

## CURS 3

### METODE NUMERICE PENTRU SISTEME DE ECUAȚII LINIARE

---

I. Metode directe: Gauss; LU; Cholesky; Cholesky – matrici bandă.

II. Metode iterative: Jacobi; Gauss-Seidel; SOR.

III. Stabilitatea soluției și condiționarea sistemului.

#### 0. INTRODUCERE

Sistemele de ecuații liniare apar în modelarea unor probleme științifice, precum și în metodele numerice de rezolvare a acestora. Exemple: sisteme de ecuații neliniare, aproximarea funcțiilor, probleme la limită pentru ecuații diferențiale, ecuații cu derivate parțiale, optimizare, etc.

Găsirea soluției unui sistem de ordin mare (sau inversarea unei matrici de ordin mare), poate fi o sarcină dificilă în practică, datorită:

(1) Numărului mare de operații aritmetice necesare pentru găsirea soluției.

(2) Erorilor propagate, într-un șir lung de operații cu numere reprezentate în virgulă flotantă.

Dificultatea (1) face nepractică rezolvarea manuală, iar (2) este inerentă în utilizarea calculatorului.

Analiza unei metode cuprinde problemele:

- Numărul de operații aritmetice necesare pentru a găsi soluția.
- Estimarea, înainte de calculare, a preciziei soluției (estimarea *a priori* a erorii).
- Verificarea preciziei soluției calculate (estimarea *a posteriori* a erorii).

Problemele practice se pot clasifica după ordinul sistemului (numărul de ecuații) și tipul matricii sistemului. Matricea sistemului se zice *densă*, dacă majoritatea elementelor sale sunt diferite de zero, și *rară* – dacă majoritatea elementelor sunt zero. Astfel, avem categoriile:

I. Sisteme de ordin moderat, cu matrice densă:

Acestea se rezolvă în memorie, și ordinul sistemului este limitat numai de memoria disponibilă. Uzual, ordinul sistemului poate ajunge la sute de ecuații. Majoritatea algoritmilor se bazează pe eliminarea Gauss, cu variante pentru clase speciale de matrici.

II. Sisteme de ordin mare, cu matrice rară:

Asemenea sisteme pot ajunge la sute de mii de ecuații. Eliminarea Gauss nu mai este convenabilă, și se utilizează *metode iterative*.

III. Sisteme de tip II – rezolvare prin utilizarea de fișiere disc

Memoria fiind insuficientă, se folosesc fișiere pe disc pentru generarea și procesarea matricii. Lucrul cu matricea sistemului se face pe ‘blocuri’: de exemplu, matricea este generată pe disc; succesiv, blocurile sunt citite, procesate în memorie și rescrise pe disc. Un exemplu este cel al sistemelor cu matrice bandă simetrică, sisteme care apar în problemele de analiză statică a structurilor ■

Un sistem de ordinul  $n$  se va nota prin:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ \dots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{cases}$$

Matriceal, sistemul se scrie:

$$\begin{bmatrix} a_{11} & \dots & \dots & a_{1n} \\ a_{21} & \dots & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & \dots & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \quad (1)$$

Sau:

$$\mathbf{Ax} = \mathbf{b}, \quad (1')$$

în care  $\mathbf{A}$  este matricea sistemului ( $n \times n$ ), iar  $\mathbf{x}$  și  $\mathbf{b}$  vectorul necunoscutelor și, respectiv, vectorul termenilor liberi (matrici coloană cu  $n$  elemente):

$$\mathbf{A} = [a_{ij}]_{i,j=1,n} \quad \mathbf{x} = [x_1 \dots x_n]^T, \quad \mathbf{b} = [b_1 \dots b_n]^T$$

## I. METODELE DIRECTE

### 1 Metoda Gauss

#### 1.1 Metoda

Sistemul dat se transformă în sistemul echivalent  $\mathbf{Ux} = \mathbf{g}$ , cu matricea  $\mathbf{U}$  superior triunghiulară, cum urmează. Sistemul dat se notează

$$\mathbf{A}^{(1)} \mathbf{x} = \mathbf{b}^{(1)}.$$

La pasul curent  $k$  ( $k = 1, 2, \dots, n-1$ ), sistemul este

$$\mathbf{A}^{(k)} \mathbf{x} = \mathbf{b}^{(k)}.$$

Se lucrează cu liniile  $i = k, \dots, n$ , din  $\mathbf{A}^{(k)}$  și  $\mathbf{b}^{(k)}$ . (Liniile  $1, \dots, k-1$ , procesate anterior, rămân neschimbate). Concret, se operează cu sub-matricea  $\mathbf{A}^{(k)}$  ( $k : n, k : n$ ) și sub-coloana  $\mathbf{b}^{(k)}$  ( $k : n$ ) a termenilor liberi.

$$\mathbf{A}^{(k)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1k}^{(1)} & \cdot & \cdot & \cdot & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2k}^{(2)} & \cdot & \cdot & \cdot & a_{2n}^{(2)} \\ \cdot & 0 & \ddots & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \dots & a_{k-1,k-1}^{(k-1)} & \cdot & \cdot & \cdot & a_{k-1,n}^{(k-1)} \\ \hline 0 & \cdot & \dots & 0 & a_{kk}^{(k)} & \cdot & a_{kj}^{(k)} & \cdot & a_{kn}^{(k)} \\ \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \dots & \cdot & a_{ik}^{(k)} & \cdot & a_{ij}^{(k)} & \cdot & a_{in}^{(k)} \\ \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & 0 & a_{nk}^{(k)} & \cdot & \cdot & \cdot & a_{nn}^{(k)} \end{bmatrix} \times (-m_{ik});$$

$$\downarrow$$

$$\oplus \leftarrow$$

$$\mathbf{b}^{(k)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ \cdot \\ b_{k-1,k-1}^{(k-1)} \\ \hline b_k^{(k)} \\ \cdot \\ b_i^{(k)} \\ \cdot \\ b_n^{(k)} \end{bmatrix} \times (-m_{ik})$$

$$\downarrow$$

$$\oplus \leftarrow$$

Elementul diagonal  $a_{kk}^{(k)}$  se numește *pivot* (la pasul  $k$ ). Presupunem  $a_{kk}^{(k)} \neq 0$  – v. mai jos [pivotare](#), și definim multiplicatorii de la pasul  $k$ :

$$m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}; \quad i = k + 1, \dots, n$$

Pentru  $i = k + 1, \dots, n$ : linia  $k$  se înmulțește cu  $-m_{ik}$ , și se adună la linia  $i$ . Rezultă elemente nule în coloana  $k$ , sub elementul diagonal  $a_{kk}^{(k)}$ . Linia  $k$  rămâne neschimbată.

Noii coeficienți și termeni liberi vor fi:

$$\left. \begin{aligned} a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} \cdot a_{kj}^{(k)}; & j &= k + 1, \dots, n \\ b_i^{(k+1)} &= b_i^{(k)} - m_{ik} \cdot b_k^{(k)} \end{aligned} \right\} \quad i = k + 1, \dots, n$$

La pasul  $n$  se notează, pentru conveniență,  $\mathbf{U} = \mathbf{A}^{(n)}$ ,  $\mathbf{g} = \mathbf{b}^{(n)}$ , și sistemul devine

$\mathbf{U}\mathbf{x} = \mathbf{g}$ . Explicit:

$$\begin{bmatrix} u_{11} & \cdot & \cdot & \cdot & \cdot & u_{1n} \\ 0 & u_{22} & \cdot & \cdot & \cdot & u_{2n} \\ & & \cdot & \cdot & \cdot & \cdot \\ 0 & & 0 & u_{kk} & \cdot & u_{kn} \\ & & & \cdot & \cdot & \cdot \\ 0 & & & & 0 & u_{nn} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ x_k \\ \cdot \\ x_n \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \\ \cdot \\ g_k \\ \cdot \\ g_n \end{bmatrix} \quad (2)$$

Observați că  $u_{kk} = a_{kk}^{(k)}$ .

