

**BLOCK DIAGONALIZATION OF NEARLY DIAGONAL
MATRICES**

E. KOVAČ STRIKO AND K. VESELIĆ

University of Zagreb, Croatia and Fernuniversität Hagen, Germany

ABSTRACT. In this paper we study the effect of block diagonalization of a nearly diagonal matrix by iterating the related Riccati equations. We show that the iteration is fast, if a matrix is diagonally dominant or scaled diagonally dominant and the block partition follows an appropriately defined spectral gap. We also show that both kinds of diagonal dominance are not destroyed after the block diagonalization.

1. INTRODUCTION

In this note we consider block analoga of plane transformations in order to block diagonalize a given matrix. Iterations of plane rotations were introduced by Jacobi [4], however it does not seem to be broadly known (see Sleipen, van der Vorst [6]) that Jacobi himself was not satisfied with the convergence of his method and tried to accelerate it by using generalized rotations which annihilate a whole off-diagonal block. Of course, this complicates the computation of rotation parameters which are easily computable only if the matrix is already diagonal enough in which case the corresponding Riccati equations are solvable by a couple of iterations.¹

We will consider the effect of a block diagonalization on matrices which are diagonally dominant i.e.

$$\|D^{-1}(A - D)\| < 1$$

where D is the diagonal part of A as well as on the scaled diagonally dominant matrices for which

$$\| |D|^{-1/2}(A - D)|D|^{-1/2} \| < 1.$$

This work was done during the first author's stay at the Fernuniversität Hagen.

¹Block rotations with the corresponding Riccati equations were also used in [5] as iterative improvement to a computed eigensolution.

(Here $\|\cdot\|$ is the spectral norm. As it is known, on symmetric matrices the scaled diagonal dominance is a weaker property than the usual scaled dominance).

In Section 2 we prove that Jacobi or Gauss - Seidel like iterations solve the Riccati equations with a number of steps independent of the dimension. This allows to find few extremal eigenpairs in some n^2 operations, both in the diagonally dominant and the scaled diagonally dominant case.

In Section 3 we study the effect of block diagonalization and show that the (scaled) diagonal dominance is preserved after the block diagonalization step. This opens the prospective to use our block transformations as an iterative correction on matrices whose eigensolution was poor because of their high condition number (see [3]).

In Section 4 we present some simple illustrative numerical examples.

2. BLOCK DIAGONALIZATION

Let $A \in C^{n \times n}$,

$$(2.1) \quad A = \begin{bmatrix} a & b \\ c & d \end{bmatrix},$$

where $a = [a_{ij}]$, $[d_{ij}]$ are square of order m , $n - m$, respectively.²

We will consider two types of block transformations

a)

$$(2.2) \quad A \rightarrow A_1 = V^{-1}AV, \quad A_1 = \begin{bmatrix} a_1 & b_1 \\ c_1 & d_1 \end{bmatrix}$$

$$(2.3) \quad V = V(u, t) = \begin{bmatrix} I & u \\ -t & I \end{bmatrix} \begin{bmatrix} (I + ut)^{-1/2} & 0 \\ 0 & (I + tu)^{-1/2} \end{bmatrix}$$

$$(2.4) \quad V^{-1} = \begin{bmatrix} (I + ut)^{-1/2} & 0 \\ 0 & (I + tu)^{-1/2} \end{bmatrix} \begin{bmatrix} I & -u \\ t & I \end{bmatrix}$$

Here the square root is taken as an analytic continuation of the binomial series (if it exists, see [9]).

b)

$$(2.5) \quad A \rightarrow A_2 = U^{-1}AU, \quad A_2 = \begin{bmatrix} a_2 & b_2 \\ c_2 & d_2 \end{bmatrix}$$

²We use lower case letters for matrix blocks to remind the connection with the classical Jacobi algorithm.

$$(2.6) \quad U = U(u, t) = \begin{bmatrix} I & u \\ -t & I \end{bmatrix} \begin{bmatrix} L_1^{-1} & 0 \\ 0 & L_2^{-1} \end{bmatrix}$$

$$(2.7) \quad U^{-1} = \begin{bmatrix} R_1^{-1} & 0 \\ 0 & R_2^{-1} \end{bmatrix} \begin{bmatrix} I & -u \\ t & I \end{bmatrix}$$

where

$$(2.8) \quad I + ut = R_1 L_1 \quad , \quad I + tu = R_2 L_2 \quad ,$$

with L_1, L_2 upper triangular, R_1, R_2 lower triangular and $\text{diag } L_i = \text{diag } R_i > 0, \quad i = 1, 2.$

These two types are the simplest block transformations which tend to the identity, if the matrix A is nearly block diagonal. The second transformation avoids the use of matrix square roots. It should, however, be noted that even in the first case the complexity in computing the square roots is tolerable, if $m \ll n$, because $I + ut, \quad I + tu$ are low - rank perturbations of identity matrices.

For instance, if $m = 1$ then in (2.3) ut is just a number and

$$(I + tu)^{-1/2} = I - \frac{tu}{1 + ut + (1 + ut)^{1/2}}$$

as is immediately seen by using the binomial formula. The variables u, t can be understood as some kind of "tangents". In [9] it was shown how to use "tan $\frac{x}{2}$ " substitution to get rid of all square roots in (2.3).

