

Mathematics, rightly viewed, possesses not only truth, but supreme beauty - a beauty cold and austere, like that of sculpture.

— BERTRAND RUSSELL

Chapter 3

System of linear algebraic equation

Topics from linear algebra form the core of numerical analysis. Almost every conceivable problem, be it curve fitting, optimization, simulation of flow sheets or simulation of distributed parameter systems requiring solution of differential equations, require at some stage the solution of a system (often a large system!) of algebraic equations. MATLAB (acronym for MATrix LABoratory) was in fact conceived as a collection of tools to aid in the interactive learning and analysis of linear systems and was derived from a well known core of linear algebra routines written in FORTRAN called LINPACK.

In this chapter we provide a quick review of concepts from linear algebra. We make frequent reference to MATLAB implementation of various concepts throughout this chapter. The reader is encouraged to try out these interactively during a MATLAB session. For a more complete treatment of topics in linear algebra see Hager (1985) and Barnett (1990). The text by Amundson (1966) is also an excellent source with specific examples drawn from Chemical Engineering. For a more rigorous, axiomatic introduction within the frame work of linear operator theory see Ramakrishna and Amundson (1985).

3.1 Matrix notation

We have already used the matrix notation to write a system of linear algebraic equations in a compact form in sections §1.3.1 and §1.3.2. While a matrix, as an object, is represented in bold face, its constituent elements are represented in index notation or as subscripted arrays in programming languages. For example the following are equivalent.

$$\mathbf{A} = [a_{ij}], \quad i = 1, \dots, m; \quad j = 1, \dots, n$$

where \mathbf{A} is an $m \times n$ matrix. a_{ij} represents an element of the matrix \mathbf{A} in row i and column j position. A vector can be thought of as an object with a single row or column. A row vector is represented by,

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]$$

while a column vector can be represented by,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

These elements can be *real* or *complex*.

Having defined objects like vectors and matrices, we can extend the notions of basic arithmetic operations between scalar numbers to higher dimensional objects like vectors and matrices. The reasons for doing so are many. It not only allows us to express a large system of equations in a compact symbolic form, but a study of the properties of such objects allows us to develop and codify very efficient ways of solving and analysing large linear systems. Packages like MATLAB and Mathematica present to us a vast array of such codified algorithms. As an engineer you should develop a conceptual understanding of the underlying principles and the skills to use such packages. But the most important task is to identify each element of a vector or a matrix, which is tied closely to the physical description of the problem.

3.1.1 Review of basic operations

The arithmetic operations are defined both in symbolic form and using index notation. The later actually provides the algorithm for implementing the rules of operation using any programming language. The syntax of these operations in MATLAB are shown with specific examples.*

*MATLAB illustrations have been tested with Ver 5.0.0.4064

The addition operation between two matrices is defined as,

$$\text{addition: } \mathbf{A} = \mathbf{B} + \mathbf{C} \quad \Rightarrow \quad a_{ij} = b_{ij} + c_{ij}$$

This implies an element-by-element addition of the matrices \mathbf{B} and \mathbf{C} . Clearly all the matrices involved must have the same dimension. Note that the addition operation is *commutative* as seen easily with its scalar counterpart. *i.e.*,

$$\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$$

Matrix addition is also *associative*, *i.e.*, independent of the order in which it is carried out, *e.g.*,

$$\mathbf{A} + \mathbf{B} + \mathbf{C} = (\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$$

The scalar multiplication of a matrix involves multiplying each element of the matrix by the scalar, *i.e.*,

$$\text{scalar multiplication: } k\mathbf{A} = \mathbf{B} \quad \Rightarrow \quad k a_{ij} = b_{ij}$$

Subtraction operation can be handled by combining addition and scalar multiplication rules as follows:

$$\text{subtraction: } \mathbf{C} = \mathbf{A} + (-1)\mathbf{B} = \mathbf{A} - \mathbf{B} \quad \Rightarrow \quad c_{ij} = a_{ij} - b_{ij}$$

The product between two matrices \mathbf{A} (of dimension $n \times m$) and \mathbf{B} (of dimension $m \times r$) is defined as,

$$\text{multiplication: } \mathbf{C} = \mathbf{A}\mathbf{B} \quad \Rightarrow \quad c_{ij} = \sum_{k=1}^m a_{ik}b_{kj} \quad \forall \quad i, j$$

and the resultant matrix has the dimension $n \times r$. The operation indicated in the index notation is carried out for each value of the free indices $i = 1 \cdots n$ and $j = 1 \cdots r$. The product is defined only if the dimensions of \mathbf{B} , \mathbf{C} are compatible - *i.e.*, number of *columns* in \mathbf{B} should equal the number of *rows* in \mathbf{C} . This implies that while the product $\mathbf{B} \mathbf{C}$ may be defined, the product $\mathbf{C} \mathbf{B}$ may not even be defined! Even when they are dimensionally compatible, in general

$$\mathbf{BC} \neq \mathbf{CB}$$

i.e., matrix multiplication is not *commutative*.

Example

Consider the matrices A , B defined below.

$$A = \begin{bmatrix} 2 & 3 & 4 \\ 1 & 3 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 2 \\ 3 & 1 \\ 4 & 1 \end{bmatrix}$$

In MATLAB they will be defined as follows:

```

»A=[2 3 4;1 3 2] % Define (2x3) matrix A. Semicolon separates rows
»B=[1 2; 3 1; 4 1]% Define (3x2) matrix B
»C=A*B           % calculate the product
C=              % Display the result
  [ 27  11 ]
  [ 18   7 ]

```

Other useful products can also be defined between vectors and matrices.

A Hadamard (or Schur) product is defined as

$$C = A \circ B \quad \Rightarrow \quad c_{ij} = a_{ij}b_{ij} \quad \forall \quad i, j$$

Obviously, the dimension of A and B should be the same.

Example

The example below illustrates the Hadamard product, called the array product in MATLAB.

```

»C=A'.*B        % Note the dimensions are made the same by transpose of A
C=             % Display the result
  [ 2  2 ]
  [ 9  3 ]
  [16  2 ]

```

A Kronecker product is defined as

$$C = A \otimes B \quad \Rightarrow \quad C = \begin{bmatrix} \mathbf{a}_{11}B & \mathbf{a}_{12}B & \cdots & \mathbf{a}_{1m}B \\ \mathbf{a}_{21}B & \mathbf{a}_{22}B & \cdots & \mathbf{a}_{2m}B \\ \vdots & & & \\ \mathbf{a}_{n1}B & \mathbf{a}_{n2}B & \cdots & \mathbf{a}_{nm}B \end{bmatrix}$$

Multiplying a scalar number by unity leaves it unchanged. Extension of this notion to matrices results in the definition of *identity* matrix,

$$I = \begin{bmatrix} \mathbf{1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \cdots & \mathbf{0} \\ \vdots & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1} \end{bmatrix} \Rightarrow \delta_{ij} = \begin{cases} \mathbf{1} & i = j \\ \mathbf{0} & i \neq j \end{cases}$$

MATLAB function that generates the identity matrix of size N is I=eye(N)

Multiplying any matrix **A** with an identity matrix **I** of appropriate dimension leaves the original matrix unchanged, *i.e.*,

$$AI = A$$

This allows us to generalize the notion of division with scalar numbers to matrices. Division operation can be thought of as the inverse of the multiplication operation. For example, given a number, say 2, we can define its inverse, **x** in such a way that the product of the two numbers produce unity. *i.e.*, $2 \times x = 1$ or $x = 2^{-1}$. In a similar way, given a matrix **A**, can we define the inverse matrix **B** such that

MATLAB function for finding the inverse of a matrix A is B=inv(A)

$$AB = I \quad \text{or} \quad B = A^{-1}$$

The task of developing an algorithm for finding the inverse of a matrix will be addressed late in this chapter.

For a square matrix, *powers* of a matrix **A** can be defined as,

$$A^2 = AA \quad A^3 = AAA = A^2A = AA^2$$

Note that $A^p A^q = A^{p+q}$ for positive integers **p** and **q**. Having extended the definition of powers, we can extend the definition of *exponential* from scalars to square matrices as follows. For a scalar **α** it is,

MATLAB operator for producing the n-th power of a matrix A is, A.^n while the syntax for producing element-by-element power is, A.^n. Make sure that you understand the difference between these two operations!

$$e^\alpha = \mathbf{1} + \alpha + \frac{\alpha^2}{2} + \cdots = \sum_{k=0}^{\infty} \frac{\alpha^k}{k!}$$

For a matrix **A** the *exponential matrix* can be defined as,

$$e^A = I + A + \frac{A^2}{2} + \cdots = \sum_{k=0}^{\infty} \frac{A^k}{k!}$$

One operation that does not have a direct counter part in the scalar world is the *transpose* of a matrix. It is defined the result of exchanging the rows and columns of a matrix, *i.e.*,

$$B = A' \Rightarrow b_{ij} = a_{ji}$$

MATLAB function exp(A) evaluates the exponential element-by-element while expm(A) evaluates the true matrix exponential.

It is easy to verify that

$$(A + B)' = A' + B'$$

Something that is not so easy to verify, nevertheless true, is

$$(AB)' = B'A'$$

3.2 Matrices with special structure

A *diagonal* matrix D has non-zero elements only along the diagonal.

$$D = \begin{bmatrix} \mathbf{d}_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{d}_{22} & \cdots & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{d}_{nn} \end{bmatrix}$$

A *lower triangular* matrix L has non-zero elements on or below the diagonal,

$$L = \begin{bmatrix} \mathbf{l}_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{l}_{21} & \mathbf{l}_{22} & \cdots & \mathbf{0} \\ \vdots & & \ddots & \mathbf{0} \\ \mathbf{l}_{n1} & \mathbf{l}_{n2} & \cdots & \mathbf{l}_{nn} \end{bmatrix}$$

A *upper triangular* matrix U has non-zero elements on or above the diagonal,

$$U = \begin{bmatrix} \mathbf{u}_{11} & \mathbf{u}_{12} & \cdots & \mathbf{u}_{1n} \\ \mathbf{0} & \mathbf{u}_{22} & \cdots & \mathbf{u}_{2n} \\ \mathbf{0} & \mathbf{0} & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{u}_{nn} \end{bmatrix}$$

A *tridiagonal* matrix T has non-zero elements on the diagonal and one off diagonal row on each side of the diagonal

$$T = \begin{bmatrix} \mathbf{t}_{11} & \mathbf{t}_{12} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{t}_{21} & \mathbf{t}_{22} & \mathbf{t}_{23} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \ddots & \ddots & \mathbf{0} \\ \vdots & \mathbf{0} & \mathbf{t}_{n-1,n-2} & \mathbf{t}_{n-1,n-1} & \mathbf{t}_{n-1,n} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{t}_{n,n-1} & \mathbf{t}_{n,n} \end{bmatrix}$$

A *sparse* matrix is a generic term to indicate those matrices without any specific structure such as above, but with a small number (typically 10 to 15 %) of non-zero elements.