Acest sistem se rezolvă prin substituție înapoi:

$$x_n = g_n / u_{nn}; \quad x_k = (g_k - \sum_{j=k+1}^n u_{kj} \cdot x_j) / u_{kk}; \quad k = n-1, n-2, \dots, 1$$

■

## 1.2 Factorizarea (descompunerea) LU

Multiplicatorii  $m_{ik}$  (unde  $i = k + 1, \dots, n$ ) din eliminarea Gauss, se pot reține pentru a rezolva  $\mathbf{Ax} = \mathbf{b}$  cu diferiți termeni liberi  $\mathbf{b}$ . Introducem atunci, matricea inferior-triunghiulară (cu 1 pe diagonală, și multiplicatorii  $m_{ik}$  în sub-coloanele  $k + 1 : n$ ):

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & \cdot & \cdot & \cdot & 0 \\ m_{21} & 1 & 0 & & & 0 \\ m_{31} & m_{32} & 1 & 0 & & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 & 0 \\ m_{n1} & m_{n2} & \cdot & \cdot & m_{n,n-1} & 1 \end{bmatrix} \quad (3)$$

Avem următoarea propoziție:

**Propoziția 1**

Dacă, în eliminarea Gauss, la fiecare pas  $k$  pivotul  $a_{kk}^{(k)} \neq 0$ , atunci  $\mathbf{A} = \mathbf{LU}$  ■

$\mathbf{L}$  este matricea inferior triunghiulară a multiplicatorilor (3), iar  $\mathbf{U}$  matricea superior triunghiulară din (2) ■

**Consecință – Calculul determinantului:**

Dacă la fiecare pas  $k$  pivotul  $a_{kk}^{(k)} \neq 0$ , atunci:

$$\det(\mathbf{A}) = u_{11} \cdot u_{22} \cdot \dots \cdot u_{nn}$$

Observați că  $u_{kk} = a_{kk}^{(k)}$ , astfel că valoarea determinantului este *produsul pivoților* ■

Într-adevăr:  $\det(\mathbf{A}) = \det(\mathbf{L}) \cdot \det(\mathbf{U})$ , iar  $\det(\mathbf{L}) = 1$ .

Pentru calculul determinantului în cazul când se schimbă linii în matricea  $\mathbf{A}$  (se pivotează), v. [mai jos](#).

**1.3 Număr de operații**

Pentru  $n = \text{''mare''}$ , numărul de operații în eliminarea Gauss este:

$$NOP_{Guass} \approx \frac{n^3}{3}$$

(Operație = adunare sau scădere; în calcule, acestea sunt urmate de obicei de o înmulțire sau împărțire.)

■

## 1.4 Pivotare

La fiecare pas  $k$ , *pivotul*  $a_{kk}^{(k)} \neq 0$  apare la numitorul multiplicatorilor  $m_{ik}$ . Astfel, avem condițiile:

### **Pivotul nu trebuie să fie nul, și nici ”foarte mic”:**

- Dacă pivotul este nul, eliminarea Gauss nu poate continua. În acest caz, matricea este singulară.
- Dacă pivotul este ”foarte mic”: rezultă multiplicatori  $m_{ik}$  mari, iar erorile de rotunjire în valoarea lui  $m_{ik}$  conduc la erori mari în calculațiile următoare. Matricea este ”aproape singulară”.

Practic, se testează dacă pivotul este mai mic decât un *prag* ales.

### V. [Exemplul de mai jos](#).

Procedeul de alegere a unui pivot  $\neq 0$  se zice *pivotare*. Aceasta constă în:

- Căutarea elementului de modul maxim, în sub-matricea cu liniile și coloanele  $k, \dots, n$ ;
- Aducerea lui pe poziția de pivot – poziția  $(k, k)$  – prin schimbări de linii, sau schimbări de linii și coloane.

Schimbarea de linii din  $\mathbf{A}$ , se face și în vectorul  $\mathbf{b}$ .

Dacă pivotul este mai mic decât pragul, eliminarea Gauss se oprește.

### **Strategii de pivotare:**

#### 1) Pivotarea parțială:

Se caută elementul de modul maxim în sub-coloana  $a(k : n, k)$ . Acesta devine pivot, aducându-l în poziția  $(k, k)$ . Fie pivotul găsit în linia  $I \geq k$ . Dacă pivotul este mai mare decât pragul și  $I > k$ , atunci se schimbă liniile  $k$  și  $I$  (în  $\mathbf{A}^{(k)}$  și în  $\mathbf{b}^{(k)}$ ).

$$\left[ \begin{array}{c|ccc} a_{kk}^{(k)} & \dots & a_{kn}^{(k)} & (k) \\ \vdots & & \vdots & \updownarrow \\ a_{lk}^{(k)} & \dots & \cdot & (I) \\ \vdots & & \vdots & \\ a_{nk}^{(k)} & \dots & a_{nn}^{(k)} & \end{array} \right]$$

$$\left[ \begin{array}{c} b_k^{(k)} \\ \vdots \\ b_I^{(k)} \\ \vdots \\ b_n^{(k)} \end{array} \right] \begin{array}{l} (k) \\ \updownarrow \\ (I) \\ \vdots \\ \end{array}$$

## 2) Pivotarea completă:

Se caută elementul de modul maxim în sub-matricea  $a(k:n, k:n)$ . Dacă maximum este atins pentru linia  $I > k$  și coloana  $J > k$ , atunci se permută liniile  $k$  și  $I$  și coloanele  $k$  și  $J$  - în  $\mathbf{A}$ , și liniile  $k$  și  $I$  - în  $\mathbf{b}$  ■

$$\left[ \begin{array}{cccc|c} a_{kk}^{(k)} & \dots & \cdot & \dots & a_{kn}^{(k)} & (k) \\ \vdots & & & & \vdots & \updownarrow \\ \cdot & & a_{IJ}^{(k)} & & \cdot & (I) \\ \vdots & & & & \vdots & \\ a_{nk}^{(k)} & \dots & \cdot & \dots & a_{nn}^{(k)} & \end{array} \right]$$

$$\left[ \begin{array}{c} b_k^{(k)} \\ \vdots \\ b_I^{(k)} \\ \vdots \\ b_n^{(k)} \end{array} \right] \begin{array}{l} (k) \\ \updownarrow \\ (I) \\ \vdots \\ \end{array}$$

$(k) \leftrightarrow (J)$

### Observații

- Prin pivotare rezultă, la fiecare pas, multiplicatori de modul  $\leq 1$ :

$$|m_{ik}| \leq 1, \quad i = k+1, \dots, n.$$

Acesta previne generarea în  $\mathbf{A}^{(k+1)}$ , de elemente de mărime foarte diferită, care ar putea duce la erori de pierdere de semnificație.

- La pivotarea completă, datorită permutării de coloane, ordinea necunoscutelor se schimbă. Ea trebuie refăcută la sfârșitul eliminării.

Teoretic, pivotarea completă duce la o precizie mai mare a soluției; în practică însă, diferența este mică față de pivotarea parțială. În plus, pivotarea completă cere mai mult timp-calculator. Din aceste motive, majoritatea algoritmilor utilizează pivotarea

parțială. În ceea ce urmează vom considera numai pivotarea parțială, numită pe scurt, pivotare ■

*Observație*

Dacă nu se pivotează, atunci avem  $\mathbf{LU} = \mathbf{A}$  (Propoziția 1).

Dacă se pivotează, atunci avem  $\mathbf{LU} = \mathbf{A}'$ , unde  $\mathbf{A}' = \mathbf{PA}$ , iar  $\mathbf{P}$  este o matrice de permutare de linii. Matricea  $\mathbf{A}'$  se obține din  $\mathbf{A}$ , permutând liniile în aceeași ordine în care se permută la pivotare. Sistemul care se rezolvă este  $\mathbf{A}'\mathbf{x} = \mathbf{b}'$ , unde  $\mathbf{b}'$  este  $\mathbf{b}$  cu liniile permutate în ordinea de la pivotare ■

**Calculul determinantului:**

Determinantul matricii  $\mathbf{A}' = \mathbf{LU}$  este  $\det(\mathbf{A}') = u_{11}u_{22} \dots u_{nn}$ .

