



# On optimal improvements of classical iterative schemes for $Z$ -matrices

D. Noutsos<sup>a,\*</sup>, M. Tzoumas<sup>b</sup>

<sup>a</sup>*Department of Mathematics, University of Ioannina, GR-451 10 Ioannina, Greece*

<sup>b</sup>*Department of Environment and Natural Resources Management, Agrinion University School, University of Ioannina, GR-301 00 Agrinion, Greece*

Received 9 September 2004; received in revised form 12 January 2005

## Abstract

Many researchers have considered preconditioners, applied to linear systems, whose matrix coefficient is a  $Z$ - or an  $M$ -matrix, that make the associated Jacobi and Gauss–Seidel methods converge asymptotically faster than the unpreconditioned ones. Such preconditioners are chosen so that they eliminate the off-diagonal elements of the same column or the elements of the first upper diagonal [Milaszewicz, LAA 93 (1987) 161–170], Gunawardena et al. [LAA 154–156 (1991) 123–143]. In this work we generalize the previous preconditioners to obtain optimal methods. “Good” Jacobi and Gauss–Seidel algorithms are given and preconditioners, that eliminate more than one entry per row, are also proposed and analyzed. Moreover, the behavior of the above preconditioners to the Krylov subspace methods is studied.

© 2005 Elsevier B.V. All rights reserved.

MSC: Primary 65F10

Keywords: Jacobi and Gauss–Seidel iterative methods; Diagonally dominant  $Z$ - and  $M$ -matrices

## 1. Introduction and preliminaries

Consider the linear system of algebraic equations

$$Ax = b, \tag{1.1}$$

\* Corresponding author.

E-mail address: [dnoutsos@cc.uoi.gr](mailto:dnoutsos@cc.uoi.gr) (D. Noutsos).

where  $A \in \mathbb{R}^{n,n}$  is an *irreducible dominant Z-matrix* with positive diagonal entries (see [4,18,20]), that is its off-diagonal elements are nonpositive, ( $a_{ij} \leq 0$ ,  $i, j = 1(1)n$ ,  $j \neq i$ ) and  $b \in \mathbb{R}^n$ . Since  $a_{ii} > 0$ ,  $i = 1(1)n$ , we may assume for simplicity that  $a_{ii} = 1$ ,  $i = 1(1)n$ . We consider the usual triangular splitting of  $A$ ,

$$A = I - L - U, \quad (1.2)$$

where  $I$  is the identity matrix and  $L$  and  $U$  are strictly lower and strictly upper triangular, respectively. Then, it is known that the iterative methods of Jacobi and Gauss–Seidel associated with (1.1) converge and by the Stein–Rosenberg theorem [18,20,4] the Gauss–Seidel method is faster than the Jacobi one.

Many researchers have considered preconditioners, applied to system (1.1) that make the associated Jacobi and Gauss–Seidel methods converge asymptotically faster than the original ones. Milaszewicz [16], basing his idea on previous ones (see, e.g., [15,5,9]), considered as a preconditioner  $P_1 \equiv I + S_1$ , the matrix that eliminates the off-diagonal elements of the  $k$ th column of  $A$ . Gunawardena et al. [6] considered as a preconditioner the matrix  $P_2 \equiv I + S_2$ , which eliminates the elements of the first upper diagonal. Kohno et al. [10] extended the main idea in [6]. Recently Li and Sun [13] extended the class of matrices considered in [10] and very recently Hadjidimos et al. [8] extended, generalized and compared the previous works. Modifications of the above preconditioners have introduced and studied by Kotakemori et al. [11] and by Niki et al. [17] (see also Li [12]). The term “preconditioning” is used in Sections 2 and 3 as in the aforementioned works, namely, to reduce the spectral radius of the iteration matrix in order to improve the convergence of the classical iterative methods. However, the same term is often used when the goal is to improve the condition number of  $A$  and hence the convergence of the Conjugate Gradient or other Krylov subspace methods and this is done in Section 4.

This work is organized as follows: In Section 2, we extend Milaszewicz’s and Gunawardena et al.’s preconditioners by giving a family of preconditioners based on the elimination of one element in each row of  $A$ , present the convergence analysis and propose the algorithm that chooses a “good” preconditioner. In Section 3 we generalize the above preconditioners by introducing the idea of eliminating more than one entry per row and perform the corresponding convergence analysis. In Section 4 we study the behavior of the proposed preconditioners when applied to Krylov subspace methods, especially to Conjugate Gradient and to restarted GMRES methods. Finally, in Section 5, numerical examples are presented in support of our theory.

## 2. Extending known preconditioners

Milaszewicz’s preconditioner [16] is based on the elimination of the entries of the  $k$ th column of  $A$ ,  $a_{ik}$ ,  $i = 1(1)n$ ,  $i \neq k$ , while Gunawardena et al.’s one [6] is based on the elimination of the entries of the first upper diagonal  $a_{i,i+1}$ ,  $i = 1(1)n - 1$ . Their common feature is that they eliminate precisely one element of  $A$  in each but one row. If we try to extend Milaszewicz’s preconditioner by eliminating an off-diagonal element of the  $k$ th row we obtain the same convergence results for the Jacobi and the Gauss–Seidel type schemes, since the spectral radii of the corresponding iteration matrices associated with  $A_1 = P_1 A$ , which are reducible, are independent of the off-diagonal elements of the first row of  $A$  (see [8]). This does not happen in the case of Gunawardena et al.’s preconditioner, which eliminates an off-diagonal element of the last row. If we choose the first element of the last row, we introduce a new preconditioner having a cyclic structure and call it *cyclic* preconditioner:  $P_3 \equiv I + S_3$ , where



and

$$\begin{aligned} H &:= (I - L)^{-1}U, \\ H' &:= (I - L - L_s + S_L)^{-1}(U + D_s + U_s - S_U), \\ \tilde{H} &:= (I - D_s - L - L_s + S_L)^{-1}(U + U_s - S_U). \end{aligned} \quad (2.8)$$

The main results of this section can now be stated and proved:

**Theorem 2.1.** (a) *Under the assumptions and the notation so far, the following hold: There exist  $y$  and  $z \in \mathbb{R}^n$ , with  $y \geq 0$  and  $z \geq 0$ , such that*

$$B'y \leq By \quad \text{and} \quad H'z \leq Hz, \quad (2.9)$$

$$\rho(\tilde{B}) \leq \rho(B') \leq \rho(B) < 1, \quad (2.10)$$

$$\rho(\tilde{H}) \leq \rho(H') \leq \rho(H) < 1, \quad (2.11)$$

$$\rho(\tilde{H}) \leq \rho(\tilde{B}), \quad \rho(H') \leq \rho(B'), \quad \rho(H) < \rho(B) < 1. \quad (2.12)$$

(Note: Equalities in (2.12) hold if and only if  $\rho(\tilde{B}) = 0$ .)

(b) *Suppose that  $A$  is irreducible. Then, the matrix  $B$  is also irreducible which implies that the first inequality in (2.9) and the middle inequality in (2.10) are strict.*

**Proof.** (a) (2.9): The matrices  $B$  and  $B'$  are related as follows:

$$B' = I - (I + S)A = I - (I + S)(I - B) = B - S(I - B). \quad (2.13)$$

For the nonnegative Jacobi iteration matrix  $B$  there exists a nonnegative vector  $y$  such that  $By = \rho(B)y$ . Postmultiplying (2.13) by  $y$  we get

$$B'y = (B - S(I - B))y = \rho(B)y - (1 - \rho(B))Sy \leq \rho(B)y = By, \quad (2.14)$$

which proves the first inequality of (2.9).

