

### 5.3 Conditioning of Linear Systems

Two linear systems that look quite similar at the first glance, may, in fact, have a very different degree of sensitivity of their solutions to the errors committed when specifying the input data. This phenomenon can be observed already for the systems  $Ax = f$  of order two:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 &= f_1, \\ a_{21}x_1 + a_{22}x_2 &= f_2. \end{aligned} \quad (5.30)$$

With no loss of generality we will assume that the coefficients of system (5.30) are normalized:  $a_{i1}^2 + a_{i2}^2 = 1$ ,  $i = 1, 2$ . Geometrically, each individual equation of system (5.30) defines a straight line on the Cartesian plane  $(x_1, x_2)$ . Accordingly, the solution of system (5.30) can be interpreted as the intersection point of these two lines as shown in Figure 5.2.

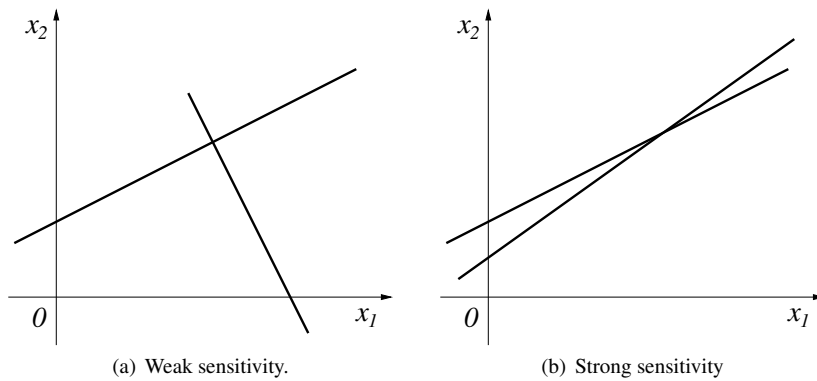


FIGURE 5.2: Sensitivity of the solution of system (5.30) to perturbations of the data.

Let us qualitatively analyze two opposite situations. The straight lines that correspond to linear equations (5.30) can intersect “almost normally,” as shown in Figure 5.2(a), or they can intersect “almost tangentially,” as shown in Figure 5.2(b). If we slightly perturb the input data, i.e., the right-hand sides  $f_i$ ,  $i = 1, 2$ , and/or the coefficients  $a_{ij}$ ,  $i, j = 1, 2$ , then each line may move parallel to itself and/or tilt. In doing so, it is clear that the intersection point (i.e., the solution) on Figure 5.2(a) will only move slightly, whereas the intersection point on Figure 5.2(b) will move much more visibly. Accordingly, one can say that in the case of Figure 5.2(a) the sensitivity of the solution to perturbations of the input data is weak, whereas in the case of Figure 5.2(b) it is strong.

Quantitatively, the sensitivity of the solution to perturbations of the input data can be characterized with the help of the so-called *condition number*  $\mu(A)$ . We will later

see that not only does the condition number determine the aforementioned sensitivity, but also that it directly influences the performance of iterative methods for solving  $\mathbf{Ax} = \mathbf{f}$ . Namely, it affects the number of iterations and as such, the number of arithmetic operations, required for finding an approximate solution of  $\mathbf{Ax} = \mathbf{f}$  within a prescribed tolerance, see Chapter 6.

### 5.3.1 Condition Number

The condition number of a linear operator  $\mathbf{A}$  acting on a normed vector space  $\mathbb{L}$  is defined as follows:

$$\mu(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|. \quad (5.31)$$

The same quantity  $\mu(\mathbf{A})$  given by formula (5.31) is also referred to as the condition number of a linear system  $\mathbf{Ax} = \mathbf{f}$ .

Recall that we have previously identified every matrix:

$$\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$$

with a linear operator acting on an  $n$ -dimensional vector space  $\mathbb{L}$  of the elements

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix} \quad (\text{for example, } \mathbb{L} = \mathbb{R}^n \text{ or } \mathbb{L} = \mathbb{C}^n). \text{ For a given } \mathbf{x} \in \mathbb{L}, \text{ the operator yields}$$

$$\mathbf{y} = \mathbf{Ax}, \text{ where } \mathbf{y} = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} \text{ is computed as follows:}$$

$$y_i = \sum_{j=1}^n a_{ij}x_j, \quad i = 1, 2, \dots, n.$$

Accordingly, the definition of the condition number  $\mu(\mathbf{A})$  by formula (5.31) also makes sense for matrices, so that one can refer to the condition number of a matrix  $\mathbf{A}$ , as well as to the condition number of a system of linear algebraic equations specified in its canonical form (5.1) rather than only in the operator form.

