

On Discussion of SOR Method for Solving Sylvester Equation

Seifedine Kadry and Zbigniew Woznicki
 Beirut-Lebanon, Allenby Street, P.O.Box: 110-435, Lebanon

Abstract: The main subject of the study is to describe the Successive Over Relaxation (SOR) method for solving the Sylvester equation $AX-XB = C$, derived by using two different separation models of boundary-value problems. As is demonstrated in several test problems, the proposed method seems to be very efficient and strongly competitive to Krylov-subspace techniques.

Key words: Successive Over Relaxation (SOR), sylvester equation

INTRODUCTION

Sylvester equations have numerous applications in control theory, signal processing, filtering, model reduction, image restoration, decoupling techniques for ordinary and partial differential equations, implementation of implicit numerical methods for ordinary differential equations and block-diagonalization of matrices (Aliev and Larin, 1998; Enright, 1978; Epton, 1980; Calvetti and Reichel, 1996; Dieci *et al.*, 1988; Golub and Van Loan, 1996; Sima, 1996).

The aim of this paper is to discuss the iterative method, introduced recently by the authors and called as SOR-like method (Woznicki and Kadry, 2003) for solving the matrix equation

$$AX-XB = C \quad (1)$$

Where $A \in \mathbb{R}^{m \times m}$, $B \in \mathbb{R}^{n \times n}$ and C , $X \in \mathbb{R}^{m \times n}$. The above equation represents the Sylvester equation, playing an essential role in control theory and when $B = -A^T$, it reduces to the well known Lyapunov equation.

The matrix Eq. 1 possesses a unique solution if and only if the matrices A and B have no common eigenvalues (Gantmacher, 1966) and by means of the Kronecker transformation, it can be equivalently written as a large linear system of the following form

$$Gx = c \quad (2)$$

with an $mn \times mn$ matrix G .

In the literature, several methods have been proposed for solving Eq. 1. A survey of solution methods, properties and applications of the Sylvester equation in control theory is presented by Datta (2004).

In the following section, the simple algorithm of SOR-like method for solving Eq. 1 is presented. The

results of numerical experiments with the solution of Sylvester equations, obtained by two different separation models of elliptic partial differential equations are reported in Section 3 for the examples taken from Hu and Reicher (1992). One of them, called the separation model A, is that used in (Hu and Reicher, 1992; Kadry, 2003) and the second, called as the separation model B, is a new one.

MATERIALS AND METHODS

Assuming that the matrix A is defined by the following decomposition:

$$A = K-L-U \quad (3)$$

Where K , L and U are nonsingular diagonal, strictly lower triangular and strictly upper triangular parts of A ; Eq. (1) can be rewritten as

$$KX = LX+UX+XB+C \quad (4)$$

Or equivalently

$$X = K^{-1} \{LX+UX+XB+C\} \quad (5)$$

We assume that the iterative scheme has the following form

$$X^{(t)} = K^{-1} \{LX^{(t)} + UX^{(t-1)} + X^{(t-1,t)}B + C\} \quad (6)$$

For $t = 1, 2, \dots$ where the term $X^{(t-1,t)}$ means that for computing the product XB , the entries of $X^{(t-1)}$ and $X^{(t)}$ are used in the current iteration t . For the acceleration of convergence in the above iterative scheme, the over-relaxation procedure can be used as follows

$$X^{(t)} = wK^{-1} \{ LX^{(t)} + UX^{(t-1)} + X^{(t-1,t)}B + C \} - (w-1)X^{(t-1)}, \quad t = 1, 2, \dots \quad (7)$$

Or written equivalently

$$X^{(t)} = [I - wK^{-1}L]^{-1} \left\{ \begin{aligned} &[(1-w)I + wK^{-1}U]X^{(t-1)} \\ &+ wK^{-1}[X^{(t-1,t)}B + C] \end{aligned} \right\}, \quad t = 1, 2, \dots \quad (8)$$

