

Mathematical Tripos Part II: Michaelmas Term 2023

Numerical Analysis – Lecture 15

4 Iterative methods for linear systems

A general *iterative* method for solving $A\mathbf{x} = \mathbf{b}$ is a rule $\mathbf{x}^{k+1} = f_k(\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^k)$. We will consider the simplest ones: *linear, one-step, stationary* iterative schemes:

$$\mathbf{x}^{k+1} = H\mathbf{x}^k + \mathbf{v}, \quad \mathbf{x}^0, \mathbf{v} \in \mathbb{R}^n. \quad (4.1)$$

Here one chooses H and \mathbf{v} so that \mathbf{x}^* , a solution of $A\mathbf{x} = \mathbf{b}$, satisfies $\mathbf{x}^* = H\mathbf{x}^* + \mathbf{v}$, i.e. it is the fixed point of the iteration (4.1) (if the scheme converges). Standard terminology:

$$\text{the iteration matrix } H, \quad \text{the error } \mathbf{e}^k := \mathbf{x}^* - \mathbf{x}^k, \quad \text{the residual } \mathbf{r}^k := A\mathbf{e}^k = \mathbf{b} - A\mathbf{x}^k.$$

For a given class of matrices A (e.g. positive definite matrices, or even a single particular matrix), we are interested in *convergent* methods, i.e. the methods such that $\mathbf{x}^k \rightarrow \mathbf{x}^* = A^{-1}\mathbf{b}$ for every starting value \mathbf{x}^0 . Subtracting $\mathbf{x}^* = H\mathbf{x}^* + \mathbf{v}$ from (4.1) we obtain

$$\mathbf{e}^{k+1} = H\mathbf{e}^k = \dots = H^{k+1}\mathbf{e}^0, \quad (4.2)$$

i.e., a method is convergent if $\mathbf{e}^k = H^k\mathbf{e}^0 \rightarrow 0$ for any $\mathbf{e}^0 \in \mathbb{R}^n$.

Scheme 4.1 (Iterative refinement) This is the scheme

$$\mathbf{x}^{k+1} = \mathbf{x}^k - S(A\mathbf{x}^k - \mathbf{b}).$$

If $S = A^{-1}$, then $\mathbf{x}^{k+1} = A^{-1}\mathbf{b} = \mathbf{x}^*$, so it is suggestive to choose S as an approximation to A^{-1} . The iteration matrix for this scheme is $H_S = I - SA$.

Scheme 4.2 (Splitting) We assume $A = B + C$ in such a way that solving a linear system with the matrix C is “easy”. We consider the scheme which can be written as $B\mathbf{x}^k + C\mathbf{x}^{k+1} = \mathbf{b}$, i.e., eliminating C

$$(A - B)\mathbf{x}^{k+1} = -B\mathbf{x}^k + \mathbf{b},$$

with the iteration matrix $H = -(A - B)^{-1}B$. Any splitting can be viewed as an iterative refinement (and vice versa) because

$$\begin{aligned} (A - B)\mathbf{x}^{k+1} = -B\mathbf{x}^k + \mathbf{b} &\Leftrightarrow (A - B)\mathbf{x}^{k+1} = (A - B)\mathbf{x}^k - (A\mathbf{x}^k - \mathbf{b}) \\ &\Leftrightarrow \mathbf{x}^{k+1} = \mathbf{x}^k - (A - B)^{-1}(A\mathbf{x}^k - \mathbf{b}), \end{aligned}$$

so we should seek a splitting such that $S = (A - B)^{-1}$ approximates A^{-1} .

Theorem 4.3 Let $H \in \mathbb{R}^{n \times n}$. Then $\lim_{k \rightarrow \infty} H^k \mathbf{z} = 0$ for any $\mathbf{z} \in \mathbb{R}^n$ if and only if $\rho(H) < 1$.