In case a) we have

$$(2.9) \quad \begin{aligned} a_1 &= (I + ut)^{-1/2}(a - uc - bt + udt)(I + ut)^{-1/2} \\ b_1 &= (I + ut)^{-1/2}(au - ud + b - ucu)(I + tu)^{-1/2} \\ c_1 &= (I + tu)^{-1/2}(ta - dt + c - tbt)(I + ut)^{-1/2} \\ d_1 &= (I + tu)^{-1/2}(tau + cu + tb + d)(I + tu)^{-1/2} \end{aligned}$$

and in case b):

$$(2.10) \quad \begin{aligned} a_2 &= R_1^{-1}(a - uc - bt + udt)L_1^{-1} \\ b_2 &= R_1^{-1}(au - ud + b - ucu)L_2^{-1} \\ c_2 &= R_2^{-1}(ta - dt + c - tbt)L_1^{-1} \\ d_2 &= R_2^{-1}(tua + cu + tb + d)L_2^{-1}. \end{aligned}$$

Block diagonalization means to solve the Riccati equations

$$(2.11) \quad ta - dt + c - tbt = 0$$

$$(2.12) \quad au - ud + b - ucu = 0 .$$

Note that both

$$V(t^H, t) \quad \text{and} \quad U(t^H, t)$$

are unitary, in this case (2.8) are Cholesky decompositions. If A is Hermitian then with t solving (2.11) the matrix $u = t^H$ solves (2.12) and vice versa.

Furthermore the Riccati equation (2.11) determines the matrix

$$S = S(t) = \begin{bmatrix} I & 0 \\ -t & I \end{bmatrix} \quad \text{with} \quad S^{-1} = \begin{bmatrix} I & 0 \\ t & I \end{bmatrix}$$

such that

$$S^{-1}AS = \begin{bmatrix} * & * \\ 0 & * \end{bmatrix}.$$

The obtained upper block triangular matrix is then block diagonalized as

$$\begin{bmatrix} I & -u \\ 0 & I \end{bmatrix} \begin{bmatrix} * & * \\ 0 & * \end{bmatrix} \begin{bmatrix} I & u \\ 0 & I \end{bmatrix} = \begin{bmatrix} * & 0 \\ 0 & * \end{bmatrix};$$

the corresponding Riccati equation is now linear.

Next we will be solving Riccati equations (2.11) and (2.12) in the case where the matrices b, c are small in some sense.

Let

$$(2.13) \quad a = \Delta_a + a^u + a^l, \quad d = \Delta_d + d^u + d^l$$

where Δ_a, a^u, a^l is the diagonal, strict upper triangular, strict lower triangular part of a , respectively (and similarly for d).

Gauss-Seidel iterates t^j for the Riccati equation (2.11) are defined by

$$(2.14) \quad t^{j+1}(\Delta_a + a^u) - (\Delta_d + d^l)t^{j+1} + t^j a^l - d^u t^j + c - t^j b t^j = 0,$$

t^0 arbitrary. Similarly, the iteration for (2.12) reads

$$(2.15) \quad (\Delta_a + a^u)u^{j+1} - u^{j+1}(\Delta_d + d^l) + a^l u^j - u^j d^u + b - u^j c u^j = 0,$$

u^0 arbitrary.

Besides (2.14) there are other iterations, defined by the leading terms

$$(2.16) \quad \begin{array}{ll} t^{j+1}(\Delta_a + a^u) & - (\Delta_d + d^u)t^{j+1} \\ t^{j+1}(\Delta_a + a^l) & - (\Delta_d + d^l)t^{j+1} \\ t^{j+1}(\Delta_a + a^l) & - (\Delta_d + d^u)t^{j+1} \\ t^{j+1}(\Delta_a + a^l) & - \Delta_d t^{j+1} \\ t^{j+1}(\Delta_a + a^u) & - \Delta_d t^{j+1} \\ t^{j+1}\Delta_a & - (\Delta_d + d^u)t^{j+1} \\ t^{j+1}\Delta_a & - (\Delta_d + d^l)t^{j+1} \\ t^{j+1}\Delta_a & - \Delta_d t^{j+1} \end{array}$$

(and similarly for the iteration (2.15) with u). The last choice is just the Jacobi iteration. They are all treated in the same way. We will now state a convergence result for the iteration (2.14).

We first consider the case in which the matrix A is nearly block-diagonal in the standard sense.

THEOREM 2.1. *Let A be as in (2.1), (2.13). Assume that*

$$(2.17) \quad 2\sqrt{\|b\|_E \|c\|_E} + \|a^u\| + \|a^l\| + \|d^u\| + \|d^l\| < \beta$$

where³

$$(2.18) \quad \beta = \min_{i,j} |a_{jj} - d_{ii}|.$$

Then the Riccati equation (2.11) has a unique solution t^* in the ball $K = K(0, r)$ with

$$r = \frac{2\|c\|_E}{\beta - \|a^u\| - \|a^l\| - \|d^u\| - \|d^l\|}.$$

Moreover, for the sequence t^j from (2.14) we have

$$\|t^* - t^j\|_E \leq \frac{k^j}{1 - k} \|t^1 - t^0\|_E,$$

with

$$k = \frac{\|a^l\| + \|d^u\| + 2r\|b\|_E}{\beta - \|a^u\| - \|d^l\|} < 1.$$

The estimates for the u -variables are obtained just by interchanging the roles of b and c . Moreover, the same conditions imply the convergence of any of the iterations in (2.16), and similarly for u^j .

