

# ON NONSYMMETRIC SADDLE POINT MATRICES THAT ALLOW CONJUGATE GRADIENT ITERATIONS

JÖRG LIESEN<sup>†</sup> AND BERESFORD N. PARLETT<sup>‡</sup>

**Abstract.** Linear systems in saddle point form are often symmetric and highly indefinite. Indefiniteness, however, is a major challenge for iterative solvers such as Krylov subspace methods. It has been noted by several authors that a simple trick, namely negating the second block row of the saddle point system, leads to an equivalent linear system with a nonsymmetric coefficient matrix  $\mathcal{A}$  whose spectrum is entirely contained in the right half plane. In this paper we study conditions so that  $\mathcal{A}$  is diagonalizable with a real and positive spectrum. These conditions are based on necessary and sufficient conditions for positive definiteness of a certain bilinear form, with respect to which  $\mathcal{A}$  is symmetric. In case the latter conditions are satisfied, there exists a well defined conjugate gradient (CG) method for solving linear systems with  $\mathcal{A}$ . We give an efficient implementation of this method, discuss practical issues such as error bounds and preconditioning, and present numerical experiments showing the effectiveness of the method.

**Key words.** saddle point problem, eigenvalues, conjugate gradient method, Stokes problem

**AMS subject classifications.** 65F15, 65N22, 65F50

**1. Introduction.** Many applications in science and engineering require solving large linear algebraic systems in saddle point form; see [4] for an extensive survey. A typical system matrix is of the form

$$(1.1) \quad \begin{bmatrix} A & B^T \\ B & -C \end{bmatrix},$$

where  $A = A^T \in \mathbb{R}^{n \times n}$  is positive definite ( $A > 0$ ),  $B \in \mathbb{R}^{m \times n}$  has rank  $r \leq m \leq n$ , and  $C = C^T \in \mathbb{R}^{m \times m}$  is positive semidefinite ( $C \geq 0$ ). The matrix in (1.1) is congruent to the block diagonal matrix  $[A \ 0; 0 \ S]$ , where  $S = -(C + BA^{-1}B^T)$ , so that  $-S = -S^T \geq 0$ . Hence this matrix is indefinite, with  $n$  positive and  $\text{rank}(S)$  negative eigenvalues. Typically,  $\text{rank}(S)$  is close to  $m$  and therefore, unless  $m$  is very small, the matrix in (1.1) is highly indefinite. This feature is a major challenge for iterative solvers such as Krylov subspace methods.

It has been noted by several authors (see [4, p. 23] for references) that, for solving linear algebraic systems, there is no harm in negating the second block row in (1.1), and using as system matrix

$$(1.2) \quad \mathcal{A} \equiv \begin{bmatrix} A & B^T \\ -B & C \end{bmatrix}.$$

Symmetry has been abandoned but what, if anything, has been gained? First, the matrix  $\mathcal{A}$  is *positive semistable*, i.e. all its eigenvalues have nonnegative real parts (even positive real parts in case  $B$  has full rank); see, e.g., [4, Theorem 3.6] for a proof. Moreover, it can be shown that there exists a “magic” hidden bilinear form with respect to which  $\mathcal{A}$  is symmetric. As shown in this paper, this bilinear form is a proper inner product, if and only if the spectra of  $A$  and  $C$  are separated, and the norm

---

<sup>†</sup>Institute of Mathematics, Technical University of Berlin, Straße des 17. Juni 136, 10623 Berlin, Germany ([liesen@math.tu-berlin.de](mailto:liesen@math.tu-berlin.de)). The work of this author was supported by the Emmy Noether-Programm of the Deutsche Forschungsgemeinschaft.

<sup>‡</sup>Department of Mathematics, University of California, Berkeley, CA 94720-3840, USA ([parlett@math.berkeley.edu](mailto:parlett@math.berkeley.edu)).

of  $B$  is small enough. When this is satisfied,  $\mathcal{A}$  is diagonalizable with nonnegative real eigenvalues. Thus, there exists a conjugate gradient (CG) method for solving linear systems with  $\mathcal{A}$ . In retrospect, the idea of negating the second block row appears to be natural. If we had  $B = 0$ , then  $\mathcal{A}$  would be positive semidefinite, while the matrix in (1.1) is indefinite. As  $\|B\|$  rises from zero, the moment comes when the eigenvalues of  $\mathcal{A}$  are no longer real and/or  $\mathcal{A}$  is no longer diagonalizable, but, until then,  $\mathcal{A}$  is a nice matrix.

Some of the theoretical results in this paper, particularly those concerning the definiteness of the bilinear form, generalize previous work in [7] and [5]. In these papers the focus is on cases with  $C = 0$ . Moreover, unlike in [7, 5], we provide an implementation of a CG method for solving linear systems with  $\mathcal{A}$ . We thereby clarify some subtleties concerning the choice of the inner product for constructing a well defined CG method for  $\mathcal{A}$ .

The paper is organized as follows. In Section 2 we analyze properties of the matrix  $\mathcal{A}$  and the bilinear form with respect to which  $\mathcal{A}$  is symmetric. In Section 3 we construct a CG method for solving linear systems with  $\mathcal{A}$ , discuss when this method is well defined, and provide an efficient implementation. In Section 4 we discuss some practical issues in the context of our CG method, including the conditioning of the CG inner product matrix, error bounds, and preconditioning. In Section 5 we present some numerical experiments that show the effectiveness of the method. Concluding remarks close the paper.

In this paper  $I$  denotes the identity matrix (of appropriate size). The bilinear form defined by a (real) symmetric matrix  $G$  is denoted by  $(u, v)_G \equiv v^T G u$ . In case  $G$  is positive definite, this bilinear form is an inner product, and the associated norm is given by  $\|u\|_G \equiv (u, u)_G^{1/2}$ . In case  $G = I$ , i.e. the Euclidean inner product and norm, we skip the index and simply write  $(u, v)$  and  $\|u\|$ . A matrix  $M$  is called symmetric with respect to a (real) symmetric matrix  $G$ , or shortly  $G$ -symmetric, if  $GM$  is symmetric, or, equivalently,  $(Mu, v)_G = (u, Mv)_G$  for all  $u, v$ . The matrix  $M$  is called  $G$ -definite, if  $GM$  is definite, i.e.  $(GMu, u) = (Mu, u)_G \neq 0$  for all  $u \neq 0$ . Of course, if  $G = I$ , we simply write symmetric and definite as usual.

**2. Analysis of the matrix  $\mathcal{A}$ .** We consider a matrix  $\mathcal{A} \in \mathbb{R}^{(n+m) \times (n+m)}$  as in (1.2), with symmetric positive definite  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{m \times n}$  of rank  $r \leq m \leq n$ , and symmetric positive semidefinite  $C \in \mathbb{R}^{m \times m}$ . As shown in [4, Theorem 3.6], the matrix  $\mathcal{A}$  is *positive semistable*, meaning that all eigenvalues of  $\mathcal{A}$  have nonnegative real parts. In case  $B$  has full rank  $m$ , the matrix  $\mathcal{A}$  is *positive stable*, i.e. all its eigenvalues have positive real parts.

We start with a useful lemma concerning this matrix.

LEMMA 2.1. *Let the matrix*

$$(2.1) \quad \mathcal{J} \equiv \begin{bmatrix} I & 0 \\ 0 & -I \end{bmatrix},$$

be conformally partitioned with  $\mathcal{A}$ . Then

$$(1) \quad \mathcal{A} \text{ is } \mathcal{J}\text{-symmetric, i.e. } \mathcal{J}\mathcal{A} = \mathcal{A}^T \mathcal{J} = (\mathcal{J}\mathcal{A})^T,$$

and, for any polynomial  $p$ ,

$$(2a) \quad p(\mathcal{A}) \text{ is } \mathcal{J}\text{-symmetric, i.e. } \mathcal{J}p(\mathcal{A}) = p(\mathcal{A}^T)\mathcal{J} = (\mathcal{J}p(\mathcal{A}))^T, \text{ and}$$

$$(2b) \quad \mathcal{A} \text{ is } \mathcal{J}p(\mathcal{A})\text{-symmetric, i.e. } (\mathcal{J}p(\mathcal{A}))\mathcal{A} = \mathcal{A}^T(p(\mathcal{A}^T)\mathcal{J}) = (\mathcal{J}p(\mathcal{A})\mathcal{A})^T.$$