Since the exact solution X satisfies the above equation, then with the error solution matrix $E^{(t)} = X - X^{(t)}$, we have

$$E^{(t)} = [I - wK^{-1}L]^{-1} \left\{ \begin{aligned} &[(1-w)I + wK^{-1}U]E^{(t-1)} \\ &+ wK^{-1}E^{(t-1,t)}B \end{aligned} \right\}, \quad t = 1, 2, \dots \quad (9)$$

When B would be the null matrix, we have the explicit form of the iteration matrix

$$\tau_{B=0} = [I - wK^{-1}L]^{-1} [(1-w)I + wK^{-1}U] \quad (10)$$

and the iterative scheme (7), representing the classical algorithm of the SOR method, is convergent for an arbitrary starting matrix $X^{(0)}$ if and only if the spectral radius of the iteration matrix $\tau_{B=0}$ is less than unity; and the error solution matrix $E^{(t)}$ can be expressed explicitly in dependence on $E^{(0)}$, that is,

$$E^{(t)} = \tau_{B=0}^t E^{(0)} \quad (11)$$

In the case when $B \neq 0$, we see that a similar relation for the error solution matrix can not be derived from Eq. 9. However, we can assume that there exists an implicit iteration matrix $\tau_{B=0}$ which form can not be expressed explicitly but we are able to compute its spectral radius according to the following equation derived from Eq. 7

$$wK^{-1} \{ LY^{(t)} + UY^{(t-1)} + Y^{(t-1,t)}B \} - (w-1)Y^{(t-1)} = \Lambda Y^{(t)}, \quad t = 1, 2, \dots \quad (12)$$

Where Λ is an eigenvalue of the implicit iteration matrix $\tau_{B=0}$ and the matrix Y play a role of an "eigenvector". When Λ is a real eigenvalue, the above equation represents the algorithm of the power method providing us the spectral radius $\rho(\tau_{B=0}) = |\Lambda|$.

As is observed in numerical experiments, the value of $\rho(\tau_{B=0})$ is frequently minimized for $0 < w < 1$.

The detailed analysis of convergence properties of the presented method and the estimation of the value of w_{best} are the subject of actual studies. In the analysis of the reliability of iterative solutions, it is convenient to consider the (true) error matrix

$$E^{(t)} = X - X^{(t)} \quad (13)$$

The inner (or pseudo-residual) error matrix

$$\delta^{(t)} = X^{(t)} - X^{(t-1)} \quad (14)$$

The relative inner error matrix computed entry-wise

$$\bar{\delta}_{j,k}^{(t)} = \frac{X_{j,k}^{(t)} - X_{j,k}^{(t-1)}}{X_{j,k}^{(t)}} \quad (15)$$

And the residual matrix

$$R^{(t)} = AX^{(t)} - X^{(t)}B - C \quad (16)$$

where X is assumed as the "exact" solution matrix.

Since the above quantities are matrices, the *Frobenius norm*

$$\|A\|_F = \left[\sum_{j=1}^m \sum_{k=1}^n |a_{j,k}|^2 \right]^{\frac{1}{2}} \quad (17)$$

is an important matrix norm for us and for needs of comparison of obtained results with those given in (Hu and Reicher, 1992), we shall use the *maximum matrix norm* defined as follows

$$\|A\|_{\max} = \max_{\substack{1 \leq j \leq m \\ 1 \leq k \leq n}} |a_{j,k}| \quad (18)$$

Most recent iterative methods terminate when the residual $R^{(t)}$ is sufficiently small and the termination test

$$\frac{\|R^{(t)}\|_F}{\|R^{(0)}\|_F} \leq \varepsilon \quad (19)$$

Called usually as a *relative residual norm*, is most commonly used criterion in Krylov subspace algorithms. However, it seems that in the case of iterative methods based on a matrix splitting, the termination test

$$\|\bar{\delta}\|_{\max} \leq \varepsilon \quad (20)$$

can be practically considered as the most useful stopping criterion independent on an used initial guess (Woznicki, 2001).

