

Numerical Analysis – Lecture 2¹

Theorem 1.7 (The Householder–John theorem) *If the real symmetric matrices A and $A - B - B^\top$ are both positive definite and B is real then the spectral radius of $H = -(A - B)^{-1}B$ is strictly less than one.*

Proof Let λ be an eigenvalue of H , so $Hv = \lambda v$ holds where $v \neq 0$ is an eigenvector. (Note that both λ and v may have nonzero imaginary parts when H is not symmetric, e.g. in the Gauss–Seidel method.) The definition of H provides $-Bv = \lambda(A - B)v$, and the value of λ is different from one because A is nonsingular. Thus we deduce

$$\bar{v}^\top Bv = \frac{\lambda}{\lambda - 1} \bar{v}^\top Av. \quad (1.2)$$

Moreover, writing $v = v_R + iv_I$, where v_R and v_I are real, we find the identity $\bar{v}^\top Av = v_R^\top Av_R + v_I^\top Av_I$, so positive definiteness implies $\bar{v}^\top Av > 0$ and $\bar{v}^\top (A - B - B^\top)v > 0$. It follows from equation (1.2), $\bar{v}^\top Bv = \lambda/(\lambda - 1)\bar{v}^\top Av$, $\bar{v}^\top B^\top v = \bar{\lambda}/(\bar{\lambda} - 1)\bar{v}^\top Av$ and from the fact that B is real that the last inequality is the condition

$$0 < \bar{v}^\top Av - \bar{v}^\top Bv - \bar{v}^\top B^\top v = \left(1 - \frac{\lambda}{\lambda - 1} - \frac{\bar{\lambda}}{\bar{\lambda} - 1}\right) \bar{v}^\top Av = \frac{(1 - |\lambda|^2)\bar{v}^\top Av}{|\lambda - 1|^2}.$$

Now $\lambda \neq 1$ implies $|\lambda - 1|^2 > 0$. Hence, recalling that $\bar{v}^\top Av > 0$, we see that $1 - |\lambda|^2$ is positive. Therefore $|\lambda| < 1$ occurs for every eigenvalue of H as required. \square

Corollary 1.8 (Application to Example 1.2) *Both Jacobi and Gauss–Seidel methods converge when A is the matrix of Example 1.2.*

Proof Positive definiteness of the symmetric matrix A has been already established in Lemma 1.3. For Jacobi’s method, $A - B - B^\top$ is the same as A except that the signs of the off-diagonal elements are reversed. Therefore the proof of Lemma 1.3 shows too that $A - B - B^\top$ is positive definite: recall that the proof depended on the *modulus* of off-diagonal elements, not on their sign! Moreover, for the Gauss–Seidel method, $A - B - B^\top$ is just the diagonal part of A , all the off-diagonal elements being zero, so this matrix is also positive definite. Therefore Theorem 1.7 implies $\rho(H) < 1$ in both cases. It follows from Revision 1.5 that the corollary is true. \square

Technique 1.9 (Relaxation). It is often possible to improve the efficiency of Method 1.4 (simple iteration) by *relaxation*. Specifically, instead of letting $(A - B)x^{(k+1)} = -Bx^{(k)} + b$, $k = 0, 1, \dots$, we let

$$(A - B)\tilde{x}^{(k+1)} = -Bx^{(k)} + b \quad \text{and} \quad x^{(k+1)} = x^{(k)} + \omega(\tilde{x}^{(k+1)} - x^{(k)}), \quad k = 0, 1, \dots,$$

where ω is a real constant called the *relaxation parameter*. Note that $\omega = 1$ corresponds to the former, “unrelaxed” iteration.

Good choice of ω leads to small spectral radius of the iteration matrix: clearly, it should be less than one, but ideally it should be the least possible: the smaller the spectral radius, the faster the iteration converges. To choose ω , we need to determine the iteration matrix \tilde{H} . First, we relate $x^{(k+1)}$ to $x^{(k)}$ by eliminating $\tilde{x}^{(k+1)}$ from the last displayed equation. Multiplying the equation for $x^{(k+1)}$ by $A - B$ we obtain

$$\begin{aligned} (A - B)x^{(k+1)} &= (A - B)[(1 - \omega)x^{(k)} + \omega\tilde{x}^{(k+1)}] \\ &= (1 - \omega)(A - B)x^{(k)} + \omega(-Bx^{(k)} + b) \\ &= [(1 - \omega)A - B]x^{(k)} + \omega b. \end{aligned}$$

¹Please email all corrections and suggestions to these notes to A. Iserles@damtp.cam.ac.uk. All handouts are available on the WWW at the URL <http://www.damtp.cam.ac.uk/user/na/PartII/Handouts.html>.