*Proof.* Item (1) follows by straightforward computation. Using (1), we see that

$$\mathcal{J}\mathcal{A}^2 = (\mathcal{J}\mathcal{A})\mathcal{A} = (\mathcal{A}^T \mathcal{J})\mathcal{A} = \mathcal{A}^T(\mathcal{J}\mathcal{A}) = (\mathcal{A}^T)^2 \mathcal{J},$$

which implies (2a) by induction. Using (1) and (2a),

$$(\mathcal{J}p(\mathcal{A}))\mathcal{A} = (\mathcal{J}\mathcal{A})p(\mathcal{A}) = \mathcal{A}^T(\mathcal{J}p(\mathcal{A})) = \mathcal{A}^T(p(\mathcal{A}^T)\mathcal{J}),$$

which proves (2b).  $\square$

Item (1) in Lemma 2.1 shows that  $\mathcal{A}$  is symmetric with respect to the symmetric *indefinite* matrix  $\mathcal{J}$ . Our goal here is to determine whether there exists a symmetric *positive definite* matrix with respect to which  $\mathcal{A}$  is symmetric. Our starting point is the relation (2b) in Lemma 2.1. It shows that  $\mathcal{A}$  is symmetric with respect to any matrix of the form  $\mathcal{J}p(\mathcal{A})$ , where  $p$  is any polynomial. As shown by item (2a), any matrix of the form  $\mathcal{J}p(\mathcal{A})$  is itself symmetric. Therefore it suffices to show conditions under which  $\mathcal{J}p(\mathcal{A})$  is positive definite. Obviously, when  $p$  is of degree zero, this cannot be satisfied. Our ansatz will therefore be a polynomial  $p$  of degree one. Without loss of generality we may take the leading coefficient of  $p$  to be equal to one, and write our polynomial in the form  $p(\zeta) = \zeta - \gamma$ , for some yet to be determined parameter  $\gamma \in \mathbb{R}$ . Hence we ask: When is

$$(2.2) \quad \mathcal{M}(\gamma) \equiv \mathcal{J}p(\mathcal{A}) = \mathcal{J}(\mathcal{A} - \gamma I) = \begin{bmatrix} A - \gamma I & B^T \\ B & \gamma I - C \end{bmatrix}$$

a positive definite matrix? The complete answer to this question is given in the following theorem.

**THEOREM 2.2.** *The symmetric matrix  $\mathcal{M}(\gamma)$  is positive definite if and only if*

$$(2.3) \quad \lambda_{\min}(A) > \gamma > \lambda_{\max}(C),$$

where  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the smallest and largest eigenvalue, respectively, and

$$(2.4) \quad \left\| (\gamma I - C)^{-1/2} B (A - \gamma I)^{-1/2} \right\| < 1.$$

*Proof.* It is easy to see that  $\mathcal{M}(\gamma) > 0$  holds only if  $A - \gamma I > 0$ , or, equivalently,  $\lambda_{\min}(A) > \gamma$ . If this holds,  $\mathcal{M}(\gamma)$  is congruent to  $(A - \gamma I) \oplus S$ , where

$$S = (\gamma I - C) - B(A - \gamma I)^{-1}B^T.$$

Therefore,  $\mathcal{M}(\gamma)$  is positive definite, if, and only if,  $\lambda_{\min}(A) > \gamma$  and  $S > 0$ . The second inequality is equivalent to

$$\gamma I - C > B(A - \gamma I)^{-1}B^T.$$

The matrix on the right hand side is positive semidefinite, which implies  $\gamma I - C > 0$ , or, equivalently,  $\gamma > \lambda_{\max}(C)$ . Finally,  $\gamma I - C > B(A - \gamma I)^{-1}B^T$  is equivalent to

$$I > \left( (\gamma I - C)^{-1/2} B (A - \gamma I)^{-1/2} \right) \left( (\gamma I - C)^{-1/2} B (A - \gamma I)^{-1/2} \right)^T,$$

which in turn is equivalent to (2.4).  $\square$

From Theorem 2.2 we can derive some useful *sufficient* conditions that make  $\mathcal{M}(\gamma)$  positive definite.

COROLLARY 2.3. *The matrix  $\mathcal{M}(\gamma)$  is symmetric positive definite when (2.3) holds, and, in addition,*

$$(2.5) \quad \|B\|^2 < (\lambda_{\min}(A) - \gamma)(\gamma - \lambda_{\max}(C)).$$

For  $\gamma = \hat{\gamma} \equiv \frac{1}{2}(\lambda_{\min}(A) + \lambda_{\max}(C))$ , the right hand side of (2.5) is maximal, and (2.5) reduces to

$$(2.6) \quad 2\|B\| < \lambda_{\min}(A) - \lambda_{\max}(C).$$

*Proof.* A simple computation shows that

$$\begin{aligned} \left\| (\gamma I - C)^{-1/2} B (A - \gamma I)^{-1/2} \right\| &\leq \|(\gamma I - C)^{-1/2}\| \|B\| \|(A - \gamma I)^{-1/2}\| \\ &= (\gamma - \lambda_{\max}(C))^{-1/2} \|B\| (\lambda_{\min}(A) - \gamma)^{-1/2}. \end{aligned}$$

Hence  $\mathcal{M}(\gamma) > 0$ , if (2.3) holds and the right hand side is less than one, which is equivalent to (2.5). The second part follows from another simple computation.  $\square$

The conditions for positive definiteness of  $\mathcal{M}(\gamma)$  yield sufficient conditions so that  $\mathcal{A}$  is diagonalizable with a positive real spectrum.

COROLLARY 2.4. *If there exists a  $\gamma \in \mathbb{R}$  so that  $\mathcal{M}(\gamma)$  is positive definite, then  $\mathcal{A}$  has a nonnegative real spectrum and a complete set of eigenvectors that are orthonormal with respect to the inner product defined by  $\mathcal{M}(\gamma)$ . In case  $B$  has full rank, the spectrum of  $\mathcal{A}$  is real and positive.*

*Proof.* We know that the matrix  $\mathcal{A}$  is  $\mathcal{M}(\gamma)$ -symmetric. If  $\mathcal{M}(\gamma)$  is positive definite, then  $\mathcal{M}(\gamma)$  defines an inner product, and hence  $\mathcal{A}$  has real eigenvalues and a complete set of eigenvectors that are orthonormal with respect to this inner product (see, e.g., [8, Chapter IX]). Moreover,  $\mathcal{A}$  is known to be positive semistable (positive stable if  $B$  has full rank), so that its eigenvalues indeed must be real and nonnegative (real and positive).  $\square$

Our results above generalize several results that appeared in the literature: Fischer et al. [7] consider  $\mathcal{A}$  with  $A = \alpha I > 0$  and  $C = 0$ . They show that  $\mathcal{A}$  has a nonnegative real spectrum when  $2\|B\| < \alpha$ , which in this case is equivalent to (2.6). Benzi and Simoncini [5, Section 3] consider (in our notation) a matrix  $\mathcal{M}(\gamma)$  with  $A = A^T > 0$  and  $C = 0$ . In [5, Proposition 3.1] they show that  $\mathcal{M}(\hat{\gamma})$  is positive definite when  $4\lambda_{\max}(BA^{-1}B^T) < \lambda_{\min}(A)$ . Note that

$$(2.7) \quad 4\lambda_{\max}(BA^{-1}B^T) = 4\|BA^{-1/2}\|^2 \leq 4\|B\|^2 \|A^{-1/2}\|^2 = \frac{4\|B\|^2}{\lambda_{\min}(A)}.$$

Condition (2.6) with  $C = 0$  is equivalent to  $4\|B\|^2/\lambda_{\min}(A) < \lambda_{\min}(A)$ . Therefore, (2.6) implies the sufficient condition given in [5, Proposition 3.1].

One of the conditions that make  $\mathcal{M}(\gamma)$  symmetric positive definite is  $\lambda_{\min}(A) > \gamma > \lambda_{\min}(C)$ , cf. (2.3). Hence, Corollary 2.4 only yields a condition for a (nonnegative) real spectrum of  $\mathcal{A}$  when the spectra of  $A$  and  $C$  are separated. This separation appears to be essential, as seen by considering a two by two matrix of the form

$$\begin{bmatrix} \alpha & \beta \\ -\beta & \eta \end{bmatrix}, \quad \alpha > 0, \eta > 0.$$

This matrix has real eigenvalues (necessarily nonnegative) only when its discriminant  $(\alpha - \eta)^2 - 4\beta^2$  is nonnegative. Thus  $\beta$  must vanish if  $\alpha = \eta$ , and, by extension to  $\mathcal{A}$ , if the spectra of  $A$  and  $C$  overlap, i.e.  $\lambda_{\max}(C) \geq \lambda_{\min}(A)$ , and  $B \neq 0$ , it is most unlikely that the spectrum of  $\mathcal{A}$  is real.

