Owe Axelsson

## Condition number estimates for elliptic difference problems with anisotropy

# CONDITION NUMBER ESTIMATES
# FOR ELLIPTIC DIFFERENCE PROBLEMS
# WITH ANISOTROPY

## AXELSSON O., NIJMEGEN, The Netherlands


### 1. Introduction

Consider the solution of a linear system $A\hat{x} = b$, where $A$ is symmetric and positive definite by an iterative method, such as a preconditioned conjugate gradient or preconditioned Chebyshev iterative method. Let $A$ be split as

$$A = D - L - L^T$$

where $D$ is the (block) diagonal of $A$ and $L$ is the strictly lower (block) triangular part of $A$.

As preconditioner, i.e. an approximation of $A$ with low computational complexity for the solution of systems with it, we shall analyse the *generalized SSOR method* (see, for instance [1], [3]), where

$$(1.1) \qquad C = (X - L)X^{-1}(X - L^T)$$

and $X$ is (block) diagonal with positive diagonal entries (or positive definite diagonal blocks) chosen as to be described below.

We have

$$C = X + LX^{-1}L^T - L - L^T$$

so

$$R \equiv C - A = X - D + LX^{-1}L^T.$$

Let $R^0$ be defined by

$$(R^0)_{i,j} = \begin{cases} 0, & i = j \\ (LX^{-1}L^T)_{i,j}, & i \neq j \end{cases}.$$

(In the block matrix case, $(R_0)_{i,j}$ denotes the $i,j$'th block of $R_0$.) Hence, $R^0$ consists of the "fill-in" entries, i.e. the entries of the matrix $LX^{-1}L^T$ which fall outside the (block) diagonal. $X$ is computed recursively from

$$(1.2) \qquad X_i = D_i - (LX^{-1}L^T)_{i,i} - \omega(R^0 e)_i \quad , \quad i = 1, 2, ...,$$

where $D_i$ is the $i$'th block of $D$, $e = (1, 1, ..., 1)^T$, and $\omega (\omega \leq 1)$ is a relaxation parameter. Note that $(R^0 e)_i$ is a scalar if $X$ and $D$ are diagonal and a diagonal matrix if $X$ and $D$ are block diagonal. Hence, in the latter case, the off diagonal entries of $X_i$ are determined so that they are equal to the corresponding entries of $D_i - (LX^{-1}L^T)_{i,i}$. Hence, $X_i$ is uniquely determined by (1.2). Note also that by choosing $\omega$ sufficiently small (even negative, if necessary) we can guarantee that $X_i$ becomes positive definite.

The method of using a relaxation parameter $\omega$ was first introduced in Axelsson and Lindskog [5] (for a more general incomplete factorization method). It follows readily that for $\omega = 1$ we have $Ce = Ae$, which is the rowsum criterion and a basis for the modified method of Gustafsson [6]. The relaxation parameter has the same effect on the spectrum of the iteration matrix $C^{-1}A$, as the use of perturbations, which latter has been used by the present author in [1] and [3].

Next we shall derive upper and lower bounds of the extreme eigenvalues of the generalized eigenvalue problem

$$\lambda C\mathbf{v} = A\mathbf{v} \tag{1.3}$$

and derive estimates of the spectral condition number of $C^{-1}A$ as a function of $\omega$.

## 2. Upper and lower bounds of the extreme eigenvalues.

To derive a lower bound note first that we have

$$\lambda C - A = (1 - \lambda)(-A) + \lambda(C - A),$$

so

$$\lambda C - A = (1 - \lambda)(-A) + \lambda R.$$

Let $\mu_i()$ denote the $i$'th eigenvalue. Then it follows by the Courant Fischer lemma (see Wilkinson [8], p.101) that for any positive $\lambda$, the $i$'th eigenvalue of $\lambda C - A$ satisfies

$$\mu_i(\lambda C - A) \leq f_i(\lambda) \equiv (1 - \lambda)\mu_i(-A) + \lambda\mu_+(R), \tag{2.1}$$

where $\mu_+(R)$ denotes the largest eigenvalue of $R$.

