting the eigenvalues of  $X \in M_n(\mathbb{R})$  by  $\lambda_i(X)$  implies

$$\det(I_n + A^{-1}UV^T) = \prod_{i=1}^n \lambda_i(I_n + A^{-1}UV^T) = \prod_{i=1}^n \{1 + \lambda_i(A^{-1}UV^T)\}.$$

The set of nonzero eigenvalues of  $A^{-1}UV^T$  and  $V^T A^{-1}U$  are identical (see (5)), thus proving the lemma.  $\square$

This lemma has a nice and for itself interesting corollary.

**COROLLARY 4.2.** *Let  $A \in M_n(\mathbb{R})$  and  $u, v \in V_n(\mathbb{R})$ . Then for nonsingular  $A$ ,*

$$(17) \quad \det(A + uv^T) = \det(A) \cdot (1 + v^T A^{-1}u).$$

*For arbitrary  $A \in M_n(\mathbb{R})$  holds ( $\text{adj}(A)$  denotes the adjoint of  $A$ ),*

$$(18) \quad \det(A + uv^T) = \det(A) + v^T \cdot \text{adj}(A) \cdot u.$$

The corollary shows that the relative change of the determinant is linear for rank-1 perturbations of the matrix. The second well-known formula follows, for example, by a continuity argument using  $A \cdot \text{adj}(A) = \det(A) \cdot I$ .

**5. Upper bounds on  $\sigma(A, E)$ .** The perturbation lemma for determinants given in §4 allows for other lower bounds on  $\sigma(A, E)$ . The first result can be found in [8], Corollary 5.1, (iii).

**THEOREM 5.1.** *Let  $A \in M_n(\mathbb{R})$  be nonsingular and  $E \in M_n(\mathbb{R})$  with  $E \geq 0$ . Then*

$$(19) \quad \sigma(A, E) \leq \frac{1}{\max_i \{|A^{-1}| \cdot E\}_{ii}},$$

where  $0^{-1}$  is interpreted as  $\infty$ .

*Proof.* Set  $\alpha := \max_i \{|A^{-1}| \cdot E\}_{ii} \neq 0$  and let  $i$  be an index, for which this maximum is achieved. Denote the  $i\nu$ -th component of  $A^{-1}$  by  $A_{i\nu}^{-1}$  and define  $u \in V_n(\mathbb{R})$  by  $u_\nu := -\alpha^{-1} \cdot \text{sign}(A_{i\nu}^{-1}) \cdot E_{\nu i}$ . Then

$$(20) \quad e_i^T \cdot A^{-1} \cdot u = -\alpha^{-1} \cdot \sum_{\nu=1}^n |A_{i\nu}^{-1}| \cdot E_{\nu i} = -1,$$

and Corollary 4.2 implies  $\det(A + u \cdot e_i^T) = 0$ . Now  $|ue_i^T| \leq \alpha^{-1} \cdot E$  yields  $\sigma(A, E) \leq \alpha^{-1}$ .  $\square$

**EXAMPLE 5.2.** *The upper bound (19) can be arbitrarily weak. Consider*

$$(21) \quad A = \begin{pmatrix} \varepsilon & 0 & 1 & 1 \\ 0 & \varepsilon & 1 & 1 \\ 1 & 1 & \varepsilon & 0 \\ 1 & 1 & 0 & \varepsilon \end{pmatrix}, E = |A| \quad \text{with} \quad |A^{-1}| \cdot |A| \approx \begin{pmatrix} 1 & 1 & 1/\varepsilon & 1/\varepsilon \\ 1 & 1 & 1/\varepsilon & 1/\varepsilon \\ 1/\varepsilon & 1/\varepsilon & 1 & 1 \\ 1/\varepsilon & 1/\varepsilon & 1 & 1 \end{pmatrix},$$

where the components of  $|A^{-1}| \cdot |A|$  are accurate up to a relative error  $\varepsilon$ . Then (19) gives  $\sigma(A, |A|) \leq 1 + 0(\varepsilon)$ . On the other hand,

$$\det(A + \varepsilon \cdot \delta A) = 0 \quad \text{for} \quad \delta A = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix},$$

showing  $\sigma(A, |A|) \leq \varepsilon$ .

In Theorem 5.1, a rank-1 perturbation is used to prove (19). In a normwise sense, the minimum distance to the nearest singular matrix is achieved by a rank-1 perturbation. This is no longer true for componentwise distances, as will be shown by the following example.

EXAMPLE 5.3. According to Corollary 4.2, the smallest  $\sigma$  such  $A + \sigma e$  is singular with  $|e| \leq \sigma \cdot |A|$  and  $\text{rank}(e) = 1$  is given by  $\sigma = |\hat{\varphi}|^{-1}$ , where  $\hat{\varphi}$  is an optimal value of the constraint optimization problem

$$\varphi(u, v) := v^T A^{-1} u = \text{Min!} \quad \text{subject to} \quad |uv^T| \leq |A|.$$

In Example 5.2, partition the vectors  $u, v \in V_4(\mathbb{R})$  into two vectors  $U_i, V_i \in V_2(\mathbb{R})$ ,  $i \in \{1, 2\}$ , either having 2 components. That means  $u = (U_1, U_2)^T$ ,  $v = (V_1, V_2)^T$ . Let

$$|uv^T| \leq |A|, \quad \text{i.e.} \quad U_i V_i^T \leq \varepsilon \cdot I \quad \text{and} \quad U_i V_j^T \leq (\mathbf{1})_{22} \quad \text{for } 1 \leq i, j \leq 2, i \neq j.$$

The large elements of  $A^{-1}$  are in the upper left and lower right 2 by 2 block:

$$A^{-1} \approx \begin{pmatrix} X & Y \\ Y & X \end{pmatrix} \quad \text{with} \quad X = \frac{1}{2\varepsilon} \cdot \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \quad \text{and} \quad Y = \frac{1}{4} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

up to a relative error of the order  $\varepsilon$ . Therefore

$$\begin{aligned} |v^T A^{-1} u| &\leq |V_1^T X U_1| + |V_2^T X U_2| + |V_1^T Y U_2| + |V_2^T Y U_1| \\ &\leq 2\varepsilon \cdot \sum |X_{ij}| + 2 \cdot \sum |Y_{ij}| \leq 6 \end{aligned}$$

Therefore, Corollary 4.2 implies that the minimum distance to the nearest singular matrix subject to rank-1 perturbations weighted by  $|A|$  is at least  $1/6$  compared to  $\sigma(A, |A|) \leq \varepsilon$ . This observation sheds light on the difficulties to calculate  $\sigma(A, E)$  or to find upper bounds for it.

One may define the rank- $k$  componentwise distance to the nearest singular matrix as follows

$$\sigma_k(A, E) := \min\{\alpha \in \mathbb{R} \mid A + \tilde{E} \text{ singular for } |\tilde{E}| \leq \alpha \cdot E \text{ and } \text{rank}(\tilde{E}) \leq k\}.$$

We use  $\text{rank}(\tilde{E}) \leq k$  because  $E$  may be rank-deficient. We have just seen in Example 5.3 that  $\sigma_2(A, E)/\sigma_1(A, E)$  may be arbitrarily small. ■

Given the lower bound (10), one may ask whether there exist finite constants  $\gamma(n) \in \mathbb{R}$  only depending on  $n$  such that

$$(22) \quad \frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \leq \frac{\gamma(n)}{\rho(|A^{-1}| \cdot E)}$$

for all nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$ . This question has been raised in [3] and answered for some classes of matrices. The main purpose of this paper is to derive bounds for  $\gamma(n)$ . This will be done by using Lemma 4.1. For this purpose we need the following result by Rohn (for notation see §1).

