

where

$$A = [a_{ij}] = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad \bar{\mathbf{x}} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix}, \quad \bar{\mathbf{b}} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \quad (3.4)$$

Here it is to be understood that i, j are integers in the range $1 \leq i \leq m$ and $1 \leq j \leq n$ with a_{ij} denoting the matrix element in the i th row and j th column. The matrix A is called a $m \times n$ (read m by n) matrix. The special $n \times 1$ matrix $\bar{\mathbf{x}}$ is called a n -dimensional column vector with x_i , $1 \leq i \leq n$, the element in the i th row. Similarly, the special $m \times 1$ matrix $\bar{\mathbf{b}}$ is called a m -dimensional column vector with b_k denoting the element in the k th row, for $1 \leq k \leq m$. We will examine direct methods of solution and iterative methods for solving such systems.

We also investigate methods for solving nonlinear systems of equations of the form

$$\begin{aligned} E_1 : & \quad f_1(x_1, x_2, \dots, x_n) = 0 \\ E_2 : & \quad f_2(x_1, x_2, \dots, x_n) = 0 \\ & \quad \vdots \\ E_m : & \quad f_m(x_1, x_2, \dots, x_n) = 0 \end{aligned} \quad (3.5)$$

where f_i , for $i = 1, \dots, m$, represent known continuous functions of the variables x_1, x_2, \dots, x_n . The problem is to find, if possible, some numerical procedure for determining the unknown quantities x_1, x_2, \dots, x_n which satisfy all of the nonlinear equations in the system of equations (3.5). Note that the system of equations (3.1) is a special case of the more general system of equations (3.5). The system of nonlinear equations (3.5) can be written in the vector form

$$\bar{\mathbf{f}}(\bar{\mathbf{x}}) = \bar{\mathbf{0}} \quad \text{where} \quad \bar{\mathbf{x}} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \text{and} \quad \bar{\mathbf{f}}(\bar{\mathbf{x}}) = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_m(x_1, x_2, \dots, x_n) \end{bmatrix} \quad (3.6)$$

with $\bar{\mathbf{0}}$ being an m -dimensional column vector of zeros.

We develop several numerical methods which are applicable for solving systems of the above types. We desire to develop numerical methods for solving the above systems of equations in cases where the numbers m and n are large. In some applied problems it is not unusual for the number of unknowns x_1, x_2, \dots, x_n

to be very large, say $n > 10^5$ or $n > 10^6$. For illustrative purposes and examples of the numerical techniques we will use much smaller values for n .

Preliminaries

In dealing with the matrix form associated with the equations (3.2) there are occasions when it is convenient to make use of special matrices and special functions associated with these matrices. Some of these special matrices and functions are as follows.

- (i) The $n \times n$ identity matrix $I = [\delta_{ij}]$, where $\delta_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$. The identity matrix has 1's down the main diagonal and 0's everywhere else. A 3×3 identity matrix has the form $I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$.
- (ii) If A is a $n \times n$ square matrix and there exists a matrix B with the property that $BA = AB = I$, then B is called the inverse of A and is written $B = A^{-1}$. That is, the inverse matrix A^{-1} has the property that $AA^{-1} = A^{-1}A = I$. If A^{-1} exists, then A is said to be nonsingular. If the matrix A does not have an inverse, then the matrix A is said to be singular. In the special case $m = n$ and A is a square matrix having an inverse, then the solution to the matrix system (3.3) can be symbolically obtained by multiplying both sides of the matrix equation (3.3) by the inverse of A to obtain

$$A^{-1}A\bar{\mathbf{x}} = A^{-1}\bar{\mathbf{b}} \quad \text{which simplifies to} \quad \bar{\mathbf{x}} = I\bar{\mathbf{x}} = A^{-1}\bar{\mathbf{b}}.$$

This technique for solving a $n \times n$ linear system of equations is only recommended in the case where n is small. This is because the number of multiply and divided operations needed to calculate A^{-1} increases like n^3 and so the inverse matrix calculation becomes very lengthy and burdensome when n is large.

