

Optional reading about gaussian elimination, LU factorization, and the condition number of linear systems

This is an excerpt from Chapter 2 of

Scientific Computing: An Introductory Survey (2nd edition)
by Michael T. Heath

Chapter 2

Systems of Linear Equations

2.1 Linear Systems

Systems of linear algebraic equations arise in almost every aspect of applied mathematics and scientific computation. Such systems often occur naturally, but they are also frequently the result of approximating nonlinear equations by linear equations or differential equations by algebraic equations. We will see many examples of such approximations throughout this book. For these reasons, the efficient and accurate solution of linear systems forms the cornerstone of many numerical methods for solving a wide variety of practical computational problems.

In matrix-vector notation, a system of linear algebraic equations has the form

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where \mathbf{A} is an $m \times n$ matrix, \mathbf{b} is a given m -vector, and \mathbf{x} is the unknown solution n -vector to be determined. Such a system of equations asks the question, “Can the vector \mathbf{b} be expressed as a linear combination of the columns of the matrix \mathbf{A} ?” If so, the coefficients of this linear combination are given by the components of the solution vector \mathbf{x} . There may or may not be a solution; and if there is a solution, it may or may not be unique. In this chapter we will consider only *square* systems, which means that $m = n$, i.e., the matrix has the same number of rows and columns. In later chapters we will consider systems where $m \neq n$.

2.1.1 Singularity and Nonsingularity

An $n \times n$ matrix \mathbf{A} is said to be *singular* if it has any one of the following equivalent properties:

1. \mathbf{A} has no inverse (i.e, there is no matrix \mathbf{M} such that $\mathbf{AM} = \mathbf{MA} = \mathbf{I}$, the identity matrix).
2. $\det(\mathbf{A}) = 0$.

3. $\text{rank}(\mathbf{A}) < n$ (the *rank* of a matrix is the maximum number of linearly independent rows or columns it contains).
4. $\mathbf{Az} = \mathbf{o}$ for some vector $\mathbf{z} \neq \mathbf{o}$.

Otherwise, the matrix is *nonsingular*. The solvability of a system of linear equations $\mathbf{Ax} = \mathbf{b}$ is determined by whether the matrix \mathbf{A} is singular or nonsingular. If the matrix \mathbf{A} is nonsingular, then its inverse, denoted by \mathbf{A}^{-1} , exists, and the system $\mathbf{Ax} = \mathbf{b}$ always has a unique solution $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ regardless of the value for \mathbf{b} . If, on the other hand, the matrix \mathbf{A} is singular, then the *number* of solutions is determined by the right-hand-side vector \mathbf{b} : for a given value of \mathbf{b} there may be no solution, but if there is a solution \mathbf{x} , so that $\mathbf{Ax} = \mathbf{b}$, then we also have $\mathbf{A}(\mathbf{x} + \gamma\mathbf{z}) = \mathbf{b}$ for any scalar γ , where the vector \mathbf{z} is as in the foregoing definition. Thus, if a singular system has a solution, then the solution cannot be unique. To summarize the possibilities, for a given matrix \mathbf{A} and right-hand-side vector \mathbf{b} , the system may have

One solution:	nonsingular
No solution:	singular
Infinitely many solutions:	singular

In two dimensions, each linear equation determines a straight line in the plane. The solution of the system is the intersection point of the two lines. If the two straight lines are not parallel, then they have a unique intersection point (the nonsingular case). If the two straight lines are parallel, then either they do not intersect at all (there is no solution) or the two lines are the same (any point along the line is a solution). In higher dimensions, each equation determines a hyperplane. In the nonsingular case, the unique solution is the intersection point of all of the hyperplanes.

Example 2.1 Singularity and Nonsingularity. The 2×2 system

$$\begin{aligned} 2x_1 + 3x_2 &= b_1, \\ 5x_1 + 4x_2 &= b_2, \end{aligned}$$

or in matrix-vector notation

$$\begin{bmatrix} 2 & 3 \\ 5 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix},$$

is nonsingular regardless of the value of b . If $\mathbf{b} = [8 \ 13]^T$, for example, then the unique solution is $\mathbf{x} = [1 \ 2]^T$.

The 2×2 system

$$\begin{bmatrix} 2 & 3 \\ 4 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

is singular regardless of the value of \mathbf{b} . With $\mathbf{b} = [4 \ 7]^T$, there is no solution. With $\mathbf{b} = [4 \ 8]^T$, then

$$\mathbf{x} = \begin{bmatrix} \gamma \\ (4 - 2\gamma)/3 \end{bmatrix}$$

is a solution for any real number γ .

2.2 Solving Linear Systems

To solve a linear system, the general strategy outlined in Section 1.1.1 suggests that we should transform the system into one whose solution is the same as that of the original system but is easier to compute. What type of transformation of a linear system leaves the solution unchanged? The answer is that we can premultiply (i.e., multiply from the left) both sides of the linear system $\mathbf{Ax} = \mathbf{b}$ by any nonsingular matrix \mathbf{M} without affecting the solution. To see why this is so, note that the solution to the linear system $\mathbf{MAx} = \mathbf{Mb}$ is given by

$$\mathbf{x} = (\mathbf{MA})^{-1}\mathbf{Mb} = \mathbf{A}^{-1}\mathbf{M}^{-1}\mathbf{Mb} = \mathbf{A}^{-1}\mathbf{b}.$$

Example 2.2 Permutations. An important example of such a transformation is the fact that the rows of \mathbf{A} and corresponding entries of \mathbf{b} can be reordered without changing the solution \mathbf{x} . This is intuitively obvious: all of the equations in the system must be satisfied simultaneously in any case, so the order in which they happen to be written down is irrelevant; they may as well have been drawn randomly from a hat. Formally, such a reordering of the rows is accomplished by premultiplying both sides of the equation by a *permutation* matrix \mathbf{P} , which is a square matrix having exactly one 1 in each row and column and zeros elsewhere (i.e., an identity matrix with its rows and columns permuted). For example,

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} b_3 \\ b_1 \\ b_2 \end{bmatrix}.$$

A permutation matrix is always nonsingular; in fact, its inverse is simply its transpose, $\mathbf{P}^{-1} = \mathbf{P}^T$ (the *transpose* of a matrix \mathbf{M} , denoted by \mathbf{M}^T , is a matrix whose columns are the rows of \mathbf{M} , that is, if $\mathbf{N} = \mathbf{M}^T$, then $n_{ij} = m_{ji}$). Thus, the reordered system can be written $\mathbf{PAx} = \mathbf{Pb}$, and the solution \mathbf{x} is unchanged.

Postmultiplying (i.e., multiplying from the right) by a permutation matrix reorders the columns of the matrix instead of the rows. Such a transformation *does* change the solution, but only in that the components of the solution are permuted. To see this, observe that the solution to the system $\mathbf{APx} = \mathbf{b}$ is given by

$$\mathbf{x} = (\mathbf{AP})^{-1}\mathbf{b} = \mathbf{P}^{-1}\mathbf{A}^{-1}\mathbf{b} = \mathbf{P}^T(\mathbf{A}^{-1}\mathbf{b}).$$

Example 2.3 Diagonal Scaling. Another simple but important type of transformation is *diagonal scaling*. Recall that a matrix \mathbf{D} is *diagonal* if $d_{ij} = 0$ for all $i \neq j$, that is, the only nonzero entries are $d_{ii}, i = 1, \dots, n$, on the *main diagonal*. Premultiplying both sides of a linear system $\mathbf{Ax} = \mathbf{b}$ by a nonsingular diagonal matrix \mathbf{D} multiplies each row of the matrix and right-hand side by the corresponding diagonal entry of \mathbf{D} , and hence is called *row scaling*. In principle, row scaling does not change the solution to the linear system, but in practice it can affect the numerical solution process and the accuracy that can be attained for a given problem, as we will see.

Column scaling—postmultiplying the matrix of a linear system by a nonsingular diagonal matrix \mathbf{D} —multiplies each column of the matrix by the corresponding diagonal entry of \mathbf{D} . Such a transformation does alter the solution, in effect changing the units in which the components of the solution are measured. The solution to the scaled system $\mathbf{AD}\mathbf{x} = \mathbf{b}$ is given by

$$\mathbf{x} = (\mathbf{AD})^{-1}\mathbf{b} = \mathbf{D}^{-1}\mathbf{A}^{-1}\mathbf{b},$$

and hence the solution to the original system is given by \mathbf{D} .