3.3 Determinant

A determinant of a square matrix is defined in such a way that a scalar value is associated with the matrix that does not change with certain row or column operations on the matrix - *i.e.*, it is one of the scalar invariants of the matrix. In the context of solving a system of linear equations, the determinant is also useful in knowing whether the system of equations is solvable uniquely. The determinant is formed by summing *all* possible products formed by choosing *one and only one* element from each row and column of the matrix. The precise definition, taken from Amundson (1966), is

$$\det(A) = |A| = \sum (-1)^h (a_{1l_1} a_{2l_2} \cdots a_{nl_n}) \quad (3.1)$$

Each term in the summation consists of a product of n elements selected such that only one element appears from each row and column. The summation involves a total of $n!$ terms accounted for as follows: for the first element l_1 in the product there are n choices, followed by $(n - 1)$ choices for the second element l_2 , $(n - 2)$ choices for the third element l_3 *etc.* resulting in a total of $n!$ choices for a particular product. Note that in this way of counting, the set of second subscripts $\{l_1, l_2, \cdots, l_n\}$ will contain all of the numbers in the range 1 to n , but they will not be in their natural order $\{1, 2, \cdots, n\}$. hence, h is the number of permutations required to arrange $\{l_1, l_2, \cdots, l_n\}$ in their natural order.

This definition is neither intuitive nor computationally efficient. But it is instructive in understanding the following properties of determinants.

1. The determinant of a diagonal matrix D , is simply the product of all the diagonal elements, *i.e.*,

$$\det(D) = \prod_{k=1}^n d_{kk}$$

This is the only product term that is non-zero in equation (3.1).

2. A little thought should convince you that it is the same for lower or upper triangular matrices as well, *viz.*

$$\det(L) = \prod_{k=1}^n l_{kk}$$

3. It should also be clear that if all the elements of any row or column are zero, then the determinant is zero.

MATLAB function for computing the determinant of a square matrix is `det(A)`

4. If every element of any row or column of a matrix is multiplied by a scalar, it is equivalent to multiplying the determinant of the original matrix by the same scalar, *i.e.*,

$$\begin{vmatrix} \mathbf{ka}_{11} & \mathbf{ka}_{12} & \cdots & \mathbf{ka}_{1n} \\ \mathbf{a}_{21} & \mathbf{a}_{22} & \cdots & \mathbf{a}_{2n} \\ \vdots & & \ddots & \vdots \\ \mathbf{a}_{n1} & \mathbf{a}_{n2} & \cdots & \mathbf{a}_{nn} \end{vmatrix} = \begin{vmatrix} \mathbf{a}_{11} & \mathbf{a}_{12} & \cdots & \mathbf{ka}_{1n} \\ \mathbf{a}_{21} & \mathbf{a}_{22} & \cdots & \mathbf{ka}_{2n} \\ \vdots & & \ddots & \vdots \\ \mathbf{a}_{n1} & \mathbf{a}_{n2} & \cdots & \mathbf{ka}_{nn} \end{vmatrix} = k \det(A)$$

5. Replacing any row (or column) of a matrix with a linear combination of that row (or column) and another row (or column) leaves the determinant unchanged.
6. A consequence of rules 3 and 5 is that if two rows (or columns) of a matrix are identical the determinant is zero.
7. If any two rows (or columns) are interchanged, it results in a sign change of the determinant.

3.3.1 Laplace expansion of the determinant

A definition of determinant that you might have seen in an earlier linear algebra course is

$$\det(A) = |A| = \begin{cases} \sum_{k=1}^n \mathbf{a}_{ik} A_{ik} & \text{for any } \mathbf{i} \\ \sum_{k=1}^n \mathbf{a}_{kj} A_{kj} & \text{for any } \mathbf{j} \end{cases} \quad (3.2)$$

where A_{ik} , called the *cofactor*, is given by,

$$A_{ik} = (-1)^{i+k} M_{ik}$$

and M_{ik} , called the *minor*, is the determinant of $(n-1) \times (n-1)$ submatrix of A obtained by deleting i th row and k th column of A . Note that the expansion in equation (3.2) can be carried out along any row \mathbf{i} or column \mathbf{j} of the original matrix A .

Example

Consider the matrix derived in Chapter 1 for the recycle example, *viz.* equation (1.8). Let us calculate the determinant of the matrix using the

Laplace expansion algorithm around the first row.

$$\begin{aligned} \det(A) &= \begin{vmatrix} 1 & 0 & 0.306 \\ 0 & 1 & 0.702 \\ -2 & 1 & 0 \end{vmatrix} \\ &= 1 \begin{vmatrix} 1 & 0.702 \\ 1 & 0 \end{vmatrix} + (-1)^{1+2} \times 0 \times \begin{vmatrix} 0 & 0.702 \\ -2 & 0 \end{vmatrix} + (-1)^{1+3} \times 0.306 \times \begin{vmatrix} 0 & 1 \\ -2 & 1 \end{vmatrix} \\ &= 1 \times (-0.702) + 0 + 0.306 \times 2 = -0.09 \end{aligned}$$

A MATLAB implementation of this will be done as follows:

```

»A=[1 0 0.306; 0 1 0.702; -2 1 0] % Define matrix A
»det(A) % calculate the determinant

```

3.4 Solving a system of linear equations

3.4.1 Cramers rule

Consider a 2×2 system of equations,

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

Direct elimination of the variable x_2 results in

$$(a_{11}a_{22} - a_{12}a_{21}) x_1 = a_{22}b_1 - a_{12}b_2$$

which can be written in an alternate form as,

$$\det(A) x_1 = \det(A(\mathbf{1}))$$

where the matrix $A(\mathbf{1})$ is obtained from A after replacing the first column with the vector \mathbf{b} . *i.e.*,

$$A(\mathbf{1}) = \begin{vmatrix} b_1 & a_{12} \\ b_2 & a_{22} \end{vmatrix}$$

This generalizes to $n \times n$ system as follows,

$$x_1 = \frac{\det(A(\mathbf{1}))}{\det(A)}, \quad \dots \quad x_k = \frac{\det(A(\mathbf{k}))}{\det(A)}, \quad \dots \quad x_n = \frac{\det(A(\mathbf{n}))}{\det(A)}.$$

where $A(\mathbf{k})$ is an $n \times n$ matrix obtained from A by replacing the k th column with the vector \mathbf{b} . It should be clear from the above that, in order to have a unique solution, the determinant of A should be *non-zero*. If the determinant is zero, then such matrices are called *singular*.

Example

Continuing with the recycle problem (equation (1.8) of Chapter 1), solution using Cramer's rule can be implemented with MATLAB as follows:

$$A \mathbf{x} = \mathbf{b} \Rightarrow \begin{bmatrix} 1 & 0 & 0.306 \\ 0 & 1 & 0.702 \\ -2 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 101.48 \\ 225.78 \\ 0 \end{bmatrix}$$

```

»A=[1 0 0.306; 0 1 0.702; -2 1 0]; % Define matrix A
»b=[101.48 225.78 0]'           % Define right hand side vector b
»A1=[b, A(:,[2 3])]            % Define A(1)
»A2=[A(:,1),b, A(:, 3)]        % Define A(2)
»A3=[A(:,[1 2]), b ]           % Define A(3)
»x(1) = det(A1)/det(A)         % solve for component x(1)
»x(2) = det(A2)/det(A)         % solve for component x(2)
»x(3) = det(A3)/det(A)         % solve for component x(3)
»norm(A*x'-b)                  % Check residual

```

3.4.2 Matrix inverse

We defined the inverse of a matrix A as that matrix B which, when multiplied by A produces the identity matrix - *i.e.*, $AB = I$; but we did not develop a scheme for finding B . We can do so now by combining Cramer's rule and Laplace expansion for a determinant as follows. Using Laplace expansion of the determinant of $A(k)$ around column k ,

$$\det A(k) = b_1 A_{1k} + b_2 A_{2k} + \cdots + b_n A_{nk} \quad k = 1, 2, \dots, n$$

where A_{ik} are the cofactors of A . The components of the solution vector, \mathbf{x} are,

$$x_1 = (b_1 A_{11} + b_2 A_{21} + \cdots + b_n A_{n1}) / \det(A)$$

$$x_j = (b_1 A_{1j} + b_2 A_{2j} + \cdots + b_n A_{nj}) / \det(A)$$

$$x_n = (b_1 A_{1n} + b_2 A_{2n} + \cdots + b_n A_{nn}) / \det(A)$$

The right hand side of this system of equations can be written as a vector matrix product as follows,

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \frac{1}{\det(A)} \begin{bmatrix} A_{11} & A_{21} & \cdots & A_{n1} \\ A_{12} & A_{22} & \cdots & A_{n2} \\ \vdots & & \ddots & \vdots \\ A_{1n} & A_{2n} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

or

$$\mathbf{x} = \mathbf{B} \mathbf{b}$$

Premultiplying the original equation $\mathbf{A} \mathbf{x} = \mathbf{b}$ by \mathbf{A}^{-1} we get

$$\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad \text{or} \quad \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$$

Comparing the last two equations, it is clear that,

$$\mathbf{B} = \mathbf{A}^{-1} = \frac{\mathbf{1}}{\det(\mathbf{A})} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{21} & \cdots & \mathbf{A}_{n1} \\ \mathbf{A}_{12} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{n2} \\ \vdots & & \ddots & \vdots \\ \mathbf{A}_{1n} & \mathbf{A}_{2n} & \cdots & \mathbf{A}_{nn} \end{bmatrix} = \frac{\mathit{adj}(\mathbf{A})}{\det(\mathbf{A})}$$

The above equation can be thought of as the definition for the *adjoint* of a matrix. It is obtained by simply replacing each element with its cofactor and then transposing the resulting matrix.

Inverse of a diagonal matrix

Inverse of a diagonal matrix, \mathbf{D} ,

$$\mathbf{D} = \begin{bmatrix} \mathbf{d}_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{d}_{22} & \cdots & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{d}_{nn} \end{bmatrix}$$

is given by,

$$\mathbf{D}^{-1} = \begin{bmatrix} \frac{1}{\mathbf{d}_{11}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \frac{1}{\mathbf{d}_{22}} & \cdots & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \frac{1}{\mathbf{d}_{nn}} \end{bmatrix}$$

It is quite easy to verify using the definition of matrix multiplication that $\mathbf{D}\mathbf{D}^{-1} = \mathbf{I}$.

Inverse of a triangular matrix

Inverse of a triangular matrix is also triangular. Suppose \mathbf{U} is a given upper triangular matrix, then the elements of $\mathbf{V} = \mathbf{U}^{-1}$, can be found sequentially in an efficient manner by simply using the definition $\mathbf{UV} =$

I. This equation, in expanded form, is

$$\begin{bmatrix} \mathbf{u}_{11} & \mathbf{u}_{12} & \cdots & \mathbf{u}_{1n} \\ \mathbf{0} & \mathbf{u}_{22} & \cdots & \mathbf{u}_{2n} \\ \mathbf{0} & \mathbf{0} & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{u}_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{v}_{11} & \mathbf{v}_{12} & \cdots & \mathbf{v}_{1n} \\ \mathbf{v}_{21} & \mathbf{v}_{22} & \cdots & \mathbf{v}_{2n} \\ & & \ddots & \vdots \\ \mathbf{v}_{n1} & \mathbf{v}_{n2} & \cdots & \mathbf{v}_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \cdots & \mathbf{0} \\ & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1} \end{bmatrix}$$

We can develop the algorithm (*i.e.*, find out the rules) by simply carrying out the matrix multiplication on the left hand side and equating it element-by-element to the right hand side. First let us convince ourself that \mathbf{V} is also upper triangular, *i.e.*,

$$\mathbf{v}_{ij} = \mathbf{0} \quad \mathbf{i} > \mathbf{j} \quad (3.3)$$

Consider the element $(\mathbf{n}, \mathbf{1})$ which is obtained by summing the product of each element of \mathbf{n} -th row of \mathbf{U} (consisting mostly of zeros!) with the corresponding element of the 1-st column of \mathbf{V} . The only non-zero term in this product is

$$\mathbf{u}_{nn}\mathbf{v}_{n1} = \mathbf{0}$$

Since $\mathbf{u}_{nn} \neq \mathbf{0}$ it is clear that $\mathbf{v}_{n1} = \mathbf{0}$. Carrying out similar arguments in a sequential manner (in the order $\{\mathbf{i} = \mathbf{n} \cdots \mathbf{j} - \mathbf{1}, \mathbf{j} = \mathbf{1} \cdots \mathbf{n}\}$ - *i.e.*, decreasing order of \mathbf{i} and increasing order of \mathbf{j}) it is easy to verify equation (3.3) and thus establish that \mathbf{V} is also upper triangular.