Determinantul matricii  $\mathbf{A}$  va fi

$$\det(\mathbf{A}) = (-1)^{n-l} \det(\mathbf{A}'),$$

unde  $n-l$  este numărul de schimbări de linii în cursul pivotării.

■

**Exemplu – Pivot ”foarte mic”**

Considerăm sistemul  $\mathbf{Ax} = \mathbf{b}$ , unde:

$$\mathbf{A} = \begin{bmatrix} 5 & 6 & 3 & 1 \\ -1 & 0 & -1 & 1 \\ 2 & 2 & 1 & 6 \\ 4 & 2 & 3 & 4 \end{bmatrix}, \quad \mathbf{b} = [1 \quad 1 \quad 1 \quad 1]^T$$

În calculul exact,  $\mathbf{A}$  este singulară. Rezolvând în simplă precizie, cu programul Gauss 2007, cu un prag de  $1E-7$ , se obțin următoarele rezultate:

Permutare, Pivot:

```

1          5.000000
2<->4    -2.800000
3         -0.2857143
4          7.7486033E-07
```



\* Solutia pentru termenii liberi nr. 1

x(1) = -0.2193944E+08

x(2) = 6452776.

x(3) = 0.2322999E+08

x(4) = 1290555.

\* Proba

Diferente (A\*x-b)

1 -4.625000

2 0.3750000

3 -1.750000

4 -4.500000

Adică, soluția nu verifică, și este inutil a o calcula.

Erorile în soluție sunt produse de pivotul de la pasul 4, care este “foarte mic”.

Cu un prag de  $1E-6$ , programul se oprește cu mesajul:

Linia 4: Pivot = 7.7486033E-07

\* Pivot < 1.00E-06 !

\* Sistemul nu se poate rezolva.

În calculul în dublă precizie, pivotul din linia 4 rezultă 0.0000000000000000.

■

### 3 Descompunerea LU

#### 3.1 Calculul direct al factorilor L și U

Problema este calculul direct al factorilor  $\mathbf{L}$  și  $\mathbf{U}$ , astfel ca  $\mathbf{A} = \mathbf{LU}$ .

Elementele lui  $\mathbf{L}$  și  $\mathbf{U}$  se găsesc din ecuațiile  $a_{ij} = \text{linia } i_{\mathbf{L}} \times \text{coloana } j_{\mathbf{U}}$ :

$$a_{ij} = \begin{bmatrix} l_{i1} & l_{i2} & \dots & l_{ii} & \vdots & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} u_{j1} \\ u_{j2} \\ \vdots \\ u_{jj} \\ \dots \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \sum_{k=1}^{\min(i,j)} l_{ik} u_{kj}$$

În particular, înmulțind linia  $i$  din  $\mathbf{L}$  cu coloana  $i$  din  $\mathbf{U}$ , rezultă:

$$a_{ii} = \sum_{k=1}^{i-1} l_{ik} u_{ki} + l_{ii} u_{ii}$$

Există o neunicitate a descompunerii  $\mathbf{LU}$ , care provine din alegerea arbitrară a lui  $l_{ii}$  sau  $u_{ii}$ ,  $i = 1, \dots, n$ . Două metode se pot considera:

$$l_{ii} = \text{arbitrar} \neq 0;$$

$$u_{ii} = \text{arbitrar} \neq 0.$$

Prezentăm, pentru exemplificare, formulele generale pentru prima metodă.

$$l_{ii} = \text{ales arbitrar (nenul)}, i = 1, 2, \dots, n$$

$$u_{ij} = \frac{a_{ij} - \sum_{k=1}^{i-1} l_{ik} \cdot u_{kj}}{l_{ii}}, \quad j = i, i+1, \dots, n \quad (\text{a})$$

Dacă  $u_{ii} \neq 0$ ,  $i = \overline{1, n-1}$ , rezultă:

$$l_{ji} = \frac{a_{ji} - \sum_{k=1}^{i-1} l_{jk} \cdot u_{ki}}{u_{ii}}, \quad j = i+1, \dots, n \quad (\text{b})$$

Elementele lui  $\mathbf{L}$  și  $\mathbf{U}$  se determină în următoarea secvență, conform schemei de mai jos (se calculează elementele din linia  $i$  din  $\mathbf{U}$ , și coloana  $i$  din  $\mathbf{L}$ ).

$$\begin{array}{c|cccc} l_{ii} & \rightarrow & u_{i,i} & u_{i,i+1} & \dots & u_{i,n} \\ \downarrow & & & & & \\ l_{i+1,i} & & & & & \\ \vdots & & & & & \\ l_{n,i} & & & & & \end{array}$$

În particular, la  $i = n$ , se calculează numai  $u_{nn}$  - din (a).

Numărul de operații este același ca în eliminarea Gauss ( $\approx n^3/3$ ).

Două metode se utilizează în practică:

$l_{ii} = 1 \dots$  Metoda *Doolittle* (descompunerea  $\mathbf{LU}$  revine la cea din eliminarea Gauss).

$u_{ii} = 1 \dots$  Metoda *Crout*.

## Posibilitatea descompunerii LU

### Propoziție

Descompunerea  $\mathbf{LU}$  există dacă și numai dacă toate sub-matricile *principale*

$\mathbf{A}(1:i, 1:i)$ ,  $i = \overline{1, n-1}$  sunt nesingulare ■

## 3.2 Pivotare la descompunerea LU

Pivotare parțială: A pivota în descompunerea  $\mathbf{LU}$  înseamnă ca, la fiecare pas  $i$ , să căutăm *pivotul*  $u_{ii}$  în coloana  $i$  a lui  $\mathbf{A}$ , în liniile  $i, i+1, \dots, n$ , și apoi să permutăm două linii.

Pivotarea trebuie făcută înainte de calculul elementelor lui  $\mathbf{U}$  și  $\mathbf{L}$ . Adică, pivotul  $u_{ii}$  trebuie calculat și testat *înainte* de aplicarea formulelor (a) și (b).

Pentru metoda Doolittle, pivotul este dat de (cf. a):

$$u_{ii} = a_{ii} - \sum_{k=1}^{i-1} l_{ik} \cdot u_{ki}$$

Așa cum s-a arătat la metoda Gauss, trebuie pivotat și dacă pivotul este foarte mic. Practic, se testează dacă pivotul este mai mic decât un *prag*.

Observația de la pivotarea în eliminarea Gauss se aplică).

Determinantul matricii **A** (în metoda Doolittle) este

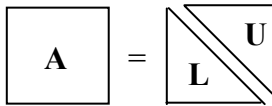
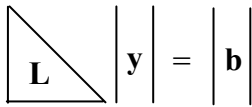
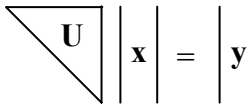
$$\det(\mathbf{A}) = (-1)^{n-l} u_{11} u_{22} \dots u_{nn},$$

unde  $n-l$  este numărul de schimbări de linii în cursul pivotării.

■

### 3.3 Metoda

Rezolvarea sistemului prin descompunerea **LU**, constă în pașii:

1. Factorizare  $\mathbf{A} = \mathbf{LU}$ , cu pivotare parțială. 
2. Rezolvarea sistemului  $\mathbf{Ly} = \mathbf{b}$ , prin substituție înainte; rezultă **y**. 
3. Rezolvarea sistemului  $\mathbf{Ux} = \mathbf{y}$ , prin substituție înapoi; rezultă **x**. 

#### Observații

- Legătura cu eliminarea Gauss este  $\mathbf{y} = \mathbf{g}$ .
- Pentru noi termeni liberi **b**, se repetă numai pașii 2, 3 ( $n^2$  operații suplimentare pentru fiecare **b**)
- Numărul de operații este același ca în eliminarea Gauss ( $\approx n^3 / 3$ ).