Note now that $\mu_i(\lambda C - A) = 0$ if and only if $\lambda$ is an eigenvalue of the generalized eigenvalue problem (1.3) and note that these eigenvalues are positive because $C$ and $A$ are both symmetric and positive definite.

If $\mu_+(R) > 0$ then there exists a zero, $\underline{\lambda}_i$ of $f_i(\lambda)$ in the interval (0,1) and we find

$$\lambda_i \geq \underline{\lambda}_i = \mu_i(A)/[\mu_i(A) + \mu_+(R)].$$

In particular, for the smallest eigenvalue we have

$$\lambda_1 \geq \underline{\lambda}_1 = \mu_1(A)/[\mu_1(A) + \mu_+(R)] \tag{2.2}$$

where we assume that the eigenvalues have been ordered in an increasing order. The method used above to derive a lower bound is based on an idea in Van der Vorst [7].

Next we shall derive two bounds for the largest eigenvalue of $C^{-1}A$. We extend then a method used by the author in [2], see also Axelson and Barker [4]. We have

$$\lambda C = [(1 - \frac{1}{\lambda})X - L + \frac{1}{\lambda}X](\frac{1}{\lambda}X)^{-1}[(1 - \frac{1}{\lambda})X - L^T + \frac{1}{\lambda}X]$$

or

$$\lambda C - A = \lambda V X^{-1} V^T + (2 - \frac{1}{\lambda})X - D$$

where $V = (1 - \frac{1}{\lambda})X - L$. Hence, since $V X^{-1} V^T$ is positive semidefinite, for any positive $\hat{\lambda}$ we find

$$(2.3) \qquad \mu_i(\hat{\lambda} C - A) \geq \mu_-((2 - \frac{1}{\lambda})X - D),$$

where $\mu_-()$ denotes the smallest eigenvalue. We shall assume that $2X - D$ is positive definite (which again can be achieved by a proper choice of $\omega$ in (1.2)). Hence, there exists a positive $\hat{\lambda}$ for which

$$\mu_-((2 - \frac{1}{\lambda})X - D) \geq 0.$$

Note now that

$$\lambda C - A = (1 - \frac{\lambda}{\hat{\lambda}})(-A) + \frac{\lambda}{\hat{\lambda}}(\hat{\lambda} C - A),$$

so, by (2.3) and the same result in Wilkinson [8] as used before, it follows that

$$\mu_i(\lambda C - A) \geq g_i(\lambda) \equiv (1 - \frac{\lambda}{\hat{\lambda}})\mu_i(-A) + \frac{\lambda}{\hat{\lambda}}\mu_-((2 - \frac{1}{\lambda})X - D).$$

When $\mu_-((2 - \frac{1}{\lambda})X - D) \geq 0$, there exists a zero, $\overline{\lambda}_i$ of $g_i(\lambda)$ in the interval $[0, \hat{\lambda}]$ and this is then an upper bound of the $i$'th eigenvalue $\lambda_i$ of $C^{-1}A$. Hence

$$\lambda_i \leq \overline{\lambda}_i = \hat{\lambda}\mu_i(A)/[\mu_i(A) + \mu_-((2 - \frac{1}{\lambda})X - D)].$$

In particular, for the largest eigenvalue we have

$$(2.4) \qquad \max_i \lambda_i \leq \hat{\lambda}/[1 + \mu_-((2 - \frac{1}{\lambda})X - D)/\max_i \mu_i(A)].$$

Next we consider an alternative upper bound for the largest eigenvalue, which is valid when $A$ is an $M$-matrix i.e. in particular requires that the off-diagonal entries of $A$ are non-positive. We have

$$\gamma A - C = (\gamma - 1)C + \gamma(A - C)$$

and for any positive $\gamma$,

$$\mu_i(\gamma A - C) \leq (\gamma - 1)\mu_i(C) + \gamma\mu_+(-R)$$

or, if $\mu_+(-R) \geq 0$,

$$\gamma_i \geq \underline{\gamma}_i = \mu_i(C)/[\mu_i(C) + \mu_+(-R)],$$

where $\gamma_i$ denotes the $i$'th eigenvalue of $A^{-1}C$.