We omit the proof of this theorem. It can be easily reconstructed from the harder, scaled diagonally case below. The conditions (2.17), (2.18) in fact imply that the spectra of a and d are sufficiently distant.

We now consider matrices with the scaled diagonal dominance i.e. we set

$$(2.19) \quad A = DA_0D = \begin{bmatrix} D_a & 0 \\ 0 & D_d \end{bmatrix} \begin{bmatrix} a_0 & b_0 \\ c_0 & d_0 \end{bmatrix} \begin{bmatrix} D_a & 0 \\ 0 & D_d \end{bmatrix}$$

³Here $\|\cdot\|_E$ is the Euclidean norm on matrices and $\|\cdot\|$ is the operator norm on matrices or matrix operators, induced by the Euclidean matrix norm.

where D is diagonal, diagonal entries of A_0 have the absolute value one and A_0 is diagonally dominant in the usual sense.

We use the substitution

$$(2.20) \quad t = D_d \tau D_a^{-1} \quad , \quad u = D_a^{-1} \nu D_d$$

such that (2.11) , (2.12) are equivalent to

$$(2.21) \quad \tau a_0 - d_0 D_d^2 \tau D_a^{-2} + c_0 - \tau b_0 D_d^2 \tau D_a^{-2} = 0,$$

$$(2.22) \quad a_0 \nu - D_a^{-2} \nu D_d^2 d_0 + b_0 - D_a^{-2} \nu D_d^2 c_0 \nu = 0.$$

Let

$$(2.23) \quad a_0 = E_a + a_0^u + a_0^l \quad , \quad d_0 = E_d + d_0^u + d_0^l$$

where E_a, a_0^u, a_0^l is the diagonal, strict upper triangular, strict lower triangular part of a_0 , respectively (and similarly for d_0). The equations (2.14), (2.15) are equivalent to

$$(2.24) \quad \tau^{j+1} (E_a + a_0^u) - (E_d + d_0^l) D_d^2 \tau^{j+1} D_a^{-2} + \tau^j a_0^l - d_0^u D_d^2 \tau^j D_a^{-2} + c_0 - \tau^j b_0 D_d^2 \tau^j D_a^{-2} = 0 ,$$

τ_0 arbitrary,

$$(2.25) \quad (E_a + a_0^u) \nu^{j+1} - D_a^{-2} \nu^{j+1} D_d^2 (E_d + d_0^l) + a_0^l \nu^j - D_a^{-2} \nu^j D_d^2 d_0^u + b_0 - D_a^{-2} \nu^j D_d^2 c_0 \nu^j = 0,$$

ν_0 arbitrary.

THEOREM 2.2. *Let A be as in (2.19) , (2.23). Assume that*

$$(2.26) \quad \alpha = \frac{\max_i |d_{ii}|}{\min_j |a_{jj}|} \leq 1$$

and

$$(2.27) \quad 2\sqrt{\alpha \|b_0\|_E \|c_0\|_E} + \|a_0^u\| + \|a_0^l\| + \alpha(\|d_0^u\| + \|d_0^l\|) < \beta_s$$

where

$$(2.28) \quad \beta_s = \min_{i,j} \left(1 - \frac{d_{ii}}{a_{jj}}\right) .$$

Then the equation (2.21) has a unique solution τ^* in the ball $K = K(0, r)$ with

$$(2.29) \quad r = \frac{2\gamma \|c_0\|_E}{1 - \gamma(\|a_0^l\| + \alpha \|d_0^u\|)} ,$$

$$(2.30) \quad \gamma = \frac{1}{\beta_s - \|a_0^u\| - \alpha \|d_0^l\|} .$$

Moreover, for the sequence τ^j from (2.24) we have

$$\|\tau^* - \tau^j\|_E \leq \frac{k^j}{1-k} \|\tau^1 - \tau^0\|_E,$$

with

$$(2.31) \quad k = \gamma(\|a_0^l\| + \alpha\|d_0^u\| + 2\alpha r\|b_0\|_E) < l.$$

Again (2.27) and (2.26) ensure a sufficient distance of the spectra of a and d but now this distance is measured in the "relative sense".

PROOF. We can write the equation (2.21) as

$$(2.32) \quad \tau(E_a + a_0^u + a_0^l) - (E_d + d_0^u + d_0^l)D_d^2 \tau D_a^{-2} + c_0 - \tau b_0 D_d^2 \tau D_a^{-2} = 0.$$

Let \mathcal{L} be the linear operator defined by

$$(2.33) \quad \mathcal{L}\tau = \tau(E_a + a_0^u) - (E_d + d_0^l)D_d^2 \tau D_a^{-2}.$$

The equation (2.32) is equivalent to

$$\mathcal{L}\tau = -\tau a_0^l + d_0^u D_d^2 \tau D_a^{-2} - c_0 + \tau b_0 D_d^2 \tau D_a^{-2}.$$

If \mathcal{L} is invertible then this equation is written as

$$\mathcal{P}(\tau) = \tau$$

with

$$\mathcal{P}(\tau) = \mathcal{L}^{-1}(-\tau a_0^l + d_0^u D_d^2 \tau D_a^{-2} - c_0 + \tau b_0 D_d^2 \tau D_a^{-2}).$$

