

where $D = \text{diag}(T_i)$, and L is a lower block bidiagonal matrix, with the main diagonal equal to the one of D , and the lower diagonal equal to the one of A . It is simple to show that $R = \text{diag}(R_i)$ [3],

$$\begin{aligned} R_1 &= T_1 - S_1 = 0, \\ R_i &= F_i(T_{i-1}^{-1} - \Sigma_{i-1})F_i, \quad i = 2, \dots, n. \end{aligned} \tag{1.4}$$

We consider the case in which the blocks T_i have the same sparsity as the diagonal blocks S_i of A . In particular, we are going to consider the standard INV and MINV preconditioner [2–4], even if the analysis could be extended to other block preconditioners.

The INV preconditioner is obtained by considering $\Sigma_{i-1} = \text{trid}(T_{i-1}^{-1})$, where $\text{trid}(T_{i-1}^{-1})$ is the tridiagonal part of T_{i-1}^{-1} . The MINV preconditioner is the modified version of INV, obtained by imposing that Σ_{i-1} and T_{i-1}^{-1} have the same row-sums [3]. We consider again the tridiagonal part of the inverse, in (1.2), but the neglected elements of the inverse are summed, on each row, to the corresponding diagonal element.

One can show [3] that the blocks T_i of INV and MINV are symmetric diagonally dominant M -matrices.

2. Vectorizing the solution of $Ms = r$

To get the preconditioned residual s , in a preconditioned conjugate gradient method, we must solve $Ms = r$, where r is the current residual:

$$\begin{aligned} T_1 y_1 &= r_1, & s_n &= y_n, \\ T_i y_i &= r_i - F_i y_{i-1}, \quad i = 2, \dots, n, & T_i (s_i - y_i) &= -F_{i+1} s_{i+1}, \quad i = n - 1, \dots, 1. \end{aligned} \tag{2.1}$$

Let us consider, now, the generic block T_i :

$$T_i = \begin{bmatrix} a_1^{(i)} & -b_2^{(i)} & & & \\ -b_2^{(i)} & \ddots & \ddots & & \\ & \ddots & \ddots & & \\ & & & -b_k^{(i)} & a_k^{(i)} \end{bmatrix}. \tag{2.2}$$

It could be factorized as $T_i = L_{T_i} D_{T_i} L_{T_i}^T$, where

$$\begin{aligned} L_{T_i} &= \begin{bmatrix} d_1^{(i)} & & & & \\ -b_2^{(i)} & \ddots & & & \\ & \ddots & \ddots & & \\ & & & -b_k^{(i)} & d_k^{(i)} \end{bmatrix}, & D_{T_i} &= \text{diag}(d_1^{(i)}, \dots, d_k^{(i)}), \\ d_1^{(i)} &= a_1^{(i)}, \\ d_j^{(i)} &= a_j^{(i)} - \frac{(b_j^{(i)})^2}{d_{j-1}^{(i)}}, \quad j = 2, \dots, k. \end{aligned} \tag{2.3}$$

If we scale T_i to get $D_{T_i} = I$, we get $T_i = (I - E_i)(I - E_i^T)$, where

$$E_i = \begin{bmatrix} 0 & & & & & \\ -\tilde{b}_2^{(i)} & 0 & & & & \\ & \ddots & \ddots & & & \\ & & \ddots & \ddots & & \\ & & & -\tilde{b}_k^{(i)} & 0 & \\ & & & & & \end{bmatrix}. \tag{2.4}$$

From the diagonal dominance of T_i , one gets $\rho(E_i) < 1$. It follows that

$$T_i^{-1} = (I - E_i^T)^{-1}(I - E_i)^{-1} = (I + E_i^T + \dots + E_i^{T^{k-1}})(I + E_i + \dots + E_i^{k-i}). \tag{2.5}$$

We can use a truncated expansion of (2.5) as approximation of T_i^{-1} :

$$T_i^{-1} \approx \tilde{T}_i^{-1} = (I + E^T + \dots + E^{T^m})(I + E + \dots + E^m), \tag{2.6}$$

where, obviously $m \ll k - 1$. It follows that the solution of the tridiagonal subsystems in (2.1) can be replaced with vectorizable operations. Now we want to get an estimate of the error due to the use of (2.6), instead of (2.5). For sake of simplicity, let us assume $F_i = -I$ (this is a common case). In such a way, instead of M we get

$$\begin{aligned} \tilde{M} &= \begin{bmatrix} \tilde{T}_1 & & -I & & & \\ -I & \tilde{T}_2 + \tilde{T}_1^{-1} & & -I & & \\ & \ddots & \ddots & \ddots & & \\ & & \ddots & \ddots & & -I \\ & & & -I & \tilde{T}_n + \tilde{T}_{n-1}^{-1} & \\ & & & & & \end{bmatrix} \\ &= M + \begin{bmatrix} \tilde{T}_1 - T_1 & & & & & \\ & \ddots & & & & \\ & & & \tilde{T}_n - T_n & & \\ & & & & & \end{bmatrix} + \begin{bmatrix} 0 & & & & & \\ \tilde{T}_1^{-1} - T_1^{-1} & & & & & \\ & \ddots & & & & \\ & & & & & \\ & & & & & \tilde{T}_{n-1}^{-1} - T_{n-1}^{-1} \end{bmatrix} \\ &= M + P_1 + P_2 = M + \tilde{R}. \end{aligned}$$