Hence, if $\mu_+(-R) > 0$, the smallest eigenvalue satisfies

$$\gamma_1 \geq 1/[1 + \mu_+(-R)/\mu_1(C)].$$

220

Since $\max_i \lambda_i = \gamma_1^{-1}$ we have then

$$(2.5) \qquad \max_i \lambda_i \leq 1 + \mu_+(-R)/\mu_1(C).$$

To estimate $\mu_1(C)$, the smallest eigenvalue of $C$, we estimate first the largest eigenvalue of $C^{-1}$, using (1.1). We find, using the property that $X^{-1}L$ has non-negative entries,

$$(2.6) \qquad \mu_1(C) = \frac{1}{\max_i \mu_i(C^{-1})} \leq \frac{1}{\max_i\{(X - L^T)^{-1}X(X - L)^{-1}\mathbf{e}\}_i}.$$

Hence, (2.5) and (2.6) show that

$$(2.7) \qquad \max_i \lambda_i \leq 1 + \mu_+(-R)\max_i\{(X - L^T)^{-1}X(X - L)^{-1}\mathbf{e}\}_i.$$

We collect the results in a theorem.

**Theorem 2.1.** Let $C$ be defined by (1.1), (1.2) and let $R = C - A$. Then

a) if $\mu_+(R) \geq 0$, the smallest eigenvalue of $C^{-1}A$ satisfies

$$\lambda_1 \geq 1/[1 + \mu_+(R)/\mu_1(A)].$$

b) If $2X - D$ is positive definite and $\hat{\lambda}$ is sufficiently small so that $(2 - \frac{1}{\hat{\lambda}})X - D$ is positive semidefinite, then

$$\max_i \lambda_i \leq \hat{\lambda}/[1 + \mu_-((2 - \frac{1}{\hat{\lambda}})X - D)/\max_i \mu_i(A)].$$

c) If $\mu_+(-R) \geq 0$ and if $A$ is an $M$-matrix, then

$$\max_i \lambda_i \leq 1 + \mu_+(-R)\max_i\{(X - L^T)^{-1}X(X - L)^{-1}\mathbf{e}\}_i.$$

**Proof.** This follows from (2.2), (2.4) and (2.7). $\qquad\qquad\square$

**Remark 2.1.** If $X$, $D$ and $2X - D$ are $M$-matrices, then

$$\mu_-((2 - \frac{1}{\hat{\lambda}})X - D) \geq \min_i\{((2 - \frac{1}{\hat{\lambda}})X - D)\mathbf{e}\}_i.$$

In particular, if $D$ is diagonal with constant diagonal, $D = dI$, then

$$\mu_-((2 - \frac{1}{\hat{\lambda}})X - D) \geq (2 - \frac{1}{\hat{\lambda}})x - d,$$

where $x$ is the smallest diagonal entry of $X$. Note that when $D$ is diagonal we can always scale $A$, i.e. consider $D^{-1/2}AD^{-1/2}$, where the scaled matrix has unit diagonal. We shall now derive an improved upper bound for the case where $\mu_-((2 - \frac{1}{\hat{\lambda}})X - D) \geq (2 - \frac{1}{\hat{\lambda}})x - d$. This will be done by finding the value of $\hat{\lambda}$ in (2.4) which minimizes the upper bound. It is readily seen that this value satisfies

$$2(1 - \frac{1}{\hat{\lambda}})\frac{x}{\mu_i} = \frac{d}{\mu_i} - 1$$

i.e.