Set $\mathcal{L} = \mathcal{L}_0 + \mathcal{L}_1$, where

$$\mathcal{L}_0\tau = \tau E_a - E_d D_d^2 \tau D_a^{-2}, \quad \mathcal{L}_1\tau = \tau a_0^u - d_0^l D_d^2 \tau D_a^{-2}.$$

Then \mathcal{L}_0 is obviously invertible and

$$\|\mathcal{L}_0^{-1}\| = \frac{1}{\beta_s}.$$

Thus

$$\mathcal{L} = (\mathcal{I} + \mathcal{L}_1\mathcal{L}_0^{-1})\mathcal{L}_0.$$

By (2.27) and

$$\|\mathcal{L}_1\| \leq \|a_0^u\| + \alpha\|d_0^l\|$$

we obtain

$$\|\mathcal{L}_1\| \|\mathcal{L}_0^{-1}\| \leq \frac{1}{\beta_s}(\|a_0^u\| + \alpha\|d_0^l\|) < 1.$$

Thus, \mathcal{L} is invertible and

$$(2.34) \quad \|\mathcal{L}^{-1}\| \leq \frac{\|\mathcal{L}_0^{-1}\|}{1 - \|\mathcal{L}_1\| \|\mathcal{L}_0^{-1}\|} \leq \gamma,$$

$$\gamma = \frac{1}{\beta_s - \|a_0^u\| - \alpha\|d_0^l\|}.$$

We will find r such that $\mathcal{P} : K \rightarrow K$, $K = K(0, r)$. Let $\tau \in K$. Then

$$\|\mathcal{P}(\tau)\|_E \leq \gamma [(\|a_0^l\| + \alpha\|d_0^u\|)\|\tau\|_E + \alpha\|b_0\|_E \|\tau\|_E^2 + \|c_0\|_E] .$$

The inclusion $\mathcal{P}(K) \subseteq K$ will follow, if

$$(2.35) \quad \alpha\gamma \|b_0\|_E r^2 - [1 - \gamma(\|a_0^l\| + \alpha\|d_0^u\|)] r + \gamma \|c_0\|_E \leq 0.$$

The assumption (2.27) implies the positivity of the discriminant of the quadratic function above. Consequently, the solutions of (2.35) make the closed interval $[r_{min}, r_{max}]$ with

$$r_{min} = \frac{2\gamma\|c_0\|_E}{g_s + \sqrt{g_s^2 - 4\alpha\gamma^2\|b_0\|_E\|c_0\|_E}} ,$$

$$r_{max} = \frac{2\gamma\|c_0\|_E}{g_s - \sqrt{g_s^2 - 4\alpha\gamma^2\|b_0\|_E\|c_0\|_E}} ,$$

with

$$g_s = 1 - \gamma(\|a_0^l\| + \alpha\|d_0^u\|) .$$

For simplicity we take

$$(2.36) \quad r = \frac{2\gamma\|c_0\|_E}{g_s} \in (r_{min}, r_{max}) .$$

We prove the contractivity. Let $\bar{\tau}, \tau \in K$, then

$$\mathcal{P}(\bar{\tau}) - \mathcal{P}(\tau) =$$

$$\mathcal{L}^{-1}[-(\bar{\tau} - \tau)a_0^l + d_0^u D_d^2(\bar{\tau} - \tau)D_a^{-2} + \bar{\tau}b_0D_d^2(\bar{\tau} - \tau)D_a^{-2} + (\bar{\tau} - \tau)b_0D_d^2 \tau D_a^{-2}] ,$$

and

$$\|\mathcal{P}(\bar{\tau}) - \mathcal{P}(\tau)\|_E \leq k\|\bar{\tau} - \tau\|_E .$$

with k from (2.31). The condition (2.27) implies $k < 1$.

Let $\tau^1 = \mathcal{P}(\tau^0)$ for any $\tau^0 \in K$,

$$\tau^j = \mathcal{P}(\tau^{j-1}) = \mathcal{L}^{-1}(-\tau^{j-1}a_0^l + d_0^u D_d^2\tau^{j-1}D_a^{-2} - c_0 + \tau^{j-1}b_0D_d^2\tau^{j-1}D_a^{-2}) .$$

By Banach fixed-point theorem [8] there is a unique solution $\tau^* \in K$ of the equation (2.21) satisfying

$$\|\tau^* - \tau^j\|_E \leq \frac{k^j}{1 - k}\|\tau^1 - \tau^0\|_E .$$

Q.E.D.

The estimates for ν are again obtained by interchanging the roles of b_0 and c_0 .⁴

Analogous results are valid for all iterations listed in (2.16), more precisely, same assumptions lead to same results, although optimal constants, like respective r_{min} and k actually vary from case to case.

⁴In fact our assumption (2.27) is made a bit stronger to accommodate both iterations.