The discriminant condition we just have derived turns out to be a special case of [5, Proposition 2.5]. In that result the matrix  $\mathcal{A}$  is assumed to be of the form (1.2) with  $C = \eta I > 0$ . If  $\mathcal{A} \in \mathbb{R}^{2 \times 2}$ , then the condition that both eigenvalues of  $\mathcal{A}$  are real is  $(\alpha + \eta)^2 \geq 4(\beta^2 + \eta)$ , or, equivalently,  $(\alpha - \eta)^2 - 4\beta^2 \geq 0$ .

Next let us consider a more representative matrix

$$\mathcal{A} = \left[ \begin{array}{ccc|cc} 1 & 0 & 0 & \beta & 0 \\ 0 & 2 & 0 & 0 & \beta \\ 0 & 0 & 3 & 0 & 0 \\ \hline -\beta & 0 & 0 & 2\eta & -\eta \\ 0 & -\beta & 0 & -\eta & 2\eta \end{array} \right], \quad \beta \neq 0, \quad \eta \geq 0.$$

Here  $\lambda_{\min}(A) = 1$ ,  $B$  has full rank,  $\|B\| = |\beta|$ , and  $\lambda_{\max}(C) = 3\eta$ . By (2.6),  $\mathcal{A}$  is diagonalizable with a positive and real spectrum when

$$2|\beta| < 1 - 3\eta.$$

For  $\eta = 1/12$  the above condition is  $|\beta| < 3/8 = 0.375$ , while MATLAB [12] shows that  $\mathcal{A}$  has five distinct real and positive eigenvalues for  $|\beta|$  as large as 0.405. For larger  $|\beta|$ ,  $\mathcal{A}$  has nonreal eigenvalues. On the other hand, if we choose  $\beta = 1/2$ , then the above condition is not satisfied for any  $\eta \geq 0$ , and indeed a MATLAB computation reveals that the matrix  $\mathcal{A}$  is not diagonalizable for  $\eta = 0$ , and has nonreal eigenvalues for  $\eta > 0$ . Therefore, the sufficient condition (2.6) cannot be relaxed, in general.

**3. Construction of a conjugate gradient (CG) method.** In this section we construct a CG method for solving linear systems with the matrix  $\mathcal{A}$ .

**3.1. The (generic) CG method.** As a starting point we consider the following (generic) statement of the CG method based on a given inner product  $(u, v)_G = v^T G u$  for solving a linear system of the form  $Mx = b$ .

ALGORITHM 3.1. (*The Conjugate Gradient (CG) method*)

*Input:* System matrix  $M$ , right hand side  $b$ , inner product matrix  $G$ , initial guess  $x_0$ .

*Initialize:*  $r_0 = b - Mx_0$ ,  $p_0 = r_0$ .

*For*  $i = 0, 1, \dots$  *until convergence:*

$$(3.1) \quad \alpha_i = \frac{(x - x_i, p_i)_G}{(p_i, p_i)_G}$$

$$(3.2) \quad x_{i+1} = x_i + \alpha_i p_i$$

$$(3.3) \quad r_{i+1} = r_i - \alpha_i M p_i$$

$$(3.4) \quad \beta_{i+1} = -\frac{(r_{i+1}, p_i)_G}{(p_i, p_i)_G}$$

$$(3.5) \quad p_{i+1} = r_{i+1} + \beta_{i+1} p_i$$

This algorithm is the unpreconditioned version of the algorithm Omin stated in [2, (4.1a)–(4.1g)]. In case  $M = M^T > 0$  and  $G = M$ , it corresponds to the classical Hestenes and Stiefel implementation of CG [10]. Its requirements and most

important properties are summarized in the following result (see [2] for proofs and further details).

**THEOREM 3.2.** *Suppose that the matrix  $G$  is symmetric positive definite and hence defines an inner product. If the matrix  $M$  is  $G$ -symmetric and  $G$ -definite, then (until convergence) the CG method stated in Algorithm 3.1 has the following properties:*

(1) *The iterates and residuals satisfy*

$$x_i \in x_0 + K_i(M, r_0) \quad \text{and} \quad r_i = b - Mx_i \in r_0 + MK_i(M, r_0),$$

where  $K_i(M, r_0) \equiv \text{span}\{r_0, Mr_0, \dots, M^{i-1}r_0\}$  is the  $i$ th Krylov subspace generated by  $M$  and  $r_0$ .

(2) *The direction vectors  $p_i \in r_0 + MK_i(M, r_0)$  satisfy*

$$(p_i, p_j)_G = 0 \quad \text{for} \quad i \neq j.$$

(3) *The error vector  $x - x_{i+1} \in (x - x_0) + K_{i+1}(M, r_0)$  satisfies*

$$(x - x_{i+1}, u)_G = 0 \quad \text{for all} \quad u \in \text{span}\{p_0, \dots, p_i\}.$$

(4) *The method is optimal in the  $G$ -norm, i.e.*

$$\begin{aligned} \|x - x_{i+1}\|_G &= \min_{u \in x_0 + K_{i+1}(M, r_0)} \|x - u\|_G \\ &= \min_{p \in \pi_{i+1}} \|p(M)(x - x_0)\|_G, \end{aligned}$$

where  $\pi_{i+1}$  denotes the set of polynomials of degree at most  $i + 1$  and with value one at the origin.

When the assumptions of Theorem 3.2 on  $M$  and  $G$  are satisfied, and hence the assertions (1)–(4) hold, we say that Algorithm 3.1 is *well defined* for  $M$  and  $G$ .

**REMARK 3.3.** In case  $M$  is  $G$ -symmetric and  $G$ -semi-definite (i.e.  $M$  is singular), Algorithm 3.1 may still be well defined. However, this will not hold globally for any right hand side  $b$ , but just for  $b \in \text{Range}(M)$ , so that a solution of the linear system exists. Details of the CG method for singular matrices are well discussed in [3, Section 11.2.8].

Note that the numerator of  $\alpha_i$  in (3.1) contains the unknown solution vector  $x$ . Therefore an essential practical requirement for the CG method, that is not mentioned in Theorem 3.2, is that the inner product matrix  $G$  must be chosen so that the scalar  $\alpha_i$  is *computable*. In the classical CG method of Hestenes and Stiefel, the system matrix  $M$  is assumed symmetric positive definite,  $G = M$ , and then

$$(x - x_i, p_i)_G = (x - x_i, p_i)_M = (M(x - x_i), p_i) = (r_i, p_i),$$

which is a computable quantity.

**3.2. A CG method for  $\mathcal{A}$ .** From now on we assume that we are given a linear system of the form  $\mathcal{A}x = b$  where  $\mathcal{A}$  is as in (1.2) with  $A = A^T > 0$ , full rank  $B$ , and  $C = C^T \geq 0$ . Generalizations of the following algorithms and results may be easily obtained for rank deficient  $B$ , cf. Remark 3.3, but here we have opted for a clean rather than the most general presentation.

A CG method for solving this system requires a symmetric positive definite matrix defining an inner product. To make this CG method well defined, the matrix  $\mathcal{A}$  and the inner product matrix should satisfy the assumptions of Theorem 3.2. Moreover, the inner product matrix must be chosen so that  $\alpha_i$  is computable. The latter requirement is *not* satisfied by the matrix  $\mathcal{M}(\gamma)$  in (2.2), since for this matrix we cannot evaluate the quantity  $(x - x_i, p_i)_{\mathcal{M}(\gamma)}$  unless we know the solution vector  $x$ . Therefore, even if  $\mathcal{M}(\gamma)$  is positive definite, this matrix cannot be used in Algorithm 3.1 to give a computable CG method for solving  $\mathcal{A}x = b$ .

REMARK 3.4. As mentioned above, our problem has been studied previously in [7] ( $\mathcal{A}$  with  $A = \alpha I > 0$ ,  $C = 0$ ) and [5] ( $\mathcal{A}$  with  $A = A^T > 0$ ,  $C = 0$ ). In our notation, it is suggested in [7, p. 532] that “the standard conjugate gradient method with the inner product defined by  $[\mathcal{M}(\gamma)]$  could be applied directly to  $[\mathcal{A}]$ ”, and in [5, p. 185], it is proposed that “the Conjugate Gradient method may be employed for solving linear systems with  $[\mathcal{A}]$  by using the inner product defined by  $[\mathcal{M}(\gamma)]$ ”. These authors, however, have not provided implementations, so it is unclear what specific methods they had in mind.