$$\hat{\lambda} = 1/\left[1 - \frac{d - \mu_i}{2x}\right]$$

and that

$$\mu_-((2 - \frac{1}{\lambda})X - D) = x - \frac{d + \mu_i}{2}$$

for this value. Hence, if $\mu_i \leq 2x - d$ we find $\mu_-() \geq 0$ and the value of $\hat{\lambda}$ found gives the smallest upper bound $\overline{\lambda}_i$ of $\lambda_i$. This upper bound is

$$\overline{\lambda}_i = \frac{4x\mu_i(A)}{[2x - d + \mu_i(A)]^2}.$$

Further, if $\overline{\lambda}_i = \hat{\lambda}$, then for any $\mu_i(A)$, when

$$(2 - \frac{1}{\lambda})x - d = 0, \text{ i.e. } \hat{\lambda} = /(2 - \frac{d}{x})$$

we find

(2.8) $$\max_i \lambda_i \leq \hat{\lambda} = 1/(2 - \frac{d}{x}).$$

This latter value is hence the best upper bound in Theorem 2.1b when $\mu_-((2 - \frac{1}{\lambda})X - D) = (2 - \frac{1}{\lambda})x - d$.

Next we consider an application of the above results to estimate the condition number of the preconditioned iteration matrix $C^{-1}A$, when $A$ is a central difference matrix.

## 3. Application for an elliptic problem with anisotropy.

Consider the selfadjoint elliptic problem $-\delta u_{xx} - u_{yy} = f$ in $[0, 1]^2$, where $\delta > 0$, $a \geq 0$ and with Dirichlet boundary conditions, discretized by central difference approximations on a uniform mesh. Using a natural ordering, one finds

$$a_{i,i-n} = -1, \ a_{i,i-1} = -\delta, \ a_{i,i} = d, \ a_{i,i+1} = -\delta, \ a_{i,i+n} = -1,$$

where $d = 2(1 + \delta)$, and $h = 1/(n + 1)$.
For the entries of $X$ we find

$$x_i = d_i - \sum l_{ij} x_j^{-1} l_{ji}^t - \omega(R^0 e)_i, \ i = 1, 2, \ldots$$

or

$$x_i = 2(1 + \delta) - \delta^2 x_{i-1}^{-1} - x_{i-n}^{-1} - \omega\delta(x_{i-m}^{-1} + x_{i-1}^{-1})$$

(apart from corrections at points next to the boundary). We see readily that as $i \to \infty$ and $h \to 0$, $x_i$ converges to a lower bound $x$, where

$$x = 2(1 + \delta) - (1 + 2w\delta + \delta^2)/x$$

or

$$x = 1 + \delta + \{2\delta(1-\omega)\}^{1/2}$$

Then

$$2x - d = 2\{2\delta(1-\omega)\}^{1/2}$$

and

$$\mu_+(R) = 2\delta(1-\omega)/x, \ \mu_+(-R) = 2\delta(1+\omega)/x \quad (h \to 0).$$

Since we require $\mu_+(R) \geq 0$ and $\mu_+(-R) \geq 0$ we shall assume that $-1 \leq \omega \leq 1$.

Since $\mu_1(A) = (1+\delta)(2\sin \pi h/2)^2$, we find from Theorem 2.1 and (2.8), with $x = 1+\delta+\{2\delta(1-\omega)\}^{1/2}$

$$\lambda_1^{-1} \leq 1 + \frac{2\delta}{1+\delta} \frac{(1-\omega)}{x} \frac{1}{(2\sin \pi h/2)^2},$$

$$\max_i \lambda_i \leq \min \left\{ 1/(2 - \frac{d}{x}), \ 1 + \frac{2\delta(1+\omega)}{x} \frac{x}{(x-(1+\delta))^2} \right\}$$

or

$$\lambda_1^{-1} \leq 1 + 2\delta \cdot \frac{1-\omega}{1+\delta+\{2\delta(1-\omega)\}^{1/2}}(\mu_1)^{-1}$$

and

$$\max_i \lambda_i \leq \min \left\{ \tfrac{1}{2} + \frac{1+\delta}{2\{2\delta(1-\omega)\}^{1/2}} , \ \frac{2}{1-\omega} \right\}.$$