For the nonnegative Gauss–Seidel iteration matrix  $H$  there exists a nonnegative vector  $z$  such that  $H z = \rho(H)z$ . Using the fact that  $H = (I - L)^{-1}U$  we have that  $(I - L)^{-1}U z = \rho(H)z$  or  $U z = \rho(H)(I - L)z$  or equivalently

$$\rho(H)z = \rho(H)Lz + Uz. \quad (2.15)$$

We rewrite now the matrix  $\tilde{A}$  of (2.2) as follows:

$$\begin{aligned} \tilde{A} &= (I + S)A = (I + S_L + S_U)(I - L - U) \\ &= I - (L - S_L + S_L L + (S_L U)_L + (S_U L)_L) \\ &\quad - (U - S_U + S_U U + (S_L U)_U + (S_U L)_U) = I - \tilde{L} - \tilde{U}, \end{aligned} \quad (2.16)$$

where by  $(Q)_L$  and by  $(Q)_U$  we have denoted the strictly lower and the strictly upper part of the matrix  $Q$ , respectively. So,

$$L = \tilde{L} + S_L - S_L L - (S_L U)_L - (S_U L)_L, \quad U = \tilde{U} + S_U - S_U U - (S_L U)_U - (S_U L)_U. \quad (2.17)$$

By substituting (2.17) in (2.15) we get

$$\rho(H)z = \rho(H)\tilde{L}z + \tilde{U}z + z', \tag{2.18}$$

where

$$z' = \rho(H)(S_L - S_L L - (S_L U)_L - (S_U L)_L)z + (S_U - S_U U - (S_L U)_U - (S_U L)_U)z. \tag{2.19}$$

If  $z' \geq 0$  then from (2.18) we get

$$\rho(H)(I - \tilde{L})z \geq \tilde{U}z \quad \text{or} \quad \rho(H)z \geq (I - \tilde{L})^{-1}\tilde{U}z \quad \text{or} \quad \rho(H)z \geq H'z \tag{2.20}$$

from which the second inequality of (2.9) follows. It remains to prove that  $z' \geq 0$ :

$$\begin{aligned} z' &= \rho(H)S_L(I - L - U)z + \rho(H)(S_L U - (S_L U)_L - (S_U L)_L)z \\ &\quad + S_U(I - L - U)z + (S_U L - (S_L U)_U - (S_U L)_U)z \\ &= S_L(\rho(H)(I - L) - \rho(H)U)z + \rho(H)((S_L U)_U - (S_U L)_L)z \\ &\quad + \frac{1}{\rho(H)} S_U(\rho(H)(I - L) - \rho(H)U)z + ((S_U L)_L - (S_L U)_U)z \\ &= (1 - \rho(H))S_L U z + \frac{1 - \rho(H)}{\rho(H)} S_U U z - (1 - \rho(H))(S_L U)_U z + (1 - \rho(H))(S_U L)_L z \\ &= (1 - \rho(H))(S_L U)_L z + \frac{1 - \rho(H)}{\rho(H)} S_U U z + (1 - \rho(H))(S_U L)_L z \geq 0. \end{aligned} \tag{2.21}$$

(a) (2.10): It is known that (see [4]) a  $Z$ -matrix  $A$  is a nonsingular  $M$ -matrix iff there exists a positive vector  $y (> 0) \in \mathbb{R}^n$  such that  $Ay > 0$ . By taking such a  $y$ , the fact that  $P = I + S \geq 0$  implies  $\tilde{A}y = PAy > 0$ . Consequently,  $\tilde{A}$ , which is a  $Z$ -matrix, is a nonsingular  $M$ -matrix. So, the last two splittings in (2.6) are regular splittings because  $M'^{-1} = I^{-1} = I \geq 0$ ,  $N' \geq 0$  and  $M''^{-1} = (I - D_s)^{-1} \geq 0$ ,  $N'' \geq 0$  and so they are convergent. Since  $M''^{-1} \geq M'^{-1}$ , it is implied (see [19]) that the left inequality in (2.10) is true. For the proof of the middle inequality in (2.10), we recall the first inequality of (2.9) which gives that  $B'y \leq \rho(B)y$ . Then, we apply Lemma 3.3 in Marek and Szyld [14] to get our assertion.

(a) (2.11): To prove the first inequality in (2.11) we use regular splittings of the matrix  $\tilde{A}$ . Specifically, consider the following splittings that define the iteration matrices in (2.8):

$$\tilde{A} = \begin{cases} M - N = (I + S)(I - L) - (I + S)U, \\ M' - N' = (I - L - L_s + S_L) - (D_s + U + U_s - S_U), \\ M'' - N'' = (I - D_s - L - L_s + S_L) - (U + U_s - S_U), \end{cases} \tag{2.22}$$

where we have used the same symbols for the two matrices of each splitting as in the case of (2.6). So, the last two splittings in (2.22) are regular splittings because  $M'^{-1} = (I - L - L_s + S_L)^{-1} = I + (L + L_s - S_L) + \dots + (L + L_s - S_L)^{n-1} \geq 0$ ,  $N' \geq 0$  and  $M''^{-1} = (I - D_s - L - L_s + S_L)^{-1} \geq 0$ ,  $N'' \geq 0$  and so they are convergent. Since  $M''^{-1} \geq M'^{-1}$ , it is implied (see [19]) that the left inequality in (2.11) is true.

To prove the second inequality of (2.11) we consider first that the Jacobi matrix  $B$  is irreducible. For the nonnegative Gauss–Seidel iteration matrix  $H$  there exists a nonnegative vector  $z$  such that

$$Hz = \rho(H)z \quad \text{or} \quad (I - L)^{-1}Uz = \rho(H)z \quad \text{or} \quad (\rho(H)L + U)z = \rho(H)z. \tag{2.23}$$

We observe here that the matrix  $\rho(H)L + U$  has the same structure as the matrix  $B$  and consequently it is also irreducible. So, from the Perron–Frobenius Theorem (see Varga [18]), the eigenvector  $z$  will be a

positive vector. Recalling relation (2.20), the following property holds: There exists a positive vector  $z$  such that  $\rho(H)z \geq H'z$ . Based on this, we can apply Lemma 3.3 in Marek and Szyld [14] to get the second inequality in (2.11). In the case where  $B$  is reducible, we consider a small number  $\varepsilon > 0$  and replace some zeros of  $B$  with  $\varepsilon$  so that the produced matrix  $B(\varepsilon)$  becomes irreducible. Then, for the associated matrices  $H(\varepsilon)$  and  $H'(\varepsilon)$  there holds:  $\rho(H'(\varepsilon)) \leq \rho(H(\varepsilon))$ . Since the spectral radius is a continuous function of the elements of the matrix, the inequality above will also hold in the limit as  $\varepsilon$  tends to zero, which is the second inequality in (2.11).