Referring to (1.3), we choose m so that $\|\tilde{R}\| \leq \|R\|$ (from now on, $\|\cdot\|$ denotes $\|\cdot\|_\infty$), that is, the error introduced by the truncation must be at most the error made in the construction of the preconditioner. It follows that we need an estimate for $\|P_1\|$, $\|P_2\|$ and $\|R\|$ (while $\|R\|$ is available immediately at run time, if the MINV preconditioner is used; this is not true for INV).

Now we show how to get an estimate for $\|P_1\|$. Observe that

$$(I + E + \dots + E^m)^{-1} = (I - E)(I - E^{m+1})^{-1} \approx (I - E)(I + E^{m+1}).$$

It follows, from (2.6),

$$\tilde{T}_i \approx T_i + (I - E_i)(E_i^{m+1} + E_i^{T^{m+1}})(I - E_i^T),$$

that is

$$\|P_1\| \approx 2\beta_1^{m+1}(1 + \beta_1)^2, \tag{2.7}$$

where

$$\beta_1 = \max_{i=1, \dots, n} \|E_i\|. \tag{2.8}$$

Now we are looking for an estimate of $\|P_2\|$. If we define the matrices

$$S_i = T_i^{-1} - \tilde{T}_i^{-1},$$

we get

$$\|S_i\| < 2\|E_i^{m+1}\|\|T_i^{-1}\| = \eta_i.$$

It follows that

$$\|P_2\| < \max_{i=1, \dots, n-1} \eta_i \leq 2\beta_2^{m+1}\gamma, \tag{2.9}$$

where

$$\beta_2 = \max_{i=1, \dots, n-1} \|E_i\|, \quad \gamma = \max_{i=1, \dots, n-1} \|T_i^{-1}\|.$$

Finally, we get

$$\|\tilde{R}\| \approx 2(\beta_1^{m+1}(1 + \beta_1)^2 + \gamma\beta_2^{m+1}) = \xi(m). \tag{2.10}$$

At last, let us obtain an estimate of $\|R\|$ for INV.

Let T_i the generic diagonal block of the INV factorization, whose structure is given in (2.2), and let $T_i^{-1} = (t_{rs}^{(i)})$. One gets [4]

$$\begin{aligned} t_{kk}^{(i)} &= 1/d_k^{(i)}, \\ t_{rk}^{(i)} &= t_{r+1,k}^{(i)}b_{r+1}^{(i)}/d_r^{(i)}, \quad r = k-1, \dots, 1, \\ &\quad \text{for } s = k-1, \dots, 1, \\ t_{ss}^{(i)} &= (1 + t_{s,s+1}^{(i)}b_{s+1}^{(i)})/d_s^{(i)}, \\ t_{sr}^{(i)} &= t_{rs}^{(i)} = t_{r+1,s}^{(i)}b_{r+1}^{(i)}/d_r^{(i)}, \quad r = s-1, \dots, 1. \end{aligned}$$

With the assumption of T_i being diagonally dominant, we get that the elements on the rows of T_i^{-1} become smaller and smaller, as we go away from the main diagonal. Moreover, as we are considering the case in which $F_i = -I$, we get (see (1.4))

$$\|R\| = \max_{i=1, \dots, n-1} \|T_i^{-1} - \text{trid}(T_i^{-1})\|.$$

It follows that

$$\delta_{\min} \sum_{j=2}^{k-1} \sigma_{\min}^j \leq \|R\| \leq 2\delta_{\max} \sum_{j=2}^{(k-1)/2} \sigma_{\max}^j,$$

where

$$\begin{aligned} \delta_{\min} &= \min_{i,j} (t_{jj}^{(i)}), & \delta_{\max} &= \max_{i,j} (t_{jj}^{(i)}), \\ \sigma_{\min} &= \min_{i,j} (b_{j+1}^{(i)}/d_j^{(i)}), & \sigma_{\max} &= \max_{i,j} (b_{j+1}^{(i)}/d_j^{(i)}), \end{aligned}$$

that is

$$r_1 = \delta_{\min}\sigma_{\min}^2 \frac{1 - \sigma_{\min}^{k-2}}{1 - \sigma_{\min}} \leq \|R\| \leq 2\delta_{\max}\sigma_{\max}^2 \frac{1 - \sigma_{\max}^{(k-2)/2}}{1 - \sigma_{\max}} = r_2. \tag{2.11}$$