The non-zero elements of \mathbf{V} can also be found in a sequential manner as follows. For each of the diagonal elements (\mathbf{i}, \mathbf{i}) summing the product of each element of \mathbf{i} -th row of \mathbf{U} with the corresponding element of the \mathbf{i} -th column of \mathbf{V} , the only non-zero term is,

$$\mathbf{v}_{ii} = \frac{\mathbf{1}}{\mathbf{u}_{ii}} \quad \mathbf{i} = \mathbf{1}, \cdots, \mathbf{n} \quad (3.4)$$

Next, for each of the upper elements (\mathbf{i}, \mathbf{j}) summing the product of each element of \mathbf{i} -th row of \mathbf{U} with the corresponding element of the \mathbf{j} -th column of \mathbf{V} , we get,

$$\mathbf{u}_{ii}\mathbf{v}_{ij} + \sum_{r=i+1}^{\mathbf{j}} \mathbf{u}_{ir}\mathbf{v}_{rj} = \mathbf{0}$$

and hence we get,

$$\mathbf{v}_{ij} = -\frac{\mathbf{1}}{\mathbf{u}_{ii}} \sum_{r=i+1}^{\mathbf{j}} \mathbf{u}_{ir}\mathbf{v}_{rj} \quad \mathbf{j} = \mathbf{2}, \cdots, \mathbf{n}; \mathbf{j} > \mathbf{i}; \mathbf{i} = \mathbf{j} - \mathbf{1}, \mathbf{1} \quad (3.5)$$

Note that equation (3.5) should be applied in a specific order, as otherwise, it may involve unknown elements \mathbf{v}_{rj} on the right hand side. First, all of the diagonal elements of \mathbf{V} (viz. \mathbf{v}_{ii}) must be calculated from equation (3.4) as they are needed on the right hand side of equation (3.5). Next the order indicated in equation (3.5), viz. increasing \mathbf{j} from $\mathbf{2}$ to \mathbf{n} and for each \mathbf{j} decreasing \mathbf{i} from $(\mathbf{j} - \mathbf{1})$ to $\mathbf{1}$, should be obeyed to avoid having unknown values appearing on the right hand side of (3.5).

A MATLAB implementation of this algorithm is shown in figure 3.1 to illustrate precisely the order of the calculations. Note that the built-in, general purpose MATLAB inverse function (viz. `inv(U)`) does not take into account the special structure of a triangular matrix and hence is computationally more expensive than the `invu` function of figure 3.1. This is illustrated with the following example.

Example

Consider the upper triangular matrix,

$$U = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 0 & 2 & 3 & 1 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

Let us find its inverse using both the built-in MATLAB function `inv(U)` and the function `invu(U)` of figure 3.1 that is applicable specifically for an upper triangular matrix. You can also compare the floating point operation count for each of the algorithm. Work through the following example using MATLAB. Make sure that the function `invu` of figure 3.1 is in the search path of MATLAB.

```
» U=[1 2 3 4; 0 2 3 1;0 0 1 2; 0 0 0 4]
```

```
U =
```

```

1      2      3      4
0      2      3      1
0      0      1      2
0      0      0      4
```

```
» flops(0) %initialize the flop count
» V=inv(U)
```

V =

```

1.0000   -1.0000         0   -0.7500
         0     0.5000   -1.5000   0.6250
         0         0     1.0000  -0.5000
         0         0         0     0.2500

```

» flops %print flop count

ans =

208

» flops(0);ch3_invu(U),flops %initialize flop count, then invert

ans =

```

1.0000   -1.0000         0   -0.7500
         0     0.5000   -1.5000   0.6250
         0         0     1.0000  -0.5000
         0         0         0     0.2500

```

ans =

57

3.4.3 Gaussian elimination

Gaussian elimination is one of the most efficient algorithms for solving a large system of linear algebraic equations. It is based on a systematic generalization of a rather intuitive elimination process that we routinely apply to a small, say, (2×2) systems. *e.g.*,

$$10x_1 + 2x_2 = 4$$

$$x_1 + 4x_2 = 3$$

From the first equation we have $x_1 = (4 - 2x_2)/10$ which is used to *eliminate* the variable x_1 from the second equation, viz. $(4 - 2x_2)/10 + 4x_2 = 3$ which is solved to get $x_2 = 0.6842$. In the second phase, the

```
function v=invu(u)
% Invert upper triangular matrix
% u - (nxn) matrix
% v - inv(a)

n=size(u,1); % get number of rows in matrix a

for i=2:n
    for j=1:i-1
        v(i,j)=0;
    end
end

for i=1:n
    v(i,i)=1/u(i,i);
end

for j=2:n
    for i=j-1:-1:1
        v(i,j) = -1/u(i,i)*sum(u(i,i+1:j)*v(i+1:j,j));
    end
end
```

Figure 3.1: MATLAB implementation of inverse of an upper triangular matrix

value of \mathbf{x}_2 is *back substituted* into the first equation and we get $\mathbf{x}_1 = \mathbf{0.2632}$. We could have reversed the order and eliminated \mathbf{x}_1 from the first equation after rearranging the second equation as $\mathbf{x}_1 = (\mathbf{3} - \mathbf{4}\mathbf{x}_2)$. Thus there are two phases to the algorithm: (a) forward elimination of one variable at a time until the last equation contains only one unknown; (b) back substitution of variables. Also, note that we have used two rules during the elimination process: (i) two equations (or two rows) can be interchanged as it is merely a matter of book keeping and it does not in any way alter the problem formulation, (ii) we can replace any equation with a linear combination of itself with another equation. A conceptual description of a naive Gaussian elimination algorithm is shown in figure 3.2. All of the arithmetic operations needed to eliminate one variable at a time are identified in the illustration. Study that carefully.

We call it a naive scheme as we have assumed that none of the diagonal elements are zero, although this is not a requirement for existence of a solution. The reason for avoiding zeros on the diagonals is to avoid division by zeros in step 2 of the illustration 3.2. If there are zeros on the diagonal, we can interchange two equations in such a way the diagonals do not contain zeros. This process is called *pivoting*. Even if we organize the equations in such a way that there are no zeros on the diagonal, we may end up with a zero on the diagonal during the elimination process (likely to occur in step 3 of illustration 3.2). If that situation arises, we can continue to exchange that particular row with another one to avoid division by zero. If the matrix is *singular* we will eventually end up with an unavoidable zero on the diagonal. This situation will arise if the original set of equations is not *linearly independent*; in other words the *rank* of the matrix \mathbf{A} is less than \mathbf{n} . Due to the finite precision of computers, the floating point operation in step 3 of illustration 3.2 will not result usually in an exact zero, but rather a very small number. Loss of precision due to round off errors is a common problem with direct methods involving large system of equations since any error introduced at one step corrupts all subsequent calculations.

A MATLAB implementation of the algorithm is given in figure 3.3 through the function `gauss.m`.

Example

Let us continue with the recycle problem (equation (1.8) of Chapter 1). First we obtain solution using the built-in MATLAB linear equation solver (viz. $\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$) and record the floating point operations (flops). Then we solve with the Gaussian elimination function `gauss` and compare the flops. Note that in order to use the naive Gaussian elimination function,

Note that it is merely for illustrating the concepts involved in the elimination process; MATLAB `backslash, \` operator provides a much more elegant solution to solve $\mathbf{Ax} = \mathbf{b}$ in the form $\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$.

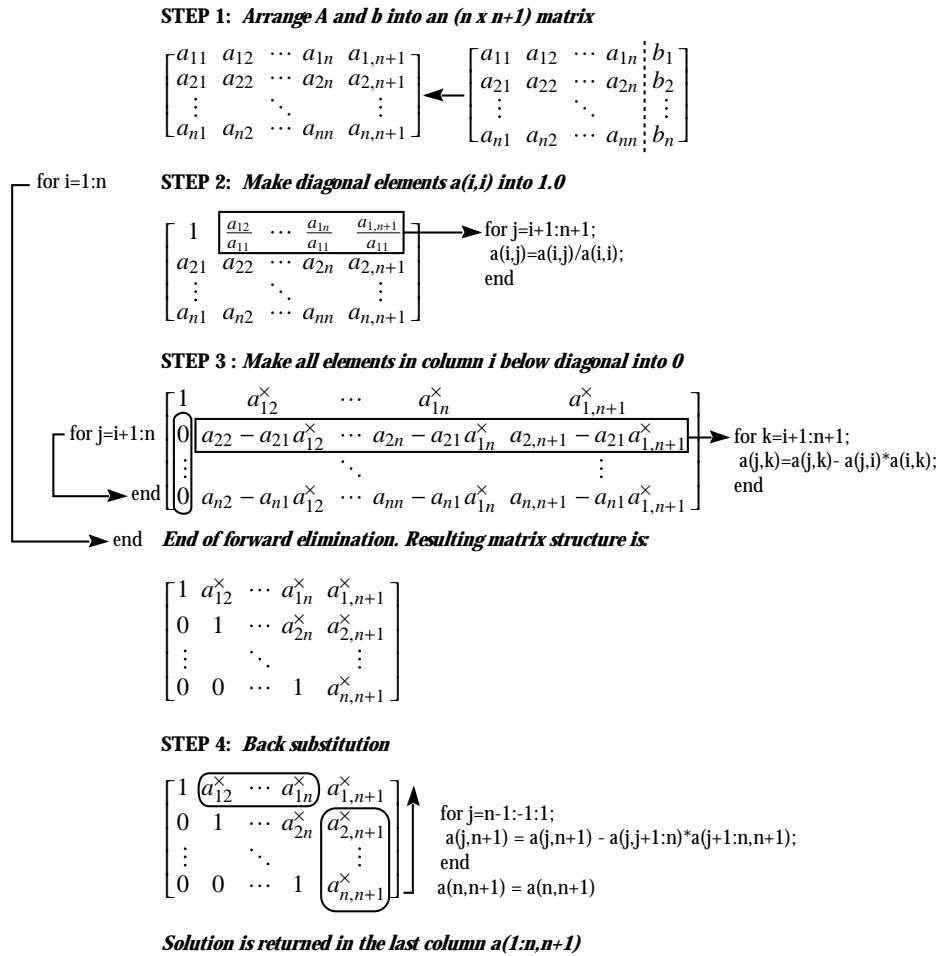


Figure 3.2: Naive Gaussian elimination scheme

```
function x=gauss(a,b)
% Naive Gaussian elimination. Cannot have zeros on diagonal
% a - (nxn) matrix
% b - column vector of length n

m=size(a,1); % get number of rows in matrix a
n=length(b); % get length of b
if (m ~= n)
    error('a and b do not have the same number of rows')
end

%Step 1: form (n,n+1) augmented matrix
a(:,n+1)=b;

for i=1:n
%Step 2: make diagonal elements into 1.0
a(i,i+1:n+1) = a(i,i+1:n+1)/a(i,i);

%Step 3: make all elements below diagonal into 0
    for j=i+1:n
        a(j,i+1:n+1) = a(j,i+1:n+1) - a(j,i)*a(i,i+1:n+1);
    end
end

%Step 4: begin back substitution
for j=n-1:-1:1
    a(j,n+1) = a(j,n+1) - a(j,j+1:n)*a(j+1:n,n+1);
end

%return solution
x=a(:,n+1)
```