THEOREM 5.4. (**Rohn**) For nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$  holds

$$\frac{1}{\max_{S_1, S_2} \rho_0(S_1 A^{-1} S_2 E)} = \sigma(A, E),$$

where  $\rho_0$  denotes the real spectral radius and the maximum is taken over all signature matrices.  $1/0$  is interpreted as  $\infty$ .

*Proof.* cf. [9]. □

We start with a theorem bounding  $\gamma(n)$  for general weight matrices  $E$ , and identify a class of matrices with  $\gamma(n) = 1$ .

THEOREM 5.5. For nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$ , the following is true.

- i) Assume a matrix  $S \in M_n(\mathbb{R})$  of rank 1 exists with

$$S_{ij} = \begin{cases} +1 & \text{if } A_{ij}^{-1} > 0 \\ -1 & \text{if } A_{ij}^{-1} < 0 \\ +1 \text{ or } -1 & \text{if } A_{ij}^{-1} = 0 \end{cases} .$$

Then (22) holds with  $\gamma(n) = 1$ .

ii) If  $0 < \eta \leq |E_{ij}| \leq \zeta$  for all  $1 \leq i, j \leq n$ , then (22) holds with  $\gamma(n) = n \cdot \zeta / \eta$ .

*Proof.* Let  $S = uv^T$  with  $u, v \in V_n(\mathbb{R})$ ,  $|u| = |v| = (\mathbf{1})$ . Defining  $S_1 = \text{diag}(u)$ ,  $S_2 = \text{diag}(v)$ , we have  $S_1 A^{-1} S_2 = |A^{-1}|$  and Rohn's characterization in Theorem 5.4 proves the first part. W.l.o.g. assume  $\sigma(A, E) < \infty$ . It is  $\eta \cdot \|A^{-1}\|_\infty \leq \max_i (|A^{-1}| \cdot E)_{ii}$  and  $\rho(|A^{-1}| \cdot E) \leq \|A^{-1}\|_\infty \cdot \|E\|_\infty \leq n \cdot \zeta \cdot \|A^{-1}\|_\infty$ . Thus, Theorem 5.1 proves the second part and therefore the theorem.  $\square$

For important classes of matrices such as nonnegative invertible matrices, among them all  $M$ -matrices, we already have a precise formula for  $\sigma(A, E)$ :

$$(23) \quad A \in M_n \mathbb{R} \text{ nonnegative invertible, } 0 \leq E \in M_n(\mathbb{R}) \Rightarrow \sigma(A, E) = \frac{1}{\rho(|A^{-1}| \cdot E)} .$$

EXAMPLE 5.6. If constants  $\gamma(n)$  with (22) exist at all, we can give a lower bound on  $\gamma(n)$  by means of the following. Define  $A \in M_n(\mathbb{R})$  by

$$(24) \quad A := \begin{pmatrix} 1 & & & & & & & s \\ 1 & 1 & & & & & & 0 \\ & 1 & 1 & & & & & \\ & & 1 & 1 & & & & \\ & & & \ddots & \ddots & & & \\ & & & & & \ddots & & \\ 0 & & & & & & 1 & \\ & & & & & & 1 & 1 \end{pmatrix} \quad \text{with } s := (-1)^{n+1} .$$

The determinant of  $A$  calculates to

$$\det(A) = \prod_{i=1}^n A_{ii} + (-1)^{n+1} \cdot \Pi_\zeta(A) = 2, \quad \text{where } \zeta = (1, \dots, n)$$

and  $\Pi_\zeta(A) = A_{12} \cdot A_{23} \cdot \dots \cdot A_{n-1,n} \cdot A_{n1}$ . If the elements of  $A$  are afflicted with relative perturbations, i.e.  $E = |A|$ , then only the 1's and  $s$  change. Therefore, any  $\tilde{A}$  with  $|\tilde{A} - A| \leq \sigma \cdot |A|$  with  $\sigma < 1$  is nonsingular, and therefore  $\sigma(A, |A|) = 1$ . On the other hand,  $|A^{-1}| \cdot |A| = (\mathbf{1})_{nn}$  and  $\rho(|A^{-1}| \cdot |A|) = n$ . This proves the following lemma.

LEMMA 5.7. If constants  $\gamma(n) \in \mathbb{R}$  with (22) for every nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$  exist at all, then  $\gamma(n) \geq n$ .

Next we show that  $\gamma(n) \leq n$  for  $E$  being of rank 1. For the proof we use Corollary 4.2, which is a consequence of Lemma 4.1 for  $k = 1$ . In the remaining part of the paper, we will extend this proof to  $k > 1$  to obtain upper bounds for  $\gamma(n)$  and for general  $A, E$ .

THEOREM 5.8. Let nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$  with  $E = uv^T$  for some  $u, v \in V_n(\mathbb{R})$ ,  $u, v \geq 0$ . Then

$$\frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \leq \frac{n}{\rho(|A^{-1}| \cdot E)} .$$

*Proof.* According to Theorem 5.4 and using (1.1),

$$(25) \quad \sigma(A, E)^{-1} = \max_{S_1, S_2} \rho_0(S_1 A^{-1} S_2 uv^T) = \max_{S_1, S_2} v^T S_1 A^{-1} S_2 u,$$

where the maximum is taken over all signature matrices  $S_1, S_2$ . For any  $i$ ,  $1 \leq i \leq n$ , we can choose appropriate signature matrices  $S_1, S_2$  such that  $v^T S_1 A^{-1} S_2 u \geq v_i \cdot (|A^{-1}| \cdot u)_i$ . Using (25) this yields

$$\sigma(A, E)^{-1} \geq \max_i v_i \cdot (|A^{-1}| \cdot u)_i .$$

On the other hand, using (1.1),

$$\rho(|A^{-1}| \cdot E) = \rho(|A^{-1}| \cdot uv^T) = v^T \cdot |A^{-1}| \cdot u \leq n \cdot \max_i v_i \cdot (|A^{-1}| \cdot u)_i. \quad \square$$

COROLLARY 5.9. For nonsingular  $A \in M_n(\mathbb{R})$  and absolute perturbations, i.e.  $E = (\mathbf{1})_{nn}$ , estimation (22) holds with  $\gamma(n) = n$ .

**6. Estimation of  $\gamma(n)$ .** To make further progress in the estimation of  $\gamma(n)$  we show that for nonsingular  $A$ ,  $\sigma(A, E)$  depends continuously on  $A$  and  $E$ . Using this we can restrict the class of matrices  $A$  and  $E$  to matrices with only nonzero components. For the proof we can hardly use a simple continuity argument on  $\rho_0(S_1 A^{-1} S_2 E)$  in connection with Theorem 5.4. This is because the search domain is restricted by  $E$ , and the (in absolute value) largest *real* eigenvalue may be multiple and become complex under arbitrarily small perturbations.

LEMMA 6.1. For nonsingular  $A \in M_n(\mathbb{R})$ ,  $\sigma(A, E)$  depends continuously on  $A$  and  $E$ .

*Proof.* For  $\sigma(A, E) = \infty$  we show that  $\sigma(\tilde{A}, \tilde{E})$  becomes unbounded for  $\tilde{A} \rightarrow A$ ,  $\tilde{E} \rightarrow E$ . A compactness and continuity argument shows that for every finite  $0 < c \in \mathbb{R}$ :

$$\forall |e| \leq c \cdot E : \quad |\det(A + e)| \geq \delta > 0.$$

For every  $\tilde{A}, \tilde{E}$  close enough to  $A, E$ , this implies  $|\det(\tilde{A} + \tilde{e})| \geq \delta/2 > 0$  for every  $|\tilde{e}| \leq c \cdot \tilde{E}$ , and hence  $\sigma(\tilde{A}, \tilde{E}) > c$ .