Thus, the iteration matrix is

$$\tilde{H} = (A - B)^{-1}[(1 - \omega)A - B] = I - \omega(A - B)^{-1}A.$$

Recall the ‘unrelaxed’ iteration matrix $H = -(A - B)^{-1}B = (A - B)^{-1}(A - B - A) = I - (A - B)^{-1}A$. Substituting $(A - B)^{-1}A = I - H$, we deduce that

$$\tilde{H} = I - \omega(I - H) = (1 - \omega)I + \omega H. \quad (1.3)$$

Suppose that $\omega \neq 0$. Then (1.3) proves that

$$\lambda \in \sigma(\tilde{H}) \quad \Leftrightarrow \quad \lambda = 1 - \omega + \omega\mu, \quad \mu \in \sigma(H),$$

where $\sigma(C)$ is the set of the eigenvalues (the *spectrum*) of the square matrix C . Therefore one may try to choose $\omega \in \mathbb{R} \setminus \{0\}$ to minimize

$$\rho(\tilde{H}) = \max\{|1 - \omega + \omega\mu| : \mu \in \sigma(H)\}.$$

In general, $\sigma(H)$ is unknown, but often we have some information about it which can be utilized to find a ‘good’ (rather than ‘best’) value of ω . For example, suppose that it is known that $\sigma(H)$ is real and resides in the interval $[\alpha, \beta]$, where $-1 < \alpha < \beta < 1$. In that case we seek ω to minimize

$$\max\{|1 - \omega + \omega\mu| : \mu \in [\alpha, \beta]\}.$$

Since maxima of the function above occur at endpoints, for optimal ω we have $|1 - \omega + \omega\alpha| = |1 - \omega + \omega\beta|$, and this is satisfied by $\omega_{\text{opt}} = 2/(2 - \alpha - \beta)$. (You can easily prove that $\omega_{\text{opt}} \in (0, 2)$.)

Approach 1.10 (An optimization calculation). We continue to assume that A is symmetric and positive definite. Therefore the quadratic function

$$F(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top A\mathbf{x} - \mathbf{b}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n, \quad (1.4)$$

is bounded below, and its least value occurs when \mathbf{x} satisfies $\nabla F(\mathbf{x}) = \mathbf{0}$, which is equivalent to \mathbf{x} being a solution of the system $A\mathbf{x} = \mathbf{b}$ of Problem 1.1. Therefore, when an iterative method generates the sequence $\mathbf{x}^{(k+1)}$, $k = 0, 1, 2, \dots$, it may be helpful to force the condition $F(\mathbf{x}^{(k+1)}) < F(\mathbf{x}^{(k)})$ for every $k \in \mathbb{Z}_+$. This remark can provide an alternative useful way of choosing ω in Technique 1.9, especially if ω is allowed to depend on k .

We now turn to algorithms of the following form. We pick any starting vector $\mathbf{x}^{(0)} \in \mathbb{R}^n$. For $k = 0, 1, 2, \dots$, the calculation stops if $\|\nabla F(\mathbf{x}^{(k)})\| = \|A\mathbf{x}^{(k)} - \mathbf{b}\|$ is acceptably small. Otherwise, a *search direction* $\mathbf{d}^{(k)}$ is generated that satisfies the *descent condition* $[\mathrm{d}F(\mathbf{x}^{(k)} + \omega\mathbf{d}^{(k)})/\mathrm{d}\omega]_{\omega=0} < 0$. Then the value of ω that minimizes $F(\mathbf{x}^{(k)} + \omega\mathbf{d}^{(k)})$, $\omega > 0$, is calculated, and we call it $\omega^{(k)}$. Finally, the k th iteration sets $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \omega^{(k)}\mathbf{d}^{(k)}$. Thus the strict inequalities $F(\mathbf{x}^{(k+1)}) < F(\mathbf{x}^{(k)})$ and $\omega^{(k)} > 0$ are achieved.

There is a convenient form of the descent condition that has been mentioned. Specifically, because the definition (1.4) implies the identity

$$F(\mathbf{x}^{(k)} + \omega\mathbf{d}^{(k)}) = F(\mathbf{x}^{(k)}) + \omega\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} + \frac{1}{2}\omega^2 \mathbf{d}^{(k)\top} A\mathbf{d}^{(k)}, \quad \omega \in \mathbb{R}, \quad (1.5)$$

where $\mathbf{g}^{(k)} = \nabla F(\mathbf{x}^{(k)})$, the search direction has to satisfy $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} < 0$, which is possible, because termination occurs when $\mathbf{g}^{(k)}$ is zero. Further, $\omega^{(k)}$ is the ω that minimizes the quadratic equation (1.5), so it has the value

$$\omega^{(k)} = -\frac{\mathbf{d}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} A\mathbf{d}^{(k)}}.$$