

Solution of Linear Systems

Peng Yu

Tel: 0755 8801 8911

Email: yup6@sustech.edu.cn

Motivation

Intersection of three planes

$$5x + y + z = 5$$

$$x + 4y + z = 4$$

$$x + y + 3z = 3.$$

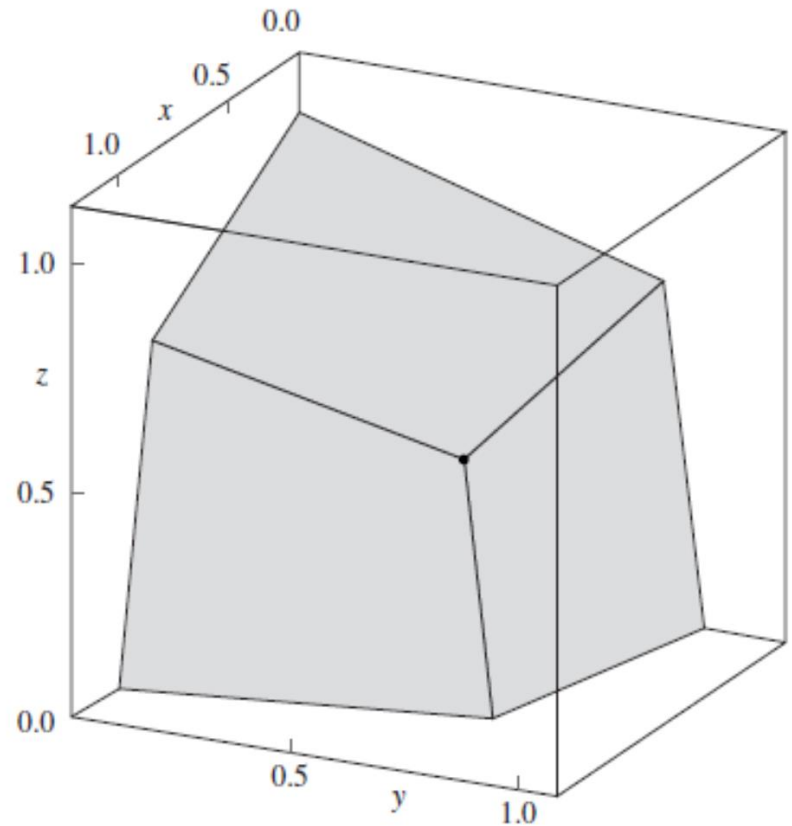


Figure 3.1 The intersection of three planes.

$$x = 0.76, \quad y = 0.68, \quad \text{and} \quad z = 0.52.$$

Solution of Linear Systems

- Introduction
- Upper-Triangular Linear Systems
- Gaussian Elimination and Pivoting
- Triangular Factorization
- Iterative Methods for Linear Systems

Introduction

Vectors and Matrices

Notations

- \mathbb{R} : real numbers; \mathbb{C} : complex numbers; \mathbb{Z} : integers;
- \mathbb{R}^n : the set of n – dimensional vectors;
- $\mathbb{R}^{m \times n}$: the set of all $m \times n$ matrices;
- $x \in \mathcal{X}$: x is a member in the set \mathcal{X} ;
- \exists : there exists;
- The matrix (vector) is denoted as a bold letters (eg. **A**)
- Define $[n] = 1, \dots, n$ and $[j:n] = j, \dots, n$.

Notations

- Vector: $\mathbf{a} = (a_1, a_2, \dots, a_n) \in \mathbb{R}^n$ and a_i is the i – th entry of \mathbf{a} .

- Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n) = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

- Matrix – vector multiplication. $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$,
 $\mathbf{Ax} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \cdots + x_n \mathbf{a}_n \in \mathbb{R}^m$.

- Matrix – matrix multiplication. $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$.
 $\mathbf{AB} = \mathbf{A}(\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_p) = (\mathbf{Ab}_1, \mathbf{Ab}_2, \dots, \mathbf{Ab}_p) \in \mathbb{R}^{m \times p}$.

Special Matrices

- Symmetric matrices \mathbb{S}^n :

$$\mathbf{A} \in \mathbb{S}^n \Leftrightarrow \mathbf{A} = \mathbf{A}^T \Leftrightarrow a_{ij} = a_{ji}, \forall i, j = [n].$$

- Lower triangular matrices \mathcal{L}

$$\mathbf{A} \in \mathcal{L} \Leftrightarrow a_{ij} = 0 \text{ if } i < j.$$

- Upper triangular matrices \mathcal{U}

$$\mathbf{A} \in \mathcal{U} \Leftrightarrow a_{ij} = 0 \text{ if } i > j \Leftrightarrow \mathbf{A}^T \in \mathcal{L}$$

- Positive semi – definite (**definite**) matrices $\mathbb{S}_+^n (\mathbb{S}_{++}^n)$.

$$\mathbf{A} \in \mathbb{S}_+^n (\mathbb{S}_{++}^n) \Leftrightarrow \mathbf{x}^T \mathbf{A} \mathbf{x} \geq (>) 0, \forall \mathbf{x} \neq 0$$

- Orthogonal matrices \mathcal{O}^n : $\mathbf{A} \in \mathcal{O}^n \Leftrightarrow \mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$.

Vector Norms

Let $\mathbf{x} \in \mathbb{R}^n$. There are several vector norms:

1. ℓ_2 -norm

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}} = (|x_1|^2 + |x_2|^2 + \cdots + |x_n|^2)^{1/2}$$

2. ℓ_1 -norm

$$\|\mathbf{x}\|_1 = (|x_1| + |x_2| + \cdots + |x_n|)$$

3. ℓ_∞ -norm

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

Proposition (Norm equivalence)

For all $x \in \mathbb{R}^n$ have

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq n\|\mathbf{x}\|_\infty$$

Basic Concepts for Vectors

- Inner product of \mathbf{x} and \mathbf{y} :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos(\theta)$$

where θ is the angle between \mathbf{x} and \mathbf{y} .

- \mathbf{x} is orthogonal to \mathbf{y} :

$$\mathbf{x}^T \mathbf{y} = 0$$

- For all $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ have:

$$\mathbf{x}^T \mathbf{A} \mathbf{y} = \mathbf{y}^T \mathbf{A}^T \mathbf{x}.$$

Matrix Norms

- Trace: $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\text{tr}(\mathbf{A}) = a_{11} + a_{22} + \cdots + a_{nn}.$$

- Frobenius norm: $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2}.$$

- Matrix norm: $\mathbf{A} \in \mathbb{R}^{m \times n}$,

$$\|\mathbf{A}\|_2 = \max_{\mathbf{x} \neq 0} \|\mathbf{A}\mathbf{x}\| / \|\mathbf{x}\|$$

Matrix Properties

■ $\text{diag}(\mathbf{A}) = (a_{11}, a_{12}, \dots, a_{nn})$ and $\text{Diag}(\mathbf{a})$ is

$$\text{Diag}(\mathbf{a}) = \begin{bmatrix} a_1 & & & \\ & a_2 & & \\ & & \ddots & \\ & & & a_n \end{bmatrix}$$

■ If $\mathbf{A} \in \mathbb{S}^n$, there exist $\mathbf{P} \in \mathcal{O}^n$ and $\mathbf{D} = \text{Diag}(\mathbf{d})$ such that $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^T$.

■ If $\mathbf{A} \in \mathbb{S}^n$, then $\|\mathbf{A}\|_2 = \lambda_{\max}(\mathbf{A})$, where $\lambda_{\max}(\mathbf{A})$ denotes the maximal eigenvalue of \mathbf{A} .

■ If $\mathbf{A} \in \mathbb{S}^n$, then $\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^T \mathbf{A}) = \sum_i d_i^2$, where d_i s are eigenvalues of \mathbf{A} .