The condition number $\mathcal{H} = \mathcal{H}(\omega) = \max_i \lambda_i/\lambda_1$ is therefore bounded above by

$$\mathcal{H}(\omega) \leq \min \left\{ \tfrac{1}{2} + \frac{1+\delta}{2\{2\delta(1-\omega)\}^{1/2}} , \ \frac{2}{1-\omega} \right\} \left[ 1 + 2\delta \frac{1-\omega}{1+\delta+\{2\delta(1-\omega)\}^{1/2}}(\mu_1)^{-1} \right]$$

or

$$\mathcal{H}(\omega) \leq \min \left\{ \frac{1}{2} \left[ \frac{1+\delta}{\{2\delta(1-\omega)\}^{1/2}} + 1 + \{2\delta(1-\omega)\}^{1/2}(\mu_1)^{-1} \right] \right.$$
$$\left. \frac{2}{1-\omega} + 4\delta \frac{1}{1+\delta+\{2\delta(1-\omega)\}^{1/2}}(\mu_1)^{-1} \right\}, -1 \leq \omega \leq 1.$$

To minimize $\mathcal{H}(\omega)$, we need to choose

$$\omega = \omega_{\text{opt}} = 1 - \frac{1+\delta}{2\delta}\mu_1(A)$$

and

$$\omega = -1,$$

respectively, for the two functions in the outer bracket.

Hence

$$\min_\omega \mathcal{H}(\omega) = \min \left\{ \mathcal{H}(\omega_{\text{opt}}), 1 + \frac{4\delta}{1+\delta+2\delta^{1/2}}(\mu_1)^{-1} \right\}$$
$$= \min \left\{ \tfrac{1}{2} + \frac{1}{2\sin \frac{\pi h}{2}}, 1 + \frac{4\delta}{(1+\delta^{1/2})^2}(\mu_1)^{-1} \right\}$$

223

and we find

$$\min_{\omega} \mathcal{H}(\omega) = \begin{cases} \frac{1}{2\sin\frac{\pi h}{2}} + \frac{1}{2}, & \delta \gtrsim \frac{1}{4}\mu_1(A)^{1/2}, \text{ for } \omega = 1 - \frac{1+\delta}{2\delta}\mu_1(A) \\ 1 + \frac{4\delta}{(1+\delta^{1/2})^2}\mu_1(A)^{-1}, & \delta \lesssim \frac{1}{4}\mu_1(A)^{1/2}, \text{ for } \omega = -1. \end{cases}$$

Note that as $\delta$ decreases, the optimal value of $\omega$ switches for $\delta \simeq \frac{1}{4}\mu_1(A)^{1/2}$ from a value slightly less than unity to the value -1.

We conclude that the spectral condition number is bounded above by

$$\frac{1}{2} + (\pi h)^{-1} \text{ for } \omega = 1 - \frac{1+\delta}{2\delta}\mu_1(A)$$

for any value of $\delta$, but for $\delta$ sufficiently small,

$$1 + \frac{4\delta}{(1+\delta^{1/2})^2}\mu_1(A)^{-1} , \text{ for } \omega = -1$$

gives a smaller upper bound.

### References.

[1] Axelsson, O., A generalized SSOR method, BIT, 12, 443-467, 1972.
[2] Axelsson, O., A class of iterative methods for finite element equations. Comput. Methods Appl. Mech. Eng. 9, 123-137, 1976.
[3] Axelsson, O., On iterative solution of elliptic difference equations on a mesh-connected array of processors, Int. J. High Speed Computing 1, 165-183, 1989.
[4] Axelsson, O., Barker, V.A., Finite Element Solution of Boundary Value Problems, Theory and Computation. Academic Press, Orlando, Fl., 1984.
[5] Axelsson, O., Lindskog, G., On the eigenvalue distribution of a class of preconditioning methods. Numer. Math. 48, 479-498, 1986.
[6] Gustafsson, I., A class of first order factorization methods, BIT 18, 142-156, 1978.
[7] Van der Vorst, H.A. The convergence behaviour of preconditioned conjugate gradient and conjugate gradient square methods, talk presented at the Conference of Preconditioned Conjugate Gradient methods, June 19-21, 1989, University of Nijmegen, The Netherlands.
[8] Wilkinson, J.H., The Algebraic Eigenvalue Problem, Clarendon Press, Oxford, 1965.