EXAMPLE 2.3. Let A in (2.1)

$$(2.37) \quad A = \begin{bmatrix} 1 & 2\varepsilon\sqrt{\varepsilon} & \varepsilon\sqrt{\varepsilon} \\ 2\varepsilon\sqrt{\varepsilon} & 4\varepsilon^2 & 2\varepsilon^2\sqrt{\varepsilon} \\ \varepsilon\sqrt{\varepsilon} & 2\varepsilon^2\sqrt{\varepsilon} & \varepsilon^2 \end{bmatrix}$$

and

$$a = \begin{bmatrix} 1 & 2\varepsilon\sqrt{\varepsilon} \\ 2\varepsilon\sqrt{\varepsilon} & 4\varepsilon^2 \end{bmatrix}$$

with $0 < \varepsilon < 1$.

Here the assumption (2.17) of Theorem 2.1 is not fulfilled.

The matrix (2.37) can be written in form (2.19) as

$$(2.38) \quad A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2\varepsilon & 0 \\ 0 & 0 & \varepsilon \end{bmatrix} \begin{bmatrix} 1 & \sqrt{\varepsilon} & \sqrt{\varepsilon} \\ \sqrt{\varepsilon} & 1 & \sqrt{\varepsilon} \\ \sqrt{\varepsilon} & \sqrt{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2\varepsilon & 0 \\ 0 & 0 & \varepsilon \end{bmatrix}$$

with

$$a_0 = \begin{bmatrix} 1 & \sqrt{\varepsilon} \\ \sqrt{\varepsilon} & 1 \end{bmatrix}.$$

Here are the assumptions of Theorem 2.2 satisfied for ε small enough. In fact, we obtain

$$\alpha = \frac{1}{4}, \quad \|b_0\|_E = \|c_0\|_E = \sqrt{2\varepsilon}, \quad \|a_0^u\| = \|a_0^l\| = \sqrt{\varepsilon}, \quad \beta_s = \frac{3}{4},$$

(d_0^u, d_0^l are not existent) and for (2.27)

$$(\sqrt{2} + 2)\sqrt{\varepsilon} < 3/4.$$

3. PRESERVING OF BLOCK DIAGONALITY

The rest of the paper gives estimates of the changes of the diagonal blocks of A after the similarity transformations (2.3) and (2.6). Due to the special type of these block transformations this change will be quadratically small thus generalizing similar properties of Jacobi - like diagonalization in the 2×2 case. Similar results hold for scaled matrices.

In case a) by (2.9) we obtain

$$a - a_1 = (I + ut)^{-1/2} [w + (uc + bt - udt)(I + ut)^{-1/2}]$$

with

$$w = (I + ut)^{1/2} a - a(I + ut)^{-1/2}.$$

For $\|ut\| < 1$ using binomial expansion we obtain

$$(3.1) \quad \|(I + ut)^{-1/2}\| \leq \frac{1}{\sqrt{1 - \|ut\|}}$$

and

$$\|w\|_E \leq \frac{\|ut\| \|a\|_E}{\sqrt{1 - \|ut\|}}.$$

Then

$$(3.2) \quad \|a - a_1\|_E \leq \frac{1}{1 - \|ut\|} (\|ut\| \|a\|_E + \|uc + bt - udt\|_E).$$

Similarly, for $\|tu\| < 1$

$$(3.3) \quad \|d - d_1\|_E \leq \frac{1}{1 - \|tu\|} (\|tu\| \|d\|_E + \|cu + tb + tua\|_E).$$

In case b) we need the following lemma.

LEMMA 3.1. *Let $\|X\|_E \leq 1/4$. Then $I + X$ can be decomposed as*

$$(3.4) \quad I + X = RL$$

R, L^T upper triangular with positive diagonal and

$$(3.5) \quad R = I + Y \quad , \quad L = I + Z$$

$$(3.6) \quad \|Y\|_E, \|Z\|_E \leq \rho_{min} = \rho(X) = \frac{2\|X\|_E}{1 + \sqrt{1 - 4\|X\|_E}}.$$

For the proof see [7], Th. 2.1.

By (2.10) we obtain

$$a - a_2 = R_1^{-1} [R_1 a - a L_1^{-1} + (uc + bt - udt) L_1^{-1}].$$

Since

$$R_1 L_1 = I + ut$$

by Lemma 3.1

$$R_1 = I + Y_1 \quad \text{and} \quad L_1^{-1} = I + Z_1$$

with

$$\|Y_1\|_E \leq \rho(ut) \leq 2\|ut\|_E \quad \text{and} \quad \|Z_1\|_E \leq \frac{\rho(ut)}{1 - \rho(ut)} \leq \frac{2\|ut\|_E}{1 - 2\|ut\|_E}$$

it follows

$$(3.7) \quad \|a - a_2\|_E \leq \|R_1^{-1}\| [(\|Y_1\| + \|Z_1\|) \|a\|_E + \|uc + bt - udt\|_E \|L_1^{-1}\|].$$

Similarly,

$$(3.8) \quad \|d - d_2\|_E \leq \|R_2^{-1}\| [(\|Y_2\| + \|Z_2\|) \|d\|_E + \|cu + tb + tua\|_E \|L_2^{-1}\|]$$

with

$$\|Y_2\|_E \leq 2\|tu\|_E \quad \text{and} \quad \|Z_2\|_E \leq \frac{2\|tu\|_E}{1 - 2\|tu\|_E}.$$

As we see in all these formulae the change of diagonal blocks is quadratical in r from Theorem 2.1.