Matrix Multiplication

Definition 3.1. If $A = [a_{ik}]_{M \times N}$ and $B = [b_{kj}]_{N \times P}$ are two matrices with the property that A has as many columns as B has rows, then the *matrix product* AB is defined to be matrix C of dimension $M \times P$:

$$(6) \quad AB = C = [c_{ij}]_{M \times P},$$

where the element c_{ij} of C is given by the dot product of the i th row of A and the j th column of B :

$$(7) \quad c_{ij} = \sum_{k=1}^N a_{ik}b_{kj} = a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{iN}b_{Nj}$$

Theorem 3.3 (Matrix Multiplication)

$$(AB)C = A(BC)$$

associativity of matrix multiplication

$$IA = AI = A$$

identity matrix

$$A(B + C) = AB + AC$$

left distributive property

$$(A + B)C = AC + BC$$

right distributive property

$$c(AB) = (cA)B = A(cB)$$

scalar associative property

Inverse of a Nonsingular Matrix

- An $N \times N$ matrix \mathbf{A} is called nonsingular or invertible if there exists an $N \times N$ matrix \mathbf{B} such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I} \quad (1)$$

- If no such matrix \mathbf{B} can be found, \mathbf{A} is said to be singular.
- When \mathbf{B} can be found, we say that \mathbf{B} is the inverse of \mathbf{A} and usually write $\mathbf{B} = \mathbf{A}^{-1}$ and use the familiar relation:

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} \quad \text{if } \mathbf{A} \text{ is nonsingular}$$

- At most one matrix \mathbf{B} can be found that satisfies relation (1). Suppose \mathbf{C} is also an inverse of \mathbf{A} , then

$$\mathbf{C} = \mathbf{IC} = (\mathbf{BA})\mathbf{C} = \mathbf{B}(\mathbf{AC}) = \mathbf{BI} = \mathbf{B}$$

Determinant of Matrix

$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & \vdots & & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NN} \end{vmatrix}.$$

If $\mathbf{A} = [a_{ij}]$ is a 1×1 matrix, we define $\det(\mathbf{A}) = a_{11}$. If $\mathbf{A} = [a_{ij}]_{N \times N}$, where $N \geq 2$, then let M_{ij} be the determinant of the $(N - 1) \times (N - 1)$ submatrix of \mathbf{A} obtained by deleting the i th row and j th column of \mathbf{A} . The determinant M_{ij} is said to be the **minor** of a_{ij} . The **cofactor** \mathbf{A}_{ij} of a_{ij} is defined as $\mathbf{A}_{ij} = (-1)^{i+j} \mathbf{M}_{ij}$. Then the determinant of an $N \times N$ matrix \mathbf{A} is given by

$$(19) \quad \det(\mathbf{A}) = \sum_{j=1}^N a_{ij} \mathbf{A}_{ij} \quad (i\text{th row expansion})$$

$$(20) \quad \det(\mathbf{A}) = \sum_{i=1}^N a_{ij} \mathbf{A}_{ij} \quad (j\text{th column expansion}).$$

Determinant of Matrix

For a 2×2 matrix

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\det A = a_{11} a_{22} - a_{12} a_{21}$$

Example 3.8. Use formula (19) with $i = 1$ and formula (20) with $j = 2$ to calculate the determinant of the matrix

$$A = \begin{bmatrix} 2 & 3 & 8 \\ -4 & 5 & -1 \\ 7 & -6 & 9 \end{bmatrix}.$$

Using formula (19) with $i = 1$, we obtain

$$\begin{aligned} \det(A) &= (2) \begin{vmatrix} 5 & -1 \\ -6 & 9 \end{vmatrix} - (3) \begin{vmatrix} -4 & -1 \\ 7 & 9 \end{vmatrix} + (8) \begin{vmatrix} -4 & 5 \\ 7 & -6 \end{vmatrix} \\ &= (2)(45 - 6) - (3)(-36 + 7) + (8)(24 - 35) \\ &= 77. \end{aligned}$$

Matrix property

Theorem 3.4. Assume that A is an $N \times N$ matrix. The following statements are equivalent.

(21) Given any $N \times 1$ matrix B , the linear system $AX = B$ has a unique solution.

(22) The matrix A is nonsingular (i.e., A^{-1} exists).

(23) The system of equations $AX = 0$ has the unique solution $X = 0$.

(24) $\det(A) \neq 0$.

Theorems 3.3 and 3.4 help relate matrix algebra to ordinary algebra. If statement (21) is true, then statement (22) together with properties (12) and (13) give the following line of reasoning:

(25) $AX = B$ implies $A^{-1}AX = A^{-1}B$, which implies $X = A^{-1}B$.

Plane Rotations

- Suppose \mathbf{U} is a positional vector $\mathbf{U} = (x, y, z)$, $\mathbf{V} = \mathbf{R}\mathbf{U}$ represents a linear transformation, spatially, rotation of the vector.

$$\mathbf{R}_x(\alpha) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\alpha) & -\sin(\alpha) \\ 0 & \sin(\alpha) & \cos(\alpha) \end{bmatrix},$$

$$\mathbf{R}_y(\beta) = \begin{bmatrix} \cos(\beta) & 0 & \sin(\beta) \\ 0 & 1 & 0 \\ -\sin(\beta) & 0 & \cos(\beta) \end{bmatrix},$$

$$\mathbf{R}_z(\gamma) = \begin{bmatrix} \cos(\gamma) & -\sin(\gamma) & 0 \\ \sin(\gamma) & \cos(\gamma) & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

Plane Rotations: An Example

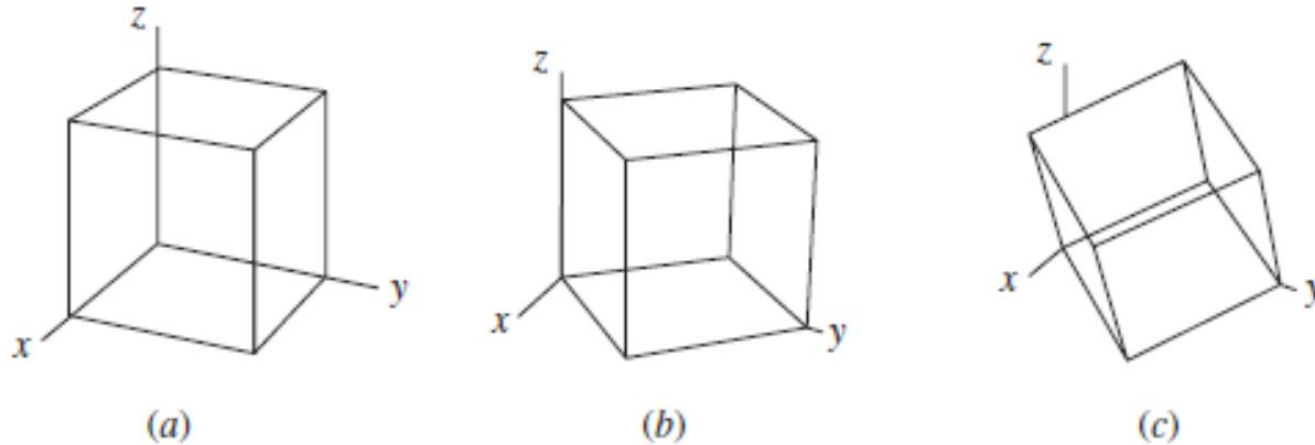


Figure 3.2 (a) The original starting cube. (b) $\mathbf{V} = \mathbf{R}_z(\pi/4)\mathbf{U}$. Rotation about the z -axis. (c) $\mathbf{W} = \mathbf{R}_y(\pi/6)\mathbf{V}$. Rotation about the y -axis.

$$\begin{aligned} \mathbf{V} &= \mathbf{R}_z\left(\frac{\pi}{4}\right)\mathbf{U} = \begin{bmatrix} \cos(\pi/4) & -\sin(\pi/4) & 0 \\ \sin(\pi/4) & \cos(\pi/4) & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} \\ &= \begin{bmatrix} 0.707107 & -0.707107 & 0.000000 \\ 0.707107 & 0.707107 & 0.000000 \\ 0.000000 & 0.000000 & 1.000000 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix}. \end{aligned} \quad \rightarrow \quad \begin{aligned} \mathbf{W} &= \mathbf{R}_y\left(\frac{\pi}{6}\right)\mathbf{V} = \begin{bmatrix} \cos(\pi/6) & 0 & \sin(\pi/6) \\ 0 & 1 & 0 \\ -\sin(\pi/6) & 0 & \cos(\pi/6) \end{bmatrix} \mathbf{V} \\ &= \begin{bmatrix} 0.866025 & 0.000000 & 0.500000 \\ 0.000000 & 1.000000 & 0.000000 \\ -0.500000 & 0.000000 & 0.866025 \end{bmatrix} \mathbf{V}. \end{aligned}$$