Assume  $\sigma := \sigma(A, E) < \infty$ . We will show that for small enough  $\varepsilon > 0$ , there exists some  $\delta > 0$  such that both of the following statements are true:

$$(26) \quad \forall e \in M_n(\mathbb{R}) : \quad |e| \leq (\sigma - \varepsilon) \cdot E \quad \Rightarrow \quad \det(A) \cdot \det(A + e) > \delta,$$

$$(27) \quad \exists e \in M_n(\mathbb{R}) : \quad |e| \leq (\sigma + \varepsilon) \cdot E \quad \text{and} \quad \det(A) \cdot \det(A + e) < -\delta.$$

(26) is seen as follows. For  $\varepsilon > 0$ , the set of matrices  $A + e$  with  $|e| \leq (\sigma - \varepsilon) \cdot E$  is nonempty and compact. Hence,  $\det(A) \cdot \det(A + e)$  achieves a minimum on this set. By definition of  $\sigma$ , this minimum is positive. To see (27), observe that  $\det(A) \cdot \det(A + e) \geq 0$  for all  $|e| \leq \sigma \cdot E$ . For any index pair  $i, j$ , the determinant  $\det(A + \varepsilon \cdot e_i e_j^T)$  depends linearly on  $\varepsilon$ . Now proceed as follows. There is some  $e$  such that  $A + e$  is singular and  $|e| = E$ . If for an index pair  $i, j$ , the determinant  $\det(A + e)$  is independent on  $e_{ij}$ , then replace  $e_{ij}$  by 0. At each step of this process,  $\det(A + e) = 0$  and  $|e| \leq E$ . The definition of  $\sigma(A, E) < \infty$  implies that during this process we must arrive at some  $e$  and an index pair  $k, l$ , such that  $\det(A + e)$  is not constant when changing  $e_{kl}$ . Then defining  $e' \in M_n(\mathbb{R})$  by  $e'_{ij} := e_{ij}$  for  $(i, j) \neq (k, l)$  and  $e'_{kl} := e_{kl} \cdot (1 + \varepsilon')$  for small  $\varepsilon' > 0$  proves (27).

Now the continuity of the determinant implies for  $\tilde{A}, \tilde{E}$  close enough to  $A, E$ ,

$$\forall |\tilde{e}| \leq (\sigma - \varepsilon) \cdot \tilde{E} : \quad \det(\tilde{A}) \cdot \det(\tilde{A} + \tilde{e}) > \delta/2 \quad \text{and}$$

$$\exists |\tilde{e}| \leq (\sigma + \varepsilon) \cdot \tilde{E} : \quad \det(\tilde{A}) \cdot \det(\tilde{A} + \tilde{e}) < -\delta/2,$$

and therefore  $\sigma(A, E) - \varepsilon < \sigma(\tilde{A}, \tilde{E}) < \sigma(A, E) + \varepsilon$ .  $\square$

COROLLARY 6.2. If (22) holds for each  $E > 0$ , then it holds for each  $E \geq 0$ .

Our goal for this chapter is to prove the following upper bound for  $\sigma(A, E)$ . The quantities  $\varphi_t$  occurring in this estimation will be quantified and estimated in §7.

PROPOSITION 6.3. Let  $A, E \in M_n(\mathbb{R})$  with  $A$  nonsingular and  $E \geq 0$  be given. Define recursively  $\varphi_1 := 1$ ,  $\varphi_2 := 1$  and  $\varphi_t \in \mathbb{R}$ ,  $2 < t \in \mathbb{N}$  to be the (unique) positive root of

$$(28) \quad P_t(x) \in \mathbb{R}[x] \quad \text{with} \quad P_t(x) := x^{t-1} - x^{t-2} - \sum_{\nu=1}^{t-1} \varphi_\nu \cdot x^{t-1-\nu}.$$

Then

$$(29) \quad \sigma(A, E) \leq \frac{n \cdot \varphi_n}{\rho(|A^{-1}| \cdot E)}.$$

Therefore, the quantities  $\gamma(n)$  defined in (5.4) satisfy

$$(30) \quad \begin{aligned} \gamma(1) &= 1, \quad \gamma(2) = 2 \quad \text{and} \\ n &\leq \gamma(n) \leq n \cdot \varphi_n. \end{aligned}$$

The proof divides into several parts and needs some preparatory lemmata. First, we will construct a specific rank- $k$  perturbation in order to be able to apply Lemma 4.1 to bound  $\gamma(n)$  for general  $A, E$ . We use the same principle as in the proof of Theorem 5.1 adapted to rank- $k$  perturbations.

LEMMA 6.4. *Let nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$  be given, and set  $C := |A^{-1}| \cdot E$ . For  $1 \leq k \leq n$  define*

$$(31) \quad i' := \begin{cases} i + 1 & \text{for } 1 \leq i < k \\ 1 & \text{for } i = k \end{cases}$$

and  $U, V \in M_{n,k}(\mathbb{R})$  by

$$U_{\nu i'} := \text{sign}(A_{i\nu}^{-1}) \cdot E_{\nu i'} \quad \text{and} \quad V_{\mu i} := \delta_{\mu i}$$

for  $1 \leq \mu, \nu \leq n, 1 \leq i \leq k$  and the Kronecker symbol  $\delta$ . Set  $\tilde{C} := V^T A^{-1} U$ . Then

- i)  $|\tilde{C}| \leq C[\omega]$  for  $\omega = (1, \dots, k)$ .
- ii)  $\tilde{C}_{ii'} = C_{ii'}$  for  $1 \leq i \leq k$ .
- iii)  $|UV^T| \leq E$ .
- iv)  $\sigma(A, E) \leq \{\rho_0(\tilde{C})\}^{-1}$ , where  $0^{-1}$  is interpreted as  $\infty$ .

*Proof.* For  $1 \leq i, j \leq k$  follows

$$|(V^T A^{-1} U)_{ij}| = \left| \sum_{\nu=1}^n \sum_{\mu=1}^n V_{\mu i} A_{\mu\nu}^{-1} U_{\nu j} \right| \leq \sum_{\nu=1}^n |A_{i\nu}^{-1}| \cdot E_{\nu j} = C_{ij},$$

and therefore  $|\tilde{C}| \leq C[\omega]$  and i). For  $1 \leq i \leq k$  holds

$$\tilde{C}_{ii'} = (V^T A^{-1} U)_{ii'} = \sum_{\nu=1}^n \sum_{\mu=1}^n V_{\mu i} A_{\mu\nu}^{-1} U_{\nu i'} = \sum_{\nu=1}^n A_{i\nu}^{-1} \cdot \text{sign}(A_{i\nu}^{-1}) \cdot E_{\nu i'} = C_{ii'}.$$

and therefore ii). For  $1 \leq \mu, \nu \leq n$  holds

$$|(UV^T)_{\nu\mu}| = \left| \sum_{i=1}^k U_{\nu i} V_{\mu i} \right|,$$

such that  $|(UV^T)_{\nu\mu}| = E_{\nu\mu}$  for  $1 \leq \mu \leq k$ , and  $|(UV^T)_{\nu\mu}| = 0$  for  $k+1 \leq \mu \leq n$ . This proves iii). For  $\lambda := \rho_0(\tilde{C}) > 0$ , it is  $\det(\lambda \cdot I - s \cdot \tilde{C}) = 0$  for  $s = -1$  or  $s = 1$ . Lemma 4.1 implies

$$\det(A - s \cdot \lambda^{-1} \cdot UV^T) = \det(A) \cdot \det(I_k - s \cdot \lambda^{-1} \cdot V^T A^{-1} U) = 0.$$

Together with iii) and the definition (0.2) of  $\sigma(A, E)$ , this proves iv) and the theorem.  $\square$

Our aim is to construct a rank- $k$  perturbation of  $A$  with large real spectral radius. Then Lemma 4.1 allows to give an upper bound on  $\sigma(A, E)$ . A first step is the following, first generalization of Theorem 5.1. It will later yield the precise value for  $\gamma(2)$ .