■

## 4 Metoda Cholesky

Metoda Cholesky se aplică pentru un sistem cu matrice *simetrică și pozitiv definită*.

$$\text{Simetrie: } \mathbf{A}^T = \mathbf{A} \Leftrightarrow a_{ij} = a_{ji}, \quad i, j = \overline{1, n}$$

$$\text{Definire pozitivă: } \forall \mathbf{x} \neq \mathbf{0} \quad \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_i \sum_j a_{ij} x_i x_j > 0$$

### 4.1 Proprietăți ale matricilor simetrice și pozitiv definite

Pentru o matrice *pozitiv definită*, au loc proprietățile:

- (1) Matricea este nesingulară.
- (2) Elementele diagonale sunt pozitive:  $a_{kk} > 0$ . Mai mult, sub-matricile principale sunt pozitiv definite.

Urmează că  $a_{11}$  poate fi luat ca pivot în eliminarea Gauss, la pasul 1.

Pentru o matrice *simetrică și pozitiv definită*, au loc și proprietățile:

- (3) Sub-matricile  $\mathbf{A}^{(k)}$  ( $k : n, k : n$ ),  $k \geq 2$ , rezultate din eliminarea Gauss sunt simetrice și pozitiv definite.

De aici, rezultă că există pivoți  $a_{kk}^{(k)} > 0$ , la fiecare pas  $k$ .

- (4) Valorile proprii ale lui  $\mathbf{A}$  sunt reale și pozitive.

Mai întâi, se arată că dacă  $\mathbf{A}$  este simetrică, valorile proprii sunt reale – v. Cap. 5.

Apoi, fie  $\lambda$  o valoare proprie și  $\mathbf{x}$  vectorul propriu asociat, atunci  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ .

Inmulțind la stânga cu  $\mathbf{x}^T$  și ținând cont de definirea pozitivă a matricii  $\mathbf{A}$ , rezultă  $\lambda\mathbf{x}^T\mathbf{x} > 0$ , din care urmează că  $\lambda$  este pozitiv.

#### Observație

Reciproc, dacă  $\mathbf{A}$  este simetrică și valorile proprii sunt pozitive, atunci  $\mathbf{A}$  este pozitiv definită. Demonstrația cere rezultate de algebră matriceală care sunt în afara scopului aceluși paragraf. V. de exemplu, R. Bellman, “Introducere în analiza matriceală”, E.T., București, 1969.

(5) Determinantul matricii este pozitiv. Mai mult, toți determinanții principali sunt pozitivi.

Întrucât determinantul matricii este  $\det(\mathbf{A}) = \lambda_1 \lambda_2 \dots \lambda_n$  - v. Cap. 5, urmează că avem  $\det(\mathbf{A}) > 0$ . A doua afirmație decurge din faptul că sub-matricile principale sunt și ele pozitiv definite – conform Proprietății 2.

*Observație*

Reciproc, dacă matricea este simetrică și toți minorii principali sunt pozitivi, matricea este pozitiv definită. Pentru demonstrație – v. Wilkinson-2.

Cu aceasta, rezultă:

*Dacă  $\mathbf{A}$  este simetrică, o condiție necesară și suficientă ca  $\mathbf{A}$  să fie pozitiv definită, este ca toți determinanții principali să fie pozitivi.*

■

## 4.2 Metoda Cholesky

*Propoziție*

Dacă  $\mathbf{A}$  este simetrică și pozitiv definită, atunci există descompunerea

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T, \quad (1)$$

unde  $\mathbf{L}$  este o matrice *inferior* triunghiulară, cu elemente diagonale *pozitive*.

■

Cu alte cuvinte, în descompunerea  $\mathbf{LU}$  se poate alege  $\mathbf{U} = \mathbf{L}^T$ . Notând

$$\mathbf{S} = \mathbf{L}^T,$$

unde  $\mathbf{S}$  este *superior* triunghiulară, avem și descompunerea

$$\mathbf{A} = \mathbf{S}^T \mathbf{S} \quad (2)$$

Cele două forme (1) și (2) reprezintă aceeași descompunere, luând ca referință una sau alta din cele două matrici triunghiulare. Ele sunt ilustrate mai jos.  $\mathbf{0}$  desemnează elementele nule.

$$\begin{array}{c} \boxed{\mathbf{A}} \end{array} = \begin{array}{c} \boxed{\begin{array}{c} \mathbf{0} \\ \mathbf{L} \end{array}} \\ \text{(S}^T\text{)} \end{array} \cdot \begin{array}{c} \boxed{\begin{array}{c} \mathbf{L}^T \\ \mathbf{0} \end{array}} \\ \text{(S)} \end{array}$$

■

**Exemplu** – Evaluarea lui  $\mathbf{L}$ , pentru  $n = 3$ :

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \cdot \begin{bmatrix} l_{11} & l_{21} & l_{31} \\ 0 & l_{22} & l_{32} \\ 0 & 0 & l_{33} \end{bmatrix}$$

Ecuții: Lucrăm cu triunghiul *inferior* al lui  $\mathbf{A}$ , luând elementele *pe coloane*:

$$a_{11} = l_{11}^2, \quad a_{21} = l_{21} \cdot l_{11}, \quad a_{31} = l_{31} \cdot l_{11}$$

$$a_{22} = l_{21}^2 + l_{22}^2, \quad a_{32} = l_{31} \cdot l_{21} + l_{32} \cdot l_{22}$$

$$a_{33} = l_{31}^2 + l_{32}^2 + l_{33}^2$$

Rezultă elementele lui  $\mathbf{L}$ , *pe coloane*:

$$l_{11} = \sqrt{a_{11}}; \quad l_{21} = \frac{a_{21}}{l_{11}}, \quad l_{31} = \frac{a_{31}}{l_{11}}.$$

$$l_{22} = (a_{22} - l_{21}^2)^{1/2}; \quad l_{32} = \frac{a_{32} - l_{31} \cdot l_{21}}{l_{22}}$$

$$l_{33} = (a_{33} - l_{31}^2 - l_{32}^2)^{1/2}$$

Analog, se poate lucra cu triunghiul superior al lui  $\mathbf{A}$  și determina elementele lui  $\mathbf{L}$ , pe linii.

*Exempu numeric:*

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 13 & 18 \\ 3 & 18 & 50 \end{bmatrix}; \quad \mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 3 & 0 \\ 3 & 4 & 5 \end{bmatrix} \quad \blacksquare$$

Pașii rezolvării sunt aceeași ca la [descompunerea LU](#), anume:

- Cu matricea  $\mathbf{L}$  (unde  $\mathbf{U} = \mathbf{L}^T$ ):

1) *Factorizare* (descompunere):  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$

2) *Substituție înainte* (rezultă  $\mathbf{y}$ ):  $\mathbf{L}\mathbf{y} = \mathbf{b}$

3) *Substituție înapoi* (rezultă  $\mathbf{x}$ ):  $\mathbf{L}^T \mathbf{x} = \mathbf{y}$

- Cu matricea  $\mathbf{S}$  (unde  $\mathbf{U} = \mathbf{S}$ ;  $\mathbf{L} = \mathbf{S}^T$ ):

1)  $\mathbf{A} = \mathbf{S}^T \mathbf{S}$

2)  $\mathbf{S}^T \mathbf{y} = \mathbf{b}$

3)  $\mathbf{S}\mathbf{x} = \mathbf{y}$

■

### 4.3 Numărul de operații

Numărul de operații (adunări/scăderi), pentru  $n = \text{'mare'}$ , este

$$NOP_{Cholesky} \approx \frac{n^3}{6}$$

Acesta este cca.  $\frac{1}{2}$  din numărul de operații cerute eliminarea Gauss (sau LU).

(În afară de adunări/scăderi, se mai calculează  $n$  rădăcini pătrate.)

Alte avantaje, în raport cu descompunerea generală LU:

- Matricea  $\mathbf{A}$  fiind simetrică, se poate stoca numai triunghiul inferior (sau superior), cerând numai  $\frac{1}{2}n(n+1)$  locații de memorie. În aceste locații se stochează matricea  $\mathbf{L}$  (sau  $\mathbf{S}$ ).
- Nu este necesară pivotarea – conform proprietăților (2, 3).