(a) (2.12): Since  $A$  is a nonsingular  $M$ -matrix, the rightmost inequality is a straightforward implication of the Stein–Rosenberg Theorem as was mentioned before. The other two inequalities in (2.12) are implied directly by the facts that  $\tilde{A}$  is a nonsingular  $M$ -matrix, and the last two pairs of splittings in (2.6) and (2.22), from which the four matrices involved,  $\tilde{H}$ ,  $\tilde{B}$ ,  $H'$ ,  $B'$ , are produced, are regular ones with  $L + L_s + U + U_s - S \geq U + U_s - S_U$  and  $D_s + L + L_s + U + U_s - S \geq D_s + U + U_s - S_U$ . It is noted that if  $\rho(\tilde{B}) = 0$  then  $\rho(\tilde{H}) = 0$  and the matrix  $\tilde{B}$  would be reducible with its *normal form* being a strictly upper triangular matrix. By the Stein–Rosenberg theorem, if  $\rho(\tilde{B}) = \rho(\tilde{H})$  it is either  $\rho(\tilde{B}) = \rho(\tilde{H}) = 0$  or  $\rho(\tilde{B}) = \rho(\tilde{H}) = 1$ . Since  $\rho(\tilde{B}) = 0$ , both spectral radii would be zero. For the second inequality of (2.12) we have that the matrix  $\tilde{B}$  has the same structure as the matrix  $B'$ . So, if  $\rho(\tilde{B}) = 0$ , the matrix  $B'$  would be reducible with its *normal form* being an upper triangular matrix. This means that in the directed graph  $\mathcal{G}(B')$  of the matrix  $B'$ , there is no strongly connected subpath except for identity paths (loops) corresponding to the nonzero diagonal elements. For the matrix  $H'$  we have that  $H' = (I - L - L_s + S_L)^{-1}(D_s + U + U_s - S_U) = (I + (L + L_s - S_L) + \dots + (L + L_s - S_L)^{n-1})(D_s + U + U_s - S_U) = D_s + [U + U_s - S_U + ((L + L_s - S_L) + \dots + (L + L_s - S_L)^{n-1})(D_s + U + U_s - S_U)]$ . The matrix in the brackets is a sum of products of nonnegative parts of  $B'$ . Therefore, if there exists a path in the graph of this matrix, then there exists also such a path, of some order, in the graph  $\mathcal{G}(B')$ . So, if there exists a strongly connected subpath in the graph of the matrix in brackets, then it will also exist such a subpath in  $\mathcal{G}(B')$ . This means that the matrix  $H'$  has also its *normal form* an upper triangular matrix, with its diagonal elements those of  $B'$ . This proves our assertion that  $\rho(B') = \rho(H')$ .

(b) (2.9): Since  $B$  is irreducible, The eigenvector  $y$ , corresponding to  $\rho(B)$ , is positive and according to the steps in the proof of (2.9) in (a) we can see that inequality (2.14) becomes a strict one.

(b) (2.10): From the inequality  $B'y < By$  we get  $B'y < \rho(B)y$ . Now we can apply Lemma 3.3 in Marek and Szyld [14] to get the strict inequality  $\rho(B') < \rho(B)$ .  $\square$

## 2.1. “Good” Jacobi and Gauss–Seidel preconditioners

It was proved that the preconditioned Jacobi and Gauss–Seidel methods converge for each choice of the matrix  $S$  and converge faster than their unpreconditioned counterparts. A question, then, arises: Is there an optimal choice for the matrix  $S$  so that the associated method will be an optimal one and if so how can one choose such a matrix  $S$ ? This question cannot be answered yet and constitutes an open problem. It seems to be difficult since we have to compare the spectral radii of  $(n - 1)^n$  different matrices. Instead, we will try to answer a simpler related question: Is there a “good” choice of the matrix  $S$  such that the associated method will be the best among many others and possibly the optimal one? To find a “good” Jacobi preconditioner or a “good” Gauss–Seidel one we will work using sufficient conditions of convergence rather than necessary and sufficient ones. So, we choose the matrix  $S$  such that to minimize the maximum norm of  $\tilde{B}$  (or of  $\tilde{H}$ ) which constitutes an upper bound for its spectral radius. In the following we give the analysis and the associated algorithm for the Jacobi method only. The analysis for

the Gauss–Seidel method is analogous and straightforward and its associated algorithm is of the same order of complexity. In the last section, we will show by numerical experiments, the efficiency of the “good” Gauss–Seidel method as well.

For the Jacobi iteration matrix,  $\tilde{B}$  to converge a sufficient condition is

$$\rho(\tilde{B}) \leq \|\tilde{B}\|_\infty \iff \rho(\tilde{B}) \leq \max_i \frac{\tilde{l}_i + \tilde{u}_i}{\tilde{d}_i} < 1, \tag{2.24}$$

where

$$\begin{aligned} \tilde{d}_i &= |\tilde{a}_{ii}| = \tilde{a}_{ii} = 1 - a_{ik_i} a_{k_i i}, \\ \tilde{l}_i &= \sum_{j=1}^{i-1} |\tilde{a}_{ij}| = - \sum_{j=1}^{i-1} \tilde{a}_{ij} = a_{ik_i} \sum_{j=1}^{i-1} a_{k_i j} - \sum_{j=1}^{i-1} a_{ij}, \\ \tilde{u}_i &= \sum_{j=i+1}^n |\tilde{a}_{ij}| = - \sum_{j=i+1}^n \tilde{a}_{ij} = a_{ik_i} \sum_{j=i+1}^n a_{k_i j} - \sum_{j=i+1}^n a_{ij}. \end{aligned} \tag{2.25}$$

For each row  $i$ , of the method we propose, we choose  $k_i$  such that all the ratios  $\frac{\tilde{l}_i + \tilde{u}_i}{\tilde{d}_i}$  are minimized and so their maximum will also be minimized. Since the choice of  $k_i$  is not unique,  $S$  is also not unique. We conjecture that since we minimize all the ratios for each row (the row sums of the nonnegative matrix  $\tilde{B}$ ), the new spectral radius will be as small as possible. We call this method the “Good” Jacobi preconditioned method and the associated preconditioner,  $I + S$ , the “Good” Jacobi preconditioner. From (2.25) we get that

$$\tilde{l}_i + \tilde{u}_i = - \sum_{j=1, j \neq i}^n \tilde{a}_{ij} = a_{ik_i} \sum_{j=1, j \neq i}^n a_{k_i j} - \sum_{j=1, j \neq i}^n a_{ij} = s_i + a_{ik_i}(1 - s_{k_i} - a_{k_i i}), \tag{2.26}$$

where  $s_i = -\sum_{j=1, j \neq i}^n a_{ij}$  is the  $i$ th row sum of  $B$ . From the nonnegativity of  $B$  and from the diagonal dominance of  $A$ ,  $0 < s_i < 1$ . So, the ratios in question are given by

$$\frac{\tilde{l}_i + \tilde{u}_i}{\tilde{d}_i} = \frac{s_i + a_{ik_i}(1 - s_{k_i} - a_{k_i i})}{1 - a_{ik_i} a_{k_i i}}. \tag{2.27}$$

To give an efficient algorithm that will choose the indices  $k_i$  and consequently the matrix  $S$  we observe that the Jacobi method requires  $\mathcal{O}(n^2)$  ops per iteration. The same number of operations is required for the multiplication  $(I + S)A$ . So, the cost of the choice of  $S$  must require at most  $\mathcal{O}(n^2)$  ops. First, we compute all the row sums  $s_i$  that require a total number of  $\mathcal{O}(n^2)$  ops. Then, we compute the ratios (2.27) for every  $i$  and every  $k_i$ . The number of ratios is  $(n - 1)n$  and so the number of required operations for each one is  $\mathcal{O}(1)$ . The number of comparisons is also  $\mathcal{O}(n^2)$  and the total cost of the algorithm is  $\mathcal{O}(n^2)$  ops per iteration. This makes it an efficient one. The previous analysis for the cost is illustrated more clearly in the following algorithm written in a pseudocode form.