Proof. 1) Let λ be an eigenvalue of (the real) H , real or complex, such that $|\lambda| = \rho(H) \geq 1$, and let \mathbf{w} be a corresponding eigenvector, i.e., $H\mathbf{w} = \lambda\mathbf{w}$. Then $H^k\mathbf{w} = \lambda^k\mathbf{w}$, and

$$\|H^k\mathbf{w}\|_\infty = |\lambda|^k \|\mathbf{w}\|_\infty \geq \|\mathbf{w}\|_\infty =: \gamma > 0. \quad (4.3)$$

If \mathbf{w} is real, we choose $\mathbf{z} = \mathbf{w}$, hence $\|H^k\mathbf{z}\|_\infty \geq \gamma$, and this cannot tend to zero.

If \mathbf{w} is complex, then $\mathbf{w} = \mathbf{u} + i\mathbf{v}$ with some real vectors \mathbf{u}, \mathbf{v} . But then at least one of the sequences $(H^k\mathbf{u}), (H^k\mathbf{v})$ does not tend to zero. For if both do, then also $H^k\mathbf{w} = H^k\mathbf{u} + iH^k\mathbf{v} \rightarrow 0$, and this contradicts (4.3).

2) Now, let $\rho(H) < 1$, and assume for simplicity that H possesses n linearly independent eigenvectors (\mathbf{w}_j) such that $H\mathbf{w}_j = \lambda_j\mathbf{w}_j$. Linear independence means that every $\mathbf{z} \in \mathbb{R}^n$ can be expressed as a linear combination of the eigenvectors, i.e., there exist $(c_j) \in \mathbb{C}$ such that $\mathbf{z} = \sum_{j=1}^n c_j\mathbf{w}_j$. Thus,

$$H^k\mathbf{z} = \sum_{j=1}^n c_j\lambda_j^k\mathbf{w}_j,$$

and since $|\lambda_j| \leq \rho(H) < 1$ we have $\lim_{k \rightarrow \infty} H^k\mathbf{z} = 0$, as required. \square

Remark 4.4 The complete proof of case (2) of Theorem 4.3 exploits the so-called Jordan normal form of the matrix H , namely $H = SJS^{-1}$, where J is a block diagonal matrix consisting of the Jordan blocks,

$$J = \begin{bmatrix} \boxed{J_1} & & & \\ & \boxed{J_2} & & \\ & & \ddots & \\ & & & \boxed{J_r} \end{bmatrix}, \quad J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}, \quad J_i \in \mathbb{R}^{n_i \times n_i}, \quad \sum_i n_i = n.$$

To prove that $J_i^k \rightarrow 0$ if $|\lambda_i| < 1$ one should split $J_i = \lambda_i I + P$, notice that $P^m = 0$ for $m \geq n_i$, and evaluate the terms of expansion $(\lambda_i I + P)^k = \sum_{m=0}^{n_i-1} \binom{k}{m} \lambda_i^{k-m} P^m$.

Applying Theorem 4.3 to the error estimate (4.2), we arrive at the following statement.

Theorem 4.5 Let \mathbf{x}^* , a solution of $A\mathbf{x} = \mathbf{b}$, satisfy $\mathbf{x}^* = H\mathbf{x}^* + \mathbf{v}$ and we are given the scheme

$$\mathbf{x}^{k+1} = H\mathbf{x}^k + \mathbf{v}, \quad \mathbf{x}^0, \mathbf{v} \in \mathbb{R}^n. \quad (4.4)$$

Then $\mathbf{x}^k \rightarrow \mathbf{x}^*$ for any choice of \mathbf{x}^0 if and only if $\rho(H) < 1$.

Note: Of course, we would like to know not just convergence but the rate of it. For example, we achieve convergence with

$$H = \begin{bmatrix} 0.99 & 10^6 \\ 0 & 0.99 \end{bmatrix},$$

but it will take quite a long time. We will discuss this topic briefly later on.

Method 4.6 (Jacobi and Gauss–Seidel) Both of these methods are versions of splitting which can be applied to any A with nonzero diagonal elements. We write A as the sum of three matrices $L_0 + D + U_0$: subdiagonal (strictly lower-triangular), diagonal and superdiagonal (strictly upper-triangular) portions of A , respectively.