RESULTS AND DISCUSSION

We shall illustrate the numerical performance of the iterative scheme presented in this paper in some examples taken from Hu and Reicher (1992).

Let $\Omega = \{(x,y) \in \mathbb{R}^2: 0 < x < 1, 0 < y < 1\}$ and let $\partial\Omega$ denotes the boundary of Ω . We shall solve the following boundary-value problem for the two-dimensional separable model of convection-diffusion equation:

$$\begin{cases} -\Delta u + 2p_1 u_x + 2p_2 u_y - 2p_3 u = F & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases} \quad (21)$$

The parameters p_1, p_2 and p_3 are nonnegative in all considered examples. The right hand-side function $F(x, y)$ is chosen so that $u(x,y) = xe^{xy} \sin\pi x \sin\pi y$ solves Eq. 21. The Laplacian in Eq. 21 is discretized by the standard five-point formula and the first-order derivatives by centered finite differences. The mesh size in both the x- and y-directions is $h = (n + 1)^{-1}$.

Define the mesh points $x_j = jh$ and $x_k = kh$ for $0 \leq j, k \leq nwe$ seek to determine values of the solution $u(x,y)$ at interior mesh points $\{(x_j, y_k)\}_{j,k=1}^n$. This yields the linear system of algebraic equations

$$Hu = c \quad (22)$$

Where an $n^2 \times n^2$ matrix H has the following form

$$H = \begin{bmatrix} H_D & H_U & & & \\ H_L & H_D & H_U & & \\ & \dots & \dots & \dots & \\ & & H_L & H_D & H_U \\ & & & H_L & H_D \end{bmatrix} \quad (23)$$

All $n \times n$ sub-matrices are nonsingular and defined as follows

$$H_D = \frac{1}{h^2} \text{tridiag}\{-1 - p_1 h, 4 - 2p_3 h^2, -1 + p_1 h\} \quad (24)$$

$$H_L = \frac{1}{h^2} \text{diag}\{1 + p_2 h\} \quad \text{and} \quad H_U = \frac{1}{h^2} \text{diag}\{1 - p_2 h\} \quad (25)$$

The above linear system can be written as a Sylvester Eq. (1) with $n \times n$ matrices A and B dependent on the assumed separation model, where the entry (j, k) of the

solution matrix $X = \tilde{U}$ of the Sylvester equation approximates $\tilde{u}(x_j, y_k)$ and $C = \{c_{j,k}\}_{j,k=1}^n$ is given by

$$c_{j,k} = F(x_j, y_k) \quad (26)$$

In the examples given in next subsections, we shall examine the behavior of the spectral Radius of the associated implicit iteration matrix $\tau_{B \neq 0}$ as function of the relaxation parameter w as well as the behavior of the following error norms versus the number of iterations

$$B = \left\| \frac{R^{(t)}}{R^{(0)}} \right\|_F, C = \|R^{(t)}\|_F, D = \|E^{(t)}\|_F \quad \text{and} \quad E = \|E^{(t)}\|_{\max} \quad (27)$$

with using the stopping test

$$A = \left\| \bar{\delta} \right\|_{\max} \leq \varepsilon = 10^{-12} \quad (28)$$

and marked in figures by the corresponding letters. The zero initial guess was used in all double precision FORTRAN computations.

Separation model A: This separation model, used in (Hu and Reicher, 1992), is represented by the following matrices

$$A = \frac{1}{h^2} \text{tridiag}\{-1 - p_1 h, 2 - p_3 h^2, -1 + p_1 h\} \quad (29)$$

and

$$B = \frac{1}{h^2} \text{tridiag}\{-1 - p_2 h, 2 - p_3 h^2, -1 - p_2 h\} \quad (30)$$

where the matrix A differs from H_D by the entries of the main diagonal and B contains the entries on sub- and supper-diagonals equal to those in H_U and H_L respectively but multiplied by -1. The diagonal entries of $A + B$ are equal to those on the diagonal of H_D .