THEOREM 6.5. *Let  $A \in M_n(\mathbb{R})$  be nonsingular and  $E \in M_n(\mathbb{R})$  with  $E \geq 0$ . For  $C := |A^{-1}| \cdot E$  holds*

$$(32) \quad \sigma(A, E) \leq \frac{1}{\max_{i,j} \sqrt{C_{ij} \cdot C_{ji}}}.$$

*Proof.* For  $i = j$ , (32) has been proved in Theorem 5.1. Reordering of indices puts the cycle  $(i, j)$ , for which the maximum in (6.7) is achieved, into the cycle  $(1, 2)$ , and Lemma 6.4 proves for  $i \neq j$  existence of a 2 by 2 matrix  $\tilde{C} = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}$  with  $0 \leq \beta = C_{ij}$ ,  $0 \leq \gamma = C_{ji}$ ,  $|\alpha| \leq C_{ii}$ ,  $|\delta| \leq C_{jj}$ , and  $\sigma(A, E) \leq \rho_0(\tilde{C})^{-1}$ . If  $|\alpha\delta| \geq \beta\gamma$ , then  $\sqrt{C_{ii} \cdot C_{jj}} \geq \sqrt{C_{ij} \cdot C_{ji}}$  and Theorem 5.1 yields (32). Otherwise,  $\det(\tilde{C}) < 0$ . The characteristic polynomial of  $\tilde{C}$  is  $\lambda^2 - \text{trace}(\tilde{C}) \cdot \lambda + \det(\tilde{C})$ , so that the eigenvalues of  $\tilde{C}$  are  $\frac{1}{2} \cdot \left\{ \text{trace}(\tilde{C}) \pm \sqrt{\text{trace}(\tilde{C})^2 - 4 \cdot \det(\tilde{C})} \right\} = \frac{1}{2} \cdot \left\{ \alpha + \delta \pm \sqrt{(\alpha - \delta)^2 + 4\beta\gamma} \right\}$  are both real. The absolute value of one of them is not less than  $\sqrt{\beta\gamma}$ , i.e.  $\sigma(A, E) \leq \rho_0(\tilde{C})^{-1} \leq (\beta\gamma)^{-1/2}$ .  $\square$

The idea of the proof of Theorem 6.5 is the following: for a given cycle of  $C$  of length 2, a suitable rank-2 perturbation of  $A$  is constructed which allows to prove an upper bound of  $\sigma(A, E)$  by using Lemma 6.4. In the following we will carry this idea to cycles of  $C$  of length  $k$ ,  $1 \leq k \leq n$ .

First, we will identify a class of matrices for which we can give explicit *lower* bounds for their real spectral radius. The class of matrices is constructed in such a way that the matrices given in Lemma 6.4 can be used to bound  $\sigma(A, E)$  from above.

LEMMA 6.6. *Let nonnegative  $C \in M_k(\mathbb{R})$  and some  $0 < a \in \mathbb{R}$  be given. Define  $\varphi_1 := 1$ ,  $\varphi_2 := 1$ , and for  $t > 2$  define recursively  $\varphi_t \in \mathbb{R}$  to be the positive zero of*

$$(33) \quad P_t(x) \in \mathbb{R}[x] \quad \text{with} \quad P_t(x) := x^{t-1} - x^{t-2} - \sum_{\nu=1}^{t-1} \varphi_\nu^\nu \cdot x^{t-1-\nu}.$$

*Suppose*

$$(34) \quad \forall 1 \leq \mu < k \quad \forall \bar{\omega} \in \Gamma_{\mu k} : \quad (\Pi_{\bar{\omega}}(C))^{1/\mu} \leq \varphi_\mu \cdot a,$$

*and for  $\omega = (1, \dots, k)$ ,*

$$(35) \quad |\Pi_\omega(C)|^{1/k} \geq \varphi_k \cdot a.$$

*Then, for  $i'$  defined as in (6.6) and every  $\tilde{C} \in M_k(\mathbb{R})$  with*

$$(36) \quad |\tilde{C}| \leq C \quad \text{and} \quad \tilde{C}_{ii'} = C_{ii'} \quad \text{for} \quad 1 \leq i \leq k,$$

*holds*

$$\rho_0(\tilde{C}) \geq a.$$

*Proof.* The **proof** divides in the following parts. First, we transform  $C$  into a standard form such that all  $C_{ii'}$  in the cycle  $(1, \dots, k)$  in (35) are equal. Second, we bound  $C$  by a circulant, show regularity of that matrix and  $\det(\tilde{C} - \lambda I) \neq 0$  for all  $0 \leq \lambda < a$ . Finally, the sign of the determinant of any  $\tilde{C}$  with (36) is determined, from which the lemma follows.

The case  $k = 1$  is trivial; for  $k = 2$  the proof of  $\rho_0(\tilde{C}) \geq a$  is included in the proof of Theorem 6.5.

Assume  $k > 2$ , and set  $b := |\Pi_\omega(C)|^{1/k}$ . Direct computation shows that any similarity transformation of  $C$  by a diagonal matrix  $D$  leaves all cycle products invariant.

Thus (34) and (35) remain valid for any diagonal  $D$  with positive diagonal entries. Define diagonal  $D \in M_k(\mathbb{R})$  by

$$f_i := b^{-1} \cdot C_{ii'} \quad \text{and} \quad D_{ii} := \prod_{\nu=i}^k f_\nu \quad \text{for} \quad 1 \leq i \leq k.$$

We show that w.l.o.g.  $C$  can be replaced by  $D^{-1}CD$ . We have  $f_i > 0$ , and (6.10) implies  $D_{11} = 1$ . It is

$$(37) \quad (D^{-1}CD)_{ii'} = \left( \prod_{\nu=i}^k f_\nu^{-1} \right) \cdot C_{ii'} \cdot \left( \prod_{\nu=i'}^k f_\nu \right) = C_{ii'} \cdot f_i^{-1} = b.$$

If  $\tilde{C} \in M_k(\mathbb{R})$  is any matrix satisfying (36), then  $|D^{-1}\tilde{C}D| \leq D^{-1}CD$ , and (37) yields  $(D^{-1}\tilde{C}D)_{ii'} = (D^{-1}CD)_{ii'} = b$ . Since the set of eigenvalues of  $\tilde{C}$  and  $D^{-1}\tilde{C}D$  are identical, we can restrict our attention to matrices  $C \in M_n(\mathbb{R})$ ,  $C \geq 0$  and

$$(38) \quad C_{ii'} = b \quad \text{for } 1 \leq i \leq k.$$

Set

$$(39) \quad C = \begin{pmatrix} c_{1,1} & b & c_{1,k-1} & \dots & c_{1,2} \\ c_{2,1} & c_{2,1} & b & c_{2,k-1} & \dots & c_{2,3} \\ & \dots & c_{3,1} & b & \dots & \\ c_{k-1,k-1} & c_{k-1,k-2} & & \dots & & b \\ b & c_{0,k-1} & & \dots & & c_{0,1} \end{pmatrix}.$$

Let  $\mu \in \mathbb{N}$ ,  $1 \leq \mu < k$  be given and define  $\bar{\omega} \in \Gamma_{\mu k}$  by  $\bar{\omega} = (1, \dots, \mu)$ . Then setting  $q := a/b$ , (34) implies

$$c_{\mu\mu} \cdot \prod_{i=1}^{\mu-1} b \leq (\varphi_\mu \cdot a)^\mu \quad \text{and therefore} \quad c_{\mu\mu} \leq b \cdot \varphi_\mu^\mu \cdot q^\mu.$$

Applying the same argument successively for  $\bar{\omega} = (i, (i+1) \bmod \mu, \dots, (i+\mu) \bmod \mu)$  yields

$$(40) \quad c_{i,\mu} \leq b \cdot \varphi_\mu^\mu \cdot q^\mu \quad \text{for all } 0 \leq i < k, 1 \leq \mu < k.$$

Therefore,

$$(41) \quad C \leq b \cdot \begin{pmatrix} c_1 & 1 & c_{k-1} & & c_2 \\ c_2 & c_1 & 1 & c_{k-1} & \dots & c_3 \\ & & c_1 & 1 & & \\ & \dots & & \dots & \dots & \\ c_{k-1} & c_{k-2} & & \dots & & 1 \\ 1 & c_{k-1} & & & & c_1 \end{pmatrix} =: b \cdot \bar{C}$$

with  $c_\mu := \varphi_\mu^\mu \cdot q^\mu$  for  $1 \leq \mu < n$ .