The composition of the two rotation is

$$\mathbf{W} = \mathbf{R}_y\left(\frac{\pi}{6}\right)\mathbf{R}_z\left(\frac{\pi}{4}\right)\mathbf{U} = \begin{bmatrix} 0.612372 & -0.612372 & 0.500000 \\ 0.707107 & 0.707107 & 0.000000 \\ -0.353553 & 0.353553 & 0.866025 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

Upper-Triangular Linear Systems

Definitions

- Upper triangular matrix
- If \mathbf{A} is an upper-triangular matrix, then $\mathbf{AX} = \mathbf{B}$ is said to be an upper triangular system of linear equations and has the form

$$a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \cdots + a_{1N-1}x_{N-1} + a_{1N}x_N = b_1$$

$$a_{22}x_2 + a_{23}x_3 + \cdots + a_{2N-1}x_{N-1} + a_{2N}x_N = b_2$$

$$a_{33}x_3 + \cdots + a_{3N-1}x_{N-1} + a_{3N}x_N = b_3$$

$$\vdots \qquad \qquad \vdots$$

$$a_{N-1N-1}x_{N-1} + a_{N-1N}x_N = b_{N-1}$$

$$a_{NN}x_N = b_N.$$

Back Substitution

- Suppose that $\mathbf{AX} = \mathbf{B}$ is an upper-triangular system. If $a_{kk} \neq 0$ for $k = 1, 2, \dots, N$, then there exists a unique solution .

Constructive Proof. The solution is easy to find. The last equation involves only x_N , so we solve it first:

$$(3) \quad x_N = \frac{b_N}{a_{NN}}.$$

Now x_N is known and it can be used in the next-to-last equation:

$$(4) \quad x_{N-1} = \frac{b_{N-1} - a_{N-1N}x_N}{a_{N-1N-1}}.$$

Now x_N and x_{N-1} are used to find x_{N-2} :

$$(5) \quad x_{N-2} = \frac{b_{N-2} - a_{N-2N-1}x_{N-1} - a_{N-2N}x_N}{a_{N-2N-2}}.$$

Once the values $x_N, x_{N-1}, \dots, x_{k+1}$ are known, the general step is

$$(6) \quad x_k = \frac{b_k - \sum_{j=k+1}^N a_{kj}x_j}{a_{kk}} \quad \text{for } k = N-1, N-2, \dots, 1.$$

An example

- The condition that $a_{kk} \neq 0$ is essential. If this requirement is not fulfilled, either no solution exists or infinitely many solutions exist.
- The system has a unique solution if and only if $\det(A) \neq 0$.
Here $\det(A) = a_{11}a_{22}\dots a_{NN}$

$$\begin{aligned}4x_1 - x_2 + 2x_3 + 3x_4 &= 20 \\0x_2 + 7x_3 - 4x_4 &= -7 \\6x_3 + 5x_4 &= 4 \\3x_4 &= 6.\end{aligned}$$



$$\begin{aligned}7x_3 - 8 &= -7 \\6x_3 + 10 &= 4.\end{aligned}$$

No solution

$$\begin{aligned}4x_1 - x_2 + 2x_3 + 3x_4 &= 20 \\0x_2 + 7x_3 - 0x_4 &= -7 \\6x_3 + 5x_4 &= 4 \\3x_4 &= 6.\end{aligned}$$



$$x_2 = 4x_1 - 16,$$

Infinitely many solutions

Back Substitution - Program

Program 3.1 (Back substitution). To solve the upper-triangular system $AX = B$ by the method of back substitution. Proceed with the method only if all the diagonal elements are nonzero. First compute $x_N = b_N/a_{NN}$ and then use the rule

$$x_k = \frac{b_k - \sum_{j=k+1}^N a_{kj}x_j}{a_{kk}} \quad \text{for } k = N-1, N-2, \dots, 1.$$

```
function X = backsub(A, B)
%Input   - A is an n × n upper-triangular nonsingular matrix
%         - B is an n × 1 matrix
%Output  - X is the solution to the linear system AX = B
%Find the dimension of B and initialize X
n = length(B);
X = zeros(n,1);

X(n) = B(n)/A(n, n);
for k = n - 1:-1:1
    X(k) = (B(k) - A(k, k + 1:n) * X(k + 1:n))/A(k, k);
end
```

Gaussian Elimination and Pivoting

Gaussian Elimination

- The goal is to construct an equivalent upper-triangular system $UX = Y$ that can be solved by the method in previous section.

Theorem 3.7 (Elementary Transformations). The following operations applied to a linear system yield an equivalent system:

- (1) **Interchanges:** The order of two equations can be changed.
- (2) **Scaling:** Multiplying an equation by a nonzero constant.
- (3) **Replacement:** An equation can be replaced by the sum of itself and a nonzero multiple of any other equation.

An example

- Find the parabola $y = A + Bx + Cx^2$ that passes through the three points $(1, 1)$, $(2, -1)$, and $(3, 1)$.

$$(1) \quad A + B + C = 1 \quad \text{at } (1,1)$$

$$(2) \quad A + 2B + 4C = -1 \quad \text{at } (2, -1)$$

$$(3) \quad A + 3B + 9C = 1 \quad \text{at } (3,1)$$



$$A + B + C = 1 \quad (a)$$

$$(2) - (1): \quad B + 3C = -2 \quad (b)$$

$$(3) - (1): \quad 2B + 8C = 0 \quad (c)$$



$$A + B + C = 1$$

$$B + 3C = -2$$

$$2C = 4 \quad (c) - (b) \times 2$$

Gaussian Elimination - Operations

- It is efficient to store all the coefficients of the linear system $AX = B$ in an array of dimension $N \times (N + 1)$.

$$(7) \quad [A|B] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1N} & b_1 \\ a_{21} & a_{22} & & a_{2N} & b_2 \\ \vdots & \vdots & & \vdots & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NN} & b_N \end{array} \right] \quad \begin{array}{l} \text{增广矩阵} \\ \text{Augmented matrix} \end{array}$$

Theorem 3.8 (Elementary Row Operations). The following operations applied to the augmented matrix (7) yield an equivalent linear system.

Interchanges: The order of two rows can be changed.

Scaling: Multiplying a row by a nonzero constant

Replacement: The row can be replaced by the sum of that row and a nonzero multiple of any other row;

$$\text{row}_r = \text{row}_r - m_{rp} \times \text{row}_p.$$

Gaussian Elimination - Operations

Definition 3.3. The number a_{rr} in the coefficient matrix A that is used to eliminate a_{kr} , where $k = r + 1, r + 2, \dots, N$, is called the r th *pivotal element* (主元), and the r th row is called the *pivot row*.

Example Express the following system in augment matrix form and find an equivalent upper-triangular system and the solution.

$$\begin{aligned}x_1 + 2x_2 + x_3 + 4x_4 &= 13 \\2x_1 + 0x_2 + 4x_3 + 3x_4 &= 28 \\4x_1 + 2x_2 + 2x_3 + x_4 &= 20 \\-3x_1 + x_2 + 3x_3 + 2x_4 &= 6.\end{aligned}$$

The augment matrix is

$$\begin{aligned}\text{pivot} &\rightarrow \\m_{21} &= 2 \\m_{32} &= 4 \\m_{41} &= -3\end{aligned} \left[\begin{array}{cccc|c} \underline{1} & 2 & 1 & 4 & 13 \\ 2 & 0 & 4 & 3 & 28 \\ 4 & 2 & 2 & 1 & 20 \\ -3 & 1 & 3 & 2 & 6 \end{array} \right].$$

Gaussian Elimination - Operations

- The first row is used to eliminate elements in the first column below the diagonal. We refer to the first row as the *pivotal row* and the element $a_{11} = 1$ is called the *pivotal element*. The values m_{k1} are the multiples of row 1 that are to be subtracted from row k for $k = 2, 3, 4$.