Example 3.1A: Let $p_1 = p_2 = p_3 = 0$ and $n = 31$. Then Eq. 21 simplifies to the boundary-value problem for the Poisson equation and the related Sylvester equation has property $B = -A$ and $A = A^T$, which is equivalent to the Lyapunov equation.

Figure 1 shows the behavior of dominant eigenvalues versus w , illustrating the behavior of the spectral radius $\rho(\tau_{B \neq 0})$ marked in Fig. 1 by a. The dominant eigenvalue is positive and a decreasing function of w for $0 < w \leq w_1$ and with $w_1 \approx 0.915$ it achieves the minimum value equal to about 0.822. For $w > w_1$ the dominant eigenvalue becomes complex and its modulus is an increasing function of w and with $w = 1$ is equal to unity. Thus in this example, the iterative scheme (7) is convergent for $0 < w < 1$.

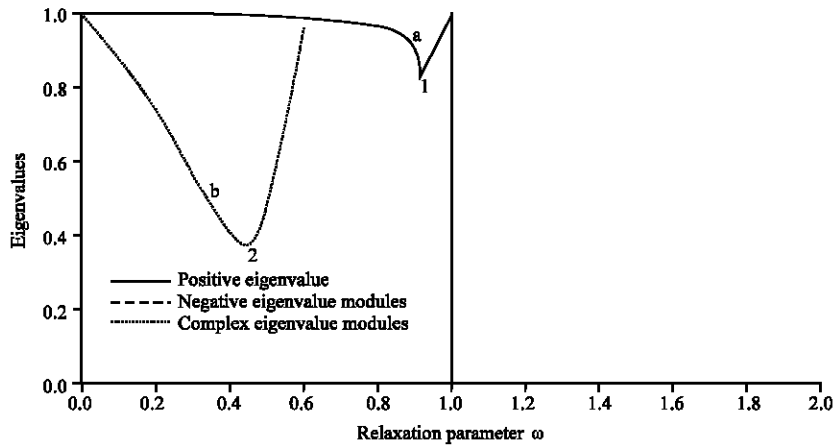


Fig. 1: (a) Example 3.1, (b) Example 3.2

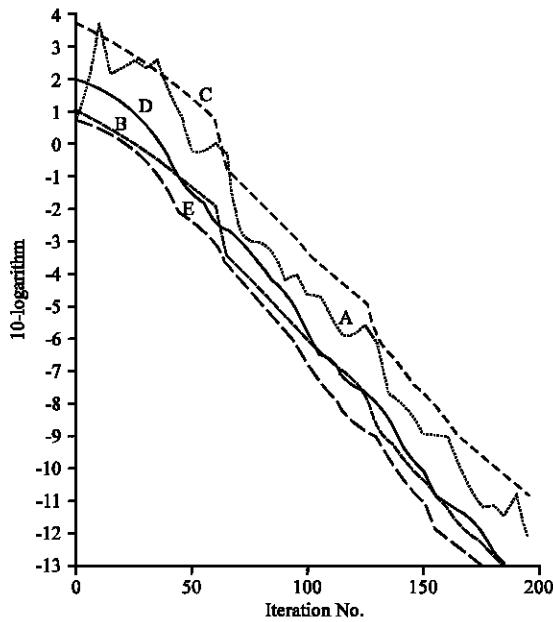


Fig. 2: Example 3.1, The behavior of error norms

The behavior of error norms is shown in Fig. 2 for the results obtained after 195 iterations with $w = 0.915$ and the stopping test (28). As can be seen, the residual norm C is satisfied with about 10^{-11} but the remaining norms are satisfied with the values equal to about 10^{-14}

Example 3.2A: Let $p_1 = 25, p_2 = 50, p_3 = 50$ and $n = 31$. Both matrices A and B are non-symmetric. The behavior of the spectral radius $\rho(\tau_{B=0})$ is marked by b in Fig. 1. The dominant eigenvalue is complex and its modulus is a decreasing function of w for $0 < w \leq w_2$ and with $w_2 \approx 0.44$ it achieves the minimum value equal to about 0.37. For $w > w_2$ its modulus is an increasing function of w and with $w \approx 0.6$