Let  $\tilde{C} \in M_k(\mathbb{R})$  with (36) be given, and let  $\lambda \in \mathbb{R}$  with  $0 \leq \lambda < a$ . Next we show that all matrices  $\tilde{C} - \lambda I$  are nonsingular. By assumption (6.11) and using (6.16),

$$(42) \quad |\tilde{C} - \lambda I| \leq C + \lambda \cdot I \leq b \cdot \bar{C} + \lambda I \quad \text{and} \quad (\tilde{C} - \lambda I)_{ii'} = \tilde{C}_{ii'} = C_{ii'} = b.$$

By (41) and (33), using  $q := a/b \leq \varphi_k^{-1}$  from (35) and  $\varphi_2 = 1$ , we have for  $k \geq 3$ ,

$$(43) \quad \begin{aligned} \lambda + b \cdot \sum_{\nu=1}^{k-1} c_\nu &< b \cdot \left\{ q + \sum_{\nu=1}^{k-1} \varphi_\nu^\nu \cdot q^\nu \right\} \leq b \cdot \left\{ \varphi_k^{-1} + \sum_{\nu=1}^{k-1} \varphi_\nu^\nu \cdot \varphi_k^{-\nu} \right\} = \\ &= b \cdot \varphi_k^{-k+1} \cdot \left\{ \varphi_k^{k-2} + \sum_{\nu=1}^{k-1} \varphi_\nu^\nu \cdot \varphi_k^{k-\nu-1} \right\} = b \cdot \varphi_k^{-k+1} \cdot \varphi_k^{k-1} = b. \end{aligned}$$

This shows that the element  $b = \tilde{C}_{ii'} = C_{ii'}$  strictly dominates the sum of the absolute values of the other components in each row of  $C + \lambda I$  and of  $\tilde{C} - \lambda I$ . That means, multiplication by a suitable permutation matrix produces a strictly diagonally dominant matrix and proves regularity of every  $\tilde{C} - \lambda I$  with  $\tilde{C}$  satisfying (36) and  $0 \leq \lambda < a$ .

We proved that for every  $\tilde{C} \in M_k(\mathbb{R})$  with (36), the determinant of  $\tilde{C} - \lambda I$  is nonzero for  $0 \leq \lambda < a$ . Therefore, the value of the characteristic polynomial  $p(\lambda) = \det(\lambda I - \tilde{C})$  of  $\tilde{C}$  has the same sign for  $0 \leq \lambda < a$ . Now  $p(\lambda) \rightarrow +\infty$  for  $\lambda \rightarrow +\infty$ . Therefore, the lemma is proved if we can show  $p(0) < 0$ , because in this case the characteristic polynomial must intersect with the real axis for some  $\lambda^* \geq a$ , thus proving  $\rho_0(\tilde{C}) \geq \lambda^* \geq a$ .

We already proved that *every* matrix  $\tilde{C}$  satisfying (36) is nonsingular. Therefore

$\text{sign}(p(0)) = \text{sign}(\det(-B))$  for every matrix  $B$  with  $|B| \leq C$  and  $B_{ii'} = C_{ii'} = b$ . Define

$$B_{ij} := \begin{cases} C_{ii'} & \text{for } j = i' \\ 0 & \text{otherwise} \end{cases}.$$

Then  $\text{sign}(\det(B)) = (-1)^{k+1}$  and therefore  $\text{sign}(p(0)) = (-1)^{2k+1} = -1$ . The theorem is proved.  $\square$

EXAMPLE 6.7. *One can show that, at least for odd  $n$ , the bounds in Lemma 6.6 are sharp in the sense that there are examples with equality in (34) and (35) such that  $\tilde{C}$  with (36) exists with  $\rho_0(\tilde{C}) = a$ . Consider*

$$C := \begin{pmatrix} a & b & c \\ c & a & b \\ b & c & a \end{pmatrix} \quad \text{with } b := \varphi_3 \cdot a \quad \text{and } c := a/\varphi_3 \quad \text{and } \tilde{C} := \begin{pmatrix} -a & b & -c \\ -c & -a & b \\ b & -c & -a \end{pmatrix}.$$

Then  $C_{11} = \varphi_1 \cdot a = a$ ,  $\sqrt{C_{12}C_{21}} = \varphi_2 \cdot a = a$  and  $(C_{12}C_{23}C_{31})^{1/3} = \varphi_3 \cdot a$ .  $\tilde{C}$  is a circulant, and its eigenvalues compute to  $P(\varepsilon^k)$ ,  $k = 0, 1, 2$  where  $\varepsilon = e^{2\pi i/3}$  and  $P(x) = bx^2 - cx - a$  (cf.[6]). It is  $b - c - a = b(1 - \varphi_2^2 q^2 - \varphi_1 q) = b \cdot \varphi_3^{-1} = a$  with  $q := a/b$ . The other two eigenvalues are complex, thus  $\rho_0(\tilde{C}) = a$ . The example extends to odd  $n \in \mathbb{N}$ .

The combination of Lemma 6.4, Theorem 6.5, and Lemma 6.6 gives the key to construct a rank- $k$  perturbation of  $A$  to achieve an upper bound for  $\sigma(A, E)$ . The following theorem is the generalization of Theorems 5.1 and 6.5 for cycles of length  $k$ ,  $1 \leq k \leq n$ .

THEOREM 6.8. *Let  $A, E \in M_n(\mathbb{R})$  with nonsingular  $A$  and  $E \geq 0$  be given and define  $C := |A^{-1}| \cdot E$ . For  $1 \leq k \leq n$  and any  $\omega \in \Gamma_{kn}$  set*

$$(44) \quad 0 \neq \tau := (\Pi_\omega(C))^{1/k}.$$

Then for  $\varphi_k$  defined as in Lemma 6.6,

$$\sigma(A, E) \leq \varphi_k / \tau.$$

In other words,  $\varphi_k$  divided by the geometric mean of the elements of any cycle of  $C$  bounds  $\sigma(A, E)$  from above.