*Algorithm of “Good” Jacobi Preconditioner*

```
for  $i = 1(1)n$ 
     $s_i = 0$ 
    for  $j = 1(1)i - 1$ 
```

```

         $s_i = s_i - a_{ij}$ 
    endfor
    for  $j = i + 1(1)n$ 
         $s_i = s_i - a_{ij}$ 
    endfor
endfor
for  $i = 1(1)n$ 
     $r = 1$ 
    for  $j = 1(1)n$ 
        if  $j \neq i$  then
             $t = \frac{s_i + a_{ij}(1 - s_j - a_{ji})}{1 - a_{ij}a_{ji}}$ 
            if  $t < r$  then
                 $r = t$ 
                 $k_i = j$ 
            endif
        endif
    endfor
endfor
End of Algorithm

```

We remark that in case there are multiple choices of the matrix  $S$ , this algorithm chooses the one with the smallest value for each  $k_i$ .

### 3. Generalized preconditioners based on multiple elimination

In this section we will generalize and extend our improved method by eliminating two or more off-diagonal elements in each row. So, the matrix  $S$ , introduced in (2.1), will have more than one elements in each row, at exactly the same positions as the elements we want to eliminate. First, we consider that in the  $i$ th row we have to eliminate the elements  $k_j$  and  $l_i$ , where  $k_i < l_i$ . For this we have to compute the elements  $s_{ik_i}$  and  $s_{il_i}$  of the matrix  $S$ . Denoting  $\tilde{A} = (I + S)A$  we have the equations:

$$\begin{aligned} \tilde{a}_{ik_i} = 0 &= a_{ik_i} + s_{ik_i} + s_{il_i}a_{lk_i} & \Leftrightarrow & \quad s_{ik_i} + s_{il_i}a_{lk_i} = -a_{ik_i} \\ \tilde{a}_{il_i} = 0 &= a_{il_i} + s_{ik_i}a_{k_l_i} + s_{il_i} & \Leftrightarrow & \quad s_{ik_i}a_{k_l_i} + s_{il_i} = -a_{il_i} \end{aligned} \quad (3.1)$$

or

$$(s_{ik_i} \quad s_{il_i}) \begin{pmatrix} 1 & a_{k_l_i} \\ a_{l_k_i} & 1 \end{pmatrix} = -(a_{ik_i} \quad a_{il_i}) \Leftrightarrow (s_{ik_i} \quad s_{il_i}) = -(a_{ik_i} \quad a_{il_i}) \begin{pmatrix} 1 & a_{k_l_i} \\ a_{l_k_i} & 1 \end{pmatrix}^{-1}. \quad (3.2)$$

We generalize the above relations by considering that  $m$  elements of the  $i$ th row are to be eliminated. For this we give the following definitions. Let  $\hat{k}_i^T = (k_{i_1} k_{i_2} \dots k_{i_m})$  be a multiindex, where the indices  $k_{i_1} < k_{i_2} < \dots < k_{i_m}$  denote the positions of the elements of row  $i$  to be eliminated. Then, we define



by  $s_{i\hat{k}_i}^T = (s_{ik_{i1}} s_{ik_{i2}} \dots s_{ik_{im}})$  the vector of nonzero off-diagonal elements of the  $i$ th row of  $S$ ,  $a_{i\hat{k}_i}^T = (a_{ik_{i1}} a_{ik_{i2}} \dots a_{ik_{im}})$  the vector of the elements of the  $i$ th row of  $A$  to be eliminated and the matrix

$$A_{\hat{k}_i} = \begin{bmatrix} 1 & a_{k_{i1}k_{i2}} & \dots & a_{k_{i1}k_{im}} \\ a_{k_{i2}k_{i1}} & 1 & \dots & a_{k_{i2}k_{im}} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k_{im}k_{i1}} & a_{k_{im}k_{i2}} & \dots & 1 \end{bmatrix}, \tag{3.3}$$

which consists of the rows and columns of  $A$  indexed by the multiindex  $\hat{k}_i$ . The matrix  $A_{\hat{k}_i}$  is a principal submatrix of the  $M$ -matrix  $A$ . So,  $A_{\hat{k}_i}$  is also an  $M$ -matrix and its inverse  $A_{\hat{k}_i}^{-1}$  is a positive matrix if  $A_{\hat{k}_i}$  is irreducible or a nonnegative one if  $A_{\hat{k}_i}$  is reducible. Therefore, relation (3.2) takes the following generalized form:

$$s_{i\hat{k}_i}^T = -a_{i\hat{k}_i}^T A_{\hat{k}_i}^{-1}. \tag{3.4}$$

After the previous notation and considerations the elements  $\tilde{a}_{ij}$  of  $\tilde{A}$  are as follows:

$$\tilde{a}_{ij} = \begin{cases} a_{ij} - a_{i\hat{k}_i}^T A_{\hat{k}_i}^{-1} a_{\hat{k}_i j} \leq 0, & j \neq i, j \notin \hat{k}_i, \\ 0, & j \in \hat{k}_i, \\ 1 - a_{i\hat{k}_i}^T A_{\hat{k}_i}^{-1} a_{\hat{k}_i i} > 0, & j = i. \end{cases} \tag{3.5}$$

The first inequality in (3.5) is obvious while the last one is to be proved.

We consider the matrix

$$\tilde{A}_{\hat{k}_i} = \begin{pmatrix} 1 & a_{i\hat{k}_i}^T \\ a_{\hat{k}_i i} & A_{\hat{k}_i} \end{pmatrix}. \tag{3.6}$$

The quantity  $1 - a_{i\hat{k}_i}^T A_{\hat{k}_i}^{-1} a_{\hat{k}_i i}$  is the Schur complement of the above matrix  $\tilde{A}_{\hat{k}_i}$ . Since the Schur complement of a nonsingular  $M$ -matrix is also a nonsingular  $M$ -matrix (e.g., see [7, p. 128]), we get that

$$1 - a_{i\hat{k}_i}^T A_{\hat{k}_i}^{-1} a_{\hat{k}_i i} > 0. \tag{3.7}$$

We define the matrix

$$D_s := \text{diag}(a_{1\hat{k}_1}^T A_{\hat{k}_1}^{-1} a_{\hat{k}_1 1}, a_{2\hat{k}_2}^T A_{\hat{k}_2}^{-1} a_{\hat{k}_2 2}, \dots, a_{n\hat{k}_n}^T A_{\hat{k}_n}^{-1} a_{\hat{k}_n n}), \tag{3.8}$$

which is the diagonal part of the matrix  $S(L + U)$ . Using the notation in (2.4) and (2.5), and considering the splittings in (2.6), as well as the associated Jacobi (2.7) and Gauss–Seidel (2.8), we can prove the theorem below, which is the generalization of Theorem 2.1.