1) *Jacobi method.* We set $A - B = D$, the diagonal part of A , and we obtain the next iteration by solving the diagonal system

$$D\mathbf{x}^{k+1} = -(L_0 + U_0)\mathbf{x}^k + \mathbf{b}, \quad H_J = -D^{-1}(L_0 + U_0).$$

2) *Gauss–Seidel method.* We take $A - B = L_0 + D = L$, the lower-triangular part of A , and we generate the sequence $(\mathbf{x}^{(k)})$ by solving the triangular system

$$(L_0 + D)\mathbf{x}^{k+1} = -U_0\mathbf{x}^k + \mathbf{b}, \quad H_{GS} = -(L_0 + D)^{-1}U_0.$$

There is no need to invert $(L_0 + D)$, we calculate the components of $\mathbf{x}^{(k+1)}$ in sequence by forward substitution:

$$a_{ii}x_i^{k+1} = -\sum_{j<i} a_{ij}x_j^{k+1} - \sum_{j>i} a_{ij}x_j^k + b_i, \quad i = 1..n.$$

As we mentioned above, the sequence \mathbf{x}^k converges to solution of $A\mathbf{x} = \mathbf{b}$ if the spectral radius of the iteration matrix, $H_J = -D^{-1}(L_0 + U_0)$ or $H_{GS} = -(L_0 + D)^{-1}U_0$, respectively, is less than one. Our next goal is to prove that this is the case for two important classes of matrices A :

- a) diagonally dominant and b) positive definite matrices.

We start with recalling the simple, but very useful Gershgorin theorem.

Revision 4.7 (Gershgorin theorem) All eigenvalues of an $n \times n$ matrix A are contained in the union of the Gershgorin discs in the complex plane:

$$\sigma(A) \subset \cup_{i=1}^n \Gamma_i, \quad \Gamma_i := \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}, \quad r_i := \sum_{j \neq i} |a_{ij}|.$$

Definition 4.8 (Strictly diagonally dominant matrices) A matrix A is called strictly diagonally dominant by rows (resp. by columns) if

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad i = 1..n \quad (\text{resp. } |a_{jj}| > \sum_{i \neq j} |a_{ij}|, \quad j = 1..n).$$

From Gershgorin theorem, it follows that strictly diagonally dominant matrices are nonsingular.

Theorem 4.9 If A is strictly diagonally dominant (either by rows or columns), then both the Jacobi and the Gauss-Seidel methods converge.

Proof. Jacobi method: We have $H_J = -D^{-1}(A - D) = I - D^{-1}A$. The diagonal elements of H_J are all zero, and the sum of the off-diagonal entries on the i 'th row is $\sum_{j \neq i} |(H_J)_{ij}| = \sum_{j \neq i} |A_{ij}|/|A_{ii}| < 1$ if A is strictly diagonally dominant by rows. Applying Gershgorin's theorem to H_J , we get that all the eigenvalues of H_J have modulus < 1 , which is what we wanted. If A is strictly diagonally dominant by columns (instead of by rows), we get that $\rho(I - AD^{-1}) < 1$ using the same argument, and use the fact that $I - D^{-1}A$ and $I - AD^{-1}$ have the same eigenvalues (since $I - D^{-1}A = D^{-1}(I - AD^{-1})D$).

Gauss-Seidel: If λ is an eigenvalue of $H_{GS} = -(L_0 + D)^{-1}U_0$, then this means that $H_{GS} - \lambda I = -(L_0 + D)^{-1}U_0 - \lambda I$ is singular, which in turn implies that $U_0 + \lambda(L_0 + D)$ is singular. It is easy to see that if $A = L_0 + D + U_0$ is strictly diagonally dominant, then the same is true for $A_\lambda = U_0 + \lambda(L_0 + D)$ for all $|\lambda| \geq 1$, and in particular A_λ is nonsingular in this case. This implies that any eigenvalue λ of $-(L_0 + D)^{-1}U_0$ must satisfy $|\lambda| < 1$. This shows convergence of the Gauss-Seidel method. (Note: a similar argument can also be used for Jacobi.) \square