Figure 3.3: MATLAB implementation of naive Gaussian elimination

we need to switch the 2nd and 3rd equations to avoid division by zero.

```

»A=[1 0 0.306;           % matrix entry continues on next two lines
   -2 1 0;
   0 1 0.702];          % Definition of matrix A complete
»b=[101.48 0 225.78]'; % Define right hand side column vector b
»flops(0);              % initialize flop count
»A\b                    % solution is : [23.8920 47.7840 253.5556]
»flops                  % examine flops for MATLAB internal solver (ans: 71)
»flops(0);              % initialize flop count
»gauss(A,b)             % solution is, of course : [23.8920 47.7840 253.5556]
»flops                  % examine flops for gauss solver (ans: 75)
»flops(0);              % initialize flop count
»inv(A)*b               % obtain solution using matrix inverse
»flops                  % examine flops for MATLAB internal solver (ans: 102)

```

Example - loss of accuracy & need for pivoting

The need for pivoting can be illustrated with the following simple example.

$$\epsilon x_1 + x_2 = 1$$

$$x_1 + x_2 = 2$$

where ϵ is a small number. In matrix form it will be,

$$\begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

In using naive Gaussian elimination without rearranging the equations, we first make the diagonal into unity, which results in

$$x_1 + \frac{1}{\epsilon} x_2 = \frac{1}{\epsilon}$$

Next we eliminate the variable x_1 from the 2nd equation which results in,

$$\left(1 - \frac{1}{\epsilon}\right) x_2 = 2 - \frac{1}{\epsilon}$$

Rearranging this and using back substitution we finally get x_2 and x_1 as,

$$x_2 = \frac{2 - \frac{1}{\epsilon}}{1 - \frac{1}{\epsilon}}$$

ϵ	Naive elimination without pivoting gauss(A,b)	Built-in MATLAB $A \setminus \mathbf{b}$
1×10^{-15}	[1 1]	[1, 1]
1×10^{-16}	[2 1]	[1, 1]
1×10^{-17}	[0 1]	[1, 1]

Table 3.1: Loss of precision and need for pivoting

$$\mathbf{x}_1 = \frac{1}{\epsilon} - \frac{\mathbf{x}_2}{\epsilon}$$

The problem in computing \mathbf{x}_1 as $\epsilon \rightarrow \mathbf{0}$ should be clear now. As ϵ crosses the threshold of finite precision of the computation (hardware or software), taking the difference of two large numbers of comparable magnitude, can result in significant loss of precision. Let us solve the problem once again after rearranging the equations as,

$$\mathbf{x}_1 + \mathbf{x}_2 = 2$$

$$\epsilon \mathbf{x}_1 + \mathbf{x}_2 = 1$$

and apply Gaussian elimination once again. Since the diagonal element in the first equation is already unity, we can eliminate \mathbf{x}_1 from the 2nd equation to obtain,

$$(1 - \epsilon)\mathbf{x}_2 = 1 - 2\epsilon \quad \text{or} \quad \mathbf{x}_2 = \frac{1 - 2\epsilon}{1 - \epsilon}$$

Back substitution yields,

$$\mathbf{x}_1 = 2 - \mathbf{x}_2$$

Both these computations are well behaved as $\epsilon \rightarrow \mathbf{0}$.

We can actually demonstrate this using the MATLAB function `gauss` shown in figure 3.3 and compare it with the MATLAB built-in function `A \setminus \mathbf{b}` which does use pivoting to rearrange the equations and minimize the loss of precision. The results are compared in table 3.1 for ϵ in the range of 10^{-15} to 10^{-17} . Since MATLAB uses double precision, this range of ϵ is the threshold for loss of precision. Observe that the naive Gaussian elimination produces incorrect results for $\epsilon < 10^{-16}$.

3.4.4 Thomas algorithm

Many problems such as the stagewise separation problem we saw in section §1.3.1 or the solution of differential equations that we will see in

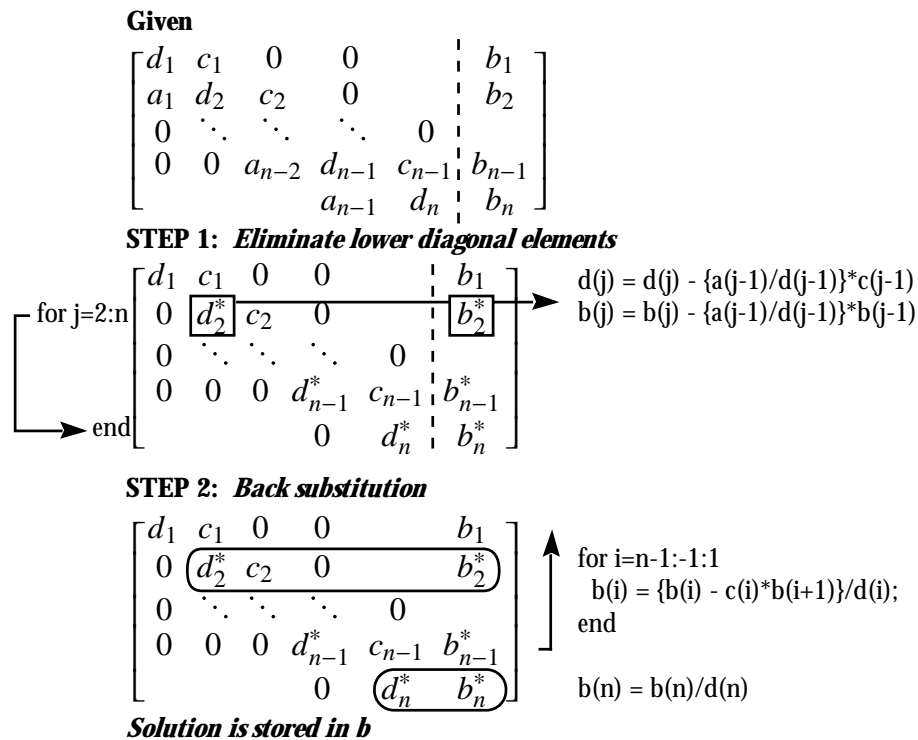


Figure 3.4: Thomas algorithm

later chapters involve solving a system of linear equations $T\mathbf{x} = \mathbf{b}$ with a *tridiagonal* matrix structure.

$$T = \begin{bmatrix} d_1 & c_1 & \mathbf{0} & \cdots & \mathbf{0} \\ a_1 & d_2 & c_2 & \cdots & \mathbf{0} \\ & & \ddots & & \\ \mathbf{0} & \mathbf{0} & a_{n-2} & d_{n-1} & c_{n-1} \\ \mathbf{0} & \cdots & \mathbf{0} & a_{n-1} & d_n \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix}$$

Since we know where the zero elements are, we do not have to carry out the elimination steps on those entries of the matrix T ; but the essential steps in the algorithm remain the same as in the Gaussian elimination scheme and are illustrated in figure 3.4. MATLAB implementation is shown in figure 3.5.

```
function x=thomas(a,b,c,d)
% Thomas algorithm for tridiagonal systems
% d - diagonal elements, n
% b - right hand side forcing term, n
% a - lower diagonal elements, (n-1)
% c - upper diagonal elements, (n-1)

na=length(a); % get length of a
nb=length(b); % get length of b
nc=length(c); % get length of c
nd=length(d); % get length of d
if (nd ~= nb | na ~= nc | (nd-1) ~= na)
    error('array dimensions not consistent')
end
n=length(d);
%Step 1: forward elimination

for i=2:n
    fctr=a(i-1)/d(i-1);
    d(i) = d(i) - fctr*c(i-1);
    b(i) = b(i) - fctr*b(i-1);
end

%Step 2: back substitution
b(n) = b(n)/d(n);
for j=n-1:-1:1
    b(j) = (b(j) - c(j)*b(j+1))/d(j);
end

%return solution
x=b;
```

Figure 3.5: MATLAB implementation of Thomas algorithm

3.4.5 Gaussian elimination - Symbolic representation

Given a square matrix A of dimension $n \times n$ it is possible to write it as the product of two matrices B and C , *i.e.*, $A = BC$. This process is called *factorization* and is in fact not at all unique - *i.e.*, there are infinitely many possibilities for B and C . This is clear with a simple counting of the unknowns - *viz.* there are $2 \times n^2$ unknown elements in B and C while only n^2 equations can be obtained by equating each element of A with the corresponding element from the product BC .

The extra degrees of freedom can be used to specify any specific structure for B and C . For example we can require $B = L$ be a lower triangular matrix and $C = U$ be an upper triangular matrix. This process is called LU factorization or decomposition. Since each triangular matrix has $n \times (n + 1)/2$ unknowns, we still have a total of $n^2 + n$ unknowns. The extra n degrees of freedom is often used in one of three ways:

- *Doolittle* method assigns the diagonal elements of L to be unity.
- *Crout* method assigns the diagonal elements of U to be unity.
- *Cholesky* method assigns the diagonal elements of L to be equal to that of U - *i.e.*, $l_{ii} = u_{ii}$.

While a simple degree of freedom analysis, indicates that it is possible to factorize a matrix into a product of lower and upper triangular matrices, it does not tell us how to find out the unknown elements.