*Proof.* Let some  $\omega \in \Gamma_{kn}$  and  $\tau$  from (44) be given and set  $a := \tau/\varphi_k$ . If for  $k = 1$  or  $k = 2$  there exists some  $\bar{\omega} \in \Gamma_{kn}$  with  $\{\Pi_{\bar{\omega}}(C)\}^{1/k} \geq \varphi_k \cdot a$ , then  $\varphi_1 = \varphi_2 = 1$  and Theorem 5.1 and Theorem 6.5 imply  $\sigma(A, E) \leq a^{-1} = \varphi_k/\tau$ . Therefore, we may assume  $\{\Pi_{\bar{\omega}}(C)\}^{1/k} < \varphi_k \cdot a$  for all  $\bar{\omega} \in \Gamma_{kn}$  and  $k \in \{1, 2\}$ . Hence, there is some  $m \in \mathbb{N}$ ,  $2 \leq m \leq k$  such that

$$\forall 1 \leq \mu < m \quad \forall \bar{\omega} \in \Gamma_{\mu n} : (\Pi_{\bar{\omega}}(C))^{1/\mu} \leq \varphi_\mu \cdot a,$$

and

$$\exists \tilde{\omega} \in \Gamma_{mn} : (\Pi_{\tilde{\omega}}(C))^{1/m} \geq \varphi_m \cdot a.$$

After suitable rearrangement of indices we may assume  $\tilde{\omega} = (1, \dots, m)$ , and Lemma 6.4 yields a matrix  $\tilde{C} \in M_m(\mathbb{R})$  with properties i) and ii) of Lemma 6.4 and  $\sigma(A, E) \leq \rho_0(\tilde{C})^{-1}$ . But Lemma 6.6 shows for all such matrices  $\rho_0(\tilde{C}) \geq a = \tau/\varphi_m$ . Regarding  $m \leq k$ , the theorem is proved if we can show

$$(45) \quad t \in \mathbb{N} \quad \Rightarrow \quad \varphi_t \leq \varphi_{t+1}.$$

We know  $\varphi_1 = \varphi_2 = 1$ , from the definition (6.8) we see  $\varphi_3 = 1 + \sqrt{2}$ , and for  $t \geq 3$

$$P_{t+1}(x) = x \cdot P_t(x) - \varphi_t^t.$$

Hence  $P_{t+1}(\varphi_t) < 0$  and  $\varphi_{t+1} > \varphi_t$ . The theorem is proved.  $\square$

Theorem 6.8 reduces the problem of finding upper bounds of  $\sigma(A, E)$  to finding proper cycles of some length  $k$  of  $|A^{-1}| \cdot E$  with large geometric mean corresponding to a suitable rank- $k$  perturbation. This is done in the following proof of Proposition 6.3.

*Proof.* of Proposition 6.3. Corollary 6.2 allows us to assume  $E > 0$ . Therefore,  $|A^{-1}| \cdot E$  is positive, and Perron-Frobenius Theory yields existence of a positive eigenvector  $x \in V_n(\mathbb{R})$  with  $|A^{-1}| \cdot E \cdot x = \rho(|A^{-1}| \cdot E) \cdot x$ ,  $\rho(|A^{-1}| \cdot E) > 0$ . Define the diagonal matrix  $D_x \in M_n(\mathbb{R})$  by  $(D_x)_{ii} := x_i$ . We may replace  $A$  by  $A \cdot D_x$  and  $E$  by  $E \cdot D_x$ , because for any nonsingular diagonal matrices  $D_1, D_2$ ,  $\sigma(A, E) = \sigma(D_1 A D_2, D_1 E D_2)$ . This is because  $|\delta A| \leq \sigma \cdot E$  iff  $|D_1 \cdot \delta A \cdot D_2| \leq \sigma \cdot |D_1 E D_2|$  and  $A + \delta A$  is singular iff  $D_1 A D_2 + D_1 \cdot \delta A \cdot D_2$  is singular (cf. [3]). Then

$$C := |(A \cdot D_x)^{-1}| \cdot E \cdot D_x = D_x^{-1} \cdot |A^{-1}| \cdot E \cdot D_x \quad \text{and} \quad C \cdot (\mathbf{1}) = \rho(|A^{-1}| \cdot E) \cdot (\mathbf{1}).$$

That means  $C$  is a multiple of a row stochastic matrix. Set  $\rho := \rho(|A^{-1}| \cdot E)$ .

Denote an index of the maximal component of  $C$  in row  $i$  by  $m_i$ . Then either  $\{m_i \mid 1 \leq i \leq n\} = \{1, \dots, n\}$  or, there is a cycle  $m_j, m_{j+1}, \dots, m_{j+k-1}, m_{j+k} = m_j$  of length  $k$ . That means, with a suitable renumbering, there is some  $k \in \mathbb{N}$ ,  $1 \leq k \leq n$  such that for the upper left  $k$  by  $k$  principal submatrix of  $C$  holds

$$(46) \quad C_{ii'} \geq \rho/n \quad \text{for} \quad 1 \leq i \leq k,$$

where  $i'$  is defined as in (6.6).

Then Theorem 6.8, (6.20) and (6.21) imply for  $\omega = (1, \dots, k)$ ,

$$\sigma(A, E) \leq \varphi_k \cdot \{\Pi_\omega(C)\}^{-1/k} \leq n \cdot \varphi_k / \rho \leq n \cdot \varphi_n / \rho. \quad \square$$

In the remaining §7, we will replace the bound (30) by giving explicit bounds for  $\gamma(n)$  only depending on  $n$ . An asymptotic bound will be given as well.

**7. Explicit bounds for  $\gamma(n)$ .** The main result in §6 is the upper bound (30) in Proposition 6.3. This bound is given in terms of  $\varphi_k$ , the positive zeros of the polynomial  $P_t$  defined in (28). In the remaining part of the paper we will give bounds on  $\gamma(n)$  showing the dependence on  $n$  by a simple function. Moreover, the asymptotic behaviour of  $\gamma(n)$  for  $n \rightarrow \infty$  is given.

The polynomials  $P_t(x) \in \mathbb{R}[x]$  defined in (28) satisfy

$$(47) \quad P_t(x) = x^{t-1} - x^{t-2} - \sum_{\nu=1}^{t-1} \varphi_\nu^\nu \cdot x^{t-1-\nu} \quad \text{and} \quad P_t(\varphi_t) = 0 \quad \text{for} \quad t > 2.$$

Therefore, for  $n \geq 3$ ,

$$(48) \quad \varphi_n^{-1} + \sum_{i=1}^{n-1} \varphi_i^i \cdot \varphi_n^{-i} = 1.$$

By (6.20),  $x + \sum_{i=1}^{n-1} \varphi_i^i \cdot x^i$  is strictly increasing for  $x > 0$ . Hence, for  $x > 0$ ,

$$(49) \quad x + \sum_{i=1}^{n-1} \varphi_i^i \cdot x^i \leq 1 \quad \text{implies} \quad x \leq \varphi_n^{-1}, \quad \text{that is} \quad \varphi_n \leq x^{-1}.$$

We are aiming on a bound of the form

$$(50) \quad \varphi_k \leq c \cdot k^\alpha$$

for some constants  $c$  and  $\alpha$ . To determine  $c$  and  $\alpha$ , we notice that if (7.4) is satisfied for  $1 \leq k < n$ , then

$$(51) \quad \sum_{i=1}^{n-1} \left(\frac{i}{n}\right)^{\alpha i} \leq 1 - c^{-1} \cdot n^{-\alpha} \quad \text{implies} \quad \varphi_n \leq c \cdot n^\alpha.$$

This is because the left hand side of (7.5) yields

$$1 \geq c^{-1} \cdot n^{-\alpha} + \sum_{i=1}^{n-1} i^{\alpha i} \cdot n^{-\alpha i} \geq (c \cdot n^\alpha)^{-1} + \sum_{i=1}^{n-1} \varphi_i^i \cdot (c \cdot n^\alpha)^{-i},$$

and (7.3) implies  $(c \cdot n^\alpha)^{-1} \leq \varphi_n^{-1}$ .