*Observație*

Rădăcinile pătrate – care consumă timp-calculator – pot fi evitate printr-o ușoară modificare a descompunerii  $\mathbf{LU}$ , și anume: căutăm o matrice inferior triunghiulară  $\tilde{\mathbf{L}}$  cu 1 pe diagonală, și o matrice diagonală  $\mathbf{D}$ , astfel ca:

$$\mathbf{A} = \tilde{\mathbf{L}}\mathbf{D}\tilde{\mathbf{L}}^T$$

Într-adevăr: dacă avem  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ , fie  $\tilde{\mathbf{L}}$  matricea obținută din  $\mathbf{L}$  punând pe diagonală  $\tilde{l}_{ii} = 1$  (în rest,  $\tilde{l}_{ij} = l_{ij}$ ), și  $\mathbf{D}' = \text{diag}(l_{ii})$ . Avem  $\mathbf{L} = \tilde{\mathbf{L}}\mathbf{D}'$ , și  $\mathbf{L}\mathbf{L}^T = \tilde{\mathbf{L}}\mathbf{D}'\mathbf{D}'\tilde{\mathbf{L}}^T$ . Astfel, matricea  $\mathbf{D} = \mathbf{D}'\mathbf{D}' = \text{diag}(l_{ii}^2)$ .

Această factorizare se poate face cu aproximativ același număr de operații ca în metoda Cholesky, și fără calculul de rădăcini pătrate ■

## 5 Matrici bandă – simetrice și pozitiv definite

### 5.1 Matrice bandă simetrică

Presupunem că matricea sistemului este simetrică și, în plus, are structura de *matrice bandă*, adică în fiecare linie elementele matricii sunt constituite din:

- Elementul diagonal, un număr de  $LIM-1$  elemente la stânga acestuia, și  $LIM-1$  elemente la dreapta.
- Celelalte elemente din linie sunt zero.

Numărul  $LIM$  va fi numit *semi-lățimea* de bandă.  $LIM$  reprezintă numărul elementelor din semi-bandă, inclusiv elementul diagonal.

*Observație*

Între elementele din bandă, pot fi și elemente nule, dar caracterul de bandă este dat de faptul că toate elementele situate în afara benzii sunt nule. În acest sens  $LIM$  poate fi considerat semi-lățimea de bandă *maximă* ■

Exemplu: Matrice  $6 \times 6$ ,  $a_{ij} = a_{ji}$ , și  $LIM = 3$ :



Stocarea se mai poate face pe linii de  $LIM$  elemente, într-un tablou  $B(n, LIM)$ .  
Liniile care conțin mai puțin de  $LIM$  elemente se completează cu elemente nule  
(exemplu: liniile 5 și 6).

### 5.3 Metoda Cholesky

Proprietatea esențială este următoarea:

Dacă matricea bandă este *simetrică și pozitiv definită*, atunci descompunerea  $LL^T$   
sau  $S^T S$  poate fi făcută lucrând *exclusiv* în bandă ■

În ceea ce urmează vom considera descompunerea  $A = S^T S$ , lucrând cu semi-banda  
superioară.

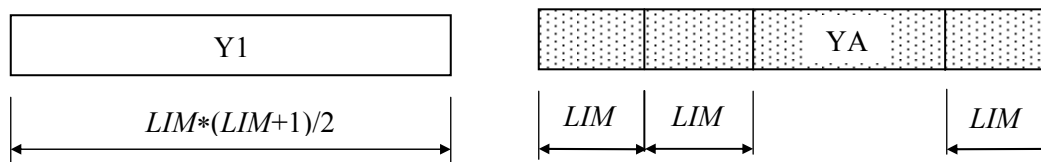
Pașii sunt cei de la Cholesky, lucrul cu [matricea S](#).

Elementele active la un pas al descompunerii, sunt conținute într-un triunghi cu  
laturile  $LIM$  – numit “triunghiul activ”. În cursul procesului, triunghiul activ coboară  
cu câte o linie în bandă.

Metoda implementată în ANA\Cholesky\_Band utilizează un vector de lucru Y, de  
dimensiune  $NY$ , unde:

$$NY0 = LIM * (LIM + 1) / 2 ; \quad NY = NY0 + n * LIM .$$

Acest vector este constituit din două părți:



Tablourile Y1 și YA

- Sub-vectorul  $Y1(1 : NY0)$ , de dimensiune  $LIM * (LIM + 1) / 2$ : servește ca front de lucru, pentru descompunerea Cholesky; în acestea se generează elementele triunghiului activ la un pas (triunghi cu laturile  $LIM$ ).

- Sub-vectorul  $YA(NY0+1:NY)$ , de dimensiune  $n * LIM$  : inițial, stochează matricea  $A$ , pe linii. În cursul descompunerii Cholesky, stochează liniile procesate ale matricii  $A$ .

## II. METODE DE ITERATIVE

### 1 Metoda JACOBI

*Exemplu numeric*

Fie sistemul:

$$8x_1 + x_2 - x_3 = 8$$

$$2x_1 + x_2 + 9x_3 = 12$$

$$x_1 - 7x_2 + 2x_3 = -4$$

Rezolvăm fiecare ecuație  $i = 1, 2, 3$ , în raport cu necunoscuta  $x_i$ , căutând ca aceasta să fie necunoscuta cu coeficientul cel mai mare din ecuație, eventual rearanjând ecuațiile. Intervertind ecuațiile 2 și 3, avem:

$$x_1 = 1 - (1/8)x_2 + (1/8)x_3$$

$$x_2 = 4/7 + (1/7)x_1 + (2/7)x_3$$

$$x_3 = 12/9 - (2/9)x_1 - (1/9)x_2$$

Sau matriceal:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 4/7 \\ 12/9 \end{bmatrix} + \begin{bmatrix} 0 & -1/8 & 1/8 \\ 1/7 & 0 & 2/7 \\ -2/9 & -1/9 & 0 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad (1)$$

care este de forma

$$\mathbf{x}^{(m+1)} = \mathbf{g} + \mathbf{M}\mathbf{x}^{(m)}$$

Se iterează, cu aproximația inițială  $\mathbf{x}^{(0)} = [1 \quad 4/7 \quad 12/9]^T$ ; testul de oprire a iterației este  $\|\mathbf{x}^{(m+1)} - \mathbf{x}^{(m)}\|_{\infty} \leq 1E-6$ . La iterația 12 rezultă soluția (1.0, 1.0, 1.0). Coordonata de modul maxim în  $\mathbf{Ax} - \mathbf{b}$  (rezidualul maxim) este 1.43E-06 (ecuația 3).

■

### Metoda

Fie sistemul dat  $\mathbf{Ax} = \mathbf{b}$ . Se rezolvă fiecare ecuație  $i$  în raport cu necunoscuta  $x_i$ .

Explicitând  $x_i$  se obține:

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

S-a presupus  $a_{ii} \neq 0$ . Iterația Jacobi va fi:

$$\begin{cases} \mathbf{x}^{(0)} = \text{arbitrar} \\ x_i^{(m+1)} = (b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(m)}) / a_{ii}, \quad i = 1, 2, \dots, n; \quad m \geq 0 \end{cases}$$

Testul de oprire a iterației este  $\|\mathbf{x}^{(m+1)} - \mathbf{x}^{(m)}\|_{\infty} \leq \text{eps}$ .

Metoda Jacobi se mai zice metoda “iterațiilor simultane”, pentru că, coordonatele  $x_i$  ale soluției  $\mathbf{x}$  se calculează independent unele de altele. (Pentru alt mod de calcul – v. metoda Metoda Gauss\_Seidel, mai jos.)

*Condiție suficientă de convergență:*

Matricea  $\mathbf{A}$  este *diagonal dominantă*, adică avem:

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|, \quad i = \overline{1, n}.$$

■

## 2 Metoda Gauss\_Seidel

Se modifică metoda Jacobi, astfel că, la calculul coordonatei  $x_i^{(m)}$  se utilizează valorile  $x_1^{(m)}, \dots, x_{i-1}^{(m)}$  deja calculate, care în general, sunt aproximații mai bune ale soluției.