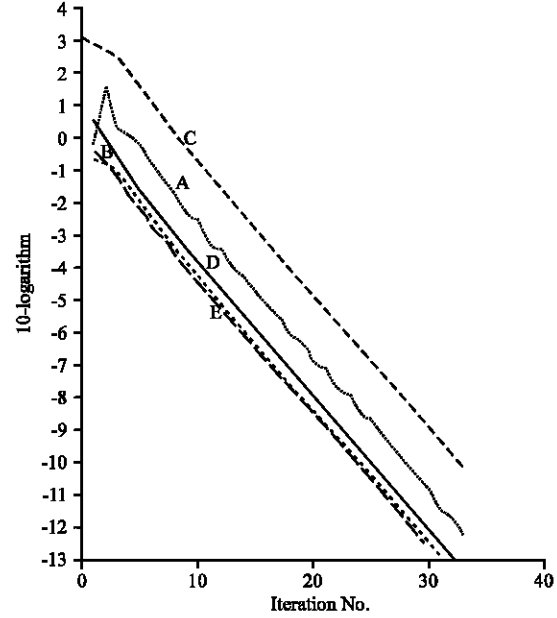


Fig. 3: Example 3.2, The behavior of error norms

is equal to unity. Thus in this example, the iterative scheme (7) is convergent for $0 < w < 0.6$.

The behavior of error norms is shown in Fig. 3 for the results obtained after 34 iterations with $w = 0.44$ and the stopping test (28). As can be seen the residual norm C is satisfied with about 10^{-10} but the remaining norms are again satisfied with the values equal to about 10^{-14} .

Example 3.3A: Let $p_1 = 50, p_2 = 100, p_3 = 50$ and $n = 63$. Both non-symmetric matrices A and B have a larger order than those in previous examples. The behavior of the spectral radius $\rho(\tau_{B=0})$ versus w , depicted in Fig. 4, is similar to that for Example 3.2A marked by b in Fig. 1 and

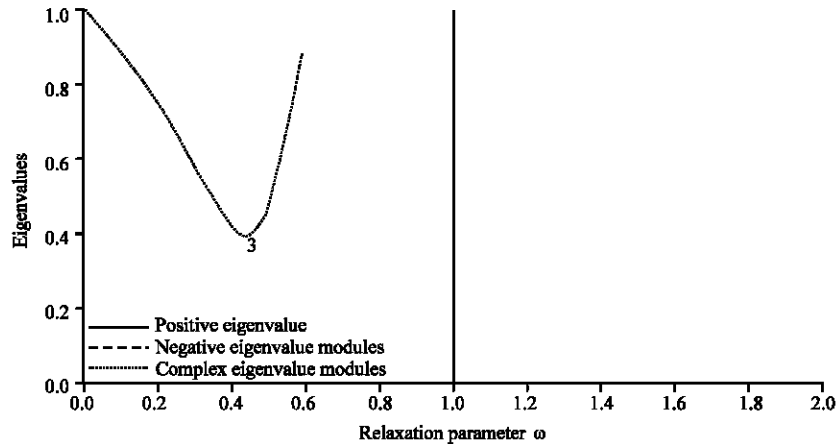


Fig. 4: Example 3.3

$\rho(\tau_{B=0})$ achieves its minimum value equal to about 0.4 with $w_3 \approx 0.45$.

The behavior of error norms, shown in Fig. 5 for the results obtained after 38 iterations with $w = 0.45$ and the stopping test (28), is similar to that for Example 3.2 depicted in Fig. 3.

Separation model B: In this separation model, we use the following matrices

$$A = \frac{1}{h^2} \text{tridiag}\{-1 - p_1 h, 4 - 2p_3 h^2, -1 + p_1 h\} \quad (31)$$

and

$$B = \frac{1}{h^2} \text{tridiag}\{-1 + p_2 h, 0, -1 - p_2 h\} \quad (32)$$

in the Sylvester Eq. (1), where now A is identical with H_D and B consists two nonzero diagonals of H_L and H_U .