To make  $\alpha_i$  computable, our choice for the inner product will be the matrix  $\mathcal{M}(\gamma)\mathcal{A}$ . From item (2b) in Lemma 2.1 we know that this matrix is symmetric. The following result gives a sufficient condition when this matrix is positive definite.

THEOREM 3.5. *If the (symmetric) matrix  $\mathcal{M}(\gamma)$  is positive definite, then for all  $k \geq 0$  the matrix  $\mathcal{M}(\gamma)\mathcal{A}^k$  is symmetric positive definite.*

*Proof.* Item (2b) in Lemma 2.1 shows that  $\mathcal{M}(\gamma)\mathcal{A} = \mathcal{A}^T\mathcal{M}(\gamma)$  and therefore, by induction,

$$\mathcal{M}(\gamma)\mathcal{A}^k = (\mathcal{A}^T)^k\mathcal{M}(\gamma) = (\mathcal{M}(\gamma)\mathcal{A}^k)^T$$

for all  $k \geq 0$ , i.e.  $\mathcal{M}(\gamma)\mathcal{A}^k$  is symmetric. Moreover,

$$\mathcal{A}^k = \mathcal{M}(\gamma)^{-1}(\mathcal{A}^T)^k\mathcal{M}(\gamma),$$

which means that  $\mathcal{A}^k$  is normal with respect to the symmetric positive definite matrix  $\mathcal{M}(\gamma)$ . A result of Givens [9, Theorem 2] implies that the  $\mathcal{M}(\gamma)$ -field of values of  $\mathcal{A}^k$ , i.e. the set of all

$$\frac{(\mathcal{A}^k u, u)_{\mathcal{M}(\gamma)}}{(u, u)_{\mathcal{M}(\gamma)}}, \quad u \neq 0,$$

is equal to the convex hull of the eigenvalues of  $\mathcal{A}^k$ . Since all eigenvalues of  $\mathcal{A}^k$  are real and positive (cf. Corollary 2.4),

$$0 < \lambda_{\min}(\mathcal{A}^k) \leq \frac{(\mathcal{A}^k u, u)_{\mathcal{M}(\gamma)}}{(u, u)_{\mathcal{M}(\gamma)}} = \frac{(\mathcal{M}(\gamma)\mathcal{A}^k u, u)}{(u, u)_{\mathcal{M}(\gamma)}} \leq \lambda_{\max}(\mathcal{A}^k) \quad \text{for all } u \neq 0.$$

Therefore  $(\mathcal{M}(\gamma)\mathcal{A}^k u, u) > 0$  for all  $u \neq 0$ , which shows that  $\mathcal{M}(\gamma)\mathcal{A}^k$  is symmetric positive definite.  $\square$

We now derive expressions for  $\alpha_i$  and  $\beta_{i+1}$  in Algorithm 3.1, in case  $\mathcal{M}(\gamma)\mathcal{A}$  is positive definite and chosen as the inner product matrix.

LEMMA 3.6. *Suppose that the (symmetric) matrix  $\mathcal{M}(\gamma)$  is positive definite. Then Algorithm 3.1 is well defined for  $M = \mathcal{A}$  and  $G = \mathcal{M}(\gamma)\mathcal{A}$ , and (until convergence) the scalars  $\alpha_i$  and  $\beta_{i+1}$ , can be computed as*

$$(3.6) \quad \alpha_i = \frac{(r_i, r_i)_{\mathcal{M}(\gamma)}}{(\mathcal{A}p_i, p_i)_{\mathcal{M}(\gamma)}},$$

$$(3.7) \quad \beta_{i+1} = \frac{(r_{i+1}, r_{i+1})_{\mathcal{M}(\gamma)}}{(r_i, r_i)_{\mathcal{M}(\gamma)}}.$$

*Proof.* Since  $\mathcal{M}(\gamma)$  is positive definite,  $\mathcal{M}(\gamma)\mathcal{A}^k$  is symmetric positive definite for all  $k \geq 0$  (cf. Theorem 3.5). In particular,  $\mathcal{M}(\gamma)\mathcal{A}$  is symmetric positive definite and hence defines an inner product. Moreover,  $\mathcal{M}(\gamma)\mathcal{A}^2 = (\mathcal{M}(\gamma)\mathcal{A})\mathcal{A}$  is symmetric positive definite, which means that  $\mathcal{A}$  is both  $\mathcal{M}(\gamma)\mathcal{A}$ -symmetric and  $\mathcal{M}(\gamma)\mathcal{A}$ -definite. Therefore, the assumptions of Theorem 3.2 are satisfied, and Algorithm 3.1 with  $M = \mathcal{A}$  and  $G = \mathcal{M}(\gamma)\mathcal{A}$  is well defined.

It is easy to see that for  $i \geq 0$  the denominator of  $\alpha_i$  in (3.1) is equal to  $(\mathcal{A}p_i, p_i)_{\mathcal{M}(\gamma)}$ . For  $i = 0$ , the numerator of  $\alpha_i$  is equal to

$$(x - x_0, p_0)_{\mathcal{M}(\gamma)\mathcal{A}} = (\mathcal{A}(x - x_0), r_0)_{\mathcal{M}(\gamma)} = (r_0, r_0)_{\mathcal{M}(\gamma)},$$

showing that (3.6) holds for  $i = 0$ . For  $i \geq 1$  we use (3.5) and the orthogonality relation in item (3) of Theorem 3.2 to obtain

$$\begin{aligned} (x - x_i, p_i)_{\mathcal{M}(\gamma)\mathcal{A}} &= (x - x_i, r_i + \beta_i p_{i-1})_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= (x - x_i, r_i)_{\mathcal{M}(\gamma)\mathcal{A}} + \beta_i (x - x_i, p_{i-1})_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= (\mathcal{A}(x - x_i), r_i)_{\mathcal{M}(\gamma)} \\ &= (r_i, r_i)_{\mathcal{M}(\gamma)}, \end{aligned}$$

which proves (3.6) for  $i \geq 1$ . Note that  $\alpha_i \neq 0$  for  $i \geq 0$ .

Next, we consider the numerator of  $\beta_{i+1}$ ,  $i \geq 0$ , in (3.4). Here we use (3.3) and again the orthogonality relation in item (3) of Theorem 3.2 to obtain

$$\begin{aligned} (r_{i+1}, p_i)_{\mathcal{M}(\gamma)\mathcal{A}} &= (x - x_{i+1}, \mathcal{A}p_i)_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= (x - x_{i+1}, \alpha_i^{-1}(r_i - r_{i+1}))_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= \alpha_i^{-1}(x - x_{i+1}, r_i)_{\mathcal{M}(\gamma)\mathcal{A}} - \alpha_i^{-1}(x - x_{i+1}, r_{i+1})_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= -\alpha_i^{-1}(r_{i+1}, r_{i+1})_{\mathcal{M}(\gamma)}. \end{aligned}$$

Therefore,

$$\beta_{i+1} = -\frac{(r_{i+1}, p_i)_{\mathcal{M}(\gamma)\mathcal{A}}}{(p_i, p_i)_{\mathcal{M}(\gamma)\mathcal{A}}} = \alpha_i^{-1} \frac{(r_{i+1}, r_{i+1})_{\mathcal{M}(\gamma)}}{(\mathcal{A}p_i, p_i)_{\mathcal{M}(\gamma)}} = \frac{(r_{i+1}, r_{i+1})_{\mathcal{M}(\gamma)}}{(r_i, r_i)_{\mathcal{M}(\gamma)}},$$

which completes the proof.  $\square$

**3.3. An efficient implementation.** In step  $i$  of Algorithm 3.1 with  $M = \mathcal{A}$  and  $G = \mathcal{M}(\gamma)\mathcal{A}$ , we can compute the scalars  $\alpha_i$  and  $\beta_{i+1}$  as shown in (3.6) and (3.7), respectively. We will now show how to replace the  $\mathcal{M}(\gamma)$ -inner products by  $\mathcal{J}$ -bilinearforms. Since  $\mathcal{M}(\gamma) = \mathcal{J}\mathcal{A} - \gamma\mathcal{J}$ , cf. (2.2), we have for all vectors  $u, v \in \mathbb{R}^{n+m}$ ,

$$(u, v)_{\mathcal{M}(\gamma)} = (\mathcal{A}u, v)_{\mathcal{J}} - \gamma(u, v)_{\mathcal{J}}.$$