Now we consider the case when u, t are the solutions of the Riccati equations in case a). From (2.12) we have

$$(3.9) \quad udt = aut + bt - uc t .$$

Inserting (3.9) in (2.9) we obtain

$$(3.10) \quad a_1 = (I + ut)^{-1/2}(a - uc)(I + ut)^{1/2} .$$

Similarly, from (2.11)

$$(3.11) \quad tau = dtu - cu + tbtu .$$

Inserting (3.11) in (2.9) we obtain

$$(3.12) \quad d_1 = (I + tu)^{-1/2}(d + tb)(I + tu)^{1/2} .$$

Now

$$a - a_1 = (I + ut)^{-1/2}[w_1 + uc(I + ut)^{1/2}]$$

with

$$w_1 = (I + ut)^{1/2}a - a(I + ut)^{1/2} .$$

For $\|ut\| < 1$ we obtain

$$(3.13) \quad \|(I + ut)^{1/2}\| \leq 1 + \frac{\|ut\|}{1 + \sqrt{1 - \|ut\|}} \leq 1 + \|ut\|$$

and

$$(3.14) \quad \|w_1\|_E \leq \frac{2\|ut\|\|a\|_E}{1 + \sqrt{1 + \|ut\|}} .$$

Using (3.1), (3.14) and (3.13) we obtain

$$(3.15) \quad \|a - a_1\|_E \leq \frac{1}{\sqrt{1 - \|ut\|}} \left[\frac{2\|ut\|\|a\|_E}{1 + \sqrt{1 + \|ut\|}} + (1 + \|ut\|)\|uc\|_E \right]$$

and for $\|tu\| < 1$

$$(3.16) \quad \|d - d_1\|_E \leq \frac{1}{\sqrt{1 - \|tu\|}} \left[\frac{2\|tu\|\|d\|_E}{1 + \sqrt{1 + \|tu\|}} + (1 + \|tu\|)\|tb\|_E \right] .$$

In case b) we have $b_2 = 0$ and $c_2 = 0$ with (2.8) we obtain

$$(3.17) \quad a_2 = R_1^{-1}(a - uc)R_1,$$

$$(3.18) \quad d_2 = R_2^{-1}(d + tb)R_2$$

By (2.8) and Lemma 3.1 we have

$$R_1 = I + Y_1 \quad , \quad R_2 = I + Y_2$$

with

$$\|Y_1\|_E \leq 2\|ut\|_E \quad , \quad \|Y_2\|_E \leq 2\|tu\|_E .$$

It follows

$$(3.19) \quad \|a - a_2\|_E \leq \|R_1^{-1}\|(2\|Y_1\|\|a\|_E + \|uc\|_E \|R_1\|) ,$$

$$(3.20) \quad \|d - d_2\|_E \leq \|R_2^{-1}\|(2\|Y_2\|\|d\|_E + \|tb\|_E \|R_2\|) .$$

Now as before we consider matrices with scaled diagonal dominance (2.19). By substituting (2.20) in (2.3) , (2.4) , (2.6) , (2.7) for

$$A_1 = V^{-1}DA_0DV$$

and

$$A_2 = U^{-1}DA_0DU$$

we obtain

$$(3.21) \quad \begin{aligned} a_1 &= D_a[(I + D_a^{-2}\nu D_d^2\tau)^{-1/2}\bar{a}_0(I + \nu D_d^2\tau D_a^{-2})^{-1/2}]D_a \\ b_1 &= D_a[(I + D_a^{-2}\nu D_d^2\tau)^{-1/2}\bar{b}_0(I + D_d^2\tau D_a^{-2}\nu)^{-1/2}]D_d \\ c_1 &= D_d[(I + \tau D_a^{-2}\nu D_d^2)^{-1/2}\bar{c}_0(I + \nu D_d^2\tau D_a^{-2})^{-1/2}]D_a \\ d_1 &= D_d[(I + \tau D_a^{-2}\nu D_d^2)^{-1/2}\bar{d}_0(I + D_d^2\tau D_a^{-2}\nu)^{-1/2}]D_d \end{aligned}$$

and

$$(3.22) \quad \begin{aligned} a_2 &= D_a(D_a^{-1}R_1^{-1}D_a\bar{a}_0D_aL_1^{-1}D_a^{-1})D_a \\ b_2 &= D_a(D_a^{-1}R_1^{-1}D_a\bar{b}_0D_dL_2^{-1}D_d^{-1})D_d \\ c_2 &= D_d(D_d^{-1}R_2^{-1}D_d\bar{c}_0D_aL_1^{-1}D_a^{-1})D_a \\ d_2 &= D_d(D_d^{-1}R_2^{-1}D_d\bar{d}_0D_dL_2^{-1}D_d^{-1})D_d \end{aligned}$$

with

$$(3.23) \quad \begin{aligned} \bar{a}_0 &= a_0 - D_a^{-2}\nu D_d^2c_0 - b_0D_d^2\tau D_a^{-2} + D_a^{-2}\nu D_d^2d_0D_d^2\tau D_a^{-2} \\ \bar{b}_0 &= a_0\nu - D_a^{-2}\nu D_d^2c_0\nu + b_0 - D_a^{-2}\nu D_d^2d_0 \\ \bar{c}_0 &= \tau a_0 + c_0 - \tau b_0D_d^2\tau D_a^{-2} - d_0D_d^2\tau D_a^{-2} \\ \bar{d}_0 &= \tau a_0\nu + c_0\nu + \tau b_0 + d_0 . \end{aligned}$$