The norms  $\|\mathbf{A}\|$  and  $\|\mathbf{A}^{-1}\|$  of the direct and inverse operators, respectively, in formula (5.31) are assumed induced by the vector norm chosen in  $\mathbb{L}$ . Consequently, the condition number  $\mu(\mathbf{A})$  also depends on the choice of the norm in  $\mathbb{L}$ . If  $\mathbf{A}$  is a matrix, and we use the maximum norm for the vectors from  $\mathbb{L}$ ,  $\|\mathbf{x}\|_\infty = \max_j |x_j|$ , then we write  $\mu = \mu_\infty(\mathbf{A})$ ; if we employ the first norm  $\|\mathbf{x}\|_1 = \sum_j |x_j|$ , then  $\mu = \mu_1(\mathbf{A})$ .

If  $\mathbb{L}$  is a Euclidean (unitary) space with the scalar product  $(\mathbf{x}, \mathbf{y})$  given, for example, by formula (5.15) or (5.16), and the corresponding Euclidean (Hermitian) norm is defined by formula (5.17):  $\|\mathbf{x}\|_2 = (\mathbf{x}, \mathbf{x})^{1/2}$ , then  $\mu = \mu_2(\mathbf{A})$ . If an alternative scalar product  $[\mathbf{x}, \mathbf{y}]_{\mathbf{B}} = (\mathbf{B}\mathbf{x}, \mathbf{y})$  is introduced in  $\mathbb{L}$  based on the original product  $(\mathbf{x}, \mathbf{y})$  and on the operator  $\mathbf{B} = \mathbf{B}^* > 0$ , see formula (5.22), and if a new norm is set up accordingly by formula (5.23):  $\|\mathbf{x}\|_{\mathbf{B}} = ([\mathbf{x}, \mathbf{y}]_{\mathbf{B}})^{1/2}$ , then the corresponding condition number is denoted by  $\mu = \mu_{\mathbf{B}}(\mathbf{A})$ .

Let us now explain the geometric meaning of the condition number  $\mu(\mathbf{A})$ . To do so, consider the set  $S \subset \mathbb{L}$  of all vectors with the norm equal to one, i.e., a unit sphere in the space  $\mathbb{L}$ . Among these vectors choose a particular two,  $\mathbf{x}_{\max}$  and  $\mathbf{x}_{\min}$ , that satisfy the following equalities:

$$\|\mathbf{A}\mathbf{x}_{\max}\| = \max_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\| \quad \text{and} \quad \|\mathbf{A}\mathbf{x}_{\min}\| = \min_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\|.$$

It is easy to see that

$$\|\mathbf{A}\| = \|\mathbf{A}\mathbf{x}_{\max}\| \quad \text{and} \quad \|\mathbf{A}^{-1}\| = \frac{1}{\|\mathbf{A}\mathbf{x}_{\min}\|}.$$

Indeed, according to formula (5.24) we can write:

$$\|\mathbf{A}\| = \max_{\mathbf{x} \in \mathbb{L}, \mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\tilde{\mathbf{x}} \in \mathbb{L}, \tilde{\mathbf{x}} \neq \mathbf{0}} \left\| \mathbf{A} \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\| = \max_{\tilde{\mathbf{x}} \in S} \|\mathbf{A}\tilde{\mathbf{x}}\| = \|\mathbf{A}\mathbf{x}_{\max}\|.$$

Likewise:

$$\|\mathbf{A}^{-1}\| = \max_{\mathbf{x} \in \mathbb{L}, \mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}^{-1}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\tilde{\mathbf{x}} \in \mathbb{L}, \tilde{\mathbf{x}} \neq \mathbf{0}} \frac{\|\tilde{\mathbf{x}}\|}{\|\mathbf{A}\tilde{\mathbf{x}}\|} = \left[ \min_{\tilde{\mathbf{x}} \in \mathbb{L}, \tilde{\mathbf{x}} \neq \mathbf{0}} \frac{\|\mathbf{A}\tilde{\mathbf{x}}\|}{\|\tilde{\mathbf{x}}\|} \right]^{-1} = \frac{1}{\|\mathbf{A}\mathbf{x}_{\min}\|}.$$