We shall consider the previous examples for this separation model, i.e.,

Example 3.1B: For $p_1 = p_2 = p_3 = 0$ and $n = 31$, Eq. 21 simplifies to the boundary-value problem for the Poisson equation and in the related Sylvester equation $A = A^T$ but it is not equivalent to the Lyapunov equation because $B \neq A$.

The behavior of dominant eigenvalues of $\rho(\tau_{B=0})$ versus w , depicted in Fig. 6, is marked by a' . It is interesting to note that this behavior of dominant eigenvalues is identical with the behavior of dominant eigenvalues of the iteration matrix \mathfrak{S}_w in SOR method applied for solving the Eq. 22 in which H is an 2-cyclic consistently ordered matrix. Thus in this case, the property of 2-cyclic consistent ordering moves to the Sylvester equation based on the separation model B so that the value of w_{opt} can be easily determined by means of the Sigma-SOR algorithm (Woznickie, 1994).

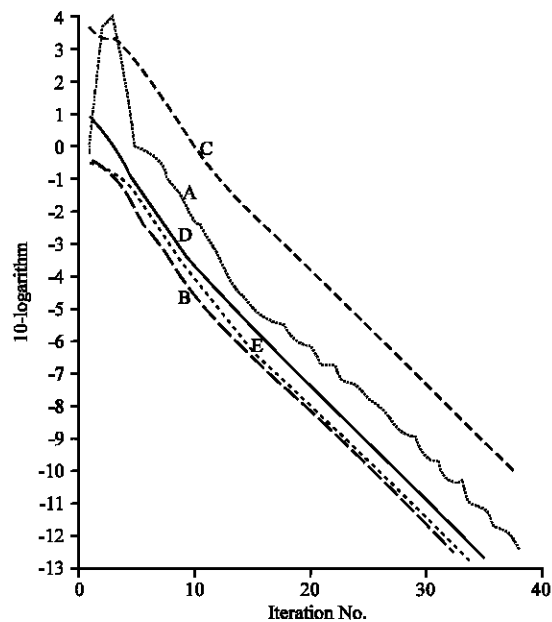


Fig. 5: Example 3.3 The behavior error norms

Example 3.2B: The values of p_1, p_2, p_3 and N are the same as those in Example 3.2A and both matrices A and B are non symmetric. The behavior of the spectral radius $\rho(\tau_{B=0})$ versus w is marked by b' in Fig. 6 The dominant eigenvalue is complex and its modulus is a decreasing function of w for $0 < w \leq w_2$ and with $w_2 \approx 0.90$ it achieves the minimum value equal to about 0.345. For $w > w_2$ its modulus is an increasing function of w and with $w \approx 1.2$ is equal to unity. Thus in this example, the iterative method (7) is convergent for $0 < w < 1.2$.

Example 3.3B: The values of p_1, p_2, p_3 and N are the same as those in Example 3.3A and both non symmetric matrices A and B have a larger order than those in two