2.2.1 Triangular Linear Systems

The next question is what type of linear system is easy to solve. Suppose there is an equation in the system $\mathbf{Ax} = \mathbf{b}$ that involves only one of the unknown solution components (i.e., only one entry in that row of \mathbf{A} is nonzero). Then that equation can easily be solved (by division) for that unknown. Now suppose there is another equation in the system that involves only two unknowns, one of which is the one already determined. By substituting the one solution component already determined into this second equation, we can then easily solve for its other unknown. If this pattern continues, with only one new unknown component arising per equation, then all of the solution components can be computed in succession. A matrix with this special property is called *triangular*, for reasons that will soon become apparent. Because triangular linear systems are easily solved by this successive substitution process, they are a suitable target in transforming a general linear system.

Although the general triangular form just described is all that is required to enable the system to be solved by successive substitution, it is convenient to define two specific triangular forms for computational purposes. A matrix \mathbf{A} is *upper triangular* if all of its entries below the main diagonal are zero (i.e., if $a_{ij} = 0$ for $i > j$). Similarly, a matrix is *lower triangular* if all of its entries above the main diagonal are zero (i.e., if $a_{ij} = 0$ for $i < j$). For an upper triangular system $\mathbf{Ax} = \mathbf{b}$, the successive substitution process is called *back-substitution* and can be expressed as follows:

$$x_n = b_n/a_{nn},$$

$$x_i = \left(b_i - \sum_{j=i+1}^n a_{ij}x_j \right) / a_{ii}, \quad i = n-1, \dots, 1.$$

Similarly, for a lower triangular system $\mathbf{Ax} = \mathbf{b}$, the successive substitution process is called *forward-substitution* and can be expressed as follows:

$$x_1 = b_1/a_{11},$$

$$x_i = \left(b_i - \sum_{j=1}^{i-1} a_{ij}x_j \right) / a_{ii}, \quad i = 2, \dots, n.$$

A matrix that is triangular in the more general sense defined earlier can be permuted into upper or lower triangular form by a suitable permutation of its rows or columns.

Example 2.4 Triangular Linear System. Consider the upper triangular linear system

$$\begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 8 \end{bmatrix}.$$

The last equation, $4x_3 = 8$, can be solved directly for $x_3 = 2$. This value can then be substituted into the second equation to obtain $x_2 = 2$, and finally both x_3 and x_2 are substituted into the first equation to obtain $x_1 = -1$.

2.2.2 Elementary Elimination Matrices

Our strategy then is to devise a nonsingular linear transformation that transforms a given general linear system into a triangular linear system that we can then solve easily by successive substitution. Thus, we need a transformation that replaces selected nonzero entries of the given matrix with zeros. This can be accomplished by taking appropriate linear combinations of the rows of the matrix, as we will now show.

Consider the 2-vector $a = [a_1 \ a_2]^T$. If $a_1 \neq 0$, then

$$\begin{bmatrix} 1 & 0 \\ -a_2/a_1 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} a_1 \\ 0 \end{bmatrix}.$$

More generally, given an n -vector a , we can annihilate *all* of its entries below the k th position, provided that $a_k \neq 0$, by the following transformation:

$$\mathbf{M}_k \mathbf{a} = \begin{bmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & -m_{k+1} & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & -m_n & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_k \\ a_{k+1} \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_k \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

where $m_i = a_i/a_k$, $i = k + 1, \dots, n$. The divisor a_k is called the *pivot*. A matrix of this form is sometimes called an *elementary elimination matrix* or *Gauss transformation*, and its effect on a vector is to add a multiple of row k to each subsequent row, with the multipliers m_i chosen so that the result in each case is zero. Note the following about these elementary elimination matrices:

1. \mathbf{M}_k is a lower triangular matrix with unit main diagonal, and hence it must be nonsingular.
2. $\mathbf{M}_k = \mathbf{I} - \mathbf{m}e_k^T$, where $\mathbf{m} = [0, \dots, 0, m_{k+1}, \dots, m_n]^T$ and e_k is the k th column of the identity matrix.
3. $\mathbf{M}_k^{-1} = \mathbf{I} + \mathbf{m}e_k^T$, which means that \mathbf{M}_k^{-1} , which we will denote by \mathbf{L}_k , is the same as \mathbf{M}_k except that the signs of the multipliers are reversed.

4. If M_j , $j > k$, is another elementary elimination matrix, with vector of multipliers \mathbf{t} , then

$$M_k M_j = \mathbf{I} - \mathbf{m}e_k^T - \mathbf{t}e_j^T + \mathbf{m}e_k^T \mathbf{t}e_j^T = \mathbf{I} - \mathbf{m}e_k^T - \mathbf{t}e_j^T,$$

since $e_k^T \mathbf{t} = \mathbf{o}$. Thus, their product is essentially their “union.” Because they have the same form, a similar result holds for the product of their inverses, $L_k L_j$. Note that the order of multiplication is significant; these results do not hold for the reverse product.

Example 2.5 Elementary Elimination Matrices. If $\mathbf{a} = [2 \ 4 \ -2]^T$, then

$$M_1 \mathbf{a} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ -2 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix}, \quad \text{and} \quad M_2 \mathbf{a} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ -2 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 0 \end{bmatrix}.$$

We also note that

$$L_1 = M_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}, \quad L_2 = M_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{1}{2} & 1 \end{bmatrix},$$

and

$$M_1 M_2 = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & \frac{1}{2} & 1 \end{bmatrix}, \quad L_1 L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & -\frac{1}{2} & 1 \end{bmatrix}.$$

2.2.3 Gaussian Elimination and LU Factorization

With elementary elimination matrices, it is a fairly simple matter to reduce a general linear system $\mathbf{Ax} = \mathbf{b}$ to upper triangular form. We first choose an elementary elimination matrix M_1 according to the recipe given in Section 2.2.2, with the first diagonal entry a_{11} as pivot, so that, when premultiplied by M_1 , the first column of \mathbf{A} becomes zero below the first row. Of course, all of the remaining columns of \mathbf{A} , as well as the right-hand-side vector \mathbf{b} , are also multiplied by M_1 , so the new system becomes $M_1 \mathbf{Ax} = M_1 \mathbf{b}$, but by our previous discussion the solution is unchanged.

Next we use the second diagonal entry as pivot to determine a second elementary elimination matrix M_2 that annihilates all of the entries of the second column of the new matrix, $M_1 \mathbf{A}$, below the second row. Again, M_2 must be applied to the entire matrix and right-hand-side vector, so that we obtain the further modified linear system $M_2 M_1 \mathbf{Ax} = M_2 M_1 \mathbf{b}$. Note that the first column of the matrix $M_1 \mathbf{A}$ is not affected by M_2 because all of its entries are zero in the relevant rows. This process is continued for each successive column until all of the subdiagonal entries of the matrix have been annihilated, so that the linear system

$$M \mathbf{Ax} = M_{n-1} \cdots M_1 \mathbf{Ax} = M_{n-1} \cdots M_1 \mathbf{b} = M \mathbf{b}$$

is upper triangular and can be solved by back-substitution to obtain the solution to the original linear system $\mathbf{Ax} = \mathbf{b}$.

The process we have just described is known as *Gaussian elimination*. It is also known as *LU factorization* or *LU decomposition* because it decomposes the matrix \mathbf{A} into a product of a unit lower triangular matrix, \mathbf{L} , and an upper triangular matrix, \mathbf{U} . To see this, recall that the product $\mathbf{L}_k\mathbf{L}_j$ is unit lower triangular if $k < j$, so that

$$\mathbf{L} = \mathbf{M}^{-1} = (\mathbf{M}_{n-1} \cdots \mathbf{M}_1)^{-1} = \mathbf{M}_1^{-1} \cdots \mathbf{M}_{n-1}^{-1} = \mathbf{L}_1 \cdots \mathbf{L}_{n-1}$$

is unit lower triangular. We have already seen that, by design, the matrix $\mathbf{U} = \mathbf{M}\mathbf{A}$ is upper triangular. Therefore, we have expressed \mathbf{A} as a product

$$\mathbf{A} = \mathbf{L}\mathbf{U},$$

where \mathbf{L} is unit lower triangular and \mathbf{U} is upper triangular. Given such a factorization, the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ can then be written as $\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{b}$ and hence can be solved by first solving the lower triangular system $\mathbf{L}\mathbf{y} = \mathbf{b}$ by forward-substitution, then the upper triangular system $\mathbf{U}\mathbf{x} = \mathbf{y}$ by back-substitution. Note that the intermediate solution \mathbf{y} is the same as the transformed right-hand-side vector, $\mathbf{M}\mathbf{b}$, in the previous formulation. Thus, Gaussian elimination and LU factorization are simply two ways of expressing the same solution process.