**Theorem 3.1.** (a) *Under the assumptions and the notation so far, there hold: There exist  $y$  and  $z \in \mathbb{R}^n$ , with  $y \geq 0$  and  $z \geq 0$ , such that*

$$B'y \leq By \quad \text{and} \quad H'z \leq Hz, \tag{3.9}$$

$$\rho(\tilde{B}) \leq \rho(B') \leq \rho(B) < 1, \tag{3.10}$$

$$\rho(\tilde{H}) \leq \rho(H') \leq \rho(H) < 1, \quad (3.11)$$

$$\rho(\tilde{H}) \leq \rho(\tilde{B}), \quad \rho(H') \leq \rho(B'), \quad \rho(H) < \rho(B) < 1. \quad (3.12)$$

(Note: Equalities in (3.12) hold if and only if  $\rho(\tilde{B}) = 0$ .)

(b) Suppose that  $A$  is irreducible. Then, the matrix  $B$  is also irreducible which implies that the first inequality in (3.9) and the middle inequality in (3.10) are strict.

**Proof.** Using the same notation, the proof of Theorem 3.1 follows step by step that of Theorem 2.1, and so, the proof of the present statement is complete.  $\square$

We remark here that our proposed multiindexed method is the most general method, among many other improved methods, based on elimination techniques. It is a generalization of the block elimination improved method proposed by Alanelli and Hadjidimos [1,2]. They have studied the block Milaszewicz's improved method, which eliminates the elements of the first  $k_1$  columns of  $A$  below the diagonal. This is precisely the multiindexed method with  $\hat{k}_i = (1 \ 2 \ 3 \ \dots \ i - 1)$ ,  $i = 1(1)k_1$  and  $\hat{k}_i = (1 \ 2 \ 3 \ \dots \ k_1)$ , otherwise.

As regards the cost of the present method we can observe that, for the construction of the matrix  $S$ , we have to solve an  $m \times m$  linear system for each row which has a cost of  $\mathcal{O}(m^3n)$  ops. The cost of the matrix–matrix product  $(I + S)A$  is  $\mathcal{O}(mn^2)$  ops. Then follows the standard iterative process of Jacobi or Gauss–Seidel method which has a cost of  $\mathcal{O}(n^2)$  ops per iteration. So, to obtain an efficient algorithm, the number  $m$  must be chosen very small and independent of the dimension  $n$ . A question which arises is: Which is the best choice of  $m$ ? Observe that by increasing  $m$  what is gained in number of iterations is lost in the construction of  $S$  and in the matrix–matrix product. Therefore there must be a golden section, which also depends on the matrix  $A$ . Another question which arises is: How can one choose the multiindices  $\hat{k}_i$ 's? Since this is difficult to answer we follow the idea we did in the “Good” Jacobi or Gauss–Seidel algorithms, where we provided search algorithms by using sufficient criteria instead of sufficient and necessary ones. In the case of  $m = 1$ , the cost of the searching algorithm is  $\mathcal{O}(n^2)$ . If we take  $m = 2$ , we have to do a double searching per row, so the cost increases to  $\mathcal{O}(n^3)$  and the algorithm becomes non-efficient. If the value of  $m$  increases further, the power of  $n$  in the cost increases too. So, the only efficient searching algorithm is the one where we search along one component of the multiindex, keeping the others fixed. As we will see in the numerical examples, in many cases, by keeping the multiindices fixed, the multiindexed preconditioner is better than the “Good” algorithm.

#### 4. Eliminated preconditioners for Krylov subspace methods

In this section we study the behavior of the proposed preconditioners when they are applied to Krylov subspace methods. First we study the conjugate gradient (CG) method when the matrix  $A$  is a real symmetric positive definite matrix. It is well known that the convergence theory of the preconditioned conjugate gradient (PCG) method holds if the preconditioner is also a symmetric and positive definite matrix. This means that, to construct our preconditioner  $I + S$ , we have to choose the entries to be eliminated in symmetric positions. We will study here the behavior of the multiindexed preconditioner, where we have fixed the multiindices by taking  $\hat{k}_1 = (2)$ ,  $\hat{k}_i = (i - 1 \ i + 1)^T$ ,  $i = 2(1)n - 1$  and  $\hat{k}_n = (n - 1)$ . In other words we eliminate the elements of the first upper and lower diagonals. So, the eliminated entries

are in symmetric positions. The study of other multiindexed preconditioners having entries in symmetric positions, could be straightforwardly generalized. First, we state and prove the following lemma which is useful for the convergence properties below:

**Lemma 4.1.** *Let  $A$  be a nonsingular  $M$ -matrix with  $a_{ii} = 1$ . Let also  $\tilde{A} = (I + S)A$  be the preconditioned matrix with  $I + S$  being the multiindexed preconditioner defined by  $\hat{k}_1 = (2)$ ,  $\hat{k}_i = (i - 1 \ i + 1)^T$ ,  $i = 2(1)n - 1$  and  $\hat{k}_n = (n - 1)$ . Then all the principal minors of  $\tilde{A}$  are positive and less than or equal to the associated principal minors of  $A$ .*

$$0 < \det(\tilde{A}_{i_1 i_2 \dots i_k}) \leq \det(A_{i_1 i_2 \dots i_k}), \tag{4.1}$$

where by the indices  $i_1 i_2 \dots i_k$  we denote the rows and the associated columns which form the principal submatrix.

**Proof.** The first inequality holds since all the principal submatrices of  $\tilde{A}$  and  $A$  are  $M$ -matrices. We will prove the second inequality by induction. First we prove that for all the principal submatrices of the preconditioner  $I + S$ , there holds

$$0 < \det((I + S)_{i_1 i_2 \dots i_k}) \leq 1.$$

Since each principal submatrix of order  $k$  can be put in the first  $k$  rows and columns by a permutation transformation of the matrix  $I + S$ , we will prove, without loss of generality, the above inequalities only for the upper left principal minors, i.e.,

$$0 < \det((I + S)_k) \leq 1,$$

where  $(I + S)_k$  is the principal submatrix, consisting of the first  $k$  rows and columns, of the matrix  $(I + S)$ . We denote also by  $A_k$  and by  $\tilde{A}_k$  the principal submatrices consisting of the first  $k$  rows and columns of the matrices  $A$  and  $\tilde{A}$ , respectively. We define now the  $k \times k$  matrix

$$A'_k = A_k + S_{k,n-k} A_{n-k,k},$$

where  $S_{k,n-k}$  is the  $k \times (n - k)$  submatrix of  $I + S$  consisting of the first  $k$  rows and the last  $n - k$  columns, while  $A_{n-k,k}$  is the  $(n - k) \times k$  submatrix of  $A$  consisting from the last  $n - k$  rows and the first  $k$  columns. From the above definition it is obvious that  $A'_k$  is the principal submatrix consisting from the first  $k$  rows and columns of the matrix  $(I + S'_{k,n-k})A$ , where  $S'_{k,n-k} = \begin{pmatrix} 0 & S_{k,n-k} \\ 0 & 0 \end{pmatrix}$ . Since  $I + S'_{k,n-k}$  is a nonnegative matrix,  $A$  is an  $M$ -matrix and  $(I + S'_{k,n-k})A$  is a  $Z$ -matrix we get that  $(I + S'_{k,n-k})A$  is also an  $M$ -matrix. Therefore its principal submatrix  $A'_k$  is an  $M$ -matrix.