Therefore,

$$(3.8) \quad (r_i, r_i)_{\mathcal{M}(\gamma)} = (\mathcal{A}r_i, r_i)_{\mathcal{J}} - \gamma(r_i, r_i)_{\mathcal{J}},$$

and, since  $\mathcal{A}$  is  $\mathcal{M}(\gamma)$ -symmetric,

$$(3.9) \quad (\mathcal{A}p_i, p_i)_{\mathcal{M}(\gamma)} = (p_i, \mathcal{A}p_i)_{\mathcal{M}(\gamma)} = (\mathcal{A}p_i, \mathcal{A}p_i)_{\mathcal{J}} - \gamma(p_i, \mathcal{A}p_i)_{\mathcal{J}}.$$

Hence to compute  $\alpha_i$  and  $\beta_{i+1}$ , the main work lies in evaluating the bilinear form  $(u, v)_{\mathcal{J}}$ , which is not more expensive than evaluating the Euclidean inner product  $(u, v)$ . Note that we need to have available both  $\mathcal{A}r_i$  and  $\mathcal{A}p_i$ . To avoid the necessity of computing both matrix-vector products in every step, we store two additional vectors, namely  $y_i = \mathcal{A}r_i$  and  $w_i = \mathcal{A}p_i$ . The former is computed by multiplying  $\mathcal{A}$  against  $r_i$ . The latter is computed via an additional recursion in the following way: Multiplying (3.5) by  $\mathcal{A}$  yields

$$\underbrace{\mathcal{A}p_{i+1}}_{=w_{i+1}} = \underbrace{\mathcal{A}r_{i+1}}_{=y_{i+1}} + \beta_{i+1} \underbrace{\mathcal{A}p_i}_{=w_i}.$$

The complete algorithm looks as follows.

ALGORITHM 3.7. (*CG method for  $\mathcal{A}$* )

*Input:* System matrix  $\mathcal{A}$ , right hand side  $b$ , real parameter  $\gamma$ , initial guess  $x_0$ .

*Initialize:*  $r_0 = b - \mathcal{A}x_0$ ,  $p_0 = r_0$ ,  $y_0 = \mathcal{A}r_0$ ,  $w_0 = y_0$

*For*  $i = 0, 1, \dots$  *until convergence:*

$$(3.10) \quad \alpha_i = \frac{(y_i, r_i)_{\mathcal{J}} - \gamma(r_i, r_i)_{\mathcal{J}}}{(w_i, w_i)_{\mathcal{J}} - \gamma(p_i, w_i)_{\mathcal{J}}}$$

$$(3.11) \quad x_{i+1} = x_i + \alpha_i p_i$$

$$(3.12) \quad r_{i+1} = r_i - \alpha_i w_i$$

$$(3.13) \quad y_{i+1} = \mathcal{A}r_{i+1}$$

$$(3.14) \quad \beta_{i+1} = \frac{(y_{i+1}, r_{i+1})_{\mathcal{J}} - \gamma(r_{i+1}, r_{i+1})_{\mathcal{J}}}{(y_i, r_i)_{\mathcal{J}} - \gamma(r_i, r_i)_{\mathcal{J}}}$$

$$(3.15) \quad p_{i+1} = r_{i+1} + \beta_{i+1} p_i$$

$$(3.16) \quad w_{i+1} = y_{i+1} + \beta_{i+1} w_i$$

Since the denominator of  $\beta_{i+1}$  is equal to the numerator of  $\alpha_i$ , this quantity only has to be evaluated once in every step. When these scalars are stored, each step of Algorithm 3.7 requires four evaluations of the bilinear form  $(u, v)_{\mathcal{J}}$ , compared to two evaluations of the inner product  $(u, v)_G$  in the (generic) CG method stated in Algorithm 3.1. In addition, Algorithm 3.7 requires two more vectors of storage and one more recurrence, namely (3.16), than Algorithm 3.1. We summarize the theoretical requirements and properties of Algorithm 3.7 in the following result.

COROLLARY 3.8. *If the (symmetric) matrix  $\mathcal{M}(\gamma)$  is positive definite, then for  $M = \mathcal{A}$  and  $G = \mathcal{M}(\gamma)\mathcal{A}$  the assumptions of Theorem 3.2 are satisfied, and Algorithm 3.7 is a well defined CG method for solving  $\mathcal{A}x = b$ .*

**4. Practical issues.** In this section we discuss several practical issues concerning our CG method in Algorithm 3.7.

**4.1. The condition number of  $\mathcal{M}(\gamma)$ .** While Algorithm 3.7 is based on the inner product defined by  $\mathcal{M}(\gamma)\mathcal{A}$ , we actually compute the scalars  $\alpha_i$  and  $\beta_{i+1}$  using the inner product defined by  $\mathcal{M}(\gamma)$ , cf. (3.6) and (3.7). Therefore it is of interest to estimate  $\kappa(\mathcal{M}(\gamma))$ , the condition number of  $\mathcal{M}(\gamma)$ , in order to assess the numerical stability of the method.

LEMMA 4.1. *Suppose that (2.3) and (2.5) hold, so that the matrix  $\mathcal{M}(\gamma)$  is symmetric positive definite. Let  $\xi \equiv (\lambda_{\min}(A) - \gamma)(\gamma - \lambda_{\max}(C)) - \|B\|^2$ , then*

$$(4.1) \quad \kappa(\mathcal{M}(\gamma)) < \frac{4}{\xi} (\lambda_{\max}(A) - \gamma)(\lambda_{\min}(A) - \gamma).$$

*Proof.* By assumption, the matrix  $\mathcal{M}(\gamma)$  permits the factorization

$$\begin{bmatrix} (A - \gamma I)^{1/2} & 0 \\ 0 & (\gamma I - C)^{1/2} \end{bmatrix} \begin{bmatrix} I & X^T \\ X & I \end{bmatrix} \begin{bmatrix} (A - \gamma I)^{1/2} & 0 \\ 0 & (\gamma I - C)^{1/2} \end{bmatrix},$$

where  $X \equiv (\gamma I - C)^{-1/2} B (A - \gamma I)^{-1/2}$ , so that

$$\|X\| \leq \frac{\|B\|}{(\lambda_{\min}(A) - \gamma)^{1/2} (\gamma - \lambda_{\max}(C))^{1/2}}.$$

Since (2.5) holds we have  $\xi > 0$ , and hence the right hand side is less than one. For any congruence  $M = F H F^T$ ,

$$\kappa(M) \leq \|F\|^2 \|H\| \|F^{-1}\|^2 \|H^{-1}\| = \kappa(F^2) \kappa(H),$$

and therefore

$$\begin{aligned} \kappa(\mathcal{M}(\gamma)) &\leq \kappa \left( \begin{bmatrix} A - \gamma I & 0 \\ 0 & \gamma I - C \end{bmatrix} \right) \kappa \left( \begin{bmatrix} I & X^T \\ X & I \end{bmatrix} \right) \\ &= \frac{\lambda_{\max}(A) - \gamma}{\gamma - \lambda_{\max}(C)} \frac{1 + \|X\|}{1 - \|X\|} \\ &= \frac{\lambda_{\max}(A) - \gamma}{\gamma - \lambda_{\max}(C)} \frac{(1 + \|X\|)^2}{1 - \|X\|^2} \\ &< 4 \frac{\lambda_{\max}(A) - \gamma}{\gamma - \lambda_{\max}(C)} \frac{(\lambda_{\min}(A) - \gamma)(\gamma - \lambda_{\max}(C))}{\xi}, \end{aligned}$$

which concludes the proof.  $\square$

The bound (4.1) indicates a relation between  $\kappa(\mathcal{M}(\gamma))$  and the sufficient condition (2.5) for positive definiteness of  $\mathcal{M}(\gamma)$ : With larger  $\xi$ , the bound on the condition number of  $\mathcal{M}(\gamma)$  becomes smaller, and vice versa.