Substituting these expressions into the definition of  $\mu(\mathbf{A})$  by means of formula (5.31), we obtain:

$$\mu(\mathbf{A}) = \frac{\max_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\|}{\min_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\|}. \quad (5.32)$$

This, in particular, implies that we always have:

$$\mu(\mathbf{A}) \geq 1.$$

According to formula (5.32), the condition number  $\mu(\mathbf{A})$  is the ratio of magnitudes of the maximally stretched unit vector to the minimally stretched (i.e., maximally shrunk) unit vector. If the operator  $\mathbf{A}$  is singular, i.e., if the inverse operator  $\mathbf{A}^{-1}$  does not exist, then we formally set  $\mu(\mathbf{A}) = \infty$ .

The geometric meaning of the quantity  $\mu(\mathbf{A})$  becomes particularly apparent in the case of an Euclidean space  $\mathbb{L} = \mathbb{R}^2$  equipped with the norm  $\|\mathbf{x}\|_2 = (\mathbf{x}, \mathbf{x})^{1/2} = \sqrt{x_1^2 + x_2^2}$ , and a self-adjoint operator  $\mathbf{A} = \mathbf{A}^*$ . In other words, let  $\mathbb{L}$  be the Euclidean plane of the variables  $x_1$  and  $x_2$ . In this case,  $S$  is a conventional unit circle:  $x_1^2 + x_2^2 = 1$ . A linear mapping by means of the operator  $\mathbf{A}$  obviously transforms this circle into an ellipse. Formula (5.32) then implies that the condition number  $\mu(\mathbf{A})$  is the ratio of the large semiaxis of this ellipse to its small semiaxis.

### THEOREM 5.3

Let  $\mathbf{A} = \mathbf{A}^*$  be an operator on the vector space  $\mathbb{L}$  self-adjoint in the sense of a given scalar product  $[\mathbf{x}, \mathbf{y}]_B$ . Then,

$$\mu_B(\mathbf{A}) = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}, \quad (5.33)$$

where  $\lambda_{\max}$  and  $\lambda_{\min}$  are the eigenvalues of  $\mathbf{A}$  with the largest and smallest moduli, respectively.

**PROOF** Let  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$  be an orthonormal basis in the  $n$ -dimensional space  $\mathbb{L}$  composed of the eigenvectors of  $\mathbf{A}$ . Orthonormality of the basis is understood in the sense of the scalar product  $[\mathbf{x}, \mathbf{y}]_{\mathbf{B}}$ . Let  $\lambda_j$  be the corresponding eigenvalues:  $\mathbf{A}\mathbf{e}_j = \lambda_j\mathbf{e}_j$ ,  $j = 1, 2, \dots, n$ , which are all real. Also assume with no loss of generality that the eigenvalues are arranged in the descending order:  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$ .

An arbitrary vector  $\mathbf{x} \in \mathbb{L}$  can be expanded with respect to this basis:

$$\mathbf{x} = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \dots + x_n\mathbf{e}_n.$$

In doing so,

$$\mathbf{A}\mathbf{x} = \lambda_1 x_1 \mathbf{e}_1 + \lambda_2 x_2 \mathbf{e}_2 + \dots + \lambda_n x_n \mathbf{e}_n,$$

and because of the orthonormality of the basis:

$$\|\mathbf{A}\mathbf{x}\|_{\mathbf{B}} = (|\lambda_1 x_1|^2 + |\lambda_2 x_2|^2 + \dots + |\lambda_n x_n|^2)^{1/2}.$$

Consequently,  $\|\mathbf{A}\mathbf{x}\|_{\mathbf{B}} \leq |\lambda_1| \|\mathbf{x}\|_{\mathbf{B}}$ , and if  $\|\mathbf{x}\|_{\mathbf{B}} = 1$ , i.e., if  $\mathbf{x} \in S$ , then  $\|\mathbf{A}\mathbf{x}\|_{\mathbf{B}} \leq |\lambda_1|$ . At the same time, for  $\mathbf{x} = \mathbf{e}_1 \in S$  we have  $\|\mathbf{A}\mathbf{e}_1\|_{\mathbf{B}} = |\lambda_1|$ . Therefore,

$$\max_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\|_{\mathbf{B}} = |\lambda_1| = |\lambda_{\max}|.$$

A similar argument immediately yields:

$$\min_{\mathbf{x} \in S} \|\mathbf{A}\mathbf{x}\|_{\mathbf{B}} = |\lambda_n| = |\lambda_{\min}|.$$

Then, formula (5.32) implies (5.33).  $\square$

Note that similarly to Theorem 5.2, the result of Theorem 5.3 can also be generalized to the case of normal matrices  $\mathbf{A}\mathbf{A}^* = \mathbf{A}^*\mathbf{A}$ .