In such a way, we get an interval estimate for $\|R\|$. Moreover, with reference to (2.9) we get

$$\gamma \leq 2\delta_{\max} \frac{1 - \sigma_{\max}^{k/2}}{1 - \sigma_{\max}} = \gamma_1.$$

If the interval $\mathcal{L}=[r_1, r_2]$ is ‘sufficiently small’, any value of m for which $\xi(m) \in \mathcal{L}$ (see (2.10)) is acceptable. Otherwise, we should remain as much as possible close to r_1 .

3. Numerical tests

We have considered three test problems. For each problem, we have compared the original preconditioner, INV or MINV, with the truncated ones, for various choices of m . Moreover, we have compared those preconditioners with the one suggested by Meurant [5], in which the inverse of the generic tridiagonal block is approximated by a band matrix, with the main seven diagonals coinciding with those of the exact inverse.

Let us briefly examine the computational cost of the considered preconditioned conjugate gradient methods, in terms of requested memory and operations (per iterate). The acronym TRUNC(m) is for the method using the truncated preconditioner, while MEUR is for the one proposed by Meurant. The acrostics MTRUNC and MMEUR are for the modified versions. N is the dimension of the problem. See Table 1.

The tests were carried out on a single processor of a CRAY X-MP/48. The language used is FORTRAN (CFT). To get the execution time (in seconds), we got, for each problem and method used, the minimum time of 50 executions.

Test problem 1. The first test problem derives from the discretization of

$$-\Delta u = f, \quad \text{on } \Omega = (0, 1) \times (0, 1),$$

$$u|_{\partial\Omega} = 0.$$

The usual 5-points scheme is used, with step size $h = (n + 1)^{-1}$, $N = n^2$, where n is the mesh size, $n = 100$.

Test problem 2. The second problem derives from the discretization of

$$-\lambda \Delta u = f \quad \text{on } \Omega = (0, 1) \times (0, 1),$$

$$u|_{\partial\Omega} = 0,$$

where

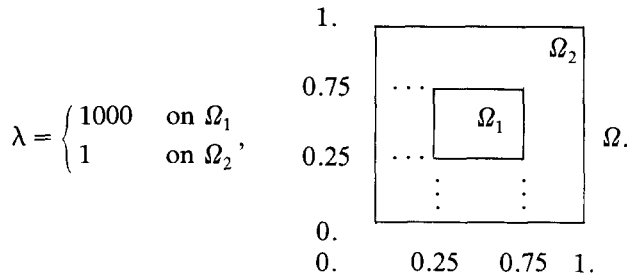


Table 1

Method	Memory	Vectorizable operations	Non vectorizable operations
INV/MINV	10N	32N	8N
TRUNC/MTRUNC(3)	11N	40N	-
TRUNC/MTRUNC(7)	12N	48N	-
TRUNC/MTRUNC(15)	13N	56N	-
MEUR/MMEUR	12N	48N	-

The usual 5-points scheme is used, with step size $h = (n + 1)^{-1}$, $N = n^2$, where n is the mesh size, $n = 100$.

In both Problems 1 and 2 we got for INV the following estimates (see (2.11) and (2.10)):

$$r_1 = 0.089, \quad r_2 = 0.4915, \quad \|R\| = 0.4915 \text{ (true value)},$$

$$\xi(3) = 0.1639, \quad \xi(7) = 0.0225, \quad \xi(16) = 4E-7.$$

For MINV we got

$$r_1 = 0.178, \quad r_2 = 1.8656, \quad \|R\| = 1.8656 \text{ (true value)},$$

$$\xi(3) = 0.4178, \quad \xi(7) = 0.0115, \quad \xi(16) = 8E-6.$$

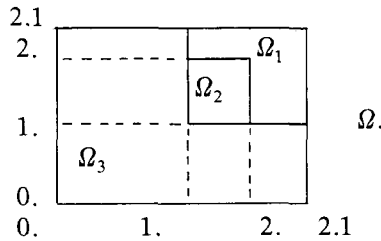
Test problem 3. The third problem derives from the discretization of

$$-\lambda \Delta u + \sigma u = \sigma \quad \text{on } \Omega = (0, 2.1) \times (0, 2.1),$$

$$\frac{\partial u}{\partial z} \Big|_{\partial\Omega} = 0,$$

where z the unit vector normal to $\partial\Omega$, and

$$\lambda = \begin{cases} 1 & \text{on } \Omega_1, \\ 2 & \text{on } \Omega_2, \\ 3 & \text{on } \Omega_3, \end{cases} \quad \sigma = \begin{cases} 0.02 & \text{on } \Omega_1, \\ 0.03 & \text{on } \Omega_2, \\ 0.05 & \text{on } \Omega_3, \end{cases}$$



The usual 5-points scheme is used, with step size $h = (n + 1)^{-1}$, $N = n^2$, where n is the mesh size, $n = 90$.