The best choice of  $\gamma$  is the one that minimizes  $\kappa(\mathcal{M}(\gamma))$ , but  $\gamma = \hat{\gamma}$  as in Corollary 2.3 is a more accessible substitute. For this choice, and the corresponding value of  $\xi = \hat{\xi}$ , (4.1) implies that

$$(4.2) \quad \kappa(\mathcal{M}(\hat{\gamma})) < \frac{(2\lambda_{\max}(A) - \lambda_{\max}(C))(\lambda_{\min}(A) - \lambda_{\max}(C))}{\hat{\xi}}.$$

In the special case  $C = 0$ , (4.2) simplifies to

$$\kappa(\mathcal{M}(\hat{\gamma})) < \frac{2}{\hat{\xi}} \lambda_{\max}(A) \lambda_{\min}(A) < \frac{8\lambda_{\max}(A)}{\lambda_{\min}(A) - 2\|B\|}.$$

In this case, Benzi and Simoncini [5, Corollary 3.2] have estimated  $\kappa(\mathcal{M}(\hat{\gamma}))$  as

$$\kappa(\mathcal{M}(\hat{\gamma})) \approx \frac{4\lambda_{\max}(A)}{\lambda_{\min}(A) - 4\lambda_{\max}(BA^{-1}B^T)}.$$

From (2.7), it is easy to see that the denominator on the right hand side is bounded from below by  $\lambda_{\min}(A) - 2\|B\|$ , so that the resulting upper bound on the right hand side corresponds to our bound up to a constant factor of two.

**4.2. Error bounds.** When the matrix  $\mathcal{M}(\gamma)$  is positive definite, the CG method in Algorithm 3.7 is optimal in the  $(\mathcal{M}(\gamma)\mathcal{A})$ -norm, see item (4) in Theorem 3.2. We know that  $\mathcal{A}$  is  $(\mathcal{M}(\gamma)\mathcal{A})$ -symmetric (cf. Theorem 3.5), and hence  $\mathcal{A}$  has a complete set of eigenvectors that are orthonormal with respect to the inner product defined by  $\mathcal{M}(\gamma)\mathcal{A}$  (cf. the proof of Corollary 2.4). Hence we may write

$$\mathcal{A} = \mathcal{Y}\Lambda\mathcal{Y}^{-1}, \quad \text{where } \mathcal{Y}^T(\mathcal{M}(\gamma)\mathcal{A})\mathcal{Y} = I.$$

Suppose that  $x_0$  is an initial guess for the solution of  $\mathcal{A}x = b$ , and write the initial error as  $x - x_0 = \mathcal{Y}v$ , for some vector  $v \in \mathbb{R}^{n+m}$ . Then the  $(\mathcal{M}(\gamma)\mathcal{A})$ -norm of the  $i$ th error satisfies (cf. item (4) in Theorem 3.2)

$$\begin{aligned} \|x - x_i\|_{\mathcal{M}(\gamma)\mathcal{A}} &= \min_{p \in \pi_i} \|p(\mathcal{A})(x - x_0)\|_{\mathcal{M}(\gamma)\mathcal{A}} \\ &= \min_{p \in \pi_i} (v^T p(\Lambda)^2 v)^{1/2} \\ (4.3) \qquad &\leq \|x - x_0\|_{\mathcal{M}(\gamma)\mathcal{A}} \min_{p \in \pi_i} \max_{\lambda \in \Lambda(\mathcal{A})} |p(\lambda)|. \end{aligned}$$

Here  $\Lambda(\mathcal{A})$  denotes the (real and positive) spectrum of  $\mathcal{A}$ . The bound (4.3) and its derivation is completely analogous to the standard convergence bound for the classical CG method in case of a symmetric positive definite system matrix  $M$  and error minimization in the  $M$ -norm. Estimation of the quantity  $\min_{p \in \pi_i} \max_{\lambda \in \Lambda(\mathcal{A})} |p(\lambda)|$  using the eigenvalue distribution of  $\mathcal{A}$  has been exhaustively done in the literature, see, e.g., [3, Chapter 13]. The important information given by (4.3) is that when the spectrum of  $\mathcal{A}$  is clustered away from the origin, then fast convergence can be expected. This gives some indication on how to choose a preconditioner.

To get a computable estimate on the error that is minimized in every step, we use that  $\mathcal{A}$  is  $(\mathcal{M}(\gamma)\mathcal{A})$ -symmetric and  $(\mathcal{M}(\gamma)\mathcal{A})$ -definite. In this case [1, Corollary 5.2] applies, and shows that

$$(4.4) \quad \left( \widehat{\kappa}(\mathcal{A})^{-1} \frac{(r_i, r_i)_{\mathcal{M}(\gamma)}}{(b, b)_{\mathcal{M}(\gamma)}} \right)^{1/2} \leq \frac{\|x - x_i\|_{\mathcal{M}(\gamma)\mathcal{A}}}{\|x\|_{\mathcal{M}(\gamma)\mathcal{A}}} \leq \left( \widehat{\kappa}(\mathcal{A}) \frac{(r_i, r_i)_{\mathcal{M}(\gamma)}}{(b, b)_{\mathcal{M}(\gamma)}} \right)^{1/2},$$

where  $(r_i, r_i)_{\mathcal{M}(\gamma)}$  is the numerator of  $\alpha_i$  (and thus available in every step at no extra cost), cf. (3.6), and

$$\widehat{\kappa}(\mathcal{A}) \equiv \frac{\max_{\lambda \in \Lambda(\mathcal{A})} \lambda}{\min_{\lambda \in \Lambda(\mathcal{A})} \lambda}.$$

Bendixon's Theorem [11, p. 69] yields

$$\min \{ \lambda_{\min}(A), \lambda_{\min}(C) \} \leq \lambda \leq \max \{ \lambda_{\max}(A), \lambda_{\max}(C) \}$$

for all  $\lambda \in \Lambda(\mathcal{A})$ , so that

$$\widehat{\kappa}(\mathcal{A}) \leq \frac{\max\{\lambda_{\max}(A), \lambda_{\max}(C)\}}{\min\{\lambda_{\min}(A), \lambda_{\min}(C)\}}.$$

Of course, when  $C$  is singular, this estimate is useless. Close estimates for the eigenvalues of  $\mathcal{A}$  may be obtained using parameters computed by the CG method itself. This gives a convergence bound that becomes tighter during the run of the method. We will not discuss this approach here, and refer the interested reader to [1, Section 7].

We next relate the  $(\mathcal{M}(\gamma)\mathcal{A})$ -norm of the error to the Euclidean norm of the residual, which is often used as a stopping criterion for the CG method (even though this quantity is not minimized, and may strongly oscillate during the iteration). Since  $\mathcal{M}(\gamma) = \mathcal{J}\mathcal{A} - \gamma\mathcal{J} > 0$ , we have  $u^T \mathcal{J}\mathcal{A}u - \gamma u^T \mathcal{J}u > 0$ , or  $-u^T \mathcal{J}\mathcal{A}u < -\gamma u^T \mathcal{J}u$ , for all vectors  $u \in \mathbb{R}^{n+m}$ , so that

$$\begin{aligned} \|x - x_i\|_{\mathcal{M}(\gamma)\mathcal{A}}^2 &= (x - x_i)^T \mathcal{M}(\gamma)\mathcal{A}(x - x_i) \\ &= (x - x_i)^T (\mathcal{A}^T \mathcal{J} - \gamma\mathcal{J})\mathcal{A}(x - x_i) \\ &= r_i^T \mathcal{J}r_i - \gamma(x - x_i)^T \mathcal{J}\mathcal{A}(x - x_i) \\ &< r_i^T \mathcal{J}r_i - \gamma^2(x - x_i)^T \mathcal{J}(x - x_i) \\ &= (r_i, r_i)_{\mathcal{J}} - \gamma^2(x - x_i, x - x_i)_{\mathcal{J}} \\ &\leq \|r_i\|^2 + \gamma^2\|x - x_i\|^2 \\ &= \|r_i\|^2 + \gamma^2\|\mathcal{A}^{-1}r_i\|^2 \\ &\leq \|r_i\|^2 \left(1 + \frac{\gamma^2}{\sigma_{\min}^2(\mathcal{A})}\right), \end{aligned}$$

where  $\sigma_{\min}(\mathcal{A})$  denotes the smallest singular value of  $\mathcal{A}$ . In particular, for  $\gamma = \hat{\gamma}$  as in Corollary 2.3,

$$\|x - x_i\|_{\mathcal{M}(\hat{\gamma})\mathcal{A}} < \|r_i\| \left(1 + \frac{(\lambda_{\min}(A) + \lambda_{\max}(C))^2}{4\sigma_{\min}^2(\mathcal{A})}\right)^{1/2}.$$

We see that if  $\mathcal{A}$  is not too ill conditioned, then the Euclidean norm of the residual gives a reasonable bound on the  $(\mathcal{M}(\gamma)\mathcal{A})$ -norm of the error.