Example 2.6 Gaussian Elimination. We illustrate Gaussian elimination by solving the linear system

$$\begin{aligned} 2x_1 + 4x_2 - 2x_3 &= 2, \\ 4x_1 + 9x_2 - 3x_3 &= 8, \\ -2x_1 - 3x_2 + 7x_3 &= 10, \end{aligned}$$

or in matrix notation

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 8 \\ 10 \end{bmatrix} = \mathbf{b}.$$

To annihilate the subdiagonal entries of the first column of \mathbf{A} , we subtract two times the first row from the second row, and add the first row to the third row:

$$\begin{aligned} \mathbf{M}_1\mathbf{A} &= \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} = \begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 1 & 5 \end{bmatrix}, \\ \mathbf{M}_1\mathbf{b} &= \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 8 \\ 10 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 12 \end{bmatrix}. \end{aligned}$$

Now to annihilate the subdiagonal entry of the second column of $\mathbf{M}_1\mathbf{A}$, we subtract the second row from the third row:

$$\mathbf{M}_2\mathbf{M}_1\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 1 & 5 \end{bmatrix} = \begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 0 & 4 \end{bmatrix},$$

$$\mathbf{M}_2\mathbf{M}_1\mathbf{b} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \\ 12 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 8 \end{bmatrix}.$$

We have therefore reduced the original system to the equivalent upper triangular system

$$\begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \\ 8 \end{bmatrix},$$

which can now be solved by back-substitution (as in Example 2.4) to obtain $\mathbf{x} = [-1 \ 2 \ 2]^T$. To write out the LU factorization explicitly, we have

$$\mathbf{L} = \mathbf{L}_1\mathbf{L}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix},$$

so that

$$\mathbf{A} = \begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & -2 \\ 0 & 1 & 1 \\ 0 & 0 & 4 \end{bmatrix} = \mathbf{LU}.$$

2.2.4 Pivoting

There is one obvious problem with the Gaussian elimination process as we have described it, as well as another, somewhat more subtle, problem. The obvious potential difficulty is that the process breaks down if the leading diagonal entry of the remaining unreduced portion of the matrix is zero at any stage, as computing the multipliers m_i for a given column requires division by the diagonal entry in that column. The solution to this problem is almost equally obvious: if the diagonal entry is zero at stage k , then interchange row k of the system with some subsequent row whose entry in column k is nonzero. As we know from Example 2.2, such an interchange does not alter the solution to the system. With a nonzero diagonal entry as pivot, the process can then proceed as usual.

But what if there is no nonzero entry on or below the diagonal in column k ? Then there is nothing to do at this stage, since all the entries to be annihilated are already zero, and we can simply move on to the next column (i.e., $\mathbf{M}_k = \mathbf{I}$). Note that this step leaves a zero on the diagonal, and hence the resulting upper triangular matrix \mathbf{U} is singular, but the LU factorization can still be completed. It does mean, however, that the subsequent back-substitution process will fail, since it requires a division by each diagonal entry of \mathbf{U} , but this is not surprising because the original matrix must have been singular anyway. A more insidious problem is that in floating-point arithmetic we may not get an exact zero, but only a very small diagonal entry, which brings us to the more subtle point.

In principle, any nonzero value will do as the pivot for computing the multipliers, but in practice the choice should be made with some care to minimize error. When the remaining portion of the matrix is multiplied by the resulting elementary elimination matrix, we should try to limit the growth of the entries of the transformed matrix in order not to

amplify rounding errors. For this reason, it is desirable for the multipliers not to exceed 1 in magnitude. This requirement can be met by choosing the entry of largest magnitude on or below the diagonal as pivot. Such a policy is called *partial pivoting*, and it is essential in practice for a numerically stable implementation of Gaussian elimination for general linear systems.

The row interchanges required by partial pivoting slightly complicate the formal description of LU factorization given earlier. In particular, each elementary elimination matrix M_k is preceded by a permutation matrix P_k that interchanges rows to bring the entry of largest magnitude into the diagonal pivot position. We still have $MA = U$, where U is upper triangular, but now

$$M = M_{n-1}P_{n-1} \cdots M_1P_1.$$

M^{-1} is still triangular in the general sense defined earlier, but because of the permutations, M^{-1} is not necessarily *lower* triangular, though we still denote it by L . Thus, “LU” factorization no longer literally means “lower times upper” triangular, but it is still equally useful for solving linear systems by successive substitution.

We note that the permutation matrix

$$P = P_{n-1} \cdots P_1$$

permutes the rows of A into the order determined by partial pivoting. An alternative interpretation, therefore, is to think of partial pivoting as a way of determining a row ordering for the system under which no pivoting would be required for numerical stability. Thus, we obtain the factorization

$$PA = LU,$$

where now L really is lower triangular. To solve the linear system $Ax = b$, we first solve the lower triangular system $Ly = Pb$ by forward-substitution, then the upper triangular system $Ux = y$ by back-substitution.

The name “partial” pivoting comes from the fact that only the current column is searched for a suitable pivot. A more exhaustive pivoting strategy is *complete pivoting*, in which the entire remaining unreduced submatrix is searched for the largest entry, which is then permuted into the diagonal pivot position. Note that this requires interchanging columns as well as rows, and hence it leads to a factorization of the form

$$PAQ = LU,$$

where L is unit lower triangular, U is upper triangular, and P and Q are permutation matrices that reorder the rows and columns, respectively, of A . To solve the linear system $Ax = b$, we first solve the lower triangular system $Ly = Pb$ by forward-substitution, then the upper triangular system $Uz = y$ by back-substitution, and finally we permute the solution components to obtain $x = Qz$. Although the numerical stability of complete pivoting is theoretically superior, it requires a much more expensive pivot search than partial pivoting. Because the numerical stability of partial pivoting is more than adequate in practice, it is almost universally used in solving linear systems by Gaussian elimination.

Since pivot selection depends on the magnitudes of individual matrix entries, the particular choice obviously depends on the scaling of the matrix. A diagonal scaling of the

matrix (recall Example 2.3) may result in a different sequence of pivots. For example, any nonzero entry in a given column can be made the largest in magnitude simply by giving that row a sufficiently heavy weighting. This does not mean that an arbitrary pivot sequence is acceptable, however: a badly skewed scaling can result in an inherently sensitive system and a correspondingly inaccurate solution. A well-formulated problem should have appropriately commensurate units for measuring the unknown variables (column scaling), and a weighting of the individual equations that properly reflects their relative importance (row scaling). It should also account for the relative accuracy of the input data. Under these circumstances, the pivoting procedure will usually produce a solution that is as accurate as the problem warrants (see Section 2.4).

Example 2.7 Pivoting. Here are some examples to illustrate the necessity of pivoting, both in theory and practice, for a stable implementation of Gaussian elimination. We first observe that the need for pivoting has nothing to do with whether the matrix is singular or nearly singular. For example, the matrix

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

is nonsingular yet has no LU factorization unless we interchange rows, whereas the singular matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

does have an LU factorization.

In practice, using finite-precision arithmetic, we must avoid not only zero pivots but also *small* pivots in order to prevent unacceptable error growth, as shown in the following example. Let

$$\mathbf{A} = \begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix},$$

where ϵ is a positive number smaller than the unit roundoff ϵ_{mach} in a given floating-point system. If we do not interchange rows, then the pivot is ϵ and the resulting multiplier is $-1/\epsilon$, so that we get the elimination matrix

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ -1/\epsilon & 1 \end{bmatrix},$$

and hence

$$\mathbf{L} = \begin{bmatrix} 1 & 0 \\ 1/\epsilon & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{U} = \begin{bmatrix} \epsilon & 1 \\ 0 & 1 - 1/\epsilon \end{bmatrix} = \begin{bmatrix} \epsilon & 1 \\ 0 & -1/\epsilon \end{bmatrix}$$

in floating-point arithmetic. But then

$$\mathbf{LU} = \begin{bmatrix} 1 & 0 \\ 1/\epsilon & 1 \end{bmatrix} \begin{bmatrix} \epsilon & 1 \\ 0 & -1/\epsilon \end{bmatrix} = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix} \neq \mathbf{A}.$$

Using a small pivot, and a correspondingly large multiplier, has caused an unrecoverable loss of information in the transformed matrix. If we interchange rows, on the other hand, then the pivot is 1 and the resulting multiplier is $-\epsilon$, so that we get the elimination matrix

$$\mathbf{M} = \begin{bmatrix} 1 & 0 \\ -\epsilon & 1 \end{bmatrix},$$

and hence

$$\mathbf{L} = \begin{bmatrix} 1 & 0 \\ \epsilon & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{U} = \begin{bmatrix} 1 & 1 \\ 0 & 1 - \epsilon \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

in floating-point arithmetic. We therefore have

$$\mathbf{LU} = \begin{bmatrix} 1 & 0 \\ \epsilon & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix},$$

which is the correct result after permutation.