Therefore, our first step is to derive upper bounds for

$$(52) \quad \sum_{i=1}^{n-1} \sigma_i \quad \text{with} \quad \sigma_i := \left(\frac{i}{n}\right)^{i\alpha}.$$

$\sigma_i$  depends on  $n$  and  $\alpha$ . We use the abbreviation  $\sigma_i$  for fixed  $n$  and  $\alpha$  and omit extra parameters for better readability. In order to estimate the sum (52), we will split it into 3 parts, which will be bounded individually. For  $i \geq 1$  holds

$$\frac{\sigma_{i+1}}{\sigma_i} = \left(\frac{i+1}{n}\right)^{(i+1)\alpha} \cdot \left(\frac{n}{i}\right)^{(i+1)\alpha} \cdot \left(\frac{i}{n}\right)^\alpha = \left(\frac{i}{n}\right)^\alpha \cdot \left(1 + \frac{1}{i}\right)^{(i+1)\alpha} > \left(\frac{i}{n}\right)^\alpha \cdot e^\alpha,$$

and therefore

$$(53) \quad \sigma_i < \left(\frac{n}{i \cdot e}\right)^\alpha \cdot \sigma_{i+1} \quad \text{for} \quad i \geq 1.$$

For all  $\beta \in \mathbb{R}$  with  $1 < \beta < e$  and  $k := \lceil \frac{n\beta}{e} \rceil$  holds  $k-1 < \frac{n\beta}{e} \leq k$ . Then (53) gives

$$\sum_{i=k}^{n-1} \sigma_i < \sigma_{n-1} + \sum_{i=k}^{n-2} \left(\frac{n}{i \cdot e}\right)^\alpha \cdot \sigma_{i+1} = \sigma_{n-1} + \sum_{i=k+1}^{n-1} \left(\frac{n}{(i-1) \cdot e}\right)^\alpha \cdot \sigma_i,$$

and therefore

$$\sigma_{n-1} - \sigma_k > \sum_{i=k+1}^{n-1} \left\{1 - \left(\frac{n}{(i-1) \cdot e}\right)^\alpha\right\} \cdot \sigma_i \geq \sum_{i=k+1}^{n-1} \left\{1 - \left(\frac{n}{k \cdot e}\right)^\alpha\right\} \cdot \sigma_i,$$

and  $\frac{n}{k \cdot e} \leq \beta^{-1}$  yields

$$\sigma_{n-1} - \sigma_k > (1 - \beta^{-\alpha}) \cdot \sum_{i=k+1}^{n-1} \sigma_i.$$

By choice,  $\beta > 1$ , and  $\alpha \geq 0$  implies  $(1 - \beta^{-\alpha})^{-1} > 1$ . Therefore,

$$(54) \quad \sum_{i=k}^{n-1} \sigma_i < (1 - \beta^{-\alpha})^{-1} \cdot \sigma_{n-1} = (1 - \beta^{-\alpha})^{-1} \cdot \left(\frac{n-1}{n}\right)^{(n-1)\alpha} =: \mu_n$$

holds for every  $\alpha \geq 0$ ,  $1 < \beta < e$  and  $k := \lceil \frac{n\beta}{e} \rceil$ . This is the first part of the sum (7.6) for a suitable  $k$  to be determined. Define

$$f(x) := \left(\frac{x}{n}\right)^{x\alpha} \quad \text{with} \quad f'(x) = \left(\frac{x}{n}\right)^{x\alpha} \cdot \left\{\alpha \cdot \ln \frac{x}{n} + \alpha\right\}.$$

For  $x > 0$ ,  $f(x)$  has exactly one minimum at  $x = \frac{n}{e}$ . Then  $f(i) = \sigma_i$  shows

$$(55) \quad \sigma_k \geq \sigma_l \quad \text{for} \quad 1 \leq k \leq l \leq \frac{n}{e}, \quad \text{and} \quad \sigma_k \leq \sigma_l \quad \text{for} \quad \frac{n}{e} \leq k \leq l \leq n-1.$$

Set  $M := \lceil n/e \rceil$ . Then  $k = \lceil \frac{n\beta}{e} \rceil$  satisfies  $M \leq k \leq n$ , and (55) implies for  $n \geq 3$ ,

$$(56) \quad \begin{aligned} \sum_{i=M}^{k-1} \sigma_i &\leq (k-M) \cdot \sigma_{k-1} < \left(\frac{n\beta}{e} + 1 - \frac{n}{e}\right) \cdot f\left(\frac{n\beta}{e}\right) < \frac{n\beta}{e} \cdot f\left(\frac{n\beta}{e}\right) \\ &\leq n \cdot \left(\frac{\beta}{e}\right)^{n\frac{\beta\alpha}{e} + 1} =: \nu_n. \end{aligned}$$

This is the second part of the sum (7.6). Finally, (7.9) implies

$$(57) \quad \sum_{i=1}^{M-1} \sigma_i < \left(\frac{1}{n}\right)^\alpha + \left(\frac{2}{n}\right)^{2\alpha} + \left(\frac{3}{n}\right)^{3\alpha} + \left(\frac{4}{n}\right)^{4\alpha} \cdot \left(\frac{n}{e}\right) =: \xi_n,$$

which is the third part of the sum (7.6). The inequalities (54), (56) and (57) together yield

$$(58) \quad n^{-\alpha} + \sum_{i=1}^{n-1} \left(\frac{i}{n}\right)^{i\alpha} \leq n^{-\alpha} + \mu_n + \nu_n + \xi_n \quad \text{for} \quad n \geq 3.$$

Next, we show that all three sequences  $\mu_n, \nu_n, \xi_n$  are decreasing for large enough  $n$ .  $(1 + \frac{1}{n})^n$  is monotonically increasing for  $n \geq 1$ , therefore for  $n \geq 2$ ,

$$\left(\frac{n+1}{n}\right)^{n\alpha} \geq \left(\frac{n}{n-1}\right)^{(n-1)\alpha} \Rightarrow \mu_{n+1} \leq \mu_n.$$

Suppose

$$(59) \quad n_0 \geq \left\{ \left( \frac{e}{\beta} \right)^{\frac{\beta\alpha}{e}} - 1 \right\}^{-1}.$$

Then for  $n \geq n_0$ ,

$$1 + \frac{1}{n} \leq \left( \frac{e}{\beta} \right)^{\frac{\beta\alpha}{e}} \Rightarrow (n+1) \cdot \left( \frac{\beta}{e} \right)^{\frac{\beta\alpha}{e}} \leq n \Rightarrow (n+1) \cdot \left( \frac{\beta}{e} \right)^{(n+1)\frac{\beta\alpha}{e}+1} \leq n \cdot \left( \frac{\beta}{\alpha} \right)^{n\frac{\beta\alpha}{e}+1},$$

and therefore  $\nu_{n+1} \leq \nu_n$  for  $n \geq n_0$  with  $n_0$  satisfying (59). Finally, for  $n \geq 1$  and  $\alpha > 0.25$ ,  $1 - 4\alpha < 0$  and therefore

$$(n+1)^{1-4\alpha} \leq n^{1-4\alpha} \Rightarrow \left( \frac{4}{n+1} \right)^{4\alpha} \cdot \left( \frac{n+1}{e} \right) \leq \left( \frac{4}{n} \right)^{4\alpha} \cdot \left( \frac{n}{e} \right) \Rightarrow \xi_{n+1} \leq \xi_n.$$

Summarizing, this proves the following lemma.