Revisiting the Gaussian elimination method from a different perspective, will show the connection between LU factorization and Gaussian elimination. Note that the algorithm outlined in section §3.4.3 is the most computationally efficient scheme for implementing Gaussian elimination. The method to be outlined below is not computationally efficient, but it is a useful conceptual aid in showing the connection between Gaussian elimination and LU factorization. Steps 2 and 3 of figure 3.2 that involve making the diagonal into unity and all the elements below the diagonal into zero is equivalent to pre-multiplying A by L_1 - *i.e.*, $L_1 A = U_1$ or,

$$\begin{bmatrix} \frac{1}{a_{11}} & \mathbf{0} & \cdots & \mathbf{0} \\ -\frac{a_{21}}{a_{11}} & \mathbf{1} & \cdots & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ -\frac{a_{n1}}{a_{11}} & \mathbf{0} & \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & a_{12}^{(1)} & \cdots & a_{1n}^{(1)} \\ \mathbf{0} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & a_{n2}^{(1)} & \cdots & a_{nn}^{(1)} \end{bmatrix}$$

Repeating this process for the 2nd row, we pre-multiply U_1 by L_2 - i.e., $L_2U_1 = U_2$ or, in expanded form,

$$\begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \frac{1}{a_{22}^{(1)}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & -\frac{a_{32}^{(1)}}{a_{22}^{(1)}} & \mathbf{1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & -\frac{a_{n2}^{(1)}}{a_{22}^{(1)}} & \mathbf{0} & \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{a}_{12}^{(1)} & \cdots & \mathbf{a}_{1n}^{(1)} \\ \mathbf{0} & \mathbf{a}_{22}^{(1)} & \cdots & \mathbf{a}_{2n}^{(1)} \\ & & \ddots & \vdots \\ \mathbf{0} & \mathbf{a}_{n2}^{(1)} & \cdots & \mathbf{a}_{nn}^{(1)} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{a}_{12}^{(1)} & \mathbf{a}_{13}^{(1)} & \cdots & \mathbf{a}_{1n}^{(1)} \\ \mathbf{0} & \mathbf{1} & \mathbf{a}_{23}^{(2)} & \cdots & \mathbf{a}_{2n}^{(2)} \\ \mathbf{0} & \mathbf{0} & \mathbf{a}_{33}^{(2)} & \cdots & \vdots \\ \mathbf{0} & \mathbf{0} & \vdots & \cdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{a}_{n3}^{(2)} & \cdots & \mathbf{a}_{nn}^{(2)} \end{bmatrix}$$

Continuing this process, we obtain in succession,

$$L_1A = U_1$$

$$L_2U_1 = U_2$$

$$L_3U_2 = U_3$$

$$L_4U_3 = U_4$$

$$L_{n-1}U_{n-2} = U_{n-1}$$

Note that each L_j is a lower triangular matrix with non-zero elements on the j -th column and unity on other diagonal elements. Eliminating all of the intermediate U_j we obtain,

$$(L_{n-1}L_{n-2} \cdots L_1)A = U_{n-1}$$

Since the product of all lower triangular matrices is yet another lower triangular matrix, we can write the above equation as,

$$LA = U$$

Also, the inverse of a lower triangular matrix is also lower triangular - i.e., L^{-1} . Hence a given square matrix A can be factored into a product of a lower and upper triangular matrix as,

$$A = L^{-1}U = LU$$

Although the development in this section provides us with an algorithm for constructing both L and U , it is quite inefficient. A more direct and efficient algorithm is developed next in section §3.4.6.

3.4.6 LU decomposition

Consider the product of L and U as shown in the expanded form below. All of the elements of L and U are unknown. By carrying out the matrix product on the left hand side and equating element-by-element to the right hand side, we can develop sufficient number of equations to find out all of the unknown elements on the left hand side. The trick, however is, (as we did with inverting a triangular matrix) to carry out the calculations in a particular sequence so that no more than one unknown appears in each equation.

$$\begin{bmatrix} l_{11} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ l_{21} & l_{22} & \mathbf{0} & \cdots & \mathbf{0} \\ l_{31} & l_{32} & l_{33} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \mathbf{0} \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{u}_{12} & \mathbf{u}_{13} & \cdots & \mathbf{u}_{1n} \\ \mathbf{0} & \mathbf{1} & \mathbf{u}_{23} & \cdots & \mathbf{u}_{2n} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \cdots & \mathbf{u}_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{11} & \mathbf{a}_{12} & \cdots & \mathbf{a}_{1n} \\ \mathbf{a}_{21} & \mathbf{a}_{22} & \cdots & \mathbf{a}_{2n} \\ \mathbf{a}_{31} & \mathbf{a}_{32} & \cdots & \mathbf{a}_{3n} \\ \vdots & \vdots & \cdots & \vdots \\ \mathbf{a}_{n1} & \mathbf{a}_{n2} & \cdots & \mathbf{a}_{nn} \end{bmatrix}$$

Let us first consider elements in column 1 of L . Carrying out the multiplication and equating we obtain,

$$l_{i1} = a_{i1} \quad i = 1, \dots, n \quad (3.6)$$

Next focusing on the elements in the first row of U we get,

$$u_{1j} = a_{1j}/l_{11} \quad j = 2, \dots, n \quad (3.7)$$

It would be inefficient to proceed to the 2nd column of L . Why?

Next we alternate between a column of L and a row of U . The general expression for any element i in column j of L is,

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \quad j = 2, \dots, n \quad i = j, \dots, n \quad (3.8)$$

Similarly the general expression for any element i in row j of U is,

$$u_{ji} = \frac{a_{ji} - \sum_{k=1}^{j-1} l_{jk} u_{ki}}{l_{jj}} \quad j = 2, \dots, n \quad i = j + 1, \dots, n \quad (3.9)$$

Equations (3.6-3.9) form the basic algorithm for LU decomposition. In order to illustrate the implementation of equations (3.6-3.9) as an algorithm, a MATLAB function called `LU.m` is shown in figure 3.6. Note that MATLAB provides a built-in function for LU decomposition called `lu(A)`.

```
function [L,U]=LU(a)
% Naive LU decomposition
% a - (nxn) matrix
% L,U - are (nxn) factored matrices
% Usage [L,U]=LU(A)

n=size(a,1); % get number of rows in matrix a

%Step 1: first column of L
L(:,1)=a(:,1);

%Step 2: first row of U
U(1,:)=a(1,:)/L(1,1);

%Step 3: Alternate between column of L and row of U
for j=2:n
    for i = j:n
        L(i,j) = a(i,j) - sum(L(i,1:j-1)'.*U(1:j-1,j));
    end
    U(j,j) = 1;
    for i=j+1:n
        U(j,i)=(a(j,i) - sum(L(j,1:j-1)'.*U(1:j-1,i) ) )/L(j,j);
    end
end
end
```

Figure 3.6: MATLAB implementation of LU decomposition algorithm

3.5 Iterative algorithms for systems of linear equations

The direct methods discussed in section §3.4 have the advantage of producing the solution in a finite number of calculations. They suffer, however, from loss of precision due to accumulated round off errors. This problem is particularly severe in large dimensional systems (more than 10,000 equations). Iterative methods, on the other hand, produce the result in an asymptotic manner by repeated application of a simple algorithm. Hence the number of floating point operations required to produce the final result cannot be known *a priori*. But they have the natural ability to eliminate errors at every step of the iteration. For an authoritative account of iterative methods for large linear systems see Young (1971).

Iterative methods rely on the concepts developed in Chapter 2. They are extended naturally from a single equation (one-dimensional system) to a system of equations (n -dimensional system). The development parallels that of section §2.7 on fixed point iterations schemes. Given an equation of the form, $\mathbf{A}\mathbf{x} = \mathbf{b}$ we can rearrange it into a form,

$$\mathbf{x}^{(p+1)} = \mathbf{G}(\mathbf{x}^{(p)}) \quad p = 0, 1, \dots \quad (3.10)$$

Here we can view the vector \mathbf{x} as a point in a n -dimensional vector space and the above equation as an iterative map that maps a point $\mathbf{x}^{(p)}$ into another point $\mathbf{x}^{(p+1)}$ in the n -dimensional vector space. Starting with an initial guess $\mathbf{x}^{(0)}$ we calculate successive iterates $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \dots$ until the sequence converges. The only difference from chapter 2 is that the above iteration is applied to a higher dimensional system of (n) equations. Note that $\mathbf{G}(\mathbf{x})$ is also vector. Since we are dealing with a linear system, \mathbf{G} will be a linear function of \mathbf{x} which is constructed from the given \mathbf{A} matrix. \mathbf{G} can typically be represented as

$$\mathbf{x}^{(p+1)} = \mathbf{G}(\mathbf{x}^{(p)}) = \mathbf{T}\mathbf{x}^{(p)} + \mathbf{c}. \quad (3.11)$$

In section §2.7 we saw that a given equation $\mathbf{f}(\mathbf{x}) = \mathbf{0}$ can be rearranged into the form $\mathbf{x} = \mathbf{g}(\mathbf{x})$ in several different ways. In a similar manner, a given equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be rearranged into the form $\mathbf{x}^{(p+1)} = \mathbf{G}(\mathbf{x}^{(p)})$ in more than one way. Different choices of \mathbf{G} results in different iterative methods. In section §2.7 we also saw that the condition for convergence of the sequence $\mathbf{x}_{i+1} = \mathbf{g}(\mathbf{x}_i)$ is $\mathbf{g}'(\mathbf{r}) < \mathbf{1}$. Recognizing that the derivative of $\mathbf{G}(\mathbf{x}^{(p)})$ with respect to $\mathbf{x}^{(p)}$ is a matrix, $\mathbf{G}' = \mathbf{T}$ a convergence condition similar to that found for the scalar case must depend on the properties of the matrix \mathbf{T} . Another way to demonstrate

this is as follows. Once the sequence $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \dots$ converges to, say, \mathbf{r} equation (3.11) becomes,

$$\mathbf{r} = T\mathbf{r} + \mathbf{c}.$$

Subtracting equation (3.11) from the above,

$$(\mathbf{x}^{(p+1)} - \mathbf{r}) = T(\mathbf{x}^{(p)} - \mathbf{r}).$$

Now, recognizing that $(\mathbf{x}^{(p)} - \mathbf{r}) = \boldsymbol{\epsilon}^{(p)}$ is a measure of the error at iteration level p , we have

$$\boldsymbol{\epsilon}^{(p+1)} = T\boldsymbol{\epsilon}^{(p)}.$$

Thus, the error at step $(p + 1)$ depend on the error at step (p) . If the matrix T has the property of amplifying the error at any step, then the iterative sequence will diverge. The property of the matrix T that determines this feature is called the *spectral radius*. The *spectral radius* is defined as the largest eigenvalue in magnitude of T . For convergence of the iterative sequence the *spectral radius* of T should be less than unity,

$$\boxed{\rho(T) < 1} \quad (3.12)$$

3.5.1 Jacobi iteration

The Jacobi iteration rearranges the given equations in the form,

$$\begin{aligned} \mathbf{x}_1^{(p+1)} &= (\mathbf{b}_1 - \mathbf{a}_{12}\mathbf{x}_2^{(p)} - \mathbf{a}_{13}\mathbf{x}_3^{(p)} - \dots - \mathbf{a}_{1n}\mathbf{x}_n^{(p)})/\mathbf{a}_{11} \\ \mathbf{x}_j^{(p+1)} &= \left[\mathbf{b}_j - \sum_{k=1}^{j-1} \mathbf{a}_{jk}\mathbf{x}_k^{(p)} - \sum_{k=j+1}^n \mathbf{a}_{jk}\mathbf{x}_k^{(p)} \right] / \mathbf{a}_{jj} \\ \mathbf{x}_n^{(p+1)} &= (\mathbf{b}_n - \mathbf{a}_{n1}\mathbf{x}_1^{(p)} - \mathbf{a}_{n2}\mathbf{x}_2^{(p)} - \dots - \mathbf{a}_{n,n-1}\mathbf{x}_{n-1}^{(p)})/\mathbf{a}_{nn} \end{aligned} \quad (3.13)$$

where the variable \mathbf{x}_j has been extracted form the j - *th* equation and expressed as a function of the remaining variables. The above set of equations can be applied repetitively to update each component of the unknown vector $\mathbf{x}=(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ provided an initial guess is known for \mathbf{x} . The above equation can be written in matrix form as,

$$L\mathbf{x}^{(p)} + D\mathbf{x}^{(p+1)} + U\mathbf{x}^{(p)} = \mathbf{b}$$

where the matrices D, L, U are defined in term of components of A as follows.

$$D = \begin{bmatrix} \mathbf{a}_{11} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{a}_{22} & \dots & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{a}_{nn} \end{bmatrix}$$