Although the foregoing example is rather extreme, the principle holds in general that larger pivots produce smaller multipliers and hence smaller errors. In particular, if the largest entry on or below the diagonal in each column is used as pivot (partial pivoting), then the multipliers are bounded in magnitude by 1. In Example 2.6, we did not use row interchanges, and some of the multipliers were greater than 1. For illustration, we now repeat that example, this time using partial pivoting. The system in Example 2.6 is

$$\begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 8 \\ 10 \end{bmatrix}.$$

The largest entry in the first column is 4, so we interchange the first two rows using the permutation matrix

$$\mathbf{P}_1 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

obtaining the permuted system

$$\mathbf{P}_1 \mathbf{A} \mathbf{x} = \begin{bmatrix} 4 & 9 & -3 \\ 2 & 4 & -2 \\ -2 & -3 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 8 \\ 2 \\ 10 \end{bmatrix} = \mathbf{P}_1 \mathbf{b}.$$

To annihilate the subdiagonal entries of the first column, we use the elimination matrix

$$\mathbf{M}_1 = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 0 & 1 \end{bmatrix},$$

obtaining the transformed system

$$\mathbf{M}_1 \mathbf{P}_1 \mathbf{A} \mathbf{x} = \begin{bmatrix} 4 & 9 & -3 \\ 0 & -\frac{1}{2} & -\frac{1}{2} \\ 0 & \frac{3}{2} & \frac{11}{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 8 \\ -2 \\ 14 \end{bmatrix} = \mathbf{M}_1 \mathbf{P}_1 \mathbf{b}.$$

The largest entry in the second column on or below the diagonal is $\frac{3}{2}$, so we interchange the last two rows using the permutation matrix

$$\mathbf{P}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix},$$

obtaining the permuted system

$$\mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} \mathbf{x} = \begin{bmatrix} 4 & 9 & -3 \\ 0 & \frac{3}{2} & \frac{11}{2} \\ 0 & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 8 \\ 14 \\ -2 \end{bmatrix} = \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{b}.$$

To annihilate the subdiagonal entry of the second column, we use the elimination matrix

$$\mathbf{M}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \frac{1}{3} & 1 \end{bmatrix},$$

obtaining the transformed system

$$\mathbf{M}_2 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{A} \mathbf{x} = \begin{bmatrix} 4 & 9 & -3 \\ 0 & \frac{3}{2} & \frac{11}{2} \\ 0 & 0 & \frac{4}{3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 8 \\ 14 \\ \frac{8}{3} \end{bmatrix} = \mathbf{M}_2 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1 \mathbf{b}.$$

We have therefore reduced the original system to an equivalent upper triangular system, which can now be solved by back-substitution to obtain the same answer as before.

To write out the LU factorization explicitly, we have

$$\begin{aligned} \mathbf{L} &= \mathbf{M}^{-1} = (\mathbf{M}_2 \mathbf{P}_2 \mathbf{M}_1 \mathbf{P}_1)^{-1} = \mathbf{P}_1^T \mathbf{L}_1 \mathbf{P}_2^T \mathbf{L}_2 = \\ & \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -\frac{1}{2} & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{1}{3} & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{3} & 1 \\ 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \end{bmatrix}, \end{aligned}$$

and hence

$$\mathbf{A} = \begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{3} & 1 \\ 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \end{bmatrix} \begin{bmatrix} 4 & 9 & -3 \\ 0 & \frac{3}{2} & \frac{11}{2} \\ 0 & 0 & \frac{4}{3} \end{bmatrix} = \mathbf{L} \mathbf{U}.$$

Note that \mathbf{L} is not lower triangular, but it is triangular in the more general sense (it is a permutation of a lower triangular matrix). Alternatively, we can take

$$\mathbf{P} = \mathbf{P}_2 \mathbf{P}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix},$$

and

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ \frac{1}{2} & -\frac{1}{3} & 1 \end{bmatrix},$$

so that

$$\mathbf{P} \mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -3 & 7 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ \frac{1}{2} & -\frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} 4 & 9 & -3 \\ 0 & \frac{3}{2} & \frac{11}{2} \\ 0 & 0 & \frac{4}{3} \end{bmatrix} = \mathbf{L} \mathbf{U},$$

where \mathbf{L} now really is lower triangular but \mathbf{A} is permuted.

As we have just seen, pivoting is generally required for Gaussian elimination to be stable. There are some classes of matrices, however, for which Gaussian elimination is stable *without* pivoting. For example, if the matrix \mathbf{A} is *diagonally dominant* by columns, which means that each diagonal entry is larger in magnitude than the sum of the magnitudes of the other entries in its column,

$$\sum_{i=1, i \neq j}^n |a_{ij}| < |a_{jj}|, \quad j = 1, \dots, n,$$

then pivoting is not required in computing its LU factorization by Gaussian elimination. If partial pivoting is used on such a matrix, then no row interchanges will actually occur. Another important class for which pivoting is not required is matrices that are symmetric and positive definite, which will be defined in Section 2.5. Avoiding an unnecessary pivot search can save a significant amount of time in computing the factorization.

2.2.5 Implementation of Gaussian Elimination

Gaussian elimination, or LU factorization, has the general form of a triple-nested loop,

```

for _____
  for _____
    for _____
       $a_{ij} = a_{ij} - (a_{ik}/a_{kk})a_{kj}$ 
    end
  end
end

```

where the indices i , j , and k of the **for** loops can be taken in any order, for a total of $3! = 6$ different ways of arranging the loops. Some of the indicated arithmetic operations can be moved outside the innermost loop for greater efficiency, depending on the specific indices involved, and additional reorderings of the operations that do not have strictly nested loops are also possible. These variations of the basic algorithm have different memory access patterns (e.g., accessing memory row-wise or column-wise), and also differ in their ability to take advantage of the architectural features of a given computer (e.g., cache, paging, vectorization, multiple processors). Thus, their performance may vary widely on a given computer or across different computers, and no single arrangement may be uniformly superior.

Numerous implementation details of the algorithm are subject to variation in this way. For example, the partial pivoting procedure we described searches along columns and interchanges rows, but alternatively, one could search along rows and interchange columns. We have also taken \mathbf{L} to have unit diagonal, but one could instead arrange for \mathbf{U} to have unit diagonal. Some of these variations of Gaussian elimination are of sufficient importance to have been given names, such as the Crout and Doolittle methods.

Although the many possible variations on Gaussian elimination may have a dramatic effect on performance, they all produce essentially the same factorization. Provided the row pivot sequence is the same, if we have two LU factorizations $\mathbf{PA} = \mathbf{LU} = \hat{\mathbf{L}}\hat{\mathbf{U}}$, then

this expression implies that $\hat{\mathbf{L}}^{-1}\mathbf{L} = \hat{\mathbf{U}}\mathbf{U}^{-1} = \mathbf{D}$ is both lower and upper triangular, and hence diagonal. If both \mathbf{L} and $\hat{\mathbf{L}}$ are assumed to be unit lower triangular, then \mathbf{D} must in fact be the identity matrix \mathbf{I} , and hence $\mathbf{L} = \hat{\mathbf{L}}$ and $\mathbf{U} = \hat{\mathbf{U}}$, so that the factorization is unique. Even without this assumption, however, we may still conclude that the LU factorization is unique up to diagonal scaling of the factors. This uniqueness is made explicit in the LDU factorization $\mathbf{PA} = \mathbf{LDU}$, where \mathbf{L} is unit lower triangular, \mathbf{U} is unit upper triangular, and \mathbf{D} is diagonal.

Storage management is another important implementation issue. The numerous matrices we considered—the elementary elimination matrices \mathbf{M}_k , their inverses \mathbf{L}_k , and the permutation matrices \mathbf{P}_k —merely describe the factorization process formally. They are not formed explicitly in an actual implementation. To conserve storage, the \mathbf{L} and \mathbf{U} factors overwrite the initial storage for the input matrix \mathbf{A} , with the transformed matrix \mathbf{U} occupying the upper triangle of \mathbf{A} (including the diagonal), and the multipliers that make up the strict lower triangle of \mathbf{L} occupying the (now zero) strict lower triangle of \mathbf{A} . The unit diagonal of \mathbf{L} need not be stored.

To minimize data movement, the row interchanges required by pivoting are not usually carried out explicitly. Instead, the rows remain in their original locations, and an auxiliary integer vector is used to keep track of the new row order. Note that a single such vector suffices, because the net effect of all of the interchanges is still just a permutation of the integers $1, \dots, n$.