Pentru [exemplul anterior](#):

$$x_1^{(1)} = 1 - (1/8)x_2^{(0)} + (1/8)x_3^{(0)}$$

$$x_2^{(1)} = 4/7 + (1/7)x_1^{(1)} + (2/7)x_3^{(0)}$$

$$x_3^{(1)} = 12/9 - (2/9)x_1^{(1)} - (1/9)x_2^{(1)}$$

Cu  $\mathbf{x}^{(0)} = [0 \ 0 \ 0]^T$  și aceeași toleranță  $1E-6$ , la iterația 9, se obține soluția  $(1.0, 1.0, 1.0)$ . Rezidualul maxim este 0.0.

■

Formulele generale ale metodei Gauss-Seidel sunt:

$$\begin{cases} \mathbf{x}^{(0)} = \text{arbitrar} \\ x_i^{(m+1)} = (b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(m)}) / a_{ii}, \quad i = 1, 2, \dots, n; \quad m \geq 0 \end{cases}$$

Testul de oprire a iterației este  $\|\mathbf{x}^{(m+1)} - \mathbf{x}^{(m)}\|_{\infty} \leq \text{eps}$ .

Metoda Gauss-Seidel se mai zice metoda “iterațiilor succesive”, pentru că la un pas  $m+1$ ,  $m \geq 0$ , de îndată ce o coordonată  $x_j^{(m+1)}$  a soluției este calculată, ea se utilizează în ecuațiile pentru coordonatele următoare  $x_i^{(m+1)}$ ,  $i > j$ .

*Condiții suficiente de convergență:*

- Matricea  $\mathbf{A}$  este *diagonal dominantă*.
- Matricea  $\mathbf{A}$  este *simetrică și pozitiv definită*.

Când ambele metode Jacobi și Gauss-Seidel converg, metoda Gauss-Seidel converge mai rapid. Pentru alte considerații privind convergența, v. Cap. 4-IV.

### 3 SOR - Metoda relaxării (Successive Over-Relaxation).

Reluăm formula metodei [Gauss-Seidel](#):

$$x_i^{(m+1)} = (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(m)}) / a_{ii}, \quad i = \overline{1, 2, \dots, n}; \quad m \geq 0$$

Formula se pune sub forma următoare, adunând și scăzând  $x_i^{(m)}$  în membrul doi (observați că acum, în a doua sumă, indicele  $j$  ia valori de la  $i$  și nu de la  $i+1$ ):

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

Termenul care se adună la  $x_i^{(m)}$  este diferența  $(x_i^{m+1} - x_i^m)$ . Metoda SOR constă în înmulțirea acestei diferențe cu un factor de accelerare (sau relaxare)  $\omega > 1$ . Întrucât  $\omega > 1$ , metoda se zice *supra-relaxare*. Alegerea lui  $\omega$  se discută mai jos. Formula metodei SOR este deci

$$x_i^{(m+1)} = x_i^{(m)} + \omega (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i}^n a_{ij} x_j^{(m)}) / a_{ii}, \quad i = \overline{1, n}$$

Explicitând  $x_i^{(m)}$  din a doua sumă, rezultă:

$$x_i^{(m+1)} = \omega (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(m)}) / a_{ii} + (1 - \omega) x_i^{(m)}, \quad i = \overline{1, n}$$

Se notează expresia din prima paranteză cu  $z_i^{(m+1)}$  (aceasta este coordonata  $i$  a iteratei  $(m+1)$  din metoda Gauss-Seidel).

Formula de iterare este echivalentă cu următoarele ecuații considerate pentru  $i = \overline{1, n}$ :

$$z_i^{(m+1)} = (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(m)}) / a_{ii},$$

$$x_i^{(m+1)} = \omega z_i^{(m+1)} + (1 - \omega) x_i^{(m)}$$

În cele mai multe cazuri, valoarea optimă a lui  $\omega$  satisface relația

$$1 < \omega < 2.$$

În practică, se poate alege  $\omega$  astfel: se utilizează valori  $\omega$  de test, pe un număr limitat de iterații; se alege ca valoare optimă acel  $\omega$ , pentru care convergența este cea mai rapidă.

### III. STABILITATEA SOLUȚIEI ȘI CONDIȚIONAREA SISTEMULUI

Se consideră sistemul de  $n$  ecuații liniare

$$\mathbf{Ax} = \mathbf{b} \quad (1)$$

cu  $\mathbf{A}$  = nesingulară. Se va analiza stabilitatea soluției sistemului, în raport cu o mică variație (perturbare), în:

- Termenii liberi  $\mathbf{b}$
- Matricea  $\mathbf{A}$  a sistemului

Prima problemă conduce la numărul de condiție al matricii  $\mathbf{A}$ . A doua problemă va conduce la o evaluare a variației soluției, în care intervine numărul de condiție.

#### 1 Perturbare în $\mathbf{b}$ . Număr de condiție

Fie sistemul (1), în care termenii liberi  $\mathbf{b}$  suferă o perturbare  $\mathbf{r}$ , devenind  $\tilde{\mathbf{b}}$ , unde

$\mathbf{r} = \tilde{\mathbf{b}} - \mathbf{b}$ , și  $\|\mathbf{r}\| \ll \|\mathbf{b}\|$ . Sistemul perturbat va fi:

$$\mathbf{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}} \quad (2)$$

Notând perturbarea soluției prin  $\mathbf{e} = \tilde{\mathbf{x}} - \mathbf{x}$ , și scăzând (1) din (2) rezultă:

$$\mathbf{A}\mathbf{e} = \mathbf{r}, \quad (3)$$

din care avem și:

$$\mathbf{e} = \mathbf{A}^{-1}\mathbf{r} \quad (3')$$

Pentru a examina stabilitatea soluției  $\mathbf{x}$ , căutăm o evaluare a raportului

$$\frac{\|\mathbf{e}\| / \|\mathbf{x}\|}{\|\mathbf{r}\| / \|\mathbf{b}\|} = \frac{\text{perturbarea relativă în } \mathbf{x}}{\text{perturbarea relativă în } \mathbf{b}}$$



Acest raport exprimă efectul perturbării în  $\mathbf{b}$  asupra soluției  $\mathbf{x}$ , și anume: După cum raportul este  $\approx 1$ , respectiv  $\gg 1$ , perturbarea relativă a soluției este de același ordin, respectiv de ordin mult mai mare, în raport cu perturbarea relativă în  $\mathbf{b}$ .

Luând norma în (3), (3') avem:

$$\|\mathbf{r}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{e}\|, \quad \|\mathbf{e}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{r}\|,$$

din care, rezultă:

$$\frac{1}{\|\mathbf{A}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{r}\|} \leq \|\mathbf{A}^{-1}\| \quad (4)$$

Înmulțind (4) cu  $\|\mathbf{b}\| / \|\mathbf{x}\|$ , rezultă:

$$\frac{1}{\|\mathbf{A}\|} \cdot \frac{\|\mathbf{b}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{e}\| / \|\mathbf{x}\|}{\|\mathbf{r}\| / \|\mathbf{b}\|} \leq \|\mathbf{A}^{-1}\| \cdot \frac{\|\mathbf{b}\|}{\|\mathbf{x}\|} \quad (5)$$

Analog, luând norma în (1) și în  $\mathbf{x} = \mathbf{A}^{-1} \cdot \mathbf{b}$ , rezultă:

$$\frac{1}{\|\mathbf{A}^{-1}\|} \leq \frac{\|\mathbf{b}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \quad (6)$$

Din (5) și (6) rezultă:

$$\frac{1}{\|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|} \leq \frac{\|\mathbf{e}\| / \|\mathbf{x}\|}{\|\mathbf{r}\| / \|\mathbf{b}\|} \leq \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| \quad (7)$$

Introducem următoarea

### **Definiție**

*Numărul de condiție* al matricii  $\mathbf{A}$ , este:

$$\text{Cond}(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| \quad \blacksquare \quad (8)$$