Note that MATLAB functions `diag`, `tril`, `triu` are useful in extracting parts of a given matrix A

$$L = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{a}_{21} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & & \ddots & \mathbf{0} \\ \mathbf{a}_{n1} & \mathbf{a}_{n2} & \cdots & \mathbf{0} \end{bmatrix} \quad U = \begin{bmatrix} \mathbf{0} & \mathbf{a}_{12} & \cdots & \mathbf{a}_{1n} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{a}_{2n} \\ \mathbf{0} & \mathbf{0} & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix}$$

which can be rearranged as,

$$\mathbf{x}^{(p+1)} = D^{-1}(\mathbf{b} - (L + U)\mathbf{x}^{(p)}) \quad (3.14)$$

and hence $\mathbf{G}(\mathbf{x}^{(p)}) = -D^{-1}(L + U)\mathbf{x}^{(p)} + D^{-1}\mathbf{b}$ and $\mathbf{G}' = \mathbf{T} = -D^{-1}(L + U)$. This method has been shown to be convergent as long as the original matrix A is diagonally dominant, *i.e.*,

An examination of equation (3.13) reveals that none of the diagonal elements can be zero. If any is found to be zero, one can easily exchange the positions of any two equations to avoid this problem. Equation (3.13) is used in actual computational implementation, while the matrix form of the equation (3.14) is useful for conceptual description and convergence analysis. Note that each element in the equation set (3.13) can be updated independent of the others in any order because the right hand side of equation (3.13) is evaluated at the p -th level of iteration. This method requires that $\mathbf{x}^{(p)}$ and $\mathbf{x}^{(p+1)}$ be stored as two separate vectors until all elements of $\mathbf{x}^{(p+1)}$ have been updated using equation (3.13). A minor variation of the algorithm which uses a new value of the element in $\mathbf{x}^{(p+1)}$ as soon as it is available is called the Gauss-Seidel method. It has the dual advantage of faster convergence than the Jacobi iteration as well as reduced storage requirement for only one array \mathbf{x} .

3.5.2 Gauss-Seidel iteration

In the Gauss-Seidel iteration we rearrange the given equations in the form,

$$\begin{aligned} \mathbf{x}_1^{(p+1)} &= (\mathbf{b}_1 - \mathbf{a}_{12}\mathbf{x}_2^{(p)} - \mathbf{a}_{13}\mathbf{x}_3^{(p)} - \cdots - \mathbf{a}_{1n}\mathbf{x}_n^{(p)})/\mathbf{a}_{11} \\ \mathbf{x}_j^{(p+1)} &= \left[\mathbf{b}_j - \sum_{k=1}^{j-1} \mathbf{a}_{jk}\mathbf{x}_k^{(p+1)} - \sum_{k=j+1}^n \mathbf{a}_{jk}\mathbf{x}_k^{(p)} \right] / \mathbf{a}_{jj} \\ \mathbf{x}_n^{(p+1)} &= (\mathbf{b}_n - \mathbf{a}_{n1}\mathbf{x}_1^{(p+1)} - \mathbf{a}_{n2}\mathbf{x}_2^{(p+1)} - \cdots - \mathbf{a}_{n,n-1}\mathbf{x}_{n-1}^{(p+1)})/\mathbf{a}_{nn} \end{aligned} \quad (3.15)$$

Observe that known values of the elements in $\mathbf{x}^{(p+1)}$ are used on the right hand side of the above equations (3.15) as soon as they are available within the same iteration. We have used the superscripts p and $(p + 1)$ explicitly in equation (3.15) to indicate where the newest values occur. In a computer program there is no need to assign separate arrays for

p and $(p + 1)$ levels of iteration. Using just a single array for \mathbf{x} will automatically propagate the newest values as soon as they are updated.

The above equation can be written symbolically in matrix form as,

$$\mathbf{L}\mathbf{x}^{(p+1)} + \mathbf{D}\mathbf{x}^{(p+1)} + \mathbf{U}\mathbf{x}^{(p)} = \mathbf{b}$$

where the matrices \mathbf{D} , \mathbf{L} , \mathbf{U} are defined as before. Factoring $\mathbf{x}^{(p+1)}$ we get,

$$\mathbf{x}^{(p+1)} = (\mathbf{L} + \mathbf{D})^{-1}(\mathbf{b} - \mathbf{U}\mathbf{x}^{(p)}) \quad (3.16)$$

and hence $\mathbf{G}(\mathbf{x}^{(p)}) = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}\mathbf{x}^{(p)} + (\mathbf{L} + \mathbf{D})^{-1}\mathbf{b}$ and $\mathbf{G}' = \mathbf{T} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}$. Thus the convergence of this scheme depends on the spectral radius of the matrix, $\mathbf{T} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}$. This method has also been shown to be convergent as long as the original matrix \mathbf{A} is diagonally dominant.

MATLAB implementation of the Gauss-Seidel algorithm is shown in figure 3.7.

3.5.3 Successive over-relaxation (SOR) scheme

The relaxation scheme can be thought of as a convergence acceleration scheme that can be applied to any of the basic iterative methods like Jacobi or Gauss-Seidel schemes. We introduce an extra parameter, ω often called the *relaxation parameter* and choose its value in such a way that we can either speed up convergence by using $\omega > 1$ (called *over-relaxation*) or in some difficult problems with poor initial guess we can attempt to enlarge the region of convergence using $\omega < 1$ (called *under-relaxation*). Let us illustrate the implementation with the Gauss-Seidel scheme. The basic Gauss-Seidel scheme is:

$$\mathbf{t} := \mathbf{x}_j^{(p+1)} = \left[\mathbf{b}_j - \sum_{k=1}^{j-1} \mathbf{a}_{jk}\mathbf{x}_k^{(p+1)} - \sum_{k=j+1}^n \mathbf{a}_{jk}\mathbf{x}_k^{(p)} \right] / \mathbf{a}_{jj} \quad (3.17)$$

Instead of accepting the value of $\mathbf{x}_j^{(p+1)}$ computed from the above formula as the current value, we store it in a temporary variable \mathbf{t} and form a better (or accelerated) estimate of $\mathbf{x}_j^{(p+1)}$ from,

$$\mathbf{x}_j^{(p+1)} = \mathbf{x}_j^{(p)} + \omega [\mathbf{t} - \mathbf{x}_j^{(p)}]$$

Observe that if $\omega = 1$, the method remains the same as Gauss-Seidel scheme. For $\omega > 1$, then the difference between two successive iterates (the term in the square brackets) is amplified and added to the current value $\mathbf{x}_j^{(p)}$.

```
function x=GS(a,b,x,tol,max)
% Gauss-Seidel iteration
% a - (nxn) matrix
% b - column vector of length n
% x - initial guess vector x
% tol - convergence tolerance
% max - maximum number of iterations
% Usage x=GS(A,b,x)

m=size(a,1); % get number of rows in matrix a
n=length(b); % get length of b
if (m ~= n)
    error('a and b do not have the same number of rows')
end
if nargin < 5, max=100; end
if nargin < 4, max=100; tol=eps; end
if nargin == 2
    error('Initial guess is required')
end
count=0;

while (norm(a*x-b) > tol & count < max),
    x(1) = ( b(1) - a(1,2:n)*x(2:n) )/a(1,1);
    for i=2:n-1
        x(i) = (b(i) - a(i,1:i-1)*x(1:i-1) - ...
                a(i,i+1:n)*x(i+1:n) )/a(i,i);
    end
    x(n) = ( b(n) - a(n,1:n-1)*x(1:n-1) )/a(n,n);
    count=count+1;
end

if (count >= max)
    fprintf(1,'Maximum iteration %3i exceeded\n',count)
    fprintf(1,'Residual is %12.5e\n ',norm(a*x-b) )
end
```

Figure 3.7: MATLAB implementation of Gauss-Seidel algorithm

The above operations can be written in symbolic matrix form as,

$$\mathbf{x}^{(p+1)} = \mathbf{x}^{(p)} + \omega[\{\mathbf{D}^{-1}(\mathbf{b} - \mathbf{L}\mathbf{x}^{(p+1)} - \mathbf{U}\mathbf{x}^{(p)})\} - \mathbf{x}^{(p)}]$$

where the term in braces represent the Gauss-Seidel scheme. After extracting $\mathbf{x}^{(p+1)}$ from the above equation, it can be cast in the standard iterative form of equation (3.11) as,

$$\mathbf{x}^{(p+1)} = (\mathbf{D} + \omega\mathbf{L})^{-1}[(\mathbf{1} - \omega)\mathbf{D} - \omega\mathbf{U}]\mathbf{x}^{(p)} + \omega(\mathbf{D} + \omega\mathbf{L})^{-1}\mathbf{b} \quad (3.18)$$

Thus the convergence of the relaxation method depends on the spectral radius of the matrix $\mathbf{T}(\omega) := (\mathbf{D} + \omega\mathbf{L})^{-1}[(\mathbf{1} - \omega)\mathbf{D} - \omega\mathbf{U}]$. Since this matrix is a function of ω we have gained a measure of control over the convergence of the iterative scheme. It has been shown (Young, 1971) that the SOR method is convergent for $0 < \omega < 2$ and that there is an optimum value of ω which results in the maximum rate of convergence. The optimum value of ω is very problem dependent and often difficult to determine precisely. For linear problems, typical values in the range of $\omega \approx 1.7 \sim 1.8$ are used.

3.5.4 Iterative refinement of direct solutions

We have seen that solutions obtained with direct methods are prone to accumulation of round-off errors, while iterative methods have the natural ability to remove such errors. In an attempt to combine the best of both worlds, one might construct an algorithm that takes the error-prone solution from a direct method as an initial guess to an iterative method and thus improve the accuracy of the solution.

Let us illustrate this concept as applied to improving the accuracy of a matrix inverse. Suppose \mathbf{B} is an approximate (error-prone) inverse of a given matrix \mathbf{A} obtained by one of the direct methods outlined in section §3.4. If \mathbf{B}_ϵ is the error in \mathbf{B} then

$$\mathbf{A}(\mathbf{B} + \mathbf{B}_\epsilon) = \mathbf{I} \quad \text{or} \quad \mathbf{A}\mathbf{B}_\epsilon = (\mathbf{I} - \mathbf{A}\mathbf{B})$$

We do not, of course, know \mathbf{B}_ϵ and our task is to attempt to estimate it approximately. Premultiplying above equation by \mathbf{B} , and recognizing $\mathbf{B}\mathbf{A} \approx \mathbf{I}$, we have

$$\mathbf{B}_\epsilon = \mathbf{B}(\mathbf{I} - \mathbf{A}\mathbf{B})$$

Observe carefully that we have used the approximation $\mathbf{B}\mathbf{A} \approx \mathbf{I}$ on the left hand side where products of order unity are involved and *not* on the right hand side where difference between numbers of order unity are

involved. Now we have an estimate of the error B_ϵ which can be added to the approximate result B to obtain,

$$B + B_\epsilon = B + B(I - AB) = B(2I - AB)$$

Hence the iterative sequence should be,

$$B^{(p+1)} = B^{(p)}(2I - AB^{(p)}) \quad p = 0, 1, \dots \quad (3.19)$$

3.6 Gram-Schmidt orthogonalization procedure

Given a set of n linearly independent vectors, $\{\mathbf{x}_i \mid i = 1, \dots, n\}$ that are not necessarily orthonormal, it is possible to produce an orthonormal set of vectors $\{\mathbf{u}_i \mid i = 1, \dots, n\}$.