2.2.6 Complexity of Solving Linear Systems

The Gaussian elimination process for computing the LU factorization requires about $n^3/3$ floating-point multiplications and a similar number of additions. Solving the resulting triangular system for a single right-hand-side vector by forward- and back-substitution requires about n^2 multiplications and a similar number of additions. Thus, as the order n of the matrix grows, the LU factorization phase becomes increasingly dominant in the cost of solving linear systems.

We can also solve a linear system by explicitly inverting the matrix so that the solution is given by $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$. But computing \mathbf{A}^{-1} is tantamount to solving n linear systems: it requires an LU factorization of \mathbf{A} followed by n forward- and back-substitutions, one for each column of the identity matrix. The total operation count is about n^3 multiplications and a similar number of additions (taking advantage of the zeros in the right-hand-side vectors for the forward-substitution).

Thus, explicit inversion is three times as expensive as LU factorization. The subsequent matrix-vector multiplication $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ to solve a linear system requires about n^2 multiplications and a similar number of additions, which is similar to the total cost of forward- and back-substitution. Hence, even for multiple right-hand-side vectors, matrix inversion is more costly than LU factorization for solving linear systems. In addition, explicit inversion gives a less accurate answer. As a simple example, if we solve the 1×1 linear system $3x = 18$ by division, we get $x = 18/3 = 6$, but explicit inversion would give $x = 3^{-1} \times 18 = 0.333 \times 18 = 5.99$ using three-digit arithmetic. In this small example, inversion requires an additional arithmetic operation and obtains a less accurate result. The disadvantages of inversion become worse as the size of the system grows.

Explicit matrix inverses often occur as a convenient notation in various formulas, but this practice does not mean that an explicit inverse is required to implement such a formula. One merely need solve a linear system with an appropriate right-hand side, which might itself be a matrix. Thus, for example, a product of the form $\mathbf{A}^{-1}\mathbf{B}$ should be computed by LU factorization of \mathbf{A} , followed by forward- and back-substitutions using each column of \mathbf{B} . It is extremely rare in practice that an explicit matrix inverse is actually needed, so whenever you see a matrix inverse in a formula, you should think “solve a system” rather than “invert a matrix.”

Another method for solving linear systems that should be avoided is *Cramer’s rule*, in which each component of the solution is computed as a ratio of determinants. Though often taught in elementary linear algebra courses, this method is astronomically expensive for full matrices of nontrivial size. Cramer’s rule is useful mostly as a theoretical tool.

2.2.7 Gauss-Jordan Elimination

The motivation for Gaussian elimination is to reduce a general matrix to triangular form, because the resulting linear system is easy to solve. Diagonal linear systems are even easier to solve, however, so diagonal form would appear to be an even more desirable target. *Gauss-Jordan elimination* is a variation of standard Gaussian elimination in which the matrix is reduced to diagonal form rather than merely to triangular form. The same type of row combinations are used to eliminate matrix entries as in standard Gaussian elimination, but they are applied to annihilate entries above as well as below the diagonal. Thus, the elimination matrix used for a given column vector a is of the form

$$\begin{bmatrix} 1 & \cdots & 0 & -m_1 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & -m_{k-1} & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & -m_{k+1} & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & -m_n & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_{k-1} \\ a_k \\ a_{k+1} \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ a_k \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

where $m_i = a_i/a_k$, $i = 1, \dots, n$. This process requires about $n^3/2$ multiplications and a similar number of additions, which is 50 percent more expensive than standard Gaussian elimination.

During the elimination phase, the same row operations are also applied to the right-hand-side vector (or vectors) of a system of linear equations. Once the elimination phase has been completed and the matrix is in diagonal form, then the components of the solution to the linear system can be computed simply by dividing each entry of the transformed right-hand side by the corresponding diagonal entry of the matrix. This computation requires a total of only n divisions, which is significantly cheaper than solving a triangular system, but not enough to make up for the more costly elimination phase. Gauss-Jordan elimination also has the numerical disadvantage that the multipliers can exceed 1 in magnitude even if pivoting is used.

Despite its higher overall cost, Gauss-Jordan elimination may be preferred in some situations because of the extreme simplicity of its final solution phase. For example, it is

occasionally advocated for implementation on parallel computers because it has a uniform workload throughout the factorization phase, and then all of the solution components can be computed simultaneously rather than one at a time as in ordinary back-substitution.

Gauss-Jordan elimination is also sometimes used to compute the inverse of a matrix explicitly, if desired. If the right-hand-side matrix is initialized to be the identity matrix \mathbf{I} and the given matrix \mathbf{A} is reduced to the identity matrix by Gauss-Jordan elimination, then the transformed right-hand-side matrix will be the inverse of \mathbf{A} . For computing the inverse, Gauss-Jordan elimination has about the same operation count as explicit inversion by Gaussian elimination followed by forward- and back-substitution.

Example 2.8 Gauss-Jordan Elimination. We illustrate Gauss-Jordan elimination by using it to compute the inverse of the matrix of Example 2.6. For simplicity, we omit pivoting. We begin with the matrix \mathbf{A} , augmented by the identity matrix \mathbf{I} as right-hand side, and repeatedly apply elimination matrices to annihilate off-diagonal entries of \mathbf{A} until we reach diagonal form, then scale by the remaining diagonal entries to produce the identity matrix on the left, and hence the inverse matrix on the right.

$$\begin{aligned} \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 2 & 4 & -2 \\ -2 & 1 & 0 & 4 & 9 & -3 \\ 1 & 0 & 1 & -2 & -3 & 7 \end{array} \right] &= \left[\begin{array}{ccc|ccc} 2 & 4 & -2 & 1 & 0 & 0 \\ 0 & 1 & 1 & -2 & 1 & 0 \\ 0 & 1 & 5 & 1 & 0 & 1 \end{array} \right], \\ \left[\begin{array}{ccc|ccc} 1 & -4 & 0 & 2 & 4 & -2 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & -1 & 1 & 0 & 1 & 5 \end{array} \right] &= \left[\begin{array}{ccc|ccc} 2 & 0 & -6 & 9 & -4 & 0 \\ 0 & 1 & 1 & -2 & 1 & 0 \\ 0 & 0 & 4 & 3 & -1 & 1 \end{array} \right], \\ \left[\begin{array}{ccc|ccc} 1 & 0 & \frac{3}{2} \\ 0 & 1 & -\frac{1}{4} \\ 0 & 0 & 1 \end{array} \right] &= \left[\begin{array}{ccc|ccc} 2 & 0 & -6 & 9 & -4 & 0 \\ 0 & 1 & 1 & -2 & 1 & 0 \\ 0 & 0 & 4 & 3 & -1 & 1 \end{array} \right] = \left[\begin{array}{ccc|ccc} 2 & 0 & 0 & \frac{27}{2} & -\frac{11}{2} & \frac{3}{2} \\ 0 & 1 & 0 & -\frac{11}{4} & \frac{5}{4} & -\frac{1}{4} \\ 0 & 0 & 4 & 3 & -1 & 1 \end{array} \right], \\ &= \left[\begin{array}{ccc|ccc} \frac{1}{2} & 0 & 0 & 2 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & \frac{1}{4} & 0 & 0 & 4 \end{array} \right] \left[\begin{array}{ccc|ccc} 2 & 0 & 0 & \frac{27}{2} & -\frac{11}{2} & \frac{3}{2} \\ 0 & 1 & 0 & -\frac{11}{4} & \frac{5}{4} & -\frac{1}{4} \\ 0 & 0 & 4 & 3 & -1 & 1 \end{array} \right] = \\ \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & \frac{27}{4} & -\frac{11}{4} & \frac{3}{4} \\ 0 & 1 & 0 & -\frac{11}{4} & \frac{5}{4} & -\frac{1}{4} \\ 0 & 0 & 1 & \frac{3}{4} & -\frac{1}{4} & \frac{1}{4} \end{array} \right], & \text{ so } A^{-1} = \frac{1}{4} \begin{bmatrix} 27 & -11 & 3 \\ -11 & 5 & -1 \\ 3 & -1 & 1 \end{bmatrix}. \end{aligned}$$

2.2.8 Solving Modified Problems

In many practical situations linear systems do not occur in isolation but as part of a sequence of related problems that change in some systematic way. For example, one may need to solve a sequence of linear systems $\mathbf{Ax} = \mathbf{b}$ having the same matrix \mathbf{A} but different right-hand sides \mathbf{b} . After having solved the initial system by Gaussian elimination, then the \mathbf{L} and \mathbf{U} factors already computed can be used to solve the additional systems by forward- and back-substitution. The factorization phase need not be repeated in solving subsequent linear systems unless the matrix changes. This procedure represents a substantial savings

in work, since additional triangular solutions cost only $\mathcal{O}(n^2)$ work, in contrast to the $\mathcal{O}(n^3)$ cost of a factorization.