$$\begin{array}{l} \text{pivot} \rightarrow \\ m_{32} = 1.5 \\ m_{42} = -1.75 \end{array} \left[\begin{array}{cccc|c} 1 & 2 & 1 & 4 & 13 \\ 0 & \underline{-4} & 2 & -5 & 2 \\ 0 & -6 & -2 & -15 & -32 \\ 0 & 7 & 6 & 14 & 45 \end{array} \right].$$



$$\begin{array}{l} \text{pivot} \rightarrow \\ m_{43} = -1.9 \end{array} \left[\begin{array}{cccc|c} 1 & 2 & 1 & 4 & 13 \\ 0 & -4 & 2 & -5 & 2 \\ 0 & 0 & \underline{-5} & -7.5 & -35 \\ 0 & 0 & 9.5 & 5.25 & 48.5 \end{array} \right].$$



$$\left[\begin{array}{cccc|c} 1 & 2 & 1 & 4 & 13 \\ 0 & -4 & 2 & -5 & 2 \\ 0 & 0 & -5 & -7.5 & -35 \\ 0 & 0 & 0 & -9 & -18 \end{array} \right].$$

Gaussian Elimination - Operations

$$AX = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2N}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} & \cdots & a_{3N}^{(1)} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{N1}^{(1)} & a_{N2}^{(1)} & a_{N3}^{(1)} & \cdots & a_{NN}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} a_{1N+1}^{(1)} \\ a_{2N+1}^{(1)} \\ a_{3N+1}^{(1)} \\ \vdots \\ a_{NN+1}^{(1)} \end{bmatrix} = B$$



$$UX = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \cdots & a_{3N}^{(3)} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & a_{NN}^{(N)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} a_{1N+1}^{(1)} \\ a_{2N+1}^{(2)} \\ a_{3N+1}^{(3)} \\ \vdots \\ a_{NN+1}^{(N)} \end{bmatrix} = Y$$

Gaussian Elimination - Operations

Step 1. Store the coefficients in the augmented matrix. The superscript on $a_{rc}^{(1)}$ means that this is the first time that a number is stored in location (r, c) :

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2N}^{(1)} & a_{2\ N+1}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} & \cdots & a_{3N}^{(1)} & a_{3\ N+1}^{(1)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_{N1}^{(1)} & a_{N2}^{(1)} & a_{N3}^{(1)} & \cdots & a_{NN}^{(1)} & a_{N\ N+1}^{(1)} \end{array} \right].$$

Step 2. If necessary, switch rows so that $a_{11}^{(1)} \neq 0$; then eliminate x_1 in rows 2 through N . In this process, m_{r1} is the multiple of row 1 that is subtracted from row r .

for $r = 2:N$

$$m_{r1} = a_{r1}^{(1)} / a_{11}^{(1)};$$

$$a_{r1}^{(2)} = 0;$$

for $c = 2:N + 1$

$$a_{rc}^{(2)} = a_{rc}^{(1)} - m_{r1} * a_{1c}^{(1)};$$

end

end

Gaussian Elimination - Operations

The new elements are written $a_{rc}^{(2)}$ to indicate that this is the second time that a number has been stored in the matrix at location (r, c) . The result after step 2 is

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} & a_{2\ N+1}^{(2)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} & \cdots & a_{3N}^{(2)} & a_{3\ N+1}^{(2)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & a_{N2}^{(2)} & a_{N3}^{(2)} & \cdots & a_{NN}^{(2)} & a_{N\ N+1}^{(2)} \end{array} \right].$$

Gaussian Elimination - Operations

Step $p + 1$. This is the general step. If necessary, switch row p with some row beneath it so that $a_{pp}^{(p)} \neq 0$; then eliminate x_p in rows $p + 1$ through N . Here m_{rp} is the multiple of row p that is subtracted from row r .

```

for  $r = p + 1:N$ 
     $m_{rp} = a_{rp}^{(p)} / a_{pp}^{(p)}$ ;
     $a_{rp}^{(p+1)} = 0$ ;
    for  $c = p + 1:N + 1$ 
         $a_{rp}^{(p+1)} = a_{rp}^{(p)} - m_{rp} * a_{pc}^{(p)}$ ;
    end
end
end

```



$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} & a_{2\ N+1}^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \cdots & a_{3N}^{(3)} & a_{3\ N+1}^{(3)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{NN}^{(N)} & a_{N\ N+1}^{(N)} \end{array} \right].$$

Pivoting to avoid $a_{pp}^{(p)}=0$

- If $a_{pp}^{(p)} = 0$, row p cannot be used to eliminate the elements in column p below the main diagonal.
- It is necessary to find row k , where $a_{kp}^{(p)} \neq 0$ and $k > p$, and then interchange row p and row k so that a nonzero pivot element is obtained.
- ***Pivoting strategy:***

If $a_{pp}^{(p)} \neq 0$, do not switch rows.

If $a_{pp}^{(p)} = 0$, locate the first row below p in which $a_{kp}^{(p)} \neq 0$ and switch rows k and p .

Pivoting to reduce error

- Because the computer uses fixed-precision arithmetic, it is possible that a small error will be introduced each time that an arithmetic operation is performed.

Example The values $x_1 = x_2 = 1.000$ are the solution to

$$1.133x_1 + 5.281x_2 = 6.414$$

$$24.14x_1 - 1.210x_2 = 22.93.$$

Use four-digit arithmetic and Gaussian elimination with trivial pivoting to find a computed approximate solution to the system.

The multiple $m_{21} = 24.14/1.133 = 21.31$ of row 1 is to be subtracted from row 2 to obtain the upper-triangular system. Using four digits in the calculations, we obtain the new coefficient

$$a_{22}^{(2)} = -1.210 - 21.31(5.281) = -1.210 - 112.5 = -113.7$$

$$a_{23}^{(2)} = 22.93 - 21.31(6.414) = 22.93 - 136.7 = -113.8.$$

The computed upper-triangular system is

$$\begin{array}{rcl} 1.133x_1 + 5.281x_2 & = & 6.414 \\ -113.7x_2 & = & -113.8. \end{array} \quad \Rightarrow \quad \begin{array}{l} x_2 = 1.001 \\ x_1 = 0.9956 \end{array}$$

Pivoting to reduce error

Example Use four-digit arithmetic and Gaussian elimination with trivial pivoting to solve the linear system

$$24.14x_1 - 1.210x_2 = 22.93$$

$$1.133x_1 + 5.281x_2 = 6.414.$$

This time $m_{21} = 1.133/24.14 = 0.04693$ is the multiple of row 1 that is to be subtracted from row 2. The new coefficients are

$$a_{22}^{(2)} = 5.281 - 0.04693(-1.210) = 5.281 + 0.05679 = 5.338$$

$$a_{23}^{(2)} = 6.414 - 0.04693(22.93) = 6.414 - 1.076 = 5.338.$$

The computed upper-triangular system is

$$24.14x_1 - 1.210x_2 = 22.93$$

$$5.338x_2 = 5.338.$$

Back substitution is used to compute $x_2 = 5.338/5.338 = 1.000$, and $x_1 = (22.91 + 1.210(1.000))/24.14 = 1.000$.

Pivoting to reduce error

- The purpose of a pivoting strategy is to move the entry of greatest magnitude to the main diagonal and then use it to eliminate the remaining entries in the column.
- The *partial pivoting* strategy: check the magnitude of all the elements in column p that lie on or below the main diagonal. Locate row k in which the element that has the largest absolute value lies,

$$|a_{kp}| = \max\{|a_{pp}|, |a_{p+1 p}|, \dots, |a_{N-1 p}|, |a_{Np}|\},$$

then switch row p with row k if $k > p$.

- Now, each of the multipliers m_{rp} for $r = p + 1, \dots, N$ will be less than or equal to 1 in absolute value

Pivoting to reduce error

- It takes a total of $(4N^3 + 9N^2 - 7N)/6$ arithmetic operations to solve an $N \times N$ system.
- The *scaled partial pivoting* strategy: search all the elements in column p that lie on or below the main diagonal for the one that is largest relative to the entries in its row.
- First search rows p through N for the largest element in magnitude in each row,

$$s_r = \max\{|a_{rp}|, |a_{r\,p+1}|, \dots, |a_{rN}|\} \text{ for } r = p, p+1, \dots, N.$$

- The pivotal row k is determined by finding

$$\frac{|a_{kp}|}{s_k} = \max \left\{ \frac{|a_{pp}|}{s_p}, \frac{|a_{p+1\,p}|}{s_{p+1}}, \dots, \frac{|a_{Np}|}{s_N} \right\}.$$

- Interchange row p and k , unless $p = k$.