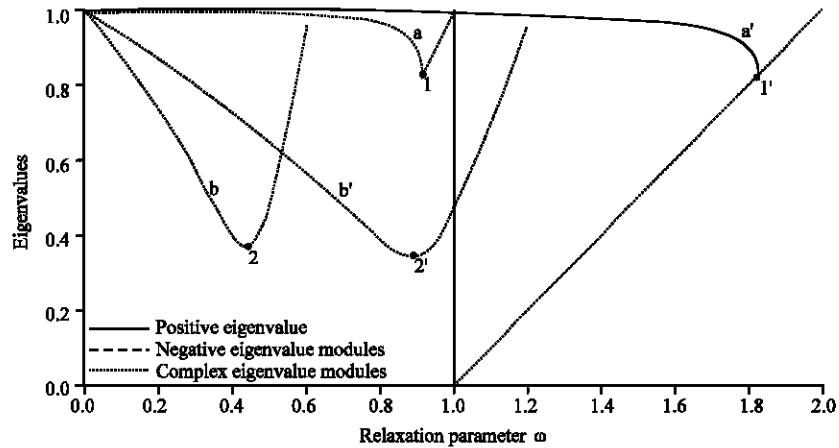


Fig. 6: a,a' Example 3.1, b,b' Example 3.2

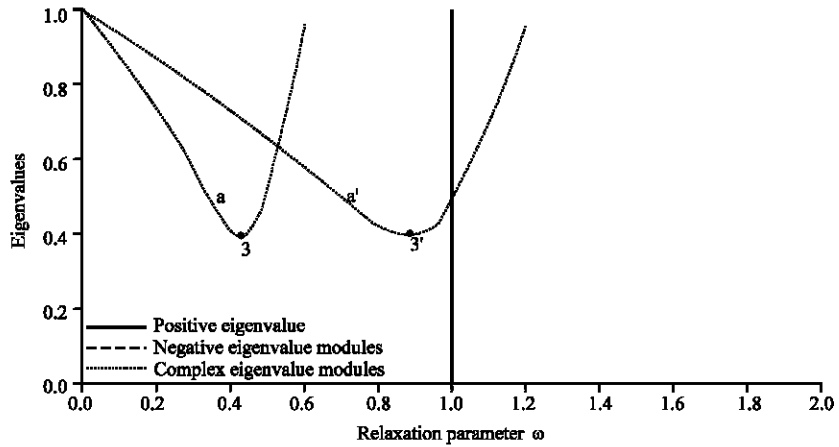


Fig. 7: Example 3.3.

previous examples. The behavior of the spectral radius $\rho(\tau_{B=0})$ versus w , marked by a' in Fig. 7, is similar to that for Example 3.2B marked by b' in Fig. 6 and $\rho(\tau_{B=0})$ achieves its minimum value equal to about 0.40 with $w_3 \approx 0.90$.

As can be seen in these figures, the minimum values of the spectral radius are the same in both separation models therefore, with using w_{best} the SOR-like method has the same convergence properties for both separation models and the behavior of error norms versus the number of iterations for the separation model B is nearly the same as that observed in Fig. 2, 3 and 5 for the separation model A.

However, the separation model B increases the range of convergence, where the SOR-like method is convergent for $0 < w < 2$ in Example 3.1B and in the case of Examples 3.2B and 3.3B for $0 < w < 1.2$, where the values of w_{best} are increased two times. This increasing the range of

convergence implies that the SOR-like method with the separation model B becomes less sensitive to the accurate estimate of w_{best} (especially in Examples 3.2B and 3.3B), which simplifies determining w_{best} and its convergence behavior is identical with that for the SOR method when is used for solving Eq. 22.

CONCLUSION

As is demonstrated in Hu and Reicher (1992), the Galerkin method gives a little better result than minimal-residual one in Examples 3.1 and 3.3 but in the case of Example 3.2, an inverse behavior is observed. In Example 3.1 the solution have been obtained after 20 iterations with the Krylov subspace dimension $m = 6$ and the relative residual norm B equal to about 10^{-3} , which corresponds to the residual norm $C = \|R^{(t=20)}\|$ in this example. The solution of Example 3.2 have been obtained

with the relative residual norm B equal to about 10^{-7} , which corresponds to the residual norm $C = \|R^{(n=20)}\|_F \approx 0.0003$. The best results are presented in Hu and Reicher (1992) for Example 3.3 and its solution were obtained with the Krylov subspace dimension $m = 4$ and the relative residual norm B equal to about 10^{-10} , corresponding to the residual norm $C = \|R^{(n=15)}\|_F \approx 2 \times 10^{-6}$. However as is shown in Figure 4 given in Hu and Reicher (1992), both minimal-residual and Galerkin methods can provide the solutions for Example 3.3 only with the true error norm $E = \|E^{(i)}\|_{\max} \approx 10^{-3}$.