In fact, in some important special cases a new factorization can be avoided even when the matrix does change. One such case that arises frequently is the addition of a matrix that is an outer product $\mathbf{u}\mathbf{v}^T$ of two nonzero vectors \mathbf{u} and \mathbf{v} . This is called a *rank-one change* because the outer product matrix $\mathbf{u}\mathbf{v}^T$ has rank one (i.e., only one linearly independent row or column), and any rank-one matrix can be expressed as such an outer product of two vectors. For example, if a single entry of the matrix \mathbf{A} changes (say the (j, k) entry changes from a_{jk} to \tilde{a}_{jk}), then the new matrix is $\mathbf{A} - \alpha \mathbf{e}_j \mathbf{e}_k^T$, where \mathbf{e}_j and \mathbf{e}_k are the corresponding columns of the identity matrix and $\alpha = a_{jk} - \tilde{a}_{jk}$.

The *Sherman-Morrison formula* gives the inverse of a matrix resulting from a rank-one change to a matrix whose inverse is already known:

$$(\mathbf{A} - \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{u}(1 - \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u})^{-1}\mathbf{v}^T\mathbf{A}^{-1},$$

where \mathbf{u} and \mathbf{v} are n -vectors. Evaluation of this formula requires only $\mathcal{O}(n^2)$ work (for matrix-vector multiplications) rather than the $\mathcal{O}(n^3)$ work normally required for inversion.

To solve a linear system $(\mathbf{A} - \mathbf{u}\mathbf{v}^T)\mathbf{x} = \mathbf{b}$ with the new matrix, we could use the foregoing formula to obtain

$$\mathbf{x} = (\mathbf{A} - \mathbf{u}\mathbf{v}^T)^{-1}\mathbf{b} = \mathbf{A}^{-1}\mathbf{b} + \mathbf{A}^{-1}\mathbf{u}(1 - \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u})^{-1}\mathbf{v}^T\mathbf{A}^{-1}\mathbf{b}.$$

We would prefer to avoid explicit inversion altogether, however. If we have an LU factorization for \mathbf{A} , then the following steps can easily be computed to obtain the solution to the modified system:

1. Solve $\mathbf{A}\mathbf{z} = \mathbf{u}$ for \mathbf{z} , so that $\mathbf{z} = \mathbf{A}^{-1}\mathbf{u}$.
2. Solve $\mathbf{A}\mathbf{y} = \mathbf{b}$ for \mathbf{y} , so that $\mathbf{y} = \mathbf{A}^{-1}\mathbf{b}$.
3. Compute $\mathbf{x} = \mathbf{y} + [(\mathbf{v}^T\mathbf{y})/(1 - \mathbf{v}^T\mathbf{z})]\mathbf{z}$.

Note that the first step is independent of \mathbf{b} and hence need not be repeated if there are multiple right-hand-side vectors \mathbf{b} . Again, this procedure requires only triangular solutions and inner products, so it requires only $\mathcal{O}(n^2)$ work and no explicit inverses.

The *Woodbury formula*, in which \mathbf{u} and \mathbf{v} become $n \times k$ matrices \mathbf{U} and \mathbf{V} , generalizes the Sherman-Morrison formula to a rank- k change in the matrix:

$$(\mathbf{A} - \mathbf{U}\mathbf{V}^T)^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} - \mathbf{V}^T\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^T\mathbf{A}^{-1}.$$

Using similar techniques, it is possible to update the factorization rather than the inverse or the solution. Caution must be exercised in using these updating formulas, however, because in general there is no guarantee of numerical stability through successive updates as the matrix changes.

Example 2.9 Rank-One Updating of Solutions. To illustrate the use of the Sherman-Morrison formula, we solve the linear system

$$\begin{bmatrix} 2 & 4 & -2 \\ 4 & 9 & -3 \\ -2 & -1 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 8 \\ 10 \end{bmatrix},$$

which is a rank-one modification of the system in Example 2.6 (only the $(3, 2)$ entry has changed). We take \mathbf{A} to be the matrix of Example 2.6, so we can use the LU factorization already computed. One way to choose the update vectors is

$$\mathbf{u} = \begin{bmatrix} 0 \\ 0 \\ -2 \end{bmatrix} \quad \text{and} \quad \mathbf{v} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix},$$

so that the matrix of the new system is $\mathbf{A} - \mathbf{u}\mathbf{v}^T$, and the right-hand-side vector \mathbf{b} has not changed. We can use the previously computed LU factorization of \mathbf{A} to solve $\mathbf{A}\mathbf{z} = \mathbf{u}$ to obtain $\mathbf{z} = [-\frac{3}{2} \quad \frac{1}{2} \quad -\frac{1}{2}]^T$, and we had already solved $\mathbf{A}\mathbf{y} = \mathbf{b}$ to obtain $\mathbf{y} = [-1 \quad 2 \quad 2]^T$. The final step is then to compute the updated solution

$$\mathbf{x} = \mathbf{y} + \frac{\mathbf{v}^T \mathbf{y}}{1 - \mathbf{v}^T \mathbf{z}} \mathbf{z} = \begin{bmatrix} -1 \\ 2 \\ 2 \end{bmatrix} + \frac{2}{1 - \frac{1}{2}} \begin{bmatrix} -\frac{3}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} = \begin{bmatrix} -7 \\ 4 \\ 0 \end{bmatrix}.$$

We have thus computed the solution to the new system without refactoring the modified matrix.

2.3 Norms and Condition Numbers

2.3.1 Vector Norms

To measure errors and sensitivity in solving linear systems, we need some notion of the “size” of vectors and matrices. The scalar concept of magnitude, modulus, or absolute value can be generalized to the concept of *norms* for vectors and matrices. Although a more general definition is possible, all of the vector norms we will use are instances of p -norms, which for an integer $p > 0$ and a vector \mathbf{x} of dimension n are defined by

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

Important special cases are as follows:

- 1-norm:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|,$$

sometimes called the Manhattan norm because in the plane it corresponds to the distance between two points as measured in “city blocks.”

- 2-norm:

$$\|\mathbf{x}\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2},$$

which corresponds to the usual notion of distance in Euclidean space.

- ∞ -norm:

$$\|\mathbf{x}\|_{\infty} = \max_{1 \leq i \leq n} |x_i|,$$

which can be viewed as a limiting case as $p \rightarrow \infty$.

All of these norms give the same results qualitatively, but in certain circumstances a particular norm may be easiest to work with analytically or computationally. The 1-norm or the ∞ -norm is usually used in analyzing the sensitivity of solutions to linear systems. We will also use the 2-norm later on in other contexts. The differences among these norms are illustrated in Fig. 2.1, which shows the unit sphere, $\{\mathbf{x}: \|\mathbf{x}\| = 1\}$, in two dimensions for each norm.

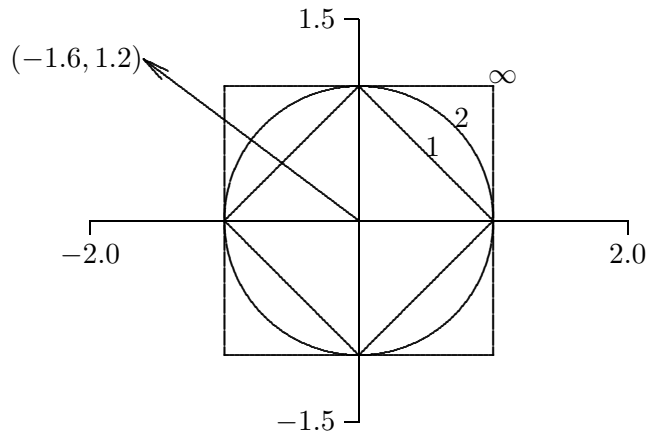


Figure 2.1: The unit sphere in various vector norms.

The norm of a vector is simply the factor by which the corresponding unit sphere must be expanded or shrunk to encompass the vector. For example, the norms have the following values for the vector shown in Fig. 2.1:

$$\left\| \begin{bmatrix} -1.6 \\ 1.2 \end{bmatrix} \right\|_1 = 2.8, \quad \left\| \begin{bmatrix} -1.6 \\ 1.2 \end{bmatrix} \right\|_2 = 2.0, \quad \left\| \begin{bmatrix} -1.6 \\ 1.2 \end{bmatrix} \right\|_{\infty} = 1.6.$$

In general, for any vector \mathbf{x} in \mathbb{R}^n , we have

$$\|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_2 \geq \|\mathbf{x}\|_{\infty}.$$

On the other hand, we also have

$$\|\mathbf{x}\|_1 \leq \sqrt{n} \|\mathbf{x}\|_2, \quad \|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_{\infty}, \quad \text{and} \quad \|\mathbf{x}\|_1 \leq n \|\mathbf{x}\|_{\infty}.$$

Thus, for a given n , any two of the norms differ by at most a constant, so they are all equivalent in the sense that if one is small, they must all be proportionally small. Hence, we can choose whichever norm is most convenient in a given context. In the remainder of this book, an appropriate subscript will be used to indicate a specific norm, when necessary, but the subscript will be omitted when it does not matter which particular norm is used.