Ill conditioning

- A matrix A is called *ill conditioned* if there exists a matrix B for which small perturbations in the coefficients of A or B will produce large changes in $X = A^{-1} B$.
- The system $AX = B$ is said to be ill conditioned when A is ill conditioned. In this case, numerical methods for computing an approximate solution are prone to have more error.
- Ill conditioning occurs when A is “nearly singular” and the determinant of A is close to zero.
- Ill conditioning can also occur in systems of two equations when two lines are nearly parallel (or in three equations when three planes are nearly parallel)

Ill conditioning: an example

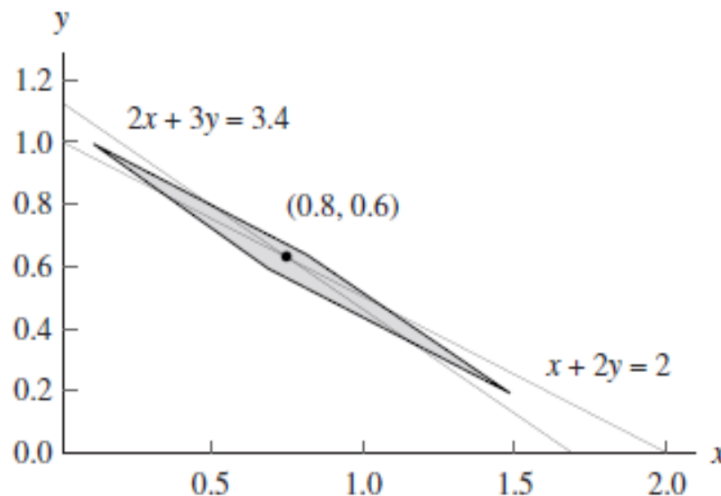
- A consequence of ill conditioning is that substitution of erroneous values may appear to be genuine solutions

(15)
$$\begin{aligned}x + 2y - 2.00 &= 0 \\2x + 3y - 3.40 &= 0.\end{aligned}$$

Substitution of $x_0 = 1.00$ and $y_0 = 0.48$ into these equations “almost produce zeros”:

$$\begin{aligned}1 + 2(0.48) - 2.00 &= 1.96 - 2.00 = -0.04 \approx 0 \\2 + 3(0.48) - 3.40 &= 3.44 - 3.40 = 0.04 \approx 0.\end{aligned}$$

- True solution $x=0.8$ and $y=0.6$.



If it is suspected that a linear system is ill conditioned, computations should be carried out in multiple precision arithmetic.

Figure 3.4 A region where two equations are “almost satisfied.”

Ill conditioning: another example

- Ill conditioning has more drastic consequences when several equations are involved.
- Consider the problem of finding the cubic polynomial $y = c_1 x^3 + c_2 x^2 + c_3 x + c_4$ that passes through the four points (2, 8), (3, 27), (4, 64), and (5, 125). Method of least squares will need to solve:

$$\begin{bmatrix} 20,514 & 4,424 & 978 & 224 \\ 4,424 & 978 & 224 & 54 \\ 978 & 224 & 54 & 14 \\ 224 & 54 & 14 & 4 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} = \begin{bmatrix} 20,514 \\ 4,424 \\ 978 \\ 224 \end{bmatrix}.$$

- A computer that carried nine digits of precision was used to compute the coefficients and obtained

$$c_1 = 1.000004, \quad c_2 = -0.000038, \quad c_3 = 0.000126, \quad \text{and} \quad c_4 = -0.000131.$$

- Furthermore, suppose that the coefficient $a_{11} = 20,514$ in the upper-left corner of the coefficient matrix is changed to the value 20,515 and the perturbed system is solved.

$$c_1 = 0.642857, \quad c_2 = 3.75000, \quad c_3 = -12.3928, \quad \text{and} \quad c_4 = 12.7500.$$

Matlab Program

Program 3.2 (Upper Triangularization Followed by Back Substitution). To construct the solution to $AX = B$, by first reducing the augmented matrix $[A|B]$ to upper-triangular form and then performing back substitution.

```
function X = uptrbk(A,B)
%Input  - A is an N x N nonsingular matrix
%        - B is an N x 1 matrix
%Output - X is an N x 1 matrix containing the solution to AX=B.
%Initialize X and the temporary storage matrix C
    [N N]=size(A);
    X=zeros(N,1);
    C=zeros(1,N+1);
%Form the augmented matrix:Aug=[A|B]
    Aug=[A B];
    for p=1:N-1
        %Partial pivoting for column p
        [Y,j]=max(abs(Aug(p:N,p)));
        %Interchange row p and j
        C=Aug(p,:);
        Aug(p,:)=Aug(j+p-1,:);
        Aug(j+p-1,:)=C;
        if Aug(p,p)==0
            'A was singular. No unique solution'
            break
        end
        %Elimination process for column p
        for k=p+1:N
            m=Aug(k,p)/Aug(p,p);
            Aug(k,p:N+1)=Aug(k,p:N+1)-m*Aug(p,p:N+1);
        end
    end
    %Back Substitution on [U|Y] using Program 3.1
    X=backsub(Aug(1:N,1:N),Aug(1:N,N+1));
```

Triangular Factorization

Triangular Factorization

- **Definition:** The nonsingular matrix A has a *triangular factorization* if it can be expressed as the product of a lower-triangular matrix L and an upper-triangular matrix U :

$$A = LU$$

- In matrix form, this is written as

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ m_{21} & 1 & 0 & 0 \\ m_{31} & m_{32} & 1 & 0 \\ m_{41} & m_{42} & m_{43} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14} \\ 0 & u_{22} & u_{23} & u_{24} \\ 0 & 0 & u_{33} & u_{34} \\ 0 & 0 & 0 & u_{44} \end{bmatrix}$$

- The condition that A is nonsingular implies that $u_{kk} \neq 0$ for all k

Solution of a Linear System

- Suppose that the coefficient matrix A for the linear system $AX = B$ has a triangular factorization; then the solution to

$$LUX = B$$

can be obtained by defining $Y = UX$ and then solving two systems:

first solve $LY = B$ for Y ; then solve $UX = Y$ for X

- In equation form, first solve the lower-triangular system to obtain Y , and use them in solving the upper-triangular system

$$\begin{array}{rclcl} y_1 & & = b_1 & & u_{11}x_1 + u_{12}x_2 + u_{13}x_3 + u_{14}x_4 = y_1 \\ m_{21}y_1 + y_2 & & = b_2 & & u_{22}x_2 + u_{23}x_3 + u_{24}x_4 = y_2 \\ m_{31}y_1 + m_{32}y_2 + y_3 & & = b_3 & \rightarrow & u_{33}x_3 + u_{34}x_4 = y_3 \\ m_{41}y_1 + m_{42}y_2 + m_{43}y_3 + y_4 & & = b_4 & & u_{44}x_4 = y_4. \end{array}$$

Example

- Solve

$$\begin{aligned}x_1 + 2x_2 + 4x_3 + x_4 &= 21 \\2x_1 + 8x_2 + 6x_3 + 4x_4 &= 52 \\3x_1 + 10x_2 + 8x_3 + 8x_4 &= 79 \\4x_1 + 12x_2 + 10x_3 + 6x_4 &= 82.\end{aligned}$$

$$A = \begin{bmatrix} 1 & 2 & 4 & 1 \\ 2 & 8 & 6 & 4 \\ 3 & 10 & 8 & 8 \\ 4 & 12 & 10 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 3 & 1 & 1 & 0 \\ 4 & 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 4 & 1 \\ 0 & 4 & -2 & 2 \\ 0 & 0 & -2 & 3 \\ 0 & 0 & 0 & -6 \end{bmatrix} = LU$$

$$\begin{aligned}y_1 &= 21 \\2y_1 + y_2 &= 52 \\3y_1 + y_2 + y_3 &= 79 \\4y_1 + y_2 + 2y_3 + y_4 &= 82\end{aligned}$$



$$Y = (21, 10, 6, -24)$$

$$\begin{aligned}x_1 + 2x_2 + 4x_3 + x_4 &= 21 \\4x_2 - 2x_3 + 2x_4 &= 10 \\-2x_3 + 3x_4 &= 6 \\-6x_4 &= -24\end{aligned}$$



$$X = (1, 2, 3, 4)$$

Triangular Factorization: An Example

- How to obtain the triangular factorization

$$\mathbf{A} = \begin{bmatrix} 4 & 3 & -1 \\ -2 & -4 & 5 \\ 1 & 2 & 6 \end{bmatrix}$$