From the results of numerical experiments presented in Sections 3.1 and 3.2, it can be concluded that the proposed SOR-like method, represented by the iterative scheme (7), is a very efficient technique for solving Sylvester equations especially with the non symmetric matrices A and B . For considered examples, the SOR-like method provides much more accurate solutions, obtained with a computational work less than a few orders in magnitude in comparison to solutions obtained by means of Krylov subspace algorithms discussed in Hu and Reicher (1992). The computational work in one iteration of the scheme (7) is roughly equivalent to that required for computing the residual matrix (16). Thus with using the stopping test (28), the solutions were obtained with the relative residual norm B , the true error norms $D = \|E^{(i)}\|_F$ and $E = \|E^{(i)}\|_{\max}$ less than 10^{-13} ; and after 195, 33 and 38 iterations (equivalent to computations of the residual matrix) for Examples 3.1, 3.2 and 3.3 respectively. From the viewpoint of simplicity in determining w_{best} , the SOR-like method with the separation model B is a more useful computational technique for solving Sylvester equations.

In Example 3.1B, the Sylvester equation preserves 2-cyclic consistent ordering property and $w_{best} = w_{opt}$ can be easily determined by means of the Sigma-SOR algorithm (Woznicki, 1994).

In the case of Examples 3.2 and 3.3, the value of w_{best} can be found experimentally by fitting a parabola curve for both separation models. Assuming that the right hand-side of a solved problem is put to zero, then with all entries of starting matrix $X^{(0)}$ equal to unity, the solution of (7) converges to the null matrix when $\rho(\tau_{B=0}) < 1$ for a given w but its rate of convergence is dependent on the used value of w . For three values of w chosen around expected w_{best} , we obtain the corresponding numbers of iterations satisfying, for instance, the stopping test $\|X^{(i)}\|_{\max} \leq 10^{-3}$ and the abscissa of the minimum of a fitted parabola provides us a quite good approximation for w_{best} .

Finally, it should mention that similar behavior of the SOR-like method was observed with solving other test examples taken from the literature.

REFERENCES

- Aliev, F.A. and V.B. Larin, 1998. Optimization of linear control systems: Analytical methods and computational algorithms, Stability and Control: Theory, Methods and Applications. Gordon and Breach.
- Calvetti, D. and L. Reichel, 1996. Application of ADI iterative methods to the restoration of noisy images. *SIAM J. Matrix Anal. Applied*, 17: 165-186.
- Datta, B.N., 2004 Numerical Methods for linear control systems, Elsevier Academic press.
- Dieci, L., M.R. Osborne and R.D. Russell, 1988. A Riccati transformation method for solving linear bvps. I: Theoretical Aspects. *SIAM J. Numer. Anal.*, 25:1055-1073.
- Enright, W.H., 1978. Improving the efficiency of matrix operations in the numerical solution of stiff ordinary differential equations. *ACM Trans. Math. Software*, 4:127-136.
- Epton, M.A., 1980. Methods for the solution of AXD-BXC = E and its application in the numerical solution of implicit ordinary differential equations. *BIT*, 20: 341-345.
- Gantmacher, F.R., 1960. The Theory of Matrices, Chelsea Publishing Company, New York.
- Golub, G.H. and C.F. Van Loan, 1996. Matrix Computations. Johns Hopkins University Press, Baltimore, (3rd Edn).
- Hu D.Y. and L. Reicher, 1992. Krylov subspace methods for the Sylvester equation, *Lin. Alg. Applied*, pp: 283-313.
- Sima, V., 1996. Algorithms for Linear-Quadratic Optimization, Pure and Applied Mathematics. Marcel Dekker, Inc., New York.
- Woznicki, Z.I., 1994. The Sigma-SOR algorithm and the optimal strategy for the SOR iterative method, *Comput. Applied Math.*, pp: 145-176.
- Woznicki, Z.I., 2001. On performance of SOR method for solving nonsymmetric linear systems, *Math. Computa.*, pp: 145-176.
- Woznicki, Z.I. and S. Kadry, 2003. A SOR-like method for solving the Sylvester equation, *Ann.* pp: 335-344.