- (iii) Associated with a $m \times n$ matrix $A = [a_{ij}]$ is the $n \times m$ transpose matrix denoted by $A^T = [a_{ji}]$ and formed by interchanging the rows and columns of the $m \times n$ matrix A . Note that column vectors such as the $\bar{\mathbf{x}}$ given in equation (3.4) can be expressed as $\bar{\mathbf{x}} = [x_1, x_2, \dots, x_n]^T$. The transpose is used to determine if a matrix is symmetric. A square matrix is said to be symmetric if $A^T = A$. The matrix product $A^T A$ always produces a square matrix. The transpose

operation satisfies the following properties.

- (i) $(AB)^T = B^T A^T$ The transpose of a product is the product of the transpose matrices in reverse order.
- (ii) $(A^{-1})^T = (A^T)^{-1}$ When A^{-1} exists, then the transpose of an inverse is the inverse of the transpose.
- (iii) $(A^T)^T = A$ The transpose of the transpose matrix returns the original matrix.

(iv) A $n \times n$ lower triangular matrix $L = [\ell_{ij}]$ has the form

$$L = \begin{bmatrix} \ell_{11} & & & & \\ \ell_{21} & \ell_{22} & & & \\ \ell_{31} & \ell_{32} & \ell_{33} & & \\ \vdots & \vdots & \vdots & \ddots & \\ \ell_{n1} & \ell_{n2} & \ell_{n3} & \cdots & \ell_{nn} \end{bmatrix} \quad \ell_{ij} = 0 \text{ for } j > i \quad (3.7)$$

with all zeros above the main diagonal. In the special case the diagonal elements of a lower triangular matrix are all 1's, then L is called a unit lower triangular matrix.

(v) A $n \times n$ upper triangular matrix $U = [u_{ij}]$ has the form

$$U = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \cdots & u_{1n} \\ & u_{22} & u_{23} & \cdots & u_{2n} \\ & & u_{33} & \cdots & u_{3n} \\ & & & \ddots & \vdots \\ & & & & u_{nn} \end{bmatrix} \quad u_{ij} = 0 \text{ for } i > j \quad (3.8)$$

with all zeros below the main diagonal. In the special case the diagonal elements of an upper triangular matrix are all 1's, then U is called a unit upper triangular matrix.

(vi) A $n \times n$ square matrix $A = [a_{ij}]$ with the property that

$$a_{ij} = \begin{cases} 0, & \text{for } i + s \leq j, \quad 1 < s < n \\ 0, & \text{for } j + t \leq i, \quad 1 < t < n \end{cases} \quad (3.9)$$

is said to be a banded matrix with band width $w = s + t - 1$. For example, the 5×5 tridiagonal matrix

$$\begin{bmatrix} a_{11} & a_{12} & 0 & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 & 0 \\ 0 & a_{32} & a_{33} & a_{34} & 0 \\ 0 & 0 & a_{43} & a_{44} & a_{45} \\ 0 & 0 & 0 & a_{54} & a_{55} \end{bmatrix} \quad (3.10)$$

- (c) $\|\bar{\mathbf{x}} + \bar{\mathbf{y}}\| \leq \|\bar{\mathbf{x}}\| + \|\bar{\mathbf{y}}\|$ for all $\bar{\mathbf{x}}, \bar{\mathbf{y}} \in R^n$.
 (d) $\|\bar{\mathbf{x}}\| = 0$ if and only if $\bar{\mathbf{x}} = 0$.