- The matrix \mathbf{L} will be constructed from an identity matrix placed at the left. For each row operation used to construct the upper-triangular matrix, the multipliers m_{ij} will be put in their proper places at the left. Start with

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 & 3 & -1 \\ -2 & -4 & 5 \\ 1 & 2 & 6 \end{bmatrix}$$

- Row 1 is used to eliminate the elements of \mathbf{A} in column 1 below a_{11} . The multiples $m_{21} = -0.5$ and $m_{31} = 0.25$ of row 1 are subtracted from rows 2 and 3, respectively. These multipliers are put in the matrix at the left and the result is

Triangular Factorization: An Example

$$A = \begin{bmatrix} 1 & 0 & 0 \\ -0.5 & 1 & 0 \\ 0.25 & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 & 3 & -1 \\ 0 & -2.5 & 4.5 \\ 0 & 1.25 & 6.25 \end{bmatrix}$$

- Row 2 is used to eliminate the elements in column 2 below the diagonal of the second factor of A in the above product. The multiple $m_{32} = -0.5$ of the second row is subtracted from row 3, and the multiplier is entered in the matrix at the left and we have the desired triangular factorization of A .

$$A = \begin{bmatrix} 1 & 0 & 0 \\ -0.5 & 1 & 0 \\ 0.25 & -0.5 & 1 \end{bmatrix} \begin{bmatrix} 4 & 3 & -1 \\ 0 & -2.5 & 4.5 \\ 0 & 0 & 8.5 \end{bmatrix}$$

Direct Factorization

- **Theorem:** Suppose that Gaussian elimination, without row interchanges, can be performed successfully to solve the general linear system $\mathbf{AX} = \mathbf{B}$. Then the matrix \mathbf{A} can be factored as the product of a lower-triangular matrix \mathbf{L} and an upper-triangular matrix \mathbf{U} :

$$\mathbf{A} = \mathbf{LU}$$

- Furthermore, \mathbf{L} can be constructed to have 1's on its diagonal and \mathbf{U} will have nonzero diagonal elements. After finding \mathbf{L} and \mathbf{U} , the solution \mathbf{X} is computed in two steps:
 1. Solve $\mathbf{LY} = \mathbf{B}$ for \mathbf{Y} using forward substitution.
 2. Solve $\mathbf{UX} = \mathbf{Y}$ for \mathbf{X} using back substitution.

Direct Factorization: Procedure

- **Step 1.** Store the coefficients in the augmented matrix.

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2N}^{(1)} & a_{2\ N+1}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} & \cdots & a_{3N}^{(1)} & a_{3\ N+1}^{(1)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_{N1}^{(1)} & a_{N2}^{(1)} & a_{N3}^{(1)} & \cdots & a_{NN}^{(1)} & a_{N\ N+1}^{(1)} \end{array} \right]$$

- **Step 2.** Eliminate x_1 in rows 2 through N and store the multiplier m_{r1} , used to eliminate x_1 in row r , in the matrix at location $(r, 1)$.

for $r = 2:N$

$$m_{r1} = a_{r1}^{(1)} / a_{11}^{(1)};$$

$$a_{r1} = m_{r1};$$

for $c = 2:N + 1$

$$a_{rc}^{(2)} = a_{rc}^{(1)} - m_{r1} * a_{1c}^{(1)};$$

end

end

Direct Factorization: Procedure

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ m_{21} & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} & a_{2\ N+1}^{(2)} \\ m_{31} & a_{32}^{(2)} & a_{33}^{(2)} & \cdots & a_{3N}^{(2)} & a_{3\ N+1}^{(2)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ m_{N1} & a_{N2}^{(2)} & a_{N3}^{(2)} & \cdots & a_{NN}^{(2)} & a_{N\ N+1}^{(2)} \end{array} \right].$$

- **Step 3.** Eliminate x_2 in rows 3 through N and store the multiplier m_{r2} , used to eliminate x_2 in row r , in the matrix at location $(r, 2)$.

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ m_{21} & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} & a_{2\ N+1}^{(2)} \\ m_{31} & m_{32} & a_{33}^{(3)} & \cdots & a_{3N}^{(3)} & a_{3\ N+1}^{(3)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ m_{N1} & m_{N2} & a_{N3}^{(3)} & \cdots & a_{NN}^{(3)} & a_{N\ N+1}^{(3)} \end{array} \right].$$

Direct Factorization: Procedure

- **Step $p + 1$** . This is the general step. Eliminate x_p in rows $p + 1$ through N and store the multipliers at the location (r, p) .

```

for  $r = p + 1:N$ 
     $m_{rp} = a_{rp}^{(p)} / a_{pp}^{(p)}$ ;
     $a_{rp} = m_{rp}$ ;
    for  $c = p + 1:N + 1$ 
         $a_{rc}^{(p+1)} = a_{rc}^{(p+1)} - m_{rp} * a_{pc}^{(p)}$ ;
    end
end

```

- *Final result*

$$\left[\begin{array}{ccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1N}^{(1)} & a_{1\ N+1}^{(1)} \\ m_{21} & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2N}^{(2)} & a_{2\ N+1}^{(2)} \\ m_{31} & m_{32} & a_{33}^{(3)} & \cdots & a_{3N}^{(3)} & a_{3\ N+1}^{(3)} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ m_{N1} & m_{N2} & m_{N3} & \cdots & a_{NN}^{(N)} & a_{N\ N+1}^{(N)} \end{array} \right].$$

Direct Factorization: Verification

- Now we verify that $\mathbf{LU} = \mathbf{A}$. Suppose $\mathbf{D} = \mathbf{LU}$ and consider the case $r \leq c$

$$d_{rc} = m_{r1}a_{1c}^{(1)} + m_{r2}a_{2c}^{(2)} + \cdots + m_{rr-1}a_{r-1c}^{(r-1)} + a_{rc}^{(r)}.$$

$$m_{r1}a_{1c}^{(1)} = a_{rc}^{(1)} - a_{rc}^{(2)},$$

$$m_{r2}a_{2c}^{(2)} = a_{rc}^{(2)} - a_{rc}^{(3)},$$

\vdots

$$m_{rr-1}a_{r-1c}^{(r-1)} = a_{rc}^{(r-1)} - a_{rc}^{(r)}.$$

$$d_{rc} = a_{rc}^{(1)} - a_{rc}^{(2)} + a_{rc}^{(2)} - a_{rc}^{(3)} + \cdots + a_{rc}^{(r-1)} - a_{rc}^{(r)} + a_{rc}^{(r)} = a_{rc}^{(1)}.$$

- The other case, $r > c$, is similar to prove

Computational Complexity

- The triangular factorization $A = LU$ requires

$$\sum_{p=1}^{N-1} (N-p)(N-p+1) = \frac{N^3 - N}{3} \quad \text{multiplications and divisions}$$

$$\sum_{p=1}^{N-1} (N-p)(N-p) = \frac{2N^3 - 3N^2 + N}{6} \quad \text{subtractions.}$$

- Finding the solution to $LUX = B$ requires
 N^2 multiplications and divisions, and $N^2 - N$ subtractions
- The bulk of the calculations lies in the triangularization.
- If the linear system is to be solved many times, with the same matrix A but with different B , it is not necessary to triangularize the matrix each time if the factors are saved. This is the reason the triangular factorization method is usually chosen over the elimination method.