Now by (2.19) and (3.21) after a lengthy calculation we obtain

$$(3.24) \quad \|D_a^{-1}(a - a_1)D_a^{-1}\|_E \leq \alpha s$$

and

$$(3.25) \quad \|D_d^{-1}(d - d_1)D_d^{-1}\|_E \leq s$$

where α is as in (2.26) and $\alpha\|\nu\|\|\tau\| < 1$ and

$$(3.26) \quad s = \frac{\|\nu\|\|\tau\|(\|a_0\|_E + \alpha\|d_0\|_E) + \|\tau\|\|b_0\|_E + \|\nu\|\|c_0\|_E}{1 - \alpha\|\nu\|\|\tau\|}.$$

Similarly, by (2.19) and (3.22) we obtain

$$(3.27) \quad \|D_a^{-1}(a - a_2)D_a^{-1}\|_E \leq \|R_1^{-1}\|[(\|Y_1\| + \|Z_1\|)\|a_0\|_E + \alpha(\alpha\|\nu\|\|\tau\|\|d_0\|_E + \|\tau\|\|b_0\|_E + \|\nu\|\|c_0\|_E)\|L_1^{-1}\|],$$

$$(3.28) \quad \|D_d^{-1}(d - d_2)D_d^{-1}\|_E \leq \|R_2^{-1}\|[(\|Y_2\| + \|Z_2\|)\|d_0\|_E + (\|\nu\|\|\tau\|\|a_0\|_E + \|\tau\|\|b_0\|_E + \|\nu\|\|c_0\|_E)\|L_2^{-1}\|]$$

where

$$\|Y_1\|_E, \|Y_2\|_E \leq 2\alpha\|\nu\|_E\|\tau\|_E \quad \text{and} \quad \|Z_1\|_E, \|Z_2\|_E \leq \frac{2\alpha\|\nu\|_E\|\tau\|_E}{1 - 2\alpha\|\nu\|_E\|\tau\|_E}.$$

The estimates (3.24), (3.25) are "scaled" analogs of (3.15), (3.16) and similarly (3.27), (3.28) are analogs of (3.7), (3.8), respectively. If $b_1 = 0$ and $c_1 = 0$ then ν, τ is the solution of (2.22), (2.21) respectively. Inserting (2.19) and (2.20) in (3.10) and (3.12) we obtain

$$a_1 = D_a(I + D_a^{-2}\nu D_d^2\tau)^{-1/2}(a_0 - D_a^{-2}\nu D_d^2c_0)(I + \nu D_d^2\tau D_a^{-2})^{1/2}D_a,$$

and

$$d_1 = D_d(I + \tau D_a^{-2}\nu D_d^2)^{-1/2}(a_0 + \tau b_0)(I + D_d^2\tau D_a^{-2}\nu)^{1/2}D_d.$$

For $\alpha\|\nu\|\|\tau\| < 1$,

$$\begin{aligned} & \|D_a^{-1}(a - a_1)D_a^{-1}\|_E \leq \\ & \leq \frac{\alpha\|\nu\|}{\sqrt{1 - \alpha\|\nu\|\|\tau\|}} \left[\frac{2\|\tau\|}{1 + \sqrt{1 + \alpha\|\nu\|\|\tau\|}} \|a_0\|_E + (1 + \alpha\|\nu\|\|\tau\|)\|c_0\|_E \right], \\ & \|D_d^{-1}(d - d_1)D_d^{-1}\|_E \leq \\ & \leq \frac{\|\tau\|}{\sqrt{1 - \alpha\|\nu\|\|\tau\|}} \left[\frac{2\alpha\|\nu\|}{1 + \sqrt{1 + \alpha\|\nu\|\|\tau\|}} \|d_0\|_E + (1 + \alpha\|\nu\|\|\tau\|)\|b_0\|_E \right]. \end{aligned}$$

Similarly, for $b_2 = 0$ and $c_2 = 0$ we obtain

$$a_2 = D_a[D_a^{-1}R_1^{-1}D_a(a_0 - D_a^{-2}\nu D_d^2c_0)(I + \nu D_d^2\tau D_a^{-2})D_aL_1^{-1}D_a^{-1}]D_a,$$

$$d_2 = D_d[D_d^{-1}R_2^{-1}D_d(d_0 + \tau b_0)(I + D_d^2\tau D_d^{-2}\nu)D_dL_2^{-1}D_d^{-1}]D_d.$$

Then

$$\|D_a^{-1}(a - a_1)D_a^{-1}\|_E \leq$$

$$\|R_1^{-1}\|[(\|Y_1\| + \|S_1\| + \|Y_1\|\|S_1\| + \alpha\|\nu\|\|\tau\|)\|a_0\|_E + \alpha\|\nu\|(1 + \alpha\|\tau\|)\|c_0\|_E]\|L_1^{-1}\|,$$

$$\|D_d^{-1}(d - d_1)D_d^{-1}\|_E \leq$$

$$\leq \|R_2^{-1}\|[(\|Y_2\| + \|S_2\| + \|Y_2\|\|S_2\| + \alpha\|\nu\|\|\tau\|)\|d_0\|_E + \|\tau\|(1 + \alpha\|\nu\|)\|b_0\|_E]\|L_2^{-1}\|$$