For any vector norm, the following important properties hold, where \mathbf{x} and \mathbf{y} are any vectors:

1. $\|\mathbf{x}\| > 0$ if $\mathbf{x} \neq \mathbf{o}$.
2. $\|\gamma\mathbf{x}\| = |\gamma| \cdot \|\mathbf{x}\|$ for any scalar γ .
3. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality).

In a more general treatment, these three properties can be taken as the *definition* of a vector norm. A useful variation on the triangle inequality is

$$\|\mathbf{x} - \mathbf{y}\| \geq \|\mathbf{x}\| - \|\mathbf{y}\|.$$

2.3.2 Matrix Norms

We also need some way to measure the size or magnitude of matrices. Again, a more general definition is possible, but all of the matrix norms we will use are defined in terms of an underlying vector norm. Specifically, given a vector norm, we define the corresponding matrix norm of a matrix \mathbf{A} as follows:

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

Such a matrix norm is said to be *subordinate* to the vector norm. Intuitively, the norm of a matrix measures the maximum stretching the matrix does to any vector, as measured in the given vector norm.

Some matrix norms are much easier to compute than others. For example, the matrix norm corresponding to the vector 1-norm is simply the maximum absolute column sum of the matrix,

$$\|\mathbf{A}\|_1 = \max_j \sum_{i=1}^n |a_{ij}|,$$

and the matrix norm corresponding to the vector ∞ -norm is simply the maximum absolute row sum of the matrix,

$$\|\mathbf{A}\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|.$$

A handy way to remember these is that the matrix norms agree with the corresponding vector norms for an $n \times 1$ matrix. Unfortunately, the matrix norm corresponding to the vector 2-norm is not so simple to compute; it turns out to be equal to the square root of the largest eigenvalue of the matrix $\mathbf{A}^T \mathbf{A}$, or, as we shall see later, the largest singular value of \mathbf{A} (see Section 4.5.2).

The matrix norms we have defined satisfy the following important properties, where \mathbf{A} and \mathbf{B} are any matrices:

1. $\|\mathbf{A}\| > 0$ if $\mathbf{A} \neq \mathbf{O}$.
2. $\|\gamma\mathbf{A}\| = |\gamma| \cdot \|\mathbf{A}\|$ for any scalar γ .
3. $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$.
4. $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|$.
5. $\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|$ for any vector \mathbf{x} .

In a more general treatment, the first three properties can be taken as the *definition* of a matrix norm. The remaining two properties, known as *submultiplicative* or *consistency* conditions, may or may not hold for these more general matrix norms, but they always hold for the matrix norms subordinate to the vector p -norms.

2.3.3 Condition Number of a Matrix

The *condition number* of a square nonsingular matrix \mathbf{A} with respect to a given norm is defined as

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|.$$

By convention, $\text{cond}(\mathbf{A}) = \infty$ if \mathbf{A} is singular. Since

$$\|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| = \left(\max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \right) \cdot \left(\min_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \right)^{-1},$$

the condition number of a matrix measures the ratio of the maximum stretching that the matrix does to any nonzero vector to the maximum shrinking. We will see in Section 2.4.2 that this concept is consistent with the general notion of condition number defined in Section 1.2.5: the condition number of the matrix bounds the ratio of the relative change in the solution of a linear system to a given relative change in the input data.

The condition number is a measure of how close a matrix is to being singular: a matrix with a large condition number (which we will quantify in Section 2.4.2) is nearly singular, whereas a matrix with a condition number close to 1 is far from being singular. Note that the determinant of a matrix is *not* a good indicator of near singularity: although a matrix \mathbf{A} is singular if $\det(\mathbf{A}) = 0$, the magnitude of a nonzero determinant, large or small, gives no information on how close to singular the matrix may be. For example, $\det(\alpha \mathbf{I}_n) = \alpha^n$, which can be arbitrarily small for $|\alpha| < 1$, yet the matrix is perfectly well-conditioned for any nonzero α , with a condition number of 1.

Some important properties of the condition number are

1. For any matrix \mathbf{A} , $\text{cond}(\mathbf{A}) \geq 1$.
2. For the identity matrix, $\text{cond}(\mathbf{I}) = 1$.
3. For any permutation matrix \mathbf{P} , $\text{cond}(\mathbf{P}) = 1$.
4. For any matrix \mathbf{A} and nonzero scalar γ , $\text{cond}(\gamma \mathbf{A}) = \text{cond}(\mathbf{A})$.
5. For any diagonal matrix $\mathbf{D} = \text{diag}(d_i)$, $\text{cond}(\mathbf{D}) = (\max |d_i|) / (\min |d_i|)$.

As we will see shortly, the usefulness of the condition number is in assessing the accuracy of solutions to linear systems. Since the definition of the condition number involves the inverse of the matrix, computing its value is obviously a nontrivial task. In fact, to compute the condition number literally would require substantially more work than solving the linear system whose accuracy is to be assessed using the condition number. In practice, therefore, the condition number is merely estimated, to perhaps within an order of magnitude, as a relatively inexpensive byproduct of the solution process.

The matrix norm $\|\mathbf{A}\|$ is easily computed as the maximum absolute column sum (or row sum, depending on the norm used). It is estimating $\|\mathbf{A}^{-1}\|$ at low cost that presents a

challenge. From the properties of norms, we know that if \mathbf{z} is the solution to $\mathbf{A}\mathbf{z} = \mathbf{y}$, then

$$\frac{\|\mathbf{z}\|}{\|\mathbf{y}\|} \leq \|\mathbf{A}^{-1}\|,$$

and the bound is achieved for some optimally chosen vector \mathbf{y} . We thus wish to pick a vector \mathbf{y} so that the ratio $\|\mathbf{z}\|/\|\mathbf{y}\|$ is as large as possible and therefore is a reasonable estimate for $\|\mathbf{A}^{-1}\|$. Finding the optimal \mathbf{y} would be prohibitively expensive, but a useful approximation can be obtained much more cheaply. One heuristic is to choose \mathbf{y} as the solution to the system $\mathbf{A}^T\mathbf{y} = \mathbf{c}$, where \mathbf{c} is a vector whose components are ± 1 , with the signs chosen successively to make the resulting \mathbf{y} as large as possible.

The motivation for this approach may not be obvious now, but it is essentially equivalent to one step of inverse iteration for computing the singular vector corresponding to the smallest singular value of \mathbf{A} (see Chapter 4). An alternative approach to condition estimation is to treat it as a convex optimization problem that can be solved very efficiently in practice using a heuristic algorithm. Most good modern software packages for solving linear systems provide an efficient and reliable condition estimator, based on a sophisticated implementation of one of the methods outlined here.

2.4 Accuracy of Solutions

2.4.1 Residual of a Solution

Intuitively, the most obvious way to check the validity of a solution is to substitute it into the equation to see how closely the two sides match. The *residual vector* of an approximate solution $\hat{\mathbf{x}}$ to the $n \times n$ linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ is defined as

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}.$$

In theory, if \mathbf{A} is nonsingular, then the error $\|\hat{\mathbf{x}} - \mathbf{x}\| = 0$ if and only if $\|\mathbf{r}\| = 0$. In practice, however, these quantities are not necessarily small simultaneously. If the computed solution $\hat{\mathbf{x}}$ exactly satisfies

$$(\mathbf{A} + \mathbf{E})\hat{\mathbf{x}} = \mathbf{b},$$

then

$$\|\mathbf{r}\| = \|\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}\| = \|\mathbf{E}\hat{\mathbf{x}}\| \leq \|\mathbf{E}\| \cdot \|\hat{\mathbf{x}}\|,$$

so that we have the inequality

$$\frac{\|\mathbf{r}\|}{\|\mathbf{A}\| \cdot \|\hat{\mathbf{x}}\|} \leq \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}$$

relating the *relative residual* to the relative change in the matrix. Thus, a large relative residual implies a large backward error in the matrix, which means that the algorithm used to compute the solution is unstable.