**4.3. Block diagonal preconditioning.** One of the most popular preconditioning techniques for saddle point systems is block diagonal preconditioning; see [4, Section 10.1] for an overview. Consider a typical saddle point matrix as in (1.1), with  $A = A^T > 0$ , full rank  $B$ , and  $C = C^T \geq 0$ . Then this matrix permits the factorization

$$(4.5) \quad \begin{bmatrix} A & B^T \\ B & -C \end{bmatrix} = \begin{bmatrix} I & 0 \\ BA^{-1} & I \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & S \end{bmatrix} \begin{bmatrix} I & A^{-1}B^T \\ 0 & I \end{bmatrix},$$

where  $S = -(C + BA^{-1}B^T)$ , so that  $-S = -S^T > 0$ . The general idea of block diagonal preconditioning is to construct a preconditioner from (an approximation of) the inverse of the block diagonal matrix  $[A \ 0; 0 \ S]$  on the right hand side of (4.5). Here we use (approximate) Cholesky factorizations of  $A$  and  $-S$ , namely  $A \approx KK^T$  and  $-S \approx LL^T$ . Let  $\alpha > 0$  be a real parameter, then

$$(4.6) \quad \begin{bmatrix} \alpha K^{-1} & 0 \\ 0 & -L^{-1} \end{bmatrix} \begin{bmatrix} A & B^T \\ B & -C \end{bmatrix} \begin{bmatrix} \alpha K^{-T} & 0 \\ 0 & L^{-T} \end{bmatrix} \equiv \begin{bmatrix} \alpha^2 \widehat{A} & \alpha \widehat{B}^T \\ -\alpha \widehat{B} & \widehat{C} \end{bmatrix} \equiv \widehat{\mathcal{A}}.$$

Clearly,  $\widehat{A} = K^{-1}AK^{-T}$  is symmetric positive definite,  $\widehat{B} = L^{-1}BK^{-T}$  has full rank, and  $\widehat{C} = L^{-1}CL^{-T}$  is symmetric positive semidefinite. We can therefore apply Corollary 2.3 to  $\widehat{A}$ , which shows that this matrix is diagonalizable with real and positive eigenvalues whenever

$$(4.7) \quad 2\alpha \|\widehat{B}\| < \alpha^2 \lambda_{\min}(\widehat{A}) - \lambda_{\max}(\widehat{C}).$$

When  $\alpha$  is chosen large enough, this condition is satisfied, even for the trivial block diagonal preconditioner with  $K = I$  and  $L = I$ , i.e. when (4.6) amounts to a scaling of the saddle point matrix. Therefore we may, at least theoretically, transform by scaling any saddle point matrix of the form (1.1) with  $A = A^T > 0$ , full rank  $B$ , and  $C = C^T \geq 0$ , into a matrix of the form (4.6) for which Algorithm 3.7 represents a well defined CG method.

Of course, if  $\|\widehat{B}\| \gg \lambda_{\min}(\widehat{A})$ , then  $\alpha$  must be chosen very large, and this might cause numerical problems, or may be incompatible with the application at hand. For example, in case of an equality-constrained optimization problem, where  $B$  represents the constraints (see [4, Section 1.1]), scaling with a very large  $\alpha$  transforms the original saddle point matrix (in the limit  $\alpha \rightarrow \infty$ ) into one that represents an unconstrained problem.

Ultimately, it depends on the properties of the specific application whether scaling can be applied or not. Nevertheless, it is instructive to consider the case of *exact* block diagonal preconditioning, i.e.  $A = KK^T$  and  $-S = LL^T$ . A simple computation shows that in this case

$$\widehat{A} = I, \quad \widehat{B} = L^{-1}BK^{-T}, \quad \widehat{C} = I - \widehat{B}\widehat{B}^T.$$

Since  $\widehat{C} \geq 0$ , we have  $I \geq \widehat{B}\widehat{B}^T$ , and thus  $\|\widehat{B}\| \leq 1$  and  $\lambda_{\max}(I - \widehat{B}\widehat{B}^T) \leq 1$ . Consequently,

$$2\alpha \|\widehat{B}\| + \lambda_{\max}(\widehat{C}) \leq 2\alpha + 1,$$

leading to the sufficient condition  $\alpha^2 > 2\alpha + 1$ , or  $\alpha^2 - 2\alpha - 1 > 0$ .

**LEMMA 4.2.** *Suppose that the matrix  $\widehat{A}$  in (4.6) is the result of exact block diagonal preconditioning as described above. If  $\alpha > 1 + \sqrt{2}$ , then  $\widehat{A}$  is diagonalizable with real and positive eigenvalues.*

In the special case  $C = 0$ , we have  $\widehat{C} = 0$ , and the condition on  $\alpha$  simplifies to  $\alpha > 2$ , which is the result shown in [5, Section 5.1].

**5. Numerical examples.** In this section we present results of numerical experiments with test problems generated by the MATLAB [12] package Incompressible Flow Iterative Solution Software (IFISS) [13]. We use the driver `stokes_testproblem` of this code with default options to set up a stabilized discretization of a Stokes equations model problem<sup>1</sup>, resulting in a linear system of the form

$$(5.1) \quad \begin{bmatrix} A & B^T \\ B & -\frac{1}{4}C \end{bmatrix} \begin{bmatrix} u \\ p \end{bmatrix} = \begin{bmatrix} f \\ g \end{bmatrix},$$

<sup>1</sup>The default parameters are: lid driven cavity; cavity type: regularized; grid parameter 4 ( $16 \times 16$  grid); uniform grid;  $Q_1 - P_0$  elements; stabilization parameter 1/4; uniform streamlines. See the IFISS user guide or [6, Chapter 5] for a detailed description of this test problem.

where  $A = A^T > 0$  is of order  $n = 578$ ,  $C = C^T \geq 0$  is of order  $m = 256$ , and, by construction,  $\text{rank}(B) = m - 2 = 254$ . To agree with the notation used in [6], we have written out the stabilization parameter  $\frac{1}{4}$  explicitly. The system matrix is of order  $n+m = 834$  and of rank  $n+m-1 = 833$ . However, the system is consistent, and the singularity of the system matrix represents no difficulty for the iterative methods considered here. Other parameters relevant for our context, and computed using MATLAB's `eig` routine, are:

$$\lambda_{\max}(A) = 3.9493, \quad \lambda_{\min}(A) = 0.0764, \quad \lambda_{\max}(\tfrac{1}{4}C) = 0.0156, \quad \lambda_{\min}(\tfrac{1}{4}C) = 0.$$

The spectra of  $A$  and  $\frac{1}{4}C$  are separated, so that (2.3) is satisfied. Since  $\|B\| = 0.2476$ , the sufficient condition (2.6) is not satisfied. Nevertheless, with

$$\hat{\gamma} = \frac{1}{2}(\lambda_{\min}(A) + \lambda_{\max}(\tfrac{1}{4}C)) = 0.0460,$$

we set up the matrix  $\mathcal{M}(\hat{\gamma})$ , which is positive definite and rather well conditioned:

$$\lambda_{\max}(\mathcal{M}(\hat{\gamma})) = 3.9191, \quad \lambda_{\min}(\mathcal{M}(\hat{\gamma})) = 0.0118, \quad \kappa(\mathcal{M}(\hat{\gamma})) = 333.3771$$

(computed using MATLAB's `eig` and `cond`). We negate the second block row in (5.1) to obtain

$$(5.2) \quad \begin{bmatrix} A & B^T \\ -B & \frac{1}{4}C \end{bmatrix} \begin{bmatrix} u \\ p \end{bmatrix} = \begin{bmatrix} f \\ -g \end{bmatrix}.$$

We now apply MATLAB's build-in MINRES algorithm to (5.1) and our Algorithm 3.7 to (5.2) (both with  $x_0 = 0$ ). The resulting convergence characteristics are shown in Fig. 5.1:

- dotted :  $\|r_i\|/\|r_0\|$  of MINRES,
- solid :  $\|r_i\|/\|r_0\|$  of Algorithm 3.7,
- dashed :  $\|x - x_i\|_{\mathcal{M}(\hat{\gamma})\mathcal{A}} / \|x - x_0\|_{\mathcal{M}(\hat{\gamma})\mathcal{A}}$  of Algorithm 3.7,
- dashed-dotted :  $(r_i, r_i)_{\mathcal{M}(\hat{\gamma})} / (b, b)_{\mathcal{M}(\hat{\gamma})}$  of Algorithm 3.7 (cf. (4.4))

(to compute the error norms,  $x^T = [u^T, p^T]$  is obtained by solving the system (5.2) using the MATLAB backslash operator). Clearly, Algorithm 3.7 is competitive with MINRES, which is optimal for the linear system (5.1) in the sense that it minimizes the Euclidean norm of the residual over the Krylov subspace generated by the system matrix and the right hand side. In fact, the convergence of the Euclidean residual norms of Algorithm 3.7 slightly outperforms those of MINRES. Note, however, that the Euclidean residual norms of Algorithm 3.7 are not monotonically decreasing; they do not satisfy a minimization property. On the other hand, the  $(\mathcal{M}(\hat{\gamma})\mathcal{A})$ -norm of the error is monotonically decreasing, and in this example it is very close to the Euclidean residual norm. Moreover, a good estimate of this norm is given by  $(r_i, r_i)_{\mathcal{M}(\hat{\gamma})} / (b, b)_{\mathcal{M}(\hat{\gamma})}$ , a quantity that is available at no additional cost during the iteration.