$$R_1 = I + Y_1 \quad , \quad R_2 = I + Y_2 \quad , \quad L_1 = I + S_1 \quad , \quad L_2 = I + S_2$$

and

$$\|Y_1\|_E, \|Y_2\|_E, \|S_1\|_E, \|S_2\|_E \leq 2\alpha\|\nu\|_E \|\tau\|_E.$$

In concluding this section we can say that in all cases the change of the diagonal elements is quadratically small in the measure of the (scaled) diagonal dominance and the same kind of diagonal dominance is shared by the transformed matrix as well. In view of the results in [1] and [2] this strongly suggests that all these transformations, at least in the Hermitian case, are candidates for an eigenvalue algorithm with high relative accuracy. The basis of our last estimates are the formulae (3.10), (3.12), (3.17) and (3.18) for the transformed diagonal blocks. These formulae are of independent interest. To illustrate this, assume that A is Hermitian and V unitary and that block diagonalization took place, i.e. $c_1 = b_1^H = 0$. Then using binomial expansion and $\|t\| = \mathcal{O}(\|c\|)$ the formula (3.10) yields

$$a_1 = a - t^H c + \frac{at^H t - t^H t a}{2} + \mathcal{O}(\|c\|^4).$$

Since a_1 has to be Hermitian we obtain

$$a_1 = a - \frac{t^H c + (t^H c)^H}{2} + \mathcal{O}(\|c\|^4),$$

i.e.

$$(3.29) \quad a_1 = \text{Hermitian part of } (a - t^H c)$$

up to a quartic error. The same is true of d_1 . A scaled variant of this approximation could be derived as well.

4. SOME NUMERICAL EXAMPLES

We have tried some typical iterations with a MATLAB - based code which just yields the number of iterations which, in turn allows to estimate the real efficiency. As predicted by our theory the number of iterations is independent of the dimensions n, m , which allows to compute few eigenpairs with the cost of some n^2 operations. Our first example were (in the MATLAB notation)

$$A = \text{diag}(1:300) + \text{rand}(300)/80$$

which looks "diagonally dominant" but is well-beyond the Gershgorin separation. For $m = 3, 5, 20, 150$ we obtained block diagonality to working accuracy after 10 Jacobi iterations or 6 – 8 Gauss-Seidel iterations (2.14), (2.15). The second example

$$A = \text{diag}(200:-1:1)(\text{eye}(200) + \text{rand}(200)/10000)\text{diag}(200:-1:1)$$

is only scaled diagonally dominant according to [1]. For $m = 2, 5, 20, 100$ we needed 10 – 13 Jacobi iterations and 9 – 12 Gauss-Seidel iterations.

To illustrate high relative accuracy of the algorithm consider the matrix

$$A = \begin{bmatrix} 1 \cdot 10^{20} & 2 & 3 & 4 \\ 2 & 4 \cdot 10^{20} & 5 & 6 \\ 3 & 5 & 7 & 8 \\ 4 & 6 & 8 & 9 \end{bmatrix}$$

whose small eigenvalues are up to a 10^{-20} error equal to those of

$$\begin{bmatrix} 7 & 8 \\ 8 & 9 \end{bmatrix}.$$

Both eigenvalues are correctly reproduced by our method (standard eigensolver from MATLAB yielded completely wrong values 0 and 10^4).

A rigorous justification of the above illustrated accuracy as well as extensive numerical experiments shall be presented elsewhere.

REFERENCES

- [1] J. Barlow, J. Demmel, Computing accurate eigensystems of scaled diagonally dominant matrices, *SIAM J. Numer. Anal.* **27** (1988) 762-791.
- [2] Demmel, J., Veselić, K., Jacobi's method is more accurate than QR, *SIAM J. Matrix Anal. Appl.*, **13** (1992) 1204-1245.
- [3] Z. Drmač, K. Veselić, On the spectra of nearly unitarily similar Hermitian matrices, *LAA* **309** (2000) 191 -215.
- [4] C.G.J. Jacobi, Über ein leichtes Verfahren die in der Theorie der Säculärstörungen vorkommenden Gleichungen numerisch aufzulösen, *Crelle's J.* **330** (1846) 51-94.
- [5] A.N. Malyshev, On iterative refinement for the spectral decomposition of Symmetric Matrices, IRISA, preprint No 628, (1992) 1-25.
- [6] G.L.G. Slepen and H.A. van der Vorst, A Jacobi-Davison iteration method for linear eigenvalue problems, *SIAM J. Matrix Anal. Appl.* **17** (1996) 401-425.

- [7] G. W. Stewart, On the perturbation of LU, Cholesky, and QR factorizations, SIAM J. Matrix Anal. Appl. **14** (1993) 551-566.
- [8] J. Stoer, R. Bulirsch, Introduction to numerical analysis, Spinger, Berlin 1980
- [9] K. Veselić, On a new class of elementary matrices, Numer. Math. **33** (1979) 173-180.

Faculty of Transport and Traffic Engineering
Vukelićeva 4
HR-10000 Zagreb
Croatia
E-mail: kovacm@fpz.hr

Lehrgebiet Mathematische Physik
Fernuniversität Hagen
Postf. 940
58084 Hagen
Germany
E-mail: kresimir.veselic@fernuni-hagen.de
Received: 12.01.1999.
Revised: 11.01.2001.