Permutation Matrices

- The $\mathbf{A} = \mathbf{LU}$ factorization in previous Theorem assumes that there are no row interchanges. It is possible that a nonsingular matrix \mathbf{A} cannot be factored directly as $\mathbf{A} = \mathbf{LU}$
- Example: Show that the following matrix cannot be factored directly as $\mathbf{A} = \mathbf{LU}$

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 6 \\ 4 & 8 & -1 \\ -2 & 3 & 5 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 6 \\ 4 & 8 & -1 \\ -2 & 3 & 5 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ m_{21} & 1 & 0 \\ m_{31} & m_{32} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix} \quad ?$$

Permutation Matrices

- Definition: An $N \times N$ **permutation matrix** \mathbf{P} is a matrix with precisely one entry whose value is 1 in each column and row, and all of whose other entries are 0. The rows of \mathbf{P} are a permutation of the rows of the identity matrix and can be written as

$$\mathbf{P} = [\mathbf{E}'_{k1} \quad \mathbf{E}'_{k2} \quad \cdots \quad \mathbf{E}'_{kN}]'.$$

The elements of $\mathbf{P} = [p_{ij}]$ have the form

$$p_{ij} = \begin{cases} 1 & j = k_i, \\ 0 & \text{otherwise.} \end{cases}$$

For example, the following 4×4 matrix is a permutation matrix,

$$\mathbf{P} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} = [\mathbf{E}'_2 \quad \mathbf{E}'_1 \quad \mathbf{E}'_4 \quad \mathbf{E}'_3]'.$$

Permutation Matrices

Theorem Suppose that $\mathbf{P} = [\mathbf{E}'_{k1} \quad \mathbf{E}'_{k2} \quad \cdots \quad \mathbf{E}'_{kN}]'$ is a permutation matrix. The product \mathbf{PA} is a new matrix whose rows consist of the rows of \mathbf{A} rearranged in the order $\text{row}_{k1}\mathbf{A}, \text{row}_{k2}\mathbf{A}, \dots, \text{row}_{kN}\mathbf{A}$.

Example Let \mathbf{A} be a 4×4 matrix and let \mathbf{P} be the permutation matrix given in (15); then \mathbf{PA} is the matrix whose rows consists of the rows of \mathbf{A} rearranged in the order $\text{row}_2\mathbf{A}, \text{row}_1\mathbf{A}, \text{row}_4\mathbf{A}, \text{row}_3\mathbf{A}$.

Computing the product, we have

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} = \begin{bmatrix} a_{21} & a_{22} & a_{23} & a_{24} \\ a_{11} & a_{12} & a_{13} & a_{14} \\ a_{41} & a_{42} & a_{43} & a_{44} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix}$$

Permutation Matrices

- **Theorem:** If P is a permutation matrix, then it is nonsingular and $P^{-1} = P^T$
- **Theorem:** If A is a nonsingular matrix, then there exists a permutation matrix P so that PA has a triangular factorization

$$PA = LU$$

- Example

$$PA = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 2 & 6 \\ 4 & 8 & -1 \\ -2 & 3 & 5 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 6 \\ -2 & 3 & 5 \\ 4 & 8 & -1 \end{bmatrix}$$

$$\begin{array}{l} \text{pivot} \rightarrow \\ m_{21} = -2 \\ m_{31} = 4 \end{array} \begin{bmatrix} 1 & 2 & 6 \\ -2 & 3 & 5 \\ 4 & 8 & -1 \end{bmatrix} \quad \Rightarrow \quad \begin{array}{l} \text{pivot} \rightarrow \\ m_{32} = 0 \end{array} \begin{bmatrix} 1 & 2 & 6 \\ 0 & 7 & 17 \\ 0 & 0 & -25 \end{bmatrix} = U$$

Extending the Gaussian Elimination Process

Theorem (Indirect Factorization: $PA = LU$). Let A be a given $N \times N$ matrix. Assume that Gaussian elimination can be performed successfully to solve the general linear system $AX = B$, but that row interchanges are required. Then there exists a permutation matrix P so that the product PA can be factored as the product of a lower-triangular matrix L and an upper-triangular matrix U :

$$PA = LU.$$

Furthermore, L can be constructed to have 1's on its main diagonal and U will have nonzero diagonal elements. The solution X is found in four steps:

1. Construct the matrices L , U , and P .
2. Compute the column vector PB .
3. Solve $LY = PB$ for Y using forward substitution.
4. Solve $UX = Y$ for X using back substitution.

Remark. Suppose that $AX = B$ is to be solved for a fixed matrix A and several different matrices B . Then step 1 is performed only once and step 2 - 4 are used to find the solution X that corresponds to B . Step 2 - 4 are a computationally efficient way to construct the solution X and require $O(N^2)$ operations instead of the $O(N^3)$ operations required by Gaussian elimination.

Matlab

The MATLAB command $[L, U, P] = \text{lu}(A)$ creates the lower-triangular matrix L , the upper-triangular matrix U (from the triangular factorization of A), and the permutation matrix P from Theorem 3.14.

Example 3.25. Use the MATLAB command $[L, U, P] = \text{lu}(A)$ on the matrix A in Example 3.22. Verify that $A = P^{-1}LU$ (equivalent to showing that $PA = LU$).

```
>> A = [1 2 6 ; 4 8 - 1; -2 3 - 5];
```

```
>> [L, U, P] = lu(A)
```

L=

```
1.0000  0      0
-0.5000 1.0000  0
0.2500  0      1.0000
```

U=

```
4.0000  8.0000 -1.0000
0       7.0000  4.5000
0       0       6.2500
```

P=

```
0 1 0
0 0 1
1 0 0
```



```
>> inv(P) * L * U
```

```
1 2 6
4 8 -1
-2 3 5
```

Program 3.3 ($PA = LU$: Factorization with Pivoting). To construct the solution to the linear system $AX = B$, where A is a nonsingular matrix.

```
function X = lufact(A,B)

%Input   - A is an N x N matrix
%         - B is an N x 1 matrix
%Output  - X is an N x 1 matrix containing the solution to AX = B.

%Initialize X, Y, the temporary storage matrix C, and the row
% permutation information matrix R
    [N,N]=size(A);
    X=zeros(N,1);
    Y=zeros(N,1);
    C=zeros(1,N);
    R=1:N;

for p=1:N-1
%Find the pivot row for column p
    [max1,j]=max(abs(A(p:N,p)));
%Interchange row p and j
    C=A(p,:);
    A(p,:)=A(j+p-1,:);
    A(j+p-1,:)=C;
    d=R(p);
    R(p)=R(j+p-1);
    R(j+p-1)=d;
end
```

```

if A(p,p)==0
    'A is singular. No unique solution'
    break
end
%Calculate multiplier and place in subdiagonal portion of A
for k=p+1:N
    mult=A(k,p)/A(p,p);
    A(k,p) = mult;
    A(k,p+1:N)=A(k,p+1:N)-mult*A(p,p+1:N);
end
end
%Solve for Y
Y(1) = B(R(1));
for k=2:N
    Y(k)= B(R(k))-A(k,1:k-1)*Y(1:k-1);
end
%Solve for X
X(N)=Y(N)/A(N,N);
for k=N-1:-1:1
    X(k)=(Y(k)-A(k,k+1:N)*X(k+1:N))/A(k,k);
end

```

Iterative Methods for Linear Systems

The goal of this section is to extend some of the iterative methods introduced in Chapter 2 to higher dimensions. We consider an extension of fixed-point iteration that applies to systems of linear equations.

Jacobi Iteration: An example

$$\begin{aligned}4x - y + z &= 7 \\4x - 7y + z &= -21 \\-2x + y + 5z &= 15.\end{aligned}$$



$$\begin{aligned}x &= \frac{7 + y - z}{4} \\y &= \frac{21 + 4x + z}{8} \\z &= \frac{15 + 2x - y}{5},\end{aligned}$$

- Consider the following iterative process

$$\begin{aligned}x_{k+1} &= \frac{7 + y_k - z_k}{4} \\y_{k+1} &= \frac{21 + 4x_k + z_k}{8} \\z_{k+1} &= \frac{15 + 2x_k - y_k}{5}.\end{aligned}$$

- Start with $\mathbf{P}_0 = (x_0, y_0, z_0)$

$$\begin{aligned}x_1 &= \frac{7 + 2 - 2}{4} = 1.75 \\y_1 &= \frac{21 + 4 + 2}{8} = 3.375 \\z_1 &= \frac{15 + 2 - 2}{5} = 3.00.\end{aligned}$$

Jacobi Iteration: An example

Table 3.2 Convergent Jacobi Iteration for the Linear System (1)

k	x_k	y_k	z_k
0	1.0	2.0	2.0
1	1.75	3.375	3.0
2	1.84375	3.875	3.025
3	1.9625	3.925	2.9625
4	1.99062500	3.97656250	3.00000000
5	1.99414063	3.99531250	3.00093750
\vdots	\vdots	\vdots	\vdots
15	1.99999993	3.99999985	2.99999993
\vdots	\vdots	\vdots	\vdots
19	2.00000000	4.00000000	3.00000000

- This process is called **Jacobi iteration** and can be used to solve certain types of linear systems.