To obtain a larger test example we again use `stokes_testproblem`, but this time we choose the grid parameter 6 in IFISS, resulting in a  $64 \times 64$  grid, and system dimensions  $n = 8450$  and  $m = 4096$ . In this case the computation of the exact eigenvalues is rather expensive, and so we only compute *estimates* in MATLAB:

$$\text{normest}(A) = 3.9965, \quad \lambda_{\min}(A) \approx 0.0048 \text{ (estimated using } \text{eigs} \text{ with } \text{maxit}=20),$$

$$\text{normest}(\tfrac{1}{4}C) = 9.7656e - 004, \quad \text{normest}(B) = 0.0625.$$

For these estimates (2.3) is satisfied, but again (2.6) is not. The run of Algorithm 3.7 is based on  $\gamma = \frac{1}{2}(0.0048 + \text{normest}(\frac{1}{4}C)) = 0.0029$ . In Fig. 5.2 we show the convergence characteristics of MINRES and Algorithm 3.7; the notation corresponds to the one of Fig. 5.1. Obviously, the qualitative behavior of the algorithms for both test problems is the same.

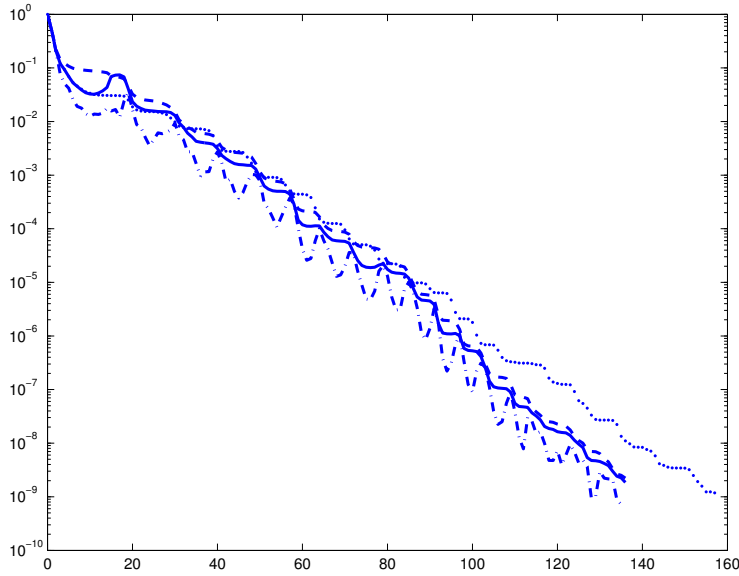


FIG. 5.1. Convergence characteristics for the default stokes\_testproblem.

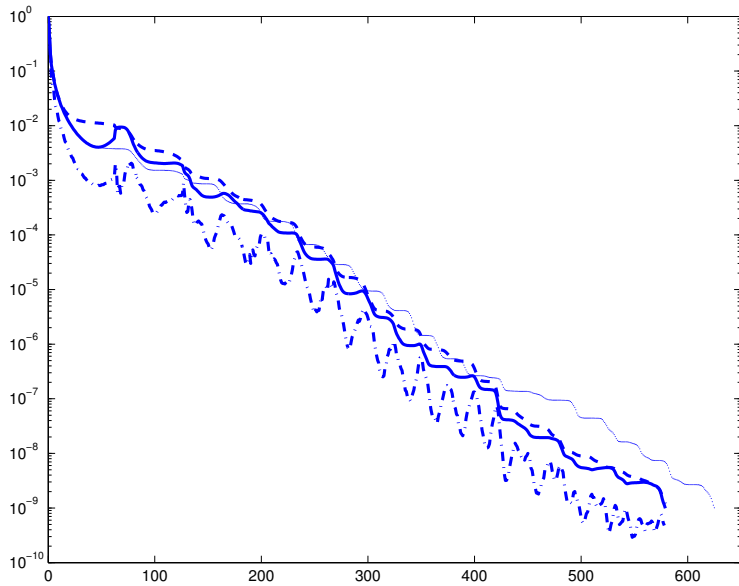


FIG. 5.2. Convergence characteristics for the larger stokes\_testproblem.

**6. Concluding remarks.** We have considered the idea of negating the second block row in a saddle point system, which leads to an unsymmetric but positive (semi)definite, rather than symmetric but indefinite system matrix. We have generalized previous results on the definiteness and conditioning of the bilinear form with respect to which the unsymmetric saddle point matrix is symmetric. In particular, we have included the case of a general positive semidefinite block  $C$ , while previous authors only considered the case  $C = 0$ .

We have derived an efficient CG method for solving the (unsymmetric positive definite) saddle point system. In numerical experiments we have seen that this method can outperform the MINRES method, which so far has been considered the standard solver for (symmetric indefinite) saddle point systems. We have discussed several practical issues concerning our new CG method, including conditioning of the inner product, error bounds, and preconditioning.

Our goal has been to present the theory and the algorithms in a clean and easily readable, rather than most general form. Many options for generalization and further analysis exist, and we hope that these will be explored in the future. In particular, further analysis is needed of the indefinite case, i.e. the practically relevant situation when  $\gamma$  has been chosen to yield an indefinite matrix  $\mathcal{M}(\gamma)$ . In very large scale applications, where only crude estimates of the relevant eigenvalues are available, such situation is not unlikely to occur.

**Acknowledgements.** Part of the work of Jörg Liesen was done during his visit of Emory University in April 2006. He thanks Michele Benzi for his kind hospitality and for very helpful discussions and suggestions. Thanks also to Valeria Simoncini and Petr Tichý for their helpful comments.

#### REFERENCES

- [1] S. F. ASHBY, M. J. HOLST, T. A. MANTEUFFEL, AND P. E. SAYLOR, *The role of the inner product in stopping criteria for conjugate gradient iterations*, BIT, 41 (2001), pp. 26–52.
- [2] S. F. ASHBY, T. A. MANTEUFFEL, AND P. E. SAYLOR, *A taxonomy for conjugate gradient methods*, SIAM J. Numer. Anal., 27 (1990), pp. 1542–1568.
- [3] O. AXELSSON, *Iterative solution methods*, Cambridge University Press, Cambridge, 1994.
- [4] M. BENZI, G. H. GOLUB, AND J. LIESEN, *Numerical solution of saddle point problems*, Acta Numer., 14 (2005), pp. 1–137.
- [5] M. BENZI AND V. SIMONCINI, *On the eigenvalues of a class of saddle point matrices*, Numer. Math., 103 (2006), pp. 173–196.
- [6] H. C. ELMAN, D. J. SILVESTER, AND A. J. WATHEN, *Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics*, Numerical Mathematics and Scientific Computation, Oxford University Press, New York, 2005.
- [7] B. FISCHER, A. RAMAGE, D. J. SILVESTER, AND A. J. WATHEN, *Minimum residual methods for augmented systems*, BIT, 38 (1998), pp. 527–543.
- [8] F. R. GANTMACHER, *The theory of matrices. Vols. 1, 2*, Chelsea Publishing Co., New York, 1959.
- [9] W. GIVENS, *Fields of values of a matrix*, Proc. Amer. Math. Soc., 3 (1952), pp. 206–209.
- [10] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, J. Research Nat. Bur. Standards, 49 (1952), pp. 409–436 (1953).
- [11] A. S. HOUSEHOLDER, *The theory of matrices in numerical analysis*, Dover Publications Inc., New York, 1975. Reprint of 1964 edition.
- [12] THE MATHWORKS COMPANY, *Matlab, version 6.5*. <http://www.mathworks.com>.
- [13] D. J. SILVESTER, H. C. ELMAN, AND A. RAMAGE, *Incompressible Flow Iterative Solution Software (IFISS), version 2.2*. <http://www.manchester.ac.uk/ifiss>.