The Euclidean norm or ℓ_2 norm is used most often and is defined

$$\|\bar{\mathbf{x}}\|_2 = \left[\sum_{i=1}^n x_i^2 \right]^{1/2} \quad (3.13)$$

This norm represents the distance of the point (x_1, x_2, \dots, x_n) from the origin. Other vector norms can be defined so long as they obey the above properties. Two other vector norms used quite frequently are the ℓ_∞ norm and ℓ_p norm defined respectively as

$$\begin{aligned} \|\bar{\mathbf{x}}\|_\infty &= \max_{1 \leq i \leq n} |x_i| \\ \|\bar{\mathbf{x}}\|_p &= \left[\sum_{i=1}^n |x_i|^p \right]^{1/p} \end{aligned} \quad (3.14)$$

(ix) A matrix norm $\|A\|$ associated with a $n \times n$ matrix $A = [a_{ij}]$ is any real-valued function $\|\cdot\|$ which satisfies the properties:

- (a) $\|A\| \geq 0$
 (b) $\|\alpha A\| \leq |\alpha| \|A\|$ where α is a scalar $\in R$
 (c) $\|A + B\| \leq \|A\| + \|B\|$, where A, B are $n \times n$ matrices.
 (d) $\|AB\| \leq \|A\| \|B\|$
 (e) $\|A\| = 0$ if and only if $a_{ij} = 0$ for all i, j values.

The quantity $\|A - B\|$ is used to measure the nearness of two $n \times n$ matrices. Let $\|\cdot\|$ denote a vector norm, then the natural matrix norm of a $n \times n$ matrix A is defined

$$\|A\| = \max_{\|\bar{\mathbf{x}}\|=1} \|A\bar{\mathbf{x}}\|. \quad (3.15)$$

For example, the ℓ_2 norm of the $n \times n$ matrix A would be represented

$$\|A\|_2 = \max_{\|\bar{\mathbf{x}}\|_2=1} \|A\bar{\mathbf{x}}\|_2 \quad (3.16)$$

and the ℓ_∞ norm of A would be represented

$$\|A\|_\infty = \max_{\|\bar{\mathbf{x}}\|_\infty=1} \|A\bar{\mathbf{x}}\|_\infty \quad (3.17)$$

It can be shown, see the Bronson reference, that if $A = (a_{ij})$ is a $n \times n$ matrix, then

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

- (x) The characteristic polynomial associated with the real $n \times n$ matrix A is defined in terms of a determinant and given by

$$p(\lambda) = \det [A - \lambda I] \quad (3.18)$$

where λ is a scalar and I is the $n \times n$ identity matrix. The values of λ which satisfy the characteristic equation $p(\lambda) = 0$ are called eigenvalues associated with the matrix A . Nonzero vectors \bar{x} with the property that $A\bar{x} = \lambda\bar{x}$ are called eigenvectors of A corresponding to the eigenvalue λ .

- (xi) The spectral radius $\rho(A)$ of the $n \times n$ matrix A is defined

$$\rho(A) = \max |\lambda| \quad (3.19)$$

where λ is an eigenvalue of A . If an eigenvalue is complex with the value $\lambda = \lambda_1 + i\lambda_2$, then $|\lambda| = \sqrt{\lambda_1^2 + \lambda_2^2}$. It can be shown that the spectral radius has the properties

$$(i) \quad \rho(A) \leq \|A\| \quad \text{for any matrix norm } \|\cdot\|.$$

$$(ii) \quad \sqrt{\rho(A^T A)} = \|A\|_2$$

- (xii) The condition number $K(A)$ associated with a nonsingular matrix A is defined

$$K(A) = \|A\| \|A^{-1}\| \quad (3.20)$$

where $\|\cdot\|$ is a natural norm. A matrix A is said to be well behaved or well-conditioned if its condition number is close to unity. It is called ill behaved or ill-conditioned if the condition number is far from unity. Great care should be taken when working with ill-conditioned matrices.

Eigenvalues and Eigenvectors

A $n \times n$ square matrix A can be used as an operator to transform a nonzero $n \times 1$ column vector \bar{x} . One can imagine an input-output system such as the one illustrated in the figure 3-1.