Numărul de condiție depinde de normă, dar este *mărginit inferior* de 1:

$$\text{Cond}(\mathbf{A}) \geq 1 \quad (9)$$

Într-adevăr, din  $\mathbf{I} = \mathbf{A}\mathbf{A}^{-1}$  rezultă:  $1 \leq \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| = \text{Cond}(\mathbf{A}) \quad \blacksquare$

Cu definiția (8), din (7) rezultă:

$$\frac{1}{\text{Cond}(\mathbf{A})} \leq \frac{\|\mathbf{e}\| / \|\mathbf{x}\|}{\|\mathbf{r}\| / \|\mathbf{b}\|} \leq \text{Cond}(\mathbf{A}) \quad (10)$$

Sau:

$$\frac{1}{\text{Cond}(\mathbf{A})} \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \text{Cond}(\mathbf{A}) \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \quad (10')$$

Din această relație rezultă că, pentru  $\|\mathbf{r}\| / \|\mathbf{b}\| = \text{'mic'}$ , avem:

- Dacă  $\text{Cond}(\mathbf{A}) \sim 1$ :  $\|\mathbf{e}\| / \|\mathbf{x}\|$  este 'mic', și  $\mathbf{A}$  este *bine-condiționată*.
- Dacă  $\text{Cond}(\mathbf{A}) \gg 1$ :  $\|\mathbf{e}\| / \|\mathbf{x}\|$  poate fi 'mare';  $\mathbf{A}$  este *rău-condiționată*.

Întrucât în Definiția (8)  $\text{Cond}(\mathbf{A})$  depinde de normă, se utilizează și un alt număr de condiție care este independent de normă. Avem  $\|\mathbf{A}\| \geq \rho(\mathbf{A})$ , unde  $\rho(\mathbf{A})$  este raza spectrală a matricii  $\mathbf{A}$ ; urmează:

$$\text{Cond}(\mathbf{A}) \geq \rho(\mathbf{A}) \cdot \rho(\mathbf{A}^{-1})$$

Sau, definind

$$\text{Cond}(\mathbf{A})_* = \rho(\mathbf{A}) \cdot \rho(\mathbf{A}^{-1}) \quad (11)$$

avem

$$\text{Cond}(\mathbf{A}) \geq \text{Cond}(\mathbf{A})_*$$

Cum valorile proprii ale matricii  $\mathbf{A}^{-1}$  sunt inversele valorilor proprii ale lui  $\mathbf{A}$ , notând cu  $\sigma(\mathbf{A})$  spectrul matricii  $\mathbf{A}$ , rezultă formula de calcul:

$$\text{Cond}(\mathbf{A})_* = \frac{\max_{\lambda \in \sigma(\mathbf{A})} |\lambda|}{\min_{\lambda \in \sigma(\mathbf{A})} |\lambda|} \quad (11')$$

*Observații asupra calculului lui  $\text{Cond}(\mathbf{A})_2$*

După definiție, avem

$$\text{Cond}(\mathbf{A})_2 = \|\mathbf{A}\|_2 \cdot \|\mathbf{A}^{-1}\|_2$$

Cu  $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ , valorile proprii ale lui  $\mathbf{B}$  sunt reale și pozitive, și avem

$$\|\mathbf{A}\|_2 = (\lambda_{B,\max})^{1/2} \quad (\text{v. 4-I, 3.3 – Observații și Cap. 5}).$$

Analog, cu  $\mathbf{D} = \mathbf{A}^{-T} \mathbf{A}^{-1}$ , avem  $\|\mathbf{A}^{-1}\|_2 = (\lambda_{D,\max})^{1/2}$ , și rezultă

$$\text{Cond}(\mathbf{A})_2 = (\lambda_{B,\max})^{1/2} (\lambda_{D,\max})^{1/2}.$$

**Putem evita inversarea matricii  $\mathbf{A}$** , conform următoarelor proprietăți:

- Valorile proprii ale lui  $\mathbf{D}$  sunt inversele valorilor proprii ale lui  $\mathbf{D}^{-1} = \mathbf{A}\mathbf{A}^T$ .

Punem  $\mathbf{C} = \mathbf{A}\mathbf{A}^T$ , și astfel,  $\lambda_{D,\max} = 1/\lambda_{C,\min}$ . Expresia anterioară devine:

$$\text{Cond}(\mathbf{A})_2 = \frac{(\lambda_{B,\max})^{1/2}}{(\lambda_{C,\min})^{1/2}}$$

- Apoi, se arată ușor că  $\mathbf{B} = \mathbf{A}^T \mathbf{A}$  și  $\mathbf{C} = \mathbf{A}\mathbf{A}^T$  au *același valori proprii* (Exercițiu).

Urmează că, în final avem:

$$\text{Cond}(\mathbf{A})_2 = \left( \frac{\lambda_{B,\max}}{\lambda_{B,\min}} \right)^{1/2}$$

Aceasta se mai scrie:  $\text{Cond}(\mathbf{A})_2 = (\text{Cond}(\mathbf{B})_*)^{1/2}$  ■

### Exemplu – 1

Fie sistemul liniar:

$$8x_1 + 9x_2 = b_1$$

$$7x_1 + 8x_2 = b_2$$

Avem:

$$\mathbf{A} = \begin{bmatrix} 8 & 9 \\ 7 & 8 \end{bmatrix}, \quad \mathbf{A}^{-1} = \begin{bmatrix} 8 & -9 \\ -7 & 8 \end{bmatrix}$$

Numerele de condiție, după diferite norme, sunt:

$$\text{Cond}(\mathbf{A})_\infty = \text{Cond}(\mathbf{A})_1 = 289,$$

$$\text{Cond}(\mathbf{A})_2 \approx 258,$$

$$\text{Cond}(\mathbf{A})_* \approx 254$$

a) Detalii pentru calculul lui  $\text{Cond}(\mathbf{A})_*$ :

Polinomul caracteristic al lui  $\mathbf{A}$ , este  $p(\lambda) = \lambda^2 - 16\lambda + 1$ , de unde  $\lambda_{1,2} = 8 \pm \sqrt{63}$ .

Urmează  $\text{Cond}(\mathbf{A})_* = \lambda_1 / \lambda_2 \approx 254$ .

b) Detalii pentru calculul lui  $\text{Cond}(\mathbf{A})_2$ :

$$\mathbf{B} = \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 8 & 7 \\ 9 & 8 \end{bmatrix} \begin{bmatrix} 8 & 9 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 113 & 128 \\ 128 & 145 \end{bmatrix}. \text{ Polinomul caracteristic al lui } \mathbf{B}, \text{ este}$$

$$p(\lambda) = \lambda^2 - 258\lambda + 1. \text{ Avem } \lambda_{1,2} = 129 \pm \sqrt{129^2 - 1}, \lambda_1 / \lambda_2 = (\lambda_1)^2. \text{ Rezultă}$$

$$\text{Cond}(\mathbf{A})_2 = \lambda_1 \approx 258.$$

Se verifică, în particular, observația anterioară:  $\mathbf{C} = \mathbf{A}\mathbf{A}^T = \begin{bmatrix} 145 & 128 \\ 128 & 113 \end{bmatrix}$ , are același

polinom caracteristic ca și  $\mathbf{B}$ .

Exercițiu: Calculați  $\text{Cond}(\mathbf{A})_2$ , direct după definiția (8).