Jacobi Iteration

- **Jacobi iteration** can be used to solve certain types of linear systems.
- The coefficient matrices for large linear systems with many variables are sparse; that is, a large percentage of the entries of the coefficient matrix are zero.
- If there is a pattern to the nonzero entries (i.e., tridiagonal systems), then an iterative process provides an efficient method for solving these large systems.
- Sometimes the Jacobi method does not work.

$$\begin{array}{rcl} -1x + y + 5z & = & 15 \\ 4x - 8y + z & = & -21 \\ 4x - y + z & = & 7. \end{array}$$



$$\begin{aligned} x &= \frac{-15 + y + 5z}{3} \\ y &= \frac{21 + 4x + z}{8} \\ z &= 7 - 4x + y, \end{aligned}$$

Jacobi Iteration

$$x_{k+1} = \frac{-15 + y_k + 5z_k}{3}$$

$$y_{k+1} = \frac{21 + 4x_k + z_k}{8}$$

$$z_{k+1} = 7 - 4x_k + y_k.$$

Table 3.3 Divergent Jacobi Iteration for the Linear System (4)

k	x_k	y_k	z_k
0	1.0	2.0	2.0
1	-1.5	3.375	5.0
2	6.6875	2.5	16.375
3	34.6875	8.015625	-17.25
4	-46.617188	17.8125	-123.73438
5	-307.929688	-36.150391	211.28125
6	502.62793	-124.929688	1202.56836
\vdots	\vdots	\vdots	\vdots

Gauss-Seidel Iteration

$$\begin{aligned} 4x - y + z &= 7 \\ 4x - 8y + z &= -21 \\ -2x + y + 5z &= 15. \end{aligned}$$



$$\begin{aligned} x_{k+1} &= \frac{7 + y_k - z_k}{4} \\ y_{k+1} &= \frac{21 + 4x_k + z_k}{8} \\ z_{k+1} &= \frac{15 + 2x_k - y_k}{5}. \end{aligned}$$

- The convergence can be speeded up.

$$\begin{aligned} x_{k+1} &= \frac{7 + y_k - z_k}{4} \\ y_{k+1} &= \frac{21 + 4x_{k+1} + z_k}{8} \\ z_{k+1} &= \frac{15 + 2x_{k+1} - y_{k+1}}{5}. \end{aligned}$$

Table 3.4 Convergent Gauss-Seidel Iteration for the System (1)

k	x_k	y_k	z_k
0	1.0	2.0	2.0
1	1.75	3.75	2.95
2	1.95	3.96875	2.98625
3	1.995625	3.99609375	2.99903125
\vdots	\vdots	\vdots	\vdots
8	1.99999983	3.99999988	2.99999996
9	1.99999998	3.99999999	3.00000000
10	2.00000000	4.00000000	3.00000000

Generalized Processes

$$\begin{array}{cccccc}
 a_{11}x_1 + a_{12}x_2 & + \cdots + a_{1j}x_j + \cdots + & a_{1N}x_N & = & b_1 \\
 a_{21}x_1 + a_{22}x_2 & + \cdots + a_{2j}x_j + \cdots + & a_{2N}x_N & = & b_2 \\
 \vdots & \vdots & \vdots & & \vdots \\
 a_{j1}x_1 + a_{j2}x_2 & + \cdots + a_{jj}x_j + \cdots + & a_{jN}x_N & = & b_j \\
 \vdots & \vdots & \vdots & & \vdots \\
 a_{N1}x_1 + a_{N2}x_2 & + \cdots + a_{Nj}x_j + \cdots + & a_{NN}x_N & = & b_N.
 \end{array}$$

Jacobi iteration:

$$(10) \quad x_j^{(k+1)} = \frac{b_j - a_{j1}x_1^{(k)} - \cdots - a_{jj-1}x_{j-1}^{(k)} - a_{jj+1}x_{j+1}^{(k)} - \cdots - a_{jN}x_N^{(k)}}{a_{jj}}$$

for $j = 1, 2, \dots, N$.

Gauss-Seidel iteration:

$$(11) \quad x_j^{(k+1)} = \frac{b_j - a_{j1}x_1^{(k+1)} - \cdots - a_{jj-1}x_{j-1}^{(k+1)} - a_{jj+1}x_{j+1}^{(k)} - \cdots - a_{jN}x_N^{(k)}}{a_{jj}}$$

for $j = 1, 2, \dots, N$.

Will the iteration converge?

Definition 3.6. A matrix A of dimension $N \times N$ is said to be *strictly diagonally dominant* provided that

$$(8) \quad |a_{kk}| > \sum_{\substack{j=1 \\ j \neq k}}^N |a_{kj}| \quad \text{for } k = 1, 2, \dots, N.$$

Theorem 3.15 (Jacobi Iteration). Suppose that A is strictly diagonally dominant matrix. Then $AX = B$ has a unique solution $X = P$. Iteration using formula (10) will produce a sequence of vectors $\{P_k\}$ that will converge to P for any choice of the starting vector P_0

- The Gauss-Seidel method will also converge when the matrix A is strictly diagonally dominant..

Convergence criteria?

- A measure of the closeness between vectors is needed so that we can determine if $\{\mathbf{P}_k\}$ is converging to \mathbf{P} . The Euclidean distance between \mathbf{P} and \mathbf{Q} is:

$$\|\mathbf{P} - \mathbf{Q}\| = \left(\sum_{j=1}^N (x_j - y_j)^2 \right)^{1/2}$$

- Euclidean distance requires considerable computing effort. Another norm is defined:

$$\|\mathbf{X}\|_1 = \sum_{j=1}^N |x_j|.$$

- The $\|\cdot\|_1$ is easier to compute and use for determining convergence in N -dimensional space.

MATLAB code

Program 3.4 (Jacobi Iteration). To solve the linear system $AX = B$ by starting with an initial guess $X = P_0$ and generating a sequence $\{P_k\}$ that converges to the solution. A sufficient condition for the method to be applicable is that A is strictly diagonally dominant.

```
function X=jacobi(A,B,P,delta, max1)
% Input  - A is an N x N nonsingular matrix
%         - B is an N x 1 matrix
%         - P is an N x 1 matrix; the initial guess
%         - delta is the tolerance for P
%         - max1 is the maximum number of iterations
% Output - X is an N x 1 matrix: the Jacobi approximation to
%         the solution of AX = B
N = length(B);
for k=1:max1
    for j=1:N
        X(j)=(B(j)-A(j,[1:j-1,j+1:N])*P([1:j-1,j+1:N]))/A(j,j);
    end
    err=abs(norm(X'-P));
    relerr=err/(norm(X)+eps);
    P=X';
    if(err<delta)|(relerr<delta)
        break
    end
end
X=X';
```


MATLAB code

Program 3.5 (Gauss-Seidel Iteration). To solve the linear system $AX = B$ by starting with the initial guess $X = P_0$ and generating a sequence $\{P_k\}$ that converges to the solution. A sufficient condition for the method to be applicable is that A is strictly diagonally dominant.

```
function X=gseid(A,B,P,delta, max1)
% Input - A is an N x N nonsingular matrix
%        - B is an N x 1 matrix
%        - P is an N x 1 matrix; the initial guess
%        - delta is the tolerance for P
%        - max1 is the maximum number of iterations
% Output - X is an N x 1 matrix: the Gauss-Seidel
%          approximation to the solution of AX = B
N = length(B);
for k=1:max1
    for j=1:N
        if j==1
            X(1)=(B(1)-A(1,2:N)*P(2:N))/A(1,1);
        elseif j==N
            X(N)=(B(N)-A(N,1:N-1)*(X(1:N-1)))'/A(N,N);
        else
            %X contains the kth approximations and P the (k-1)st
            X(j)=(B(j)-A(j,1:j-1)*X(1:j-1)'
                -A(j,j+1:N)*P(j+1:N))/A(j,j);
        end
    end
    err=abs(norm(X'-P));
    relerr=err/(norm(X)+eps);
    P=X';
    if(err<delta)|(relerr<delta)
        break
    end
end
X=X';
```