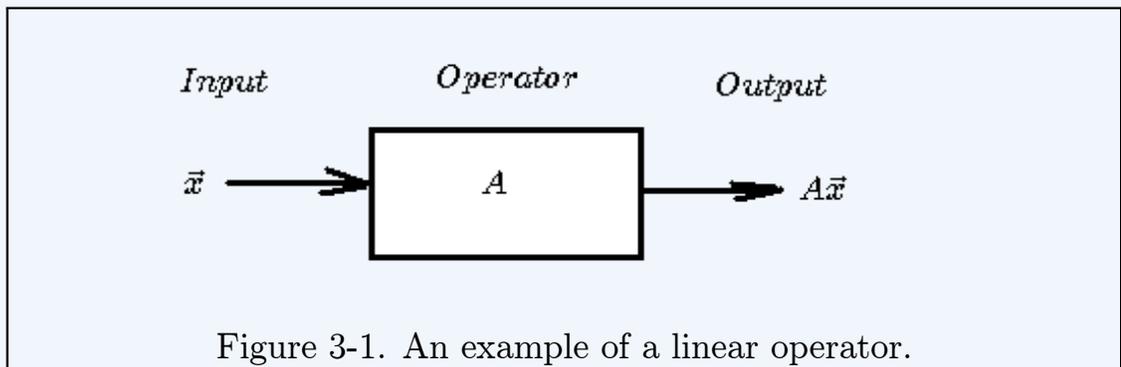


Figure 3-1. An example of a linear operator.

Those nonzero vectors \vec{x} having the special property that the output is proportional to the input must satisfy

$$A\vec{x} = \lambda\vec{x} \quad (3.21)$$

where λ is a scalar proportionality constant. The special nonzero vectors \vec{x} satisfying equation (3.21) are called eigenvectors and the corresponding scalars λ are called eigenvalues. The equation (3.21) is equivalent to the homogeneous system

$$A\vec{x} - \lambda\vec{x} = (A - \lambda I)\vec{x} = \vec{0} \quad (3.22)$$

where I is the $n \times n$ identity matrix. The equation (3.22) has a nonzero solution if and only if

$$\det(A - \lambda I) = |A - \lambda I| = 0. \quad (3.23)$$

The equation (3.23) when expanded is a polynomial equation in λ of degree n having the form

$$C(\lambda) = |A - \lambda I| = (-\lambda)^n + c_{n-1}(-\lambda)^{n-1} + \cdots + c_1(-\lambda) + c_0 = 0 \quad (3.24)$$

which is called the characteristic equation associated with the square matrix A . The solutions of equation (3.24) give the eigenvalues associated with the $n \times n$ square matrix A . For any given eigenvalue λ , the matrix $A - \lambda I$ is a singular matrix such that the homogeneous equations (3.22) produce a nonzero eigenvector \vec{x} . Note that if \vec{x} is an eigenvector, then any nonzero constant times \vec{x} is also an eigenvector.

Example 3-1. (Eigenvalues and Eigenvectors)

Find the eigenvalues and eigenvectors associated with the matrix

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

Solution: We construct the homogeneous system

$$(A - \lambda I)\vec{x} = \vec{0} \quad \text{or} \quad \begin{bmatrix} 1 - \lambda & 0 \\ 1 & 1 - \lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}. \quad (3.25)$$

The characteristic equation is found to be

$$C(\lambda) = \begin{vmatrix} 1 - \lambda & 0 \\ 1 & 1 - \lambda \end{vmatrix} = (1 - \lambda)^2 = 0.$$

The eigenvalues are $\lambda_1 = 1$ and $\lambda_2 = 1$. For $\lambda = 1$, the equation (3.25) reduces to

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Therefore, $\vec{x} = k \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ is an eigenvector for any nonzero constant k . ■

Elementary Row Operations

To solve the system of equations (3.1) one is allowed to perform any of the following elementary row operations on the system of equations.

