

## Mathematical Tripos Part II: Michaelmas Term 2024

## Numerical Analysis – Lecture 16

## 4 Iterative methods for linear algebraic systems

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

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

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

the *iteration matrix*  $H$ , the *error*  $e^k := x^* - x^k$ , the *residual*  $r^k := Ae^k = b - Ax^k$ .

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  $x^k \rightarrow x^* = A^{-1}b$  for every starting value  $x^0$ . Subtracting  $x^* = Hx^* + v$  from (4.1) we obtain

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

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

**Scheme 4.1 (Iterative refinement)** This is the scheme

$$x^{k+1} = x^k - S(Ax^k - b).$$

If  $S = A^{-1}$ , then  $x^{k+1} = A^{-1}b = x^*$ , 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)** This is the scheme

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

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

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

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

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

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

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

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

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

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

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

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