For this problem the value of β_1 (see (2.8)) in the relation (2.7) is obtained in correspondence of the last element of the last block. We get

$$\|P_1\| = \|T_n - \tilde{T}_n\| \approx \|I + E_n\|^2 \|E_n^{m+1}\|$$

where

$$E_n = \begin{bmatrix} 0 & & & & \\ \tilde{\beta}_1 & 0 & & & \\ & \ddots & \ddots & & \\ & & \tilde{\beta}_1 & 0 & \\ & & & \beta_1 & 0 \end{bmatrix},$$

and (see (2.9) and (2.4))

$$\tilde{\beta}_1 = \max\left(\beta_2, \max_{i=2, \dots, k} \{\tilde{b}_i^{(n)}\}\right), \quad \beta_1 - \tilde{\beta}_1 = \epsilon > 0.$$

It follows that, instead of (2.7), we get

$$\|P_1\| \approx 2\tilde{\beta}_1\beta_1(1 + \beta_1)^2.$$

Table 2
Test problem 1

Method	Speedup	Time	Iterates
INV	1.0	0.656	28
TRUNC(3)	3.2	0.206	31
TRUNC(7)	3.0	0.220	28
TRUNC(15)	2.6	0.248	28
MEUR	2.9	0.228	30
MINV	1.0	0.464	20
MTRUNC(3)	3.3	0.142	22
MTRUNC(7)	2.5	0.184	21
MTRUNC(15)	2.5	0.183	20
MMEUR	2.8	0.167	22

Table 3
Test problem 2

Method	Speedup	Time	Iterates
INV	1.0	0.712	30
TRUNC(3)	3.1	0.227	34
TRUNC(7)	3.0	0.236	30
TRUNC(15)	2.7	0.266	30
MEUR	2.8	0.252	33
MINV	1.0	0.440	19
MTRUNC(3)	2.7	0.163	24
MTRUNC(7)	2.6	0.171	20
MTRUNC(15)	2.5	0.176	19
MMEUR	2.4	0.183	23

We obtain for INV,

$$r_1 = 0.2457, \quad r_2 = 2.7676, \quad \|R\| = 0.8446 \text{ (true value),}$$

$$\xi(3) = 9.5059, \quad \xi(7) = 2.5945, \quad \xi(16) = 0.1933.$$

Observe that, even if $\xi(3) \notin [r_1, r_2]$, the truncated preconditioner TRUNC(3) is effective, as we are going to see.

For each method, and each problem, we give the speedup with respect to the original preconditioner (INV/MINV), the execution time, and the number of iterations to get convergence, see Tables 2–4. The stopping criterion used is $\|r_i\|_2 / \|r_0\|_2 < 10^{-6}$, where r_i is the residual at the i th step, and r_0 the initial residual. The initial point is $x_0 = 0$.

Table 4
Test problem 3

Method	Speedup	Time	Iterates
INV	1.0	1.349	69
TRUNC(3)	3.0	0.448	79
TRUNC(7)	2.9	0.460	69
TRUNC(15)	2.5	0.530	69
MEUR	2.4	0.569	80

References

- [1] O. Axelson, A survey of preconditioned iterative methods for linear systems of algebraic equations, *BIT* **25** (1985) 166–187.
- [2] L. Brugnano and M. Marrone, Metodi dei gradienti coniugati: Teoria, Tecniche di preconditionamento, ed Implementazione su calcolatori vettoriali, Quaderno n.ro 1/89 del Dipartimento di Matematica dell' Università di Bari, Italy.
- [3] P. Concus, G.H. Golub and G. Meurant, Block preconditioning for the conjugate gradient method, *SIAM J. Sci. Statist. Comput.* **6** (1) (1985) 220–252.
- [4] P. Concus and G. Meurant, On computing the INV block preconditionings for the conjugate gradient method, *BIT* **26** (1986) 493–504.
- [5] G. Meurant, Multitasking the conjugate gradient method on the CRAY X-MP/48, *Parallel Comput.* **5** (1987) 267–280.
- [6] H.A. Van Der Vorst, A vectorizable variant of some ICCG methods, *SIAM J. Sci. Statist. Comput.* **3** (3) (1982) 350–356.
- [7] R.S. Varga, *Matrix Iterative Analysis* (Prentice-Hall, Englewood Cliffs, NJ, 1962).