LEMMA 7.1. Define  $\varphi_1 := 1$ ,  $\varphi_2 := 1$  and recursively  $\varphi_n$  to be the positive zero of  $P_n(x)$  given in (47). Let constants  $c, \alpha \in \mathbb{R}$ ,  $\alpha \geq \ln 2$  and  $3 \leq n_0 \in \mathbb{N}$  be given with  $\varphi_n \leq c \cdot n^\alpha$  for  $n < n_0$ . If a constant  $\beta \in \mathbb{R}$ ,  $1 < \beta < e$  exists such that (59) is satisfied and  $\mu_n, \nu_n, \xi_n$  defined in (54), (56), (57) satisfy

$$(60) \quad n^{-\alpha} + \mu_n + \nu_n + \xi_n \leq 1 \quad \text{for } n = n_0,$$

then

$$\varphi_n \leq c \cdot n^\alpha \quad \text{for all } n \in \mathbb{N}.$$

*Proof.* (50) is satisfied for  $1 \leq k < n$ , and (58) and (7.14) prove the left hand side of (51) for  $n = n_0$ , and therefore (7.4) for  $k = n$ . The quantities  $n^{-\alpha}$ ,  $\mu_n$ ,  $\nu_n$  and  $\xi_n$  are decreasing for increasing  $n$ . Thus, (7.14) and therefore (7.4) is valid for all  $n \geq n_0$ . By assumption,  $\varphi_n \leq c \cdot n^\alpha$  for  $n < n_0$  as well.  $\square$

For example, for  $\beta := 2.697$ ,  $\alpha := 0.7$  and  $n_0 := 3000$ , one checks by explicit calculation  $\varphi_n \leq 2.321 \cdot n^\alpha$  for  $1 \leq n \leq n_0$ . The lower bound (59) for  $n_0$  is less than 183,  $\mu_n < 0.992$ ,  $\nu_n < 0.0003$ ,  $\xi_n < 0.0038$ , and  $n^{-\alpha} < 0.0038$  for  $n = n_0$ . This proves the following result.

COROLLARY 7.2. For all  $n \geq 1$ ,  $\varphi_n \leq 2.321 \cdot n^{0.7}$ . The difference  $2.321 \cdot n^{0.7} - \varphi_n$  is less than 2.8 for  $1 \leq n < 20$ , and less than 2.0 for  $20 \leq n \leq 2000$ .

Summarizing, Corollary 7.2, Proposition 6.3, and Lemma 5.7 prove the following result.

PROPOSITION 7.3. Let  $A, E \in M_n(\mathbb{R})$  with nonsingular  $A$  and  $E \geq 0$  be given. Then for all  $n \geq 1$

$$\frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \leq \frac{\gamma(n)}{\rho(|A^{-1}| \cdot E)},$$

with

$$n \leq \gamma(n) \leq 2.321 \cdot n^{1.7}.$$

The lower bound for  $\gamma(n)$  is sharp.

Finally, we will show the asymptotic behaviour of upper bounds for  $\gamma(n)$ . Let  $\alpha := \ln(2 + 2\eta)$ ,  $\eta > 0$ . For any  $1 < \beta < e$  and  $n \rightarrow \infty$ ,

$$n^{-\alpha} \rightarrow 0, \quad \mu_n \rightarrow (1 - \beta^{-\alpha})^{-1} \cdot e^{-\alpha}, \quad \nu_n \rightarrow 0 \quad \text{and} \quad \xi_n \rightarrow 0.$$

For  $\ln \beta := (2 + \eta)/(2 + 2\eta)$ , a short computation yields

$$(1 - \beta^{-\alpha})^{-1} \cdot e^{-\alpha} = \frac{2 + \eta}{2 + 4\eta + 2\eta^2} < 1.$$

Hence, for this  $\beta$  and large enough  $n_0$ , (59) holds and

$$n^{-\alpha} + \mu_n + \nu_n + \xi_n < 1 \quad \text{for all } n \geq n_0.$$

Therefore, for large enough  $c$  with  $\varphi_n \leq c \cdot n^\alpha$  for  $n < n_0$ , Lemma 7.1 implies that  $\varphi_n \leq c \cdot n^\alpha$  for all  $n \in \mathbb{N}$ . Using  $\alpha > \ln 2$  proves the following.

PROPOSITION 7.4. *Let  $\gamma(n)$  be defined as follows:*

$$\gamma(n) := \inf\{\sigma(A, E) \cdot \rho(|A^{-1}| \cdot E) \mid A \in M_n(\mathbb{R}) \text{ nonsingular and } 0 \leq E \in M_n(\mathbb{R})\}.$$

*Then  $\gamma(n)$  is finite for all  $n \in \mathbb{N}$ . Moreover, for any  $\varepsilon > 0$  there exists some  $n_0 \in \mathbb{N}$  such that for all  $n \geq n_0$  holds*

$$(61) \quad n \leq \gamma(n) \leq n^{1+\ln 2+\varepsilon}.$$

*The lower bound in (61) is sharp. †)*

In his paper [3], Demmel showed that for the Bauer-Skeel condition number  $\kappa(A, E) := \||A^{-1}| \cdot E\|$  with any  $p$ -norm,  $1 \leq p \leq \infty$ , there holds

$$\frac{1}{\rho(|A^{-1}| \cdot E)} = \frac{1}{\min_D \kappa(AD, ED)},$$

where the minimum is taken over all diagonal  $D$ . Thus, Proposition 7.3 and Proposition 7.4 prove that the componentwise relative distance to the nearest singular matrix for any weight matrix  $E \geq 0$  is not too far from the reciprocal of the smallest condition number achievable by column scaling. The evidence presented in this paper leads us to the following conjecture.

CONJECTURE 7.5. *For all nonsingular  $A \in M_n(\mathbb{R})$  and  $0 \leq E \in M_n(\mathbb{R})$  holds*

$$(62) \quad \frac{1}{\rho(|A^{-1}| \cdot E)} \leq \sigma(A, E) \leq \frac{n}{\rho(|A^{-1}| \cdot E)}.$$

*If the conjecture is true, Lemma 5.7 shows that it is sharp.*

**Acknowledgments.** The author wishes to thank Jiri Rohn for many helpful comments.

#### REFERENCES

- [1] F. BAUER, *Optimally scaled matrices*, Numerische Mathematik 5, (1963), pp. 73–87.
- [2] L. COLLATZ, *Einschließungssatz für die charakteristischen Zahlen von Matrizen*, Math. Z., 48 (1942), pp. 221–226.
- [3] J. DEMMEL, *The Componentwise Distance to the Nearest Singular Matrix*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 10–19.
- [4] G. ENGEL AND H. SCHNEIDER, *The Hadamard-Fischer Inequality for a Class of Matrices Defined by Eigenvalue Monotonicity*, Linear and Multilinear Algebra 4, (1976), pp. 155 – 176.
- [5] R. HORN AND C. JOHNSON, *Matrix Analysis*, Cambridge University Press, 1985.
- [6] M. MARCUS AND H. MINC, *A survey of matrix theory and matrix inequalities*, Dover publications, New York, 1992.
- [7] S. POLJAK AND J. ROHN, *Checking Robust Nonsingularity Is NP-Hard*, Math. of Control, Signals, and Systems 6, (1993), pp. 1–9.
- [8] J. ROHN, *Nearness of Matrices to Singularity*, KAM Series on Discrete Mathematics and Combinatorics, (1988).
- [9] ———, *Systems of Linear Interval Equations*, Linear Algebra Appl. 126, (1989), pp. 39–78.
- [10] ———, *Interval Matrices: Singularity and Real Eigenvalues*, SIAM J. Matrix Anal. Appl. 14, (1993), pp. 82–91.
- [11] S. RUMP, *Estimation of the Sensitivity of Linear and Nonlinear Algebraic Problems*, Linear Algebra and its Applications 153, (1991), pp. 1–34.
- [12] ———, *Verification Methods for Dense and Sparse Systems of Equations*, in Topics in Validated Computations — Studies in Computational Mathematics, J. Herzberger, ed., Elsevier, Amsterdam, 1994, pp. 63–136.

---

†)Note added in proof: In the meantime it has been shown by the author that  $n \leq \gamma(n) \leq (3 + 2\sqrt{2}) \cdot n$ .

[13] G. STEWART AND J. SUN, *Matrix Perturbation Theory*, Academic Press, 1990.