From the symmetric structure of the matrix  $I + S$  and since  $s_{ij} \geq 0$ ,  $i, j = 1, 2, \dots, n$ ,  $i \neq j$ , it is easily checked that if we use the Gauss elimination process step by step, the diagonal elements do not increase. So, all the first principal minors would be less than or equal to 1. We will prove in the sequel that these principal minors could not be less than 0. We will prove inequality (4.13) for the first principal minors, i.e.

$$0 < \det(\tilde{A}_k) \leq \det(A_k), \quad k = 1, 2, \dots, n.$$

The proof for any other determinant is given directly by a similar analysis which becomes a little more complicated. From the notation above, we have that

$$\tilde{A}_k = (I + S)_k A_k + S_{k,n-k} A_{n-k,k} = (I + S)_k A_k + \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix},$$

since the matrix  $S_{k,n-k}$  has only one entry different from zero in the position  $(k, k + 1)$ . We have denoted by  $a_{k+1;1:k}^T$  the row vector consisting of the first  $k$  entries of the  $(k + 1)$ st row of the matrix  $A$  and by  $0_{k-1,k}$  the  $(k - 1) \times k$  zero matrix. The above equality is written as

$$(I + S)_k A_k = \tilde{A}_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} = \left( I_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} \tilde{A}_k^{-1} \right) \tilde{A}_k. \quad (4.2)$$

Since  $\tilde{A}$  is an M-matrix,  $\tilde{A}_k$  is also an M-matrix and consequently  $\tilde{A}_k^{-1}$  is a nonnegative matrix. So,  $-s_{k,k+1} a_{k+1;1:k}^T \tilde{A}_k^{-1}$  is a nonnegative row vector. Therefore, the matrix

$$\left( I_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} \tilde{A}_k^{-1} \right)$$

is a lower triangular matrix with all the diagonal entries equal to one except the last one which is greater than or equal to 1. This means that

$$\det \left( \left( I_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} \tilde{A}_k^{-1} \right) \right) \geq 1.$$

By taking determinants of the matrices in (4.2) we get

$$\det((I + S)_k) \det(A_k) = \det \left( I_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} \tilde{A}_k^{-1} \right) \det(\tilde{A}_k).$$

Since  $A_k$  and  $\tilde{A}_k$  are both M-matrices, the associated determinants are both positive. So, from the above equality it has been proved that  $\det((I + S)_k) > 0$ . From the fact that  $\det((I + S)_k) \leq 1$  and

$$\det \left( I_k - \begin{pmatrix} 0_{k-1,k} \\ s_{k,k+1} a_{k+1;1:k}^T \end{pmatrix} \tilde{A}_k^{-1} \right) \geq 1$$

the above equality gives us the inequality we had to prove:

$$\det(\tilde{A}_k) \leq \det(A_k). \quad \square$$

#### 4.1. Preconditioned conjugate gradient method

**Theorem 4.1.** Let  $A$  be an irreducible and symmetric positive definite M-matrix with  $a_{ii} = 1$  and with its eigenvalues  $\lambda_i \in [a, b]$ ,  $i = 1, 2, \dots, n$ . Let also  $\tilde{A} = (I + S)A$  be the preconditioned matrix with  $I + S$  being the multiindexed preconditioner defined by  $\hat{k}_1 = (2)$ ,  $\hat{k}_i = (i - 1 \ i + 1)^T$ ,  $i = 2(1)n - 1$

and  $\hat{k}_n = n - 1$ , which is symmetric and positive definite. Then the eigenvalues  $\tilde{\lambda}_i, i = 1, 2, \dots, n$ , of the preconditioned matrix  $\tilde{A}$  are clustered in the interval  $[\tilde{a}, \tilde{b}]$  except for a possible small number of them that are greater than  $\tilde{b}$ , where  $\tilde{a} > a$  and  $\tilde{b} < b$ . Consequently, the Preconditioned Conjugate Gradient method converges faster than the unpreconditioned one.

**Proof.** It is well known that the characteristic polynomial of a matrix  $A$  is given by

$$P_n(\lambda) = \lambda^n - \sum_{i=1}^n a_{ii} \lambda^{n-1} + \dots + (-1)^k \sum_{i_1 i_2 \dots i_k} \det(A_{i_1 i_2 \dots i_k}) \lambda^{n-k} + \dots + (-1)^n \det(A). \quad (4.3)$$

In terms of the eigenvalues the same polynomial is given by

$$P_n(\lambda) = \lambda^n - \left( \sum_{i=1}^n \lambda_i \right) \lambda^{n-1} + \dots + (-1)^k \left( \sum_{i_1 i_2 \dots i_k} \lambda_{i_1} \lambda_{i_2} \dots \lambda_{i_k} \right) \lambda^{n-k} + \dots + (-1)^n \lambda_1 \lambda_2 \dots \lambda_n. \quad (4.4)$$

So,

$$\sum_{i_1 i_2 \dots i_k} \lambda_{i_1} \lambda_{i_2} \dots \lambda_{i_k} = \sum_{i_1 i_2 \dots i_k} \det(A_{i_1 i_2 \dots i_k}), \quad k = 1, 2, \dots, n, \quad (4.5)$$

where the sums have been taken over all the combinations of  $n$  over  $k$ . By taking the same arguments of the preconditioned matrix  $\tilde{A}$  we take

$$\sum_{i_1 i_2 \dots i_k} \tilde{\lambda}_{i_1} \tilde{\lambda}_{i_2} \dots \tilde{\lambda}_{i_k} = \sum_{i_1 i_2 \dots i_k} \det(\tilde{A}_{i_1 i_2 \dots i_k}), \quad k = 1, 2, \dots, n. \quad (4.6)$$

From Lemma 4.1 we have that

$$0 < \det(\tilde{A}_{i_1 i_2 \dots i_k}) \leq \det(A_{i_1 i_2 \dots i_k}), \quad \forall i_1 i_2 \dots i_k, \quad k = 1, 2, \dots, n. \quad (4.7)$$

This inequality and equalities (4.5), (4.6) give us the main inequality which compares the eigenvalues of the preconditioned matrix with those of the unpreconditioned one:

$$\sum_{i_1 i_2 \dots i_k} \tilde{\lambda}_{i_1} \tilde{\lambda}_{i_2} \dots \tilde{\lambda}_{i_k} \leq \sum_{i_1 i_2 \dots i_k} \lambda_{i_1} \lambda_{i_2} \dots \lambda_{i_k}, \quad k = 1, 2, \dots, n. \quad (4.8)$$

We observe that if for at least one combination  $i_1 i_2 \dots i_k$  inequality (4.7) is strict, then so is inequality (4.20). It is easily seen that there are values of  $k$  such that an inequality exists. Obvious values of such  $k$ 's are  $k = 1, k = 2$  or  $k = n$ .