- (i) An equation in row i can be multiplied by a nonzero constant α . That is, equation E_i is replaced by the equation αE_i . This is denoted by the notation $(\alpha E_i \rightarrow E_i)$ and is read, "The constant α times equation E_i replaces the equation E_i ."
- (ii) Equation E_j can be replaced by a multiple of equation E_i added to equation E_j . This can be expressed using the above notation as $(\alpha E_i + E_j \rightarrow E_j)$ where $i \neq j$.
- (iii) Any two equations can be interchanged. This is denoted by the notation $(E_i \leftrightarrow E_j)$ where $i \neq j$.

Example 3-2. (Elementary row matrices.)

Row operations performed upon the identity matrix I produces elementary row matrices E . Consider the 3×3 identity matrix $I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$, then some examples of elementary row matrices are:

- (i) Interchanging rows 2 and 3 gives

the elementary row matrix $E_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$

- (ii) Multiplying row 3 by the scalar 5 gives

the elementary row matrix $E_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 5 \end{bmatrix}$

- (iii) Multiplying row 1 by 6 and adding the result

to row 2 gives the elementary row matrix $E_3 = \begin{bmatrix} 1 & 0 & 0 \\ 6 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

Note that if $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$ is a 3×3 matrix, then

$$E_1 A = \begin{bmatrix} d & e & f \\ a & b & c \\ g & h & i \end{bmatrix}, \quad E_2 A = \begin{bmatrix} a & b & c \\ d & e & f \\ 5g & 5h & 5i \end{bmatrix}, \quad E_3 A = \begin{bmatrix} a & b & c \\ d + 6a & e + 6b & f + 6c \\ g & h & i \end{bmatrix}$$

Observe that the elementary row operations recorded in the matrices E_1, E_2, E_3 have been applied to the matrix A . ■

A shorthand notation for recording the row operations performed upon a linear system of equations is to write down the coefficient matrix A of the linear system $A\bar{x} = \bar{b}$ and then append to the right-hand side of A the column vector \bar{b} . The resulting array is called an augmented matrix and is written

$$[A|\bar{b}] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & b_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} & b_m \end{array} \right].$$

Our objective is to perform row operations upon the resulting array and try to reduce the array A to an upper triangular form. The row operations are then recorded in the augmented column vector. For example, consider the 2×3 augmented array

$$\left[\begin{array}{cc|c} a_{11} & a_{12} & b_1 \\ a_{21} & a_{22} & b_2 \end{array} \right],$$

with a_{11} nonzero, where we multiply the first row by $-a_{21}/a_{11}$ and add the result to row 2 to obtain the triangular system

$$\left[\begin{array}{cc|c} a_{11} & a_{12} & b_1 \\ a_{21} & a_{22} & b_2 \end{array} \right] \xrightarrow{E_1: -\frac{a_{21}}{a_{11}}E_1 + E_2 \rightarrow E_2} \left[\begin{array}{cc|c} a_{11} & a_{12} & b_1 \\ 0 & c_{22} & d_2 \end{array} \right]$$

where $c_{22} = a_{22} - a_{21}a_{12}/a_{11}$ and $d_2 = b_2 - a_{21}b_1/a_{11}$. The nonzero element a_{11} is called a pivot element in the diagonalization process. The resulting upper triangular system can then be solved by back substitution methods.

Gaussian Elimination

The Gaussian elimination method reduces a matrix to upper triangular (or lower triangular) and then uses back substitution to solve for the unknowns. The method is illustrated using an example.

Example 3-3. (Gaussian elimination method.)

Solve the system of equations

$$\begin{aligned} E_1: & \quad 3x_1 - 2x_2 + x_3 = 8 \\ E_2: & \quad 4x_1 + x_2 - 3x_3 = 3 \\ E_3: & \quad x_1 + 5x_2 - 4x_3 = 5 \end{aligned} \tag{3.26}$$