We begin by normalizing the \mathbf{x}_1 vector using the norm $\|\mathbf{x}_1\| = \sqrt{\mathbf{x}_1 \cdot \mathbf{x}_1}$ and call it \mathbf{u}_1 .

$$\mathbf{u}_1 = \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|}$$

Subsequently we construct other vectors orthogonal to \mathbf{u}_1 and normalize each one. For example we construct \mathbf{u}'_2 by subtracting \mathbf{u}_1 from \mathbf{x}_2 in such a way that \mathbf{u}'_2 contains no components of \mathbf{u}_1 - i.e.,

$$\mathbf{u}'_2 = \mathbf{x}_2 - c_0 \mathbf{u}_1$$

In the above c_0 is to be found in such a way that \mathbf{u}'_2 is orthogonal to \mathbf{u}_1 .

$$\mathbf{u}_1^T \cdot \mathbf{u}'_2 = 0 = \mathbf{u}_1^T \cdot \mathbf{x}_2 - c_0 \quad \text{OR} \quad c_0 = \mathbf{u}_1^T \cdot \mathbf{x}_2$$

Similarly we have,

$$\mathbf{u}'_3 = \mathbf{x}_3 - c_1 \mathbf{u}_1 - c_2 \mathbf{u}_2$$

Requiring orthogonality with respect to both \mathbf{u}_1 and \mathbf{u}_2

$$\mathbf{u}_1^T \cdot \mathbf{u}'_3 = 0 = \mathbf{u}_1^T \cdot \mathbf{x}_3 - c_1 \quad \text{OR} \quad c_1 = \mathbf{u}_1^T \cdot \mathbf{x}_3$$

$$\mathbf{u}_2^T \cdot \mathbf{u}'_3 = 0 = \mathbf{u}_2^T \cdot \mathbf{x}_3 - c_2 \quad \text{OR} \quad c_2 = \mathbf{u}_2^T \cdot \mathbf{x}_3$$

$$\mathbf{u}'_3 = \mathbf{x}_3 - (\mathbf{u}_1^T \cdot \mathbf{x}_3) \mathbf{u}_1 - (\mathbf{u}_2^T \cdot \mathbf{x}_3) \mathbf{u}_2$$

In general we have,

$$\mathbf{u}'_s = \mathbf{x}_s - \sum_{j=1}^{s-1} (\mathbf{u}_j^T \cdot \mathbf{x}_s) \mathbf{u}_j \quad \mathbf{u}_s = \frac{\mathbf{u}'_s}{\|\mathbf{u}'_s\|} \quad s = 2, \dots, n; \quad (3.20)$$

The Gram-Schmidt algorithm is illustrated in figure 3.8. Note that MATLAB has a built-in function $Q = \text{orth}(A)$ which produces an orthonormal set from A . Q spans the same space as A and the number of columns in Q is the rank of A .

```
function u=GS_orth(x)
%Gram-Schmidt procedure applied on input x
% x is a matrix of nxn.
% Each x(:,j) represents a vector
n=length(x);
if rank(x) ~= n
    error('vectors in x are not linearly independent')
end
u(:,1) = x(:,1)/norm(x(:,1));
for s=2:n
    sum = zeros(n,1);
    for j=1:s-1
        sum = sum + (u(:,j)'*x(:,s)) * u(:,j);
    end
    uprime = x(:,s) - sum;
    u(:,s) = uprime/norm(uprime);
end
```

Figure 3.8: Illustration of Gram-Schmidt algorithm

3.7 The eigenvalue problem

A square matrix, A , when operated on certain vectors, called eigenvectors, \mathbf{x} , leaves the vector unchanged except for a scaling factor, λ . This fact can be represented as,

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

The problem of finding such eigenvectors and eigenvalues is addressed by rewriting equation (3.21) as,

$$(A - \lambda I)\mathbf{x} = \mathbf{0}$$

which represents a set of homogeneous equations that admit non-trivial solutions only if

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

i.e., only certain values of λ will make the above determinant zero. This requirement produces an n -th degree polynomial in λ , called the *characteristic polynomial*. Fundamental results from algebra tell us that this polynomial will have exactly n roots, $\{\lambda_j | j = 1 \cdot \cdot \cdot n\}$ and corresponding to each root, λ_j we can determine an eigenvector, \mathbf{x}_j from

$$(A - \lambda_j I)\mathbf{x}_j = \mathbf{0}.$$

Note that if, \mathbf{x}_j satisfies the above equation, then $a\mathbf{x}_j$ will also satisfy the same equation. Viewed alternately, since the $\det(A - \lambda_j I)$ is zero, \mathbf{x}_j can be determined only up to an unknown constant - *i.e.*, only the direction of the eigenvector can be determined, its magnitude being arbitrary.

The MATLAB built-in function `[V,D]=eig(A)` computes all of the eigenvalues and the associated eigenvectors.

Example

Determine the eigenvalues and eigenvectors of the matrix A defined below.

$$A = \begin{bmatrix} 2 & 4 & 3 \\ 4 & 4 & 3 \\ 3 & 3 & 5 \end{bmatrix}$$

A typical MATLAB session follows.

```

» A=[2 4 3;4 4 3;3 3 5] %Define the matrix

A =

     2     4     3
     4     4     3
     3     3     5

%Compute the eigenvalues (in D) and the eigenvectors (in V)
» [V,D]=eig(A)

V =

     0.8197    -0.2674     0.5066
    -0.5587    -0.5685     0.6039
    -0.1265     0.7780     0.6154

D =

    -1.1894         0         0
         0     1.7764         0
         0         0    10.4130

%compute the coefficients of the characteric polynomial (in c)
» c=poly(A)

c =

     1.0000    -11.0000     4.0000    22.0000

%compute the eigenvalues (in r) and compare against the marix D
» r=roots(c)

r =

    10.4130
     1.7764
    -1.1894

```

Note that in the above example, the coefficients of the characteristic polynomial, contained in the vector c , are defined as follows.

$$p_n(\lambda) = c_1\lambda^n + c_2\lambda^{n-1} + \cdots + c_{n+1}$$

3.7.1 Left and right eigenvectors

The eigenvector \mathbf{x} , defined in equation (3.21), is called the right eigenvector, since the matrix product operation is performed from the right. Consider the operation with a vector \mathbf{y}

$$\mathbf{y}'\mathbf{A} = \mathbf{y}'\lambda \quad (3.22)$$

Since equation (3.22) can be written as

$$\mathbf{y}'(\mathbf{A} - \lambda\mathbf{I}) = \mathbf{o}$$

the criterion for admitting nontrivial solutions is still the same, viz.

$$\det(\mathbf{A} - \lambda\mathbf{I}) = \mathbf{o}.$$

Thus the eigenvalues are the same as those defined for the right eigenvector. If \mathbf{A} is symmetric, then taking the transpose of equation (3.22) leads to equation (3.21) pointing out that the distinction between the left and right eigenvectors disappear. However for a nonsymmetric matrix, \mathbf{A} there is a set of left $\{\mathbf{x}_i | i = 1 \cdots n\}$ and right $\{\mathbf{y}_j | j = 1 \cdots n\}$ eigenvectors that are distinct and form a *bi-orthogonal* set as shown below.

3.7.2 Bi-orthogonality

As long as the eigenvalues are distinct, the left and right eigenvectors form a bi-orthogonal set. To demonstrate this, take the left product of equation (3.21) with \mathbf{y}'_j and the right product of equation (3.22) with \mathbf{x}_i . These lead respectively to

$$\mathbf{y}'_j\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{y}'_j\mathbf{x}_i$$

$$\mathbf{y}'_j\mathbf{A}\mathbf{x}_i = \lambda_j\mathbf{y}'_j\mathbf{x}_i$$

Subtracting one from the other, we get,

$$\mathbf{o} = (\lambda_i - \lambda_j)\mathbf{y}'_j\mathbf{x}_i.$$

Since λ_i and λ_j are distinct, we get

$$\mathbf{y}'_j\mathbf{x}_i = \mathbf{o} \quad \forall i \neq j$$

which is the condition for bi-orthogonality.

3.7.3 Power iteration

A simple algorithm to find the largest eigenvalue and its associated eigenvector of a matrix \mathbf{A} is the power iteration scheme. Conceptually the algorithm works as follows. Starting with any arbitrary vector, $\mathbf{z}^{(0)}$, produce a sequence of vectors from

$$\begin{aligned}\mathbf{y}^{(p+1)} &= \mathbf{A}\mathbf{z}^{(p)} \\ \mathbf{y}^{(p+1)} &= \mathbf{k}^{(p+1)}\mathbf{z}^{(p+1)} \quad \mathbf{p} = \mathbf{0}, \mathbf{1}, \dots\end{aligned}$$

The second step above is merely a scaling operation to keep the vector length $\mathbf{z}^{(p)}$ bounded. The vector $\mathbf{z}^{(p)}$ converges to the eigenvector \mathbf{x}_1 corresponding to the largest eigenvalue as $\mathbf{p} \rightarrow \infty$.

The reason why it works can be understood, by beginning with a spectral representation of the arbitrary vector, $\mathbf{z}^{(0)}$ and tracing the effect of each operation, *i.e.*,

$$\begin{aligned}\mathbf{z}^{(0)} &= \sum_i \alpha_i \mathbf{x}_i \\ \mathbf{y}^{(1)} &= \mathbf{A}\mathbf{z}^{(0)} = \sum_i \alpha_i \mathbf{A}\mathbf{x}_i = \sum_i \alpha_i \lambda_i \mathbf{x}_i \\ \mathbf{y}^{(2)} &= \mathbf{A}\mathbf{z}^{(1)} = \mathbf{A} \frac{\mathbf{y}^{(1)}}{\mathbf{k}^{(1)}} = \frac{\mathbf{1}}{\mathbf{k}^{(1)}} \sum_i \alpha_i \lambda_i \mathbf{A}\mathbf{x}_i = \frac{\mathbf{1}}{\mathbf{k}^{(1)}} \sum_i \alpha_i \lambda_i^2 \mathbf{x}_i\end{aligned}$$

Repeating this process we get,

$$\mathbf{y}^{(p+1)} = \mathbf{A}\mathbf{z}^{(p)} = \frac{\mathbf{1}}{\prod_j^p \mathbf{k}^{(j)}} \sum_i \alpha_i \lambda_i^{p+1} \mathbf{x}_i$$

Factoring the largest eigenvalue, λ_1 ,

$$\mathbf{y}^{(p+1)} = \mathbf{A}\mathbf{z}^{(p)} = \frac{\lambda_1^{p+1}}{\prod_j^p \mathbf{k}^{(j)}} \sum_i \alpha_i \left(\frac{\lambda_i}{\lambda_1}\right)^{p+1} \mathbf{x}_i$$

Since $(\lambda_i/\lambda_1) < \mathbf{1}$ for $i > \mathbf{1}$ only the first term in the summation survives as $\mathbf{p} \rightarrow \infty$. Thus $\mathbf{y}^{(p)} \rightarrow \mathbf{x}_1$, the eigenvector corresponding to the largest eigenvalue.

Several other features also become evident from the above analysis.