On the other hand, from the inequalities (3.10) of Theorem 3.1 and from the irreducibility of  $A$  we have that  $\rho(B') < \rho(B)$  or  $\rho(I - \tilde{A}) < \rho(I - A)$  which means that

$$\tilde{\lambda}_{\min} > \lambda_{\min}. \quad (4.9)$$

Inequalities (4.20) and (4.21) imply that, when we apply the improved process to the matrix  $A$ , although the minimum eigenvalue increases, some sums of products of them do not increase while some others strictly decrease. This observation and the strong numerical evidence suggest that, although the new largest eigenvalue might be outside the interval  $[a, b]$ , the majority of the remaining eigenvalues must lie (are clustered) in a smaller interval  $[\tilde{a}, \tilde{b}] \subset [a, b]$ . By the known Axelsson's Theorem [3] concerning the convergence of the Conjugate Gradient method, we have that the Preconditioned Conjugate Gradient method converges faster than the unpreconditioned one.  $\square$

We remark that we could have had the same result as in Theorem 4.1 if we had proved the above inequalities only for exact values of  $k$  ( $k = 1, 2, n$ ) for which the proof could have been very easy. We think that this general proof makes our result stronger. It is also remarked that the same result could be proved, in the more general case of the choice of the multiindex. This will be shown by experiments, in the next section, when we solve the Helmholtz equation.

#### 4.2. Preconditioned GMRES method

**Theorem 4.2.** *Let  $A$  be an irreducible  $M$ -matrix with  $a_{ii} = 1$  and let its eigenvalues  $\lambda_i \in \mathcal{S}$ ,  $i = 1, 2, \dots, n$ , where  $\mathcal{S}$  is a bounded convex region in the complex plane with  $a = \min_{x \in \mathcal{S}} \operatorname{Re}(x) > 0$ . Let also  $\tilde{A} = (I + S)A$  be the preconditioned matrix with  $I + S$  being the multiindexed preconditioner defined by  $\hat{k}_1 = (2)$ ,  $\hat{k}_i = (i - 1 \ i + 1)^T$ ,  $i = 2(1)n - 1$  and  $\hat{k}_n = n - 1$ . Then, under the assumption that the matrices  $V$  and  $\tilde{V}$  of the eigenvectors of  $A$  and of  $\tilde{A}$ , respectively have both small enough euclidian condition number, the eigenvalues  $\tilde{\lambda}_i$ ,  $i = 1, 2, \dots, n$ , of the preconditioned matrix  $\tilde{A}$  are clustered in a convex region  $\tilde{\mathcal{S}} \subset \mathcal{S}$  with  $\tilde{a} = \min_{x \in \tilde{\mathcal{S}}} \operatorname{Re}(x) > a$ , except for a possible small number of them that are outside  $\tilde{\mathcal{S}}$ , but with real parts no less than  $\tilde{a}$ . Consequently, the preconditioned GMRES method converges faster than the unpreconditioned one.*

**Proof.** The proof is straightforward and the same as the one of Theorem 4.1. All relations (4.2)–(4.21) hold with the exception that now  $\lambda_i, \tilde{\lambda}_i$ ,  $i = 1, 2, \dots, n$ , are complex numbers with positive real parts. Since  $B'$  and  $B$  are nonnegative matrices,  $\lambda_{\min}$  and  $\tilde{\lambda}_{\min}$  are real positive numbers with (4.21) holding. It is obvious that the sums of products in (4.20) decrease when the moduli of the eigenvalues decrease. Since the smaller of them increase, most of the larger moduli should decrease. So, the eigenvalues are clustered in a shorter region  $\tilde{\mathcal{S}}$  which belongs in  $\mathcal{S}$ , except for a small number of eigenvalues that lie outside  $\tilde{\mathcal{S}}$ . Let  $c$  be the modulus of the center of the region in which most of the eigenvalues lie and  $s$  be the radius of the same region. It is known that the convergence rate of GMRES method depends on the ratio  $\frac{s}{c}$ , in the case where the euclidian condition number of the matrix of the eigenvectors is small enough. For the eigenvalues that lie outside, a constant number of GMRES additional iterations are performed. From the analysis above, we have no information for the centers of the regions. Since the smallest moduli increase and the largest ones decrease, the centers should be almost the same. For the radii, it is obvious that the radius of  $\tilde{\mathcal{S}}$  is smaller than the one of  $\mathcal{S}$ , due to the better clustering. So, the preconditioned GMRES method converges faster than the unpreconditioned one, and the proof is complete.  $\square$

The same remarks, as before, could be stated here. In the next section we show the validity of the above result by solving the Convection–Diffusion equation.

### 5. Numerical examples

For 10,000 randomly generated nonsingular  $M$ -matrices for  $n = 10, 20$  and  $50$  we have determined the spectral radii of the iteration matrices of all the classical methods mentioned previously. Below, we present two tables for the Jacobi and Gauss–Seidel methods, respectively (Tables 1 and 2). Each number in the tables represents the percentage of the case worked out, the method in the first column is better than the one in the head of the table, where  $M, G, C, \text{“G”}$  and  $M_2$  denote a method with the Milaszewicz’s, the Gunawardena et al.’s, the Cyclic, the “Good” and the multiindexed, with  $m = 2$ , preconditioners, respectively. For the multiindexed preconditioner we have fixed the multiindices by taking  $\hat{k}_1 = (2\ n)^T$ ,  $\hat{k}_i = (i - 1\ i + 1)^T$ ,  $i = 2(1)n - 1$  and  $\hat{k}_n = (1\ n - 1)^T$ . In other words we have eliminated the elements corresponding to the Cyclic preconditioner and the symmetrically placed elements.

We remark that for the Jacobi method and for  $n$  large enough, the multiindexed improved method is 100% better than all the others. The “Good” preconditioner is 100% better than the remaining ones, the cyclic preconditioner is 100% better than that of Gunawardena et al.’s, while the last two preconditioners tend to be equivalent, regarding their performance, to the one of Milaszewicz’s, as  $n$  increases. For the Gauss–Seidel method we can see that the “Good” preconditioner is better than all the others, then follow the multiindexed, the cyclic, the Gunawardena et al.’s and finally the Milaszewicz’s preconditioner. At this point, another question is raised: Is the “Good” preconditioner indeed better than the multiindexed one? The answer is *no*! It depends on the choice of the multiindices. In this example we have chosen one element over the diagonal and one under it. We observe that the elimination of the over-diagonal

Table 1  
Jacobi method

	$n = 10$				$n = 20$				$n = 50$			
	$M$	$G$	$C$	“G”	$M$	$G$	$C$	“G”	$M$	$G$	$C$	“G”
$M$		60.88	48.8	6.19		58.01	48.94	1.35		54.61	48.98	0
$G$	39.12		0	1.05	41.99		0	0.06	45.39		0	0
$C$	51.2	100		2.88	51.06	100		0.19	51.02	100		0
“G”	93.81	98.95	97.12		98.65	99.94	99.81		100	100	100	
$M_2$	99.83	100	100	97.23	99.99	100	100	99.42	100	100	100	99.96

Table 2  
Gauss–Seidel method

	$n = 10$				$n = 20$				$n = 50$			
	$M$	$G$	$C$	“G”	$M$	$G$	$C$	“G”	$M$	$G$	$C$	“G”
$M$		2.54	0.44	1.08		0.1	0.02	0		0	0	0
$G$	97.46		0	36.62	99.9		0	25.49	100		0	5.11
$C$	99.56	100		56.63	99.98	100		36.34	100	100		7.13
“G”	98.92	63.38	43.37		100	74.51	63.66		100	94.89	92.87	
$M_2$	99.75	100	100	68.39	100	100	100	44.8	100	100	100	8.95

Table 3

Preconditioned conjugate gradient method, for two-dimension Helmholtz equation

$n$	$I$	$S_2$	$S_4$
16	51	45	27
32	100	86	52
64	191	163	100
128	374	306	194

elements play the most important role for the Gauss–Seidel method than that of the under-diagonal ones. Since we have chosen one fixed over-diagonal element per row, for the multiindexed preconditioner, while for the “Good” one we have searched one element per row (we guess that most of them are over-diagonal elements), the last preconditioner becomes better than the first one. It has been checked that the multiindexed preconditioner ( $m = 2$ ) is 100% better than all the others for the Gauss–Seidel method if we choose both elements to be over-diagonal ones for each row. So, the numerical examples confirm the theoretical results for the proposed improved techniques.