### 5.3.2 Characterization of a Linear System by Means of Its Condition Number

#### THEOREM 5.4

Let  $\mathbb{L}$  be a normed vector space (e.g.,  $\mathbb{L} = \mathbb{R}^n$  or  $\mathbb{L} = \mathbb{C}^n$ ), and let  $\mathbf{A} : \mathbb{L} \mapsto \mathbb{L}$  be a non-singular linear operator. Consider a system of linear algebraic equations:

$$\mathbf{A}\mathbf{x} = \mathbf{f}, \tag{5.34}$$

where  $\mathbf{x} \in \mathbb{L}$  is the vector of unknowns and  $\mathbf{f} \in \mathbb{L}$  is a given right-hand side. Let  $\Delta\mathbf{f} \in \mathbb{L}$  be a perturbation of the right-hand side that leads to the perturbation  $\Delta\mathbf{x} \in \mathbb{L}$  of the solution so that:

$$\mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{f} + \Delta\mathbf{f}. \quad (5.35)$$

Then the relative error of the solution  $\|\Delta\mathbf{x}\|/\|\mathbf{x}\|$  satisfies the inequality:

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \mu(\mathbf{A}) \frac{\|\Delta\mathbf{f}\|}{\|\mathbf{f}\|}, \quad (5.36)$$

where  $\mu(\mathbf{A})$  is the condition number of the operator  $\mathbf{A}$ , see formula (5.31). Moreover, there are particular  $\mathbf{f} \in \mathbb{L}$  and  $\Delta\mathbf{f} \in \mathbb{L}$  for which (5.36) transforms into a precise equality.

**PROOF** Formulae (5.34) and (5.35) imply that  $\mathbf{A}\Delta\mathbf{x} = \Delta\mathbf{f}$ , and consequently,  $\Delta\mathbf{x} = \mathbf{A}^{-1}\Delta\mathbf{f}$ . Let us also employ the expression  $\mathbf{A}\mathbf{x} = \mathbf{f}$  that defines the original system itself. Then,

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} = \frac{\|\mathbf{A}^{-1}\Delta\mathbf{f}\|}{\|\mathbf{x}\|} = \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \frac{\|\mathbf{A}^{-1}\Delta\mathbf{f}\|}{\|\Delta\mathbf{f}\|} \frac{\|\Delta\mathbf{f}\|}{\|\mathbf{A}\mathbf{x}\|} = \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \frac{\|\mathbf{A}^{-1}\Delta\mathbf{f}\|}{\|\Delta\mathbf{f}\|} \frac{\|\Delta\mathbf{f}\|}{\|\mathbf{f}\|}. \quad (5.37)$$

According to the definition of the operator norm, see formula (5.24), we have:

$$\frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \quad \text{and} \quad \frac{\|\mathbf{A}^{-1}\Delta\mathbf{f}\|}{\|\Delta\mathbf{f}\|} \leq \|\mathbf{A}^{-1}\|. \quad (5.38)$$

Combining formulae (5.37) and (5.38), we obtain for any  $\mathbf{f} \in \mathbb{L}$  and  $\Delta\mathbf{f} \in \mathbb{L}$ :

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta\mathbf{f}\|}{\|\mathbf{f}\|} = \mu(\mathbf{A}) \frac{\|\Delta\mathbf{f}\|}{\|\mathbf{f}\|}, \quad (5.39)$$

which means that inequality (5.36) holds.

Furthermore, if  $\Delta\mathbf{f}$  is the specific vector from the space  $\mathbb{L}$  for which

$$\frac{\|\mathbf{A}^{-1}\Delta\mathbf{f}\|}{\|\Delta\mathbf{f}\|} = \|\mathbf{A}^{-1}\|,$$

and  $\mathbf{f} = \mathbf{A}\mathbf{x}$  is the element of  $\mathbb{L}$  for which

$$\frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \|\mathbf{A}\|,$$

then expression (5.37) coincides with inequalities (5.39) and (5.36) that transform into precise equalities for these particular  $\mathbf{f}$  and  $\Delta\mathbf{f}$ .  $\square$