- The convergence rate will be faster, if the largest eigenvalue is well separated from the remaining eigenvalues.
- If the largest eigenvalue occurs with a multiplicity of two, then the above sequence will not converge, or rather converge to a subspace spanned by the eigenvectors corresponding to the eigenvalues that occur with multiplicity.

```

function [Rayleigh_Q,V] = Power(T,MaxIt)
%Power iteration to find the largest e.value of T

if nargin < 2,
    MaxIt = 100;
end

n=length(T);           %Find the size of T
z_old=rand(n,1);       %Generate random vector, z
z_new=z_old/norm(z_old); %Scale it
Check_sign=1; count = 0; %Initialize
while (norm(z_old-z_new) > 1.e-10 & count < MaxIt),
    count = count + 1;
    z_old=z_new*Check_sign;
    z_new=T*z_new;
    z_new=z_new/norm(z_new);
    Check_sign=sign((z_new'*(T*z_new))/(z_new'*z_new));
end
if (count >= MaxIt)
    error('Power iteration failed to converge')
end
%Compute the Rayliegh quotient
V=z_new;
Rayleigh_Q=(z_new'*(T*z_new))/(z_new'*z_new);

```

Figure 3.9: MATLAB implementation of power iteration

- If the initial guess vector does not contain any component in the direction of \mathbf{x}_1 , then the sequence will converge to the next largest eigenvalue. This can be achieved by making the guess vector, $\mathbf{z}^{(0)}$, orthogonal to the known eigenvector, \mathbf{x}_1 .

3.7.4 Inverse iteration

To find the smallest eigenvalue of a given a matrix \mathbf{A} , the power iteration could still be applied, but on the inverse of matrix \mathbf{A} . Consider the original eigenvalue problem

$$\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i.$$

Premultiply by A^{-1} to get,

$$A^{-1}A\mathbf{x}_i = \lambda_i A^{-1}\mathbf{x}_i$$

which can be rewritten as,

$$\frac{1}{\lambda_i}\mathbf{x}_i = A^{-1}\mathbf{x}_i.$$

Hence it is clear that the eigenvalues, μ_i of A^{-1} viz.

$$\mu_i\mathbf{x}_i = A^{-1}\mathbf{x}_i$$

are related to eigenvalues, λ_i of A through,

$$\lambda_i = \frac{1}{\mu_i}.$$

Although the illustration below uses the inverse of the matrix, in reality there is no need to find A^{-1} since power iteration only requires computation of a new vector \mathbf{y} from a given vector \mathbf{z} using,

$$\mathbf{y}^{(p+1)} = A^{-1}\mathbf{z}^{(p)}.$$

This can be done most effectively by solving the linear system

$$A\mathbf{y}^{(p+1)} = \mathbf{z}^{(p)}$$

using LU factorization, which needs to be done only once.

3.7.5 Shift-Inverse iteration

To find the eigenvalue closest to a selected point, σ , the power iteration could be applied to the matrix $(A - \sigma I)^{-1}$, which is equivalent to solving the eigenvalue problem

$$(A - \sigma I)^{-1}\mathbf{x}_i = \mu_i\mathbf{x}_i.$$

This can be rewritten as,

$$\frac{1}{\mu_i}\mathbf{x}_i = (A - \sigma I)\mathbf{x}_i = A\mathbf{x}_i - \sigma\mathbf{x}_i = (\lambda_i - \sigma)\mathbf{x}_i.$$

Hence it is clear that the eigenvalues, μ_i of $(A - \sigma I)^{-1}$ are related to eigenvalues, λ_i of A through,

$$\lambda_i = \frac{1}{\mu_i} + \sigma.$$

As in the previous section, there is no need to find the inverse since power iteration only requires computation of a new vector \mathbf{y} from a given vector \mathbf{z} using,

$$\mathbf{y}^{(p+1)} = (\mathbf{A} - \sigma\mathbf{I})^{-1}\mathbf{z}^{(p)}.$$

This can be done most effectively by solving the linear system

$$(\mathbf{A} - \sigma\mathbf{I})\mathbf{y}^{(p+1)} = \mathbf{z}^{(p)}$$

using LU factorization of $(\mathbf{A} - \sigma\mathbf{I})$, which needs to be done only once.

Example

Determine the largest and smallest eigenvalues of the matrix T defined below using the power iteration.

$$T = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -2 & 2 & -1 & 0 \\ 0 & -2 & 2 & -1 \\ 0 & 0 & -2 & 2 \end{bmatrix}$$

A typical MATLAB session follows.

```
» T=trid(2,-1,-2,4) %generate the matrix
```

```
T =
```

```

     2    -1     0     0
    -2     2    -1     0
     0    -2     2    -1
     0     0    -2     2
```

```
»[d,v]=ch3_poweri(T) %apply power iteration
```

```
d =
```

```

    4.2882
```

```
v =
```

```

   -0.2012
    0.4603
   -0.6510
```

```
0.5690

»max(eig(T)) %Find the largest e-value from ‘eig’ function

ans =

    4.2882

%Find the smallest e. value
»mu=ch3_poweri(inv(T));
»lambda=1/mu

lambda =

   -0.2882
»min(eig(T))

ans =

   -0.2882

»%Find the e.value closest to 2.5
»sigma=2.5;
»mu=ch3_poweri(inv(T-sigma*eye(4)))

mu =

    2.6736

»lambda=sigma+1/mu

lambda =

    2.8740
%verfiy using ‘eig’
»eig(T)

ans =

   -0.2882
    4.2882
    2.8740
```

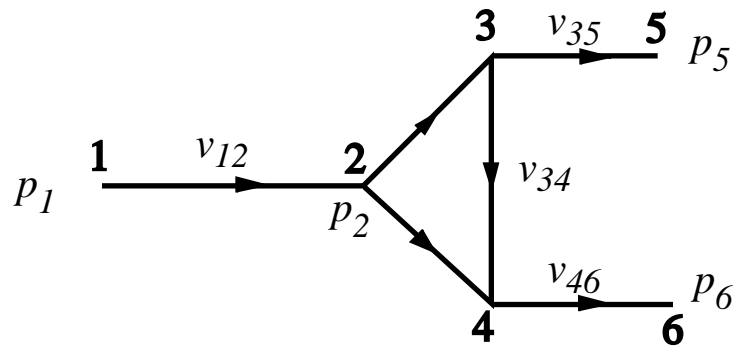


Figure 3.10: Laminar flow in a pipe network

1.1260

3.7.6 Wei-Prater analysis of a reaction system

3.8 Singular value decomposition

3.9 Genaralized inverse

3.10 Software tools

3.10.1 Lapack, Eispack library

3.10.2 MATLAB

3.10.3 Mathematica

3.11 Exercise problems

3.11.1 Laminar flow through a pipeline network

Consider laminar flow through the network shown in figure 3.10. The governing equations are the pressure drop equations for each pipe element $i - j$ and the mass balance equation at each node.

The pressure drop between nodes i and j is given by,

$$p_i - p_j = \alpha_{ij} v_{ij} \quad \text{where} \quad \alpha_{ij} = \frac{32\mu l_{ij}}{d_{ij}^2} \quad (3.23)$$

The mass balance at node 2 is given, for example by,

$$d_{12}^2 v_{12} = d_{23}^2 v_{23} + d_{24}^2 v_{24} \quad (3.24)$$

Similar equations apply at nodes 3 and 4. Let the unknown vector be

$$\underline{x} = [p_2 \ p_3 \ p_4 \ v_{12} \ v_{23} \ v_{24} \ v_{34} \ v_{35} \ v_{46}]$$

There will be six momentum balance equations, one for each pipe element, and three mass balance (for incompressible fluids volume balance) equations, one at each node. Arrange them as a system of nine equations in nine unknowns and solve the resulting set of equations. Take the viscosity of the fluid, $\mu = 0.1 \text{ Pa} \cdot \text{s}$. The dimensions of the pipes are given below.

Table 1

Element no	12	23	24	34	35	46
d_{ij} (m)	0.1	0.08	0.08	0.10	0.09	0.09
l_{ij} (m)	1000	800	800	900	1000	1000

- a) Use MATLAB to solve this problem for the specified pressures of $p_1 = 300 \text{ kPa}$ and $p_5 = p_6 = 100 \text{ kPa}$. You need to assemble the system of equations in the form $A \mathbf{x} = \mathbf{b}$. Report flops. When reporting flops, report only for that particular operation - *i.e.*, initialize the counter using `flops(0)` before every operation.
- Compute the determinant of A . Report Flops.
 - Compute the LU factor of A using built-in function `lu`. Report flops. What is the structure of L ? Explain. The function `LU` provided in the lecture notes will fail on this matrix. Why?
 - Compute the solution using `inv(A)*b`. Report flops.
 - Compute the rank of A . Report Flops.
 - Since A is sparse (*i.e.*, mostly zeros) we can avoid unnecessary operations, by using sparse matrix solvers. MATLAB Ver4.0 (not 3.5) provides such a facility. Sparse matrices are stored using triplets (i, j, s) where (i, j) identifies the non-zero entry in the matrix and s its corresponding value. The MATLAB function `find(A)` examines A and returns the triplets. Use,

» **[ii,jj,s]=find(A)**

Then construct the sparse matrix and store it in **S** using

» **S=sparse(ii,jj,s)**

Then to solve using the sparse solver and keep flops count, use

» **flops(0); x = S\b; flops**

Compare flops for solution by full and sparse matrix solution. To graphically view the structure of the sparse matrix, use

» **spy(S)**

Remember that you should have started MATLAB under X-windows for any graphics display of results!

- Compute the determinant of the sparse matrix, **S** (should be the same as the full matrix!). Report and compare flops.

- b) Find out the new velocity and pressure distributions when **p₆** is changed to **150kPa**.
- c) Suppose the forcing (column) vector in part (a) is **b₁** and that in part (b) is **b₂**, report and explain the difference in flops for the following two ways of obtaining the two solutions using the sparse matrix. Note that in the first case both solutions are obtained simultaneously.

» **b = [b₁, b₂];flops(0);x = S\b,flops**
 » **flops(0);x1 = S\b1, x2 = S\b2,flops**

Repeat the above experiment with the full matrix, **A** and report flops.

- d) Comment on how you would adopt the above problem formulation if a valve on line **34** is shut so that there is no flow in that line **34**.

3.11.2 Gram-Schmidt procedure

Consider the vectors

$$\mathbf{x}_1 = [1234] \quad (3.25)$$

$$\mathbf{x}_2 = [1212] \quad (3.26)$$

$$\mathbf{x}_3 = [2131] \quad (3.27)$$

$$\mathbf{x}_4 = [2154] \quad (3.28)$$

- Check if they form a linearly independent set.
- Construct an orthonormal set using the Gram-Schmidt procedure

3.11.3 Power iteration

Consider the matrix, A given below.

$A =$

$$\begin{array}{cccc} 2.0 + 4.0i & -1.0 - 2.0i & 0 & 0 \\ -2.0 - 4.0i & 2.0 + 4.0i & -1.0 - 2.0i & 0 \\ 0 & -2.0 - 4.0i & 2.0 + 4.0i & -1.0 - 2.0i \\ 0 & 0 & -2.0 - 4.0i & 2.0 + 4.0i \end{array}$$

Use the power iteration algorithm given in Fig. 3.9 to solve the following problems.

- Find the largest eigenvalue. [Ans: $4.2882 + 8.5765i$]
- Find the smallest eigenvalue. [Ans: $-0.2882 - 0.5765i$]
- Find the eigenvalue closest to $\sigma = (1 + 2i)$. [Ans: $1.1260 + 2.2519i$]