For the Conjugate Gradient method we solved the two-dimensional Helmholtz equation

$$-\Delta u + u = f$$

in the unit square  $\Omega$  with Dirichlet boundary conditions. We have discretized the unit square by taking the  $n \times n$  uniform grid, then we approximated the above problem by finite differences. We solved the  $n^2 \times n^2$  irreducible positive definite system yielded by the CG method, by the PCG method, with improved two-indexed preconditioner  $S_2$ , as was described in Theorem 4.1, and by the PCG method, with improved four-indexed preconditioner  $S_4$  which is based on the elimination of the first and  $n$ th diagonals over and under the main diagonal. As a stopping criterion we chose  $\|r\|_2 \leq 0.5 \times 10^{-6}$ , where  $r$  is the residual vector. In Table 3 the numbers of required iterations, for the CG and for the two PCG methods, are given for various values of  $n$ . The efficiency of the proposed method is clearly shown.

For the GMRES method we solved the two-dimensional convection–diffusion equation

$$-\Delta u + \frac{\partial u}{\partial x} + 2 \frac{\partial u}{\partial y} = f$$

in the unit square  $\Omega$  with Dirichlet boundary conditions. We followed the same technique to approximate the above problem. Then, we solved the  $n^2 \times n^2$  irreducible  $M$ -system yielded by the GMRES method, by the PGMRES( $S_2$ ) method and by the PGMRES( $S_4$ ) one. The last two PGMRES methods are based on the improved two-indexed  $S_2$  and four-indexed  $S_4$  preconditioners, respectively, as they were described previously. The same stopping criterion was chosen and the restarted technique of GMRES was run. In the first part of Table 4 the number of  $m = n$  iterations was chosen to restart the method while in the second,  $m = \frac{n}{2}$ . For each method, a pair of integers is given, where the first one corresponds to the number of completed outer iterations while the second, to the number of the inner iterations, required in the last uncompleted outer one. We also counted the number of the required operations in each case. So, in Table 4 we give them, for PGMRES methods, in percentages with respect to the GMRES method, for each case.



Table 4  
Restarted preconditioned GMRES method, for two-dimension convection–diffusion equation

$m = n$									
$n$	GMRES		PGMRES( $S_2$ )		ops%	PGMRES( $S_4$ )		ops%	
4	6	0	4	0	93.95	2	0	55.71	
8	5	3	3	7	97.42	2	0	59.37	
16	4	9	3	5	92.10	1	14	59.75	
32	4	1	2	25	78.67	1	26	56.15	
64	3	42	2	28	71.96	1	45	53.07	
128	3	83	2	35	64.77	1	65	41.75	
$m = n/2$									
4	17	0	10	0	86.62	6	2	70.61	
8	14	3	9	3	94.14	5	3	66.23	
16	12	5	7	6	83.08	4	6	59.91	
32	12	4	7	11	79.03	4	2	48.96	
64	11	20	7	17	75.89	3	29	43.78	
128	10	51	7	1	72.63	3	54	41.55	

## Acknowledgements

Warm thanks are due to the referees who made useful suggestions and provided some more references that helped us to improve the quality of the paper.

## References

- [1] M. Alanelli, Block elementary block elimination preconditioners for the numerical solution of non-singular and singular linear systems, Master Dissertation (in Greek), Department of Mathematics, University of Crete, Heraklion, Greece, 2001.
- [2] M. Alanelli, A. Hadjidimos, Block gauss elimination followed by a classical iterative method for the solution of linear systems, JCAM 163 (2004) 381–400.
- [3] O. Axelsson, G. Lindskog, On the rate of convergence of the preconditioned conjugate gradient method, Numer. Math. 48 (1986) 499–523.
- [4] A. Berman, R.J. Plemmons, Nonnegative Matrices in the Mathematical Sciences, Classics in Applied Mathematics, SIAM, Philadelphia, 1994.
- [5] R.E. Funderlic, R.J. Plemmons, LU decomposition of M-matrices by elimination without pivoting, LAA 41 (1981) 99–110.
- [6] A.D. Gunawardena, S.K. Jain, L. Snyder, Modified iterative methods for consistent linear systems, LAA 154–156 (1991) 123–143.
- [7] R.A. Horn, C.R. Jonson, Topics in Matrix Analysis, Cambridge University Press, Cambridge, 1991 (Also: second ed. Revised and Expanded, Springer, Berlin, 2000.).
- [8] A. Hadjidimos, D. Noutsos, M. Tzoumas, More on modifications and improvements of classical iterative schemes for Z-matrices, LAA 364 (2003) 253–279.
- [9] M.L. Juncosa, T.W. Mulliken, On the increase of convergence rates of relaxation procedures for elliptic partial differential equations, J. Assoc. Comput. Math. 7 (1960) 29–36.

- [10] T. Kohno, H. Kotakemori, H. Niki, M. Usui, Improving the Gauss–Seidel Method for Z-Matrices, *LAA* 267 (1997) 113–123.
- [11] H. Kotakemori, K. Harada, M. Morimoto, H. Niki, A comparison theorem for the iterative method with the preconditioner (I+Smax), *JCAM* 145 (2002) 373–378.
- [12] W. Li, Comparison results for solving preconditioned linear systems, *JCAM* 176 (2005) 319–329.
- [13] W. Li, W. Sun, Modified Gauss–Seidel type methods and Jacobi type methods for Z-matrices, *LAA* 317 (2000) 227–240.
- [14] I. Marek, D.B. Szyld, Comparison theorems for weak splittings of bounded operators, *Numer. Math.* 58 (1990) 387–397.
- [15] J.P. Milaszewicz, On modified Jacobi linear operators, *LAA* 51 (1983) 127–136.
- [16] J.P. Milaszewicz, Improving Jacobi and Gauss–Seidel iterations, *LAA* 93 (1987) 161–170.
- [17] H. Niki, K. Harada, M. Morimoto, M. Sakakihara, The survey of preconditioners used for accelerating the rate of convergence in the Gauss–Seidel method, *JCAM* 164–165 (2004) 587–600.
- [18] R.S. Varga, *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1962 (Also: second ed., Revised and Expanded, Springer, Berlin, 2000.).
- [19] Z. Woźnicki, Nonnegative splitting theory, *Jap. J. Ind. Appl. Math.* 11 (1994) 289–342.
- [20] D.M. Young, *Iterative Solution of Large Linear Systems*, Academic Press, New York, 1971.