But how large is $\|\mathbf{E}\|$ likely to be in practice? Wilkinson [273] showed that for LU factorization by Gaussian elimination, a bound of the form

$$\frac{\|\mathbf{E}\|}{\|\mathbf{A}\|} \leq \rho n \epsilon_{\text{mach}}$$

holds, where ρ , called the *growth factor*, is the ratio of the largest entry of \mathbf{U} to the largest entry of \mathbf{A} . Without pivoting, ρ can be arbitrarily large, and hence Gaussian elimination without pivoting is unstable, as we have already seen. With partial pivoting, the growth factor can still be as large as 2^{n-1} (since in the worst case the size of the entries can double at each stage of elimination), but such behavior is extremely rare. In practice, there is little or no growth in the size of the entries, so that

$$\frac{\|\mathbf{E}\|}{\|\mathbf{A}\|} \approx n \epsilon_{\text{mach}}.$$

This relation means that solving a linear system by Gaussian elimination with partial pivoting followed by back-substitution almost always yields a very small relative residual, regardless of how ill-conditioned the system may be. Thus, a small relative residual is not necessarily a good indicator that a computed solution is close to the “true” solution unless the system is well-conditioned.

Complete pivoting yields an even smaller growth factor, in both theory and practice, but the additional margin of stability it provides is usually not worth the extra expense.

Example 2.10 Small Residual. Using three-digit decimal arithmetic to solve the system

$$\begin{bmatrix} 0.641 & 0.242 \\ 0.321 & 0.121 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.883 \\ 0.442 \end{bmatrix},$$

Gaussian elimination with partial pivoting yields the triangular system

$$\begin{bmatrix} 0.641 & 0.242 \\ 0 & -0.000242 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.883 \\ -0.000383 \end{bmatrix},$$

and back-substitution then gives the computed solution

$$\mathbf{x} = \begin{bmatrix} 0.782 \\ 1.58 \end{bmatrix}.$$

The exact residual for this solution is

$$\mathbf{r} = \mathbf{b} - \mathbf{Ax} = \begin{bmatrix} -0.000622 \\ -0.000202 \end{bmatrix},$$

which is as small as we can expect using only three-digit arithmetic. Yet the exact solution for this system is easily seen to be

$$\mathbf{x} = \begin{bmatrix} 1.00 \\ 1.00 \end{bmatrix},$$

so that the error is almost as large as the solution. The cause of this phenomenon is that the matrix is very nearly singular (its condition number is more than 4000). The division that determines x_2 is between two quantities that are both on the order of rounding error, and hence the result is essentially arbitrary. Yet, by design, when this arbitrary value for x_2 is then substituted into the first equation, a value for x_1 is computed so that the first equation is satisfied. Thus, we get a small residual, but a poor solution.

2.4.2 Estimating Accuracy

In addition to being a reliable indicator of near singularity, the condition number also provides a quantitative estimate for the error in the computed solution to a linear system, as we will now see. Let \mathbf{x} be the solution to the nonsingular linear system $\mathbf{Ax} = \mathbf{b}$, and let $\hat{\mathbf{x}}$ be the solution to the system $\mathbf{A}\hat{\mathbf{x}} = \mathbf{b} + \Delta\mathbf{b}$ with a perturbed right-hand side. If we define $\Delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x}$, then we have

$$\mathbf{b} + \Delta\mathbf{b} = \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{Ax} + \mathbf{A}\Delta\mathbf{x}.$$

Since $\mathbf{Ax} = \mathbf{b}$, we must have $\mathbf{A}\Delta\mathbf{x} = \Delta\mathbf{b}$, and hence $\Delta\mathbf{x} = \mathbf{A}^{-1}\Delta\mathbf{b}$. Now

$$\mathbf{b} = \mathbf{Ax} \Rightarrow \|\mathbf{b}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|,$$

and

$$\Delta\mathbf{x} = \mathbf{A}^{-1}\Delta\mathbf{b} \Rightarrow \|\Delta\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\Delta\mathbf{b}\|,$$

which, upon using the definition $\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|$, yields the estimate

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|}.$$

Thus, the condition number of the matrix determines the possible relative change in the solution due to a given relative change in the right-hand-side vector, regardless of the algorithm used to compute the solution (compare with the general notion of condition number defined in Section 1.2.5). A similar result holds for relative changes in the entries of the matrix \mathbf{A} . If $\mathbf{Ax} = \mathbf{b}$ and

$$(\mathbf{A} + \mathbf{E})\hat{\mathbf{x}} = \mathbf{b},$$

then

$$\mathbf{x} - \hat{\mathbf{x}} = \mathbf{A}^{-1}(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = \mathbf{A}^{-1}\mathbf{E}\hat{\mathbf{x}},$$

so that

$$\|\Delta\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{E}\| \cdot \|\hat{\mathbf{x}}\|,$$

which yields the estimate

$$\frac{\|\Delta\mathbf{x}\|}{\|\hat{\mathbf{x}}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}.$$

As an alternative to the algebraic derivations just given, calculus can be used to estimate the sensitivity of linear systems. Introducing the real-valued parameter t , we define $\mathbf{A}(t) = \mathbf{A} + t\mathbf{E}$ and $\mathbf{b}(t) = \mathbf{b} + t\Delta\mathbf{b}$, and consider the solution $\mathbf{x}(t)$ to the linear system $\mathbf{A}(t)\mathbf{x}(t) = \mathbf{b}(t)$. Differentiating this equation with respect to t , we obtain

$$\mathbf{A}'(t)\mathbf{x}(t) + \mathbf{A}(t)\mathbf{x}'(t) = \mathbf{b}'(t),$$

so that we have

$$\mathbf{x}'(t) = -\mathbf{A}^{-1}(t)\mathbf{A}'(t)\mathbf{x}(t) + \mathbf{A}^{-1}(t)\mathbf{b}'(t),$$

and hence, evaluating at $t = 0$ and taking norms,

$$\frac{\|\mathbf{x}'\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}'\| + \|\mathbf{A}^{-1}\| \cdot \frac{\|\mathbf{b}'\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \left(\frac{\|\mathbf{A}'\|}{\|\mathbf{A}\|} + \frac{\|\mathbf{b}'\|}{\|\mathbf{b}\|} \right).$$

Thus, we again see that the relative change in the solution is bounded by the condition number times the relative change in the problem data.

A geometric interpretation in two dimensions of these sensitivity results is that if the two straight lines defined by the two equations are nearly parallel, then their point of intersection is not sharply defined if the lines are a bit fuzzy because of rounding errors or other sources of error. If, on the other hand, the lines are far from parallel, say nearly perpendicular, then their intersection is relatively sharply defined. These two cases are illustrated in Fig. 2.2, where the dashed lines indicate the region of uncertainty for each solid line, so that the intersection point in each case could be anywhere within the shaded parallelogram. Thus, a large condition number is associated with a large uncertainty in the solution.

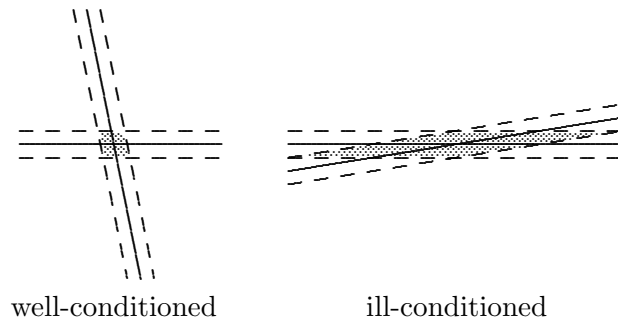


Figure 2.2: Well-conditioned and ill-conditioned linear systems.

To summarize, if the input data are accurate to machine precision, then a reasonable estimate for the relative error in the computed solution to a linear system is given by

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} \approx \text{cond}(\mathbf{A}) \epsilon_{\text{mach}}.$$

One simple way of interpreting these results is that the computed solution loses about $\log_{10}(\text{cond}(\mathbf{A}))$ decimal digits of accuracy relative to the accuracy of the input. In Example 2.10, for instance, with a condition number greater than 10^3 , we lost all of the three-digit precision available and obtained an arbitrary solution.

Before leaving the subject of assessing accuracy in terms of condition numbers, note these two caveats:

- The foregoing analysis estimates the relative error in the *largest* components of the solution vector. The *relative* error in the smaller components can be much larger, because a vector norm is dominated by the largest components of a vector. Componentwise error bounds can be obtained but are somewhat more complicated to compute, and we will not pursue this topic. Componentwise bounds are of particular interest when the system is poorly scaled.
- The condition number of a matrix is affected by the scaling of the matrix (recall Example 2.3). A large condition number can result simply from poor scaling, as well as from near singularity. Rescaling the matrix can help the former, but not the latter (see Section 2.4.3).