■

Întrucât numărul de condiție este mare, aceasta arată că sistemul este sensibil la mici schimbări în  $\mathbf{b}$ . Într-adevăr, fie de exemplu, termenii liberi dați și perturbați, și soluțiile respective, cum urmează:

$$\mathbf{b} = \begin{bmatrix} 17 \\ 15 \end{bmatrix}; \quad \tilde{\mathbf{b}} = \begin{bmatrix} 16.9 \\ 15.1 \end{bmatrix}; \quad \mathbf{r} = \tilde{\mathbf{b}} - \mathbf{b} = \begin{bmatrix} -0.1 \\ 0.1 \end{bmatrix}.$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}; \quad \tilde{\mathbf{x}} = \begin{bmatrix} -0.7 \\ 2.5 \end{bmatrix}; \quad \mathbf{e} = \tilde{\mathbf{x}} - \mathbf{x} = \begin{bmatrix} -1.7 \\ 1.5 \end{bmatrix}.$$

Cu acestea, rezultă:

$$\frac{\|\mathbf{r}\|_\infty}{\|\mathbf{b}\|_\infty} = \frac{0.1}{17}, \quad \frac{\|\mathbf{e}\|_\infty}{\|\mathbf{x}\|_\infty} = \frac{1.7}{1} = 1.7, \quad \frac{\|\mathbf{e}\|_\infty / \|\mathbf{x}\|_\infty}{\|\mathbf{r}\|_\infty / \|\mathbf{b}\|_\infty} = \frac{1.7}{0.1/17} = 289.$$

Se observă că schimbări “mici” în  $\mathbf{b}$  produc schimbări ”mari” în  $\mathbf{x}$ : perturbarea relativă în  $\mathbf{x}$  este de 289 ori mai mare decât perturbarea relativă în  $\mathbf{b}$ .

În particular, în acest exemplu, marginea superioară din (10) este atinsă – conform  $Cond(\mathbf{A})_{\infty} = 289$ .

Un asemenea sistem se zice **rău-condiționat**. Se observă că numerele  $Cond(\mathbf{A})$  și  $Cond(\mathbf{A})_*$  sunt o măsură bună pentru condiționarea sistemului.

### Exemplu – 2

Există sisteme care nu sunt rău-condiționate, deși  $Cond(\mathbf{A})$  este mare. De exemplu, să considerăm următoarea matrice, în care  $m$  este un întreg  $\geq 1$ :

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 10^{-m} \end{bmatrix}; \quad \mathbf{A}^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 10^m \end{bmatrix}.$$

Avem  $Cond(\mathbf{A})_{1,\infty} = Cond(\mathbf{A})_* = 10^m$ . Totuși sistemul este bine-condiționat.

Valoarea mare a numărului de condiție se elimină prin scalarea matricii  $\mathbf{A}$ .

### Exemplu – 3: Interpretare geometrică a condiționării sistemului

Considerăm sistemul următor:

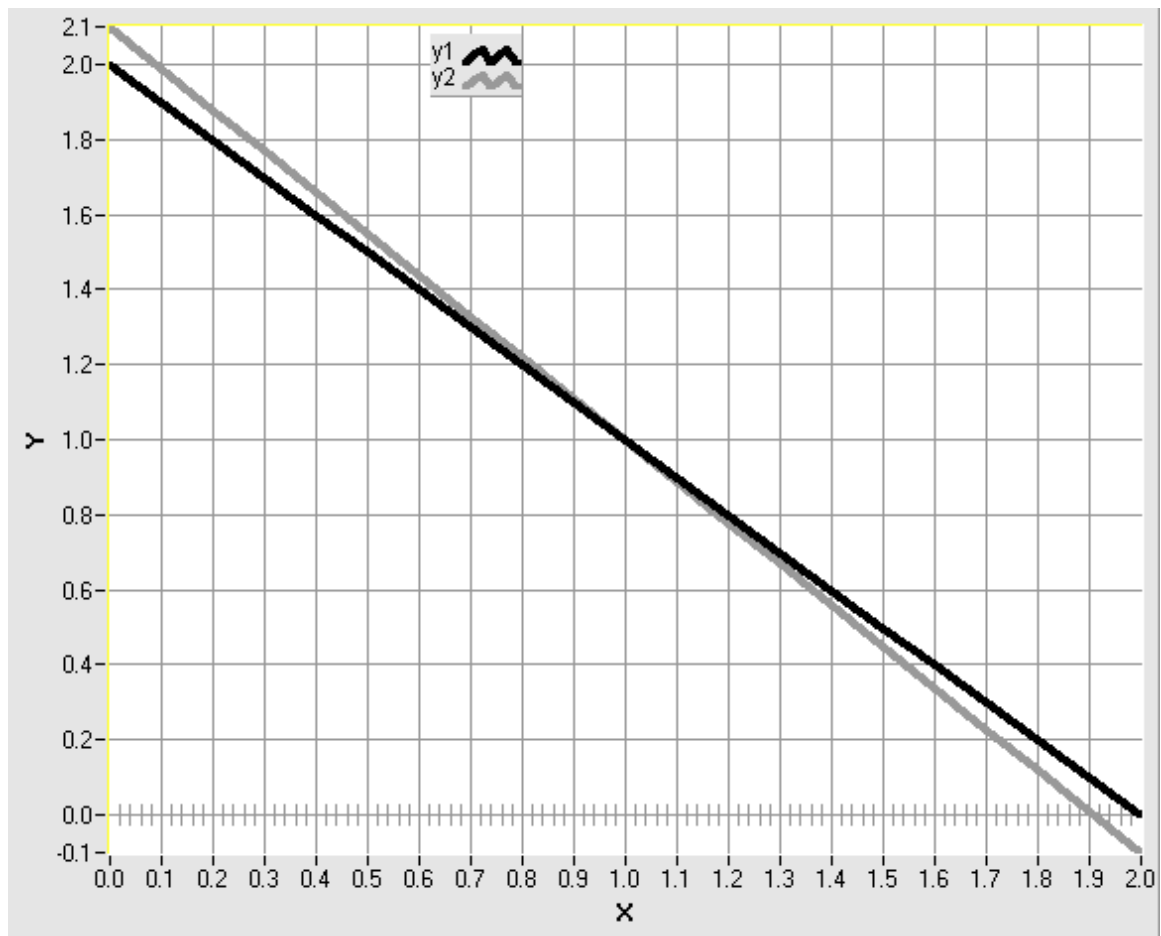
$$x + y = 2$$

$$1.1x + y = 2.1$$

care are soluția  $x = y = 1.0$ . Numerele de condiție sunt:

$$Cond(\mathbf{A})_1 = Cond(\mathbf{A})_{\infty} = 44.10, \quad Cond(\mathbf{A})_2 \approx 42.08, \quad Cond(\mathbf{A})_* \approx 41.98,$$

astfel că sistemul este rău-condiționat. Rezolvarea sistemului revine la găsirea intersecției dreptelor reprezentate de cele două ecuații.



Sistem rău condiționat: Interpretare geometrică

Pantele celor două drepte sunt respectiv  $-1.0$  și  $-1.1$ , adică apropiate. Dacă considerăm o incertitudine în coeficienții sistemului, banda de incertitudine a rădăcinii va fi ‘largă’ – și cu atât mai largă cu cât pantele sunt mai apropiate. Vezi Figura de mai sus: incertitudinea în valorile  $y(x)$  este reprezentată de linia ‘groasă’ a graficului. În schimb, dacă pantele sunt diferite, intersecția dreptelor este netă și banda de incertitudine este ‘îngustă’.

#### Exemplu – 4: Matricea Hilbert

Un exemplu de matrice rău-condiționată este următoarea matrice numită matricea Hilbert:

$$\mathbf{H}_n = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \cdots & \frac{1}{n+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n} & \frac{1}{n+1} & \frac{1}{n+2} & \cdots & \frac{1}{2n-1} \end{bmatrix}$$

Inversa matricii  $\mathbf{H}_n$  este cunoscută analitic, și anume: punând  $\mathbf{H}_n^{-1} = [\alpha_{ij}^{(n)}]$ , avem –  
v. Atkinson (1978), Ralston & Rabinowitz (1978):

$$\alpha_{ij}^{(n)} = (-1)^{i+j} \frac{(n+i-1)!(n+j-1)!}{(i+j-1)[(i-1)!(j-1)!]^2 (n-i)!(n-j)!}, \quad \overline{i, j} = \overline{1, n}$$

Numărul de condiție al matricii  $\mathbf{H}_n$  crește cu  $n$ , matricea fiind cu atât mai rău condiționată cu cât  $n$  este mai mare. Exemple:

$n$	$\text{Cond}(\mathbf{H}_n)_1$
3	7.48 E+02
4	2.84 E+04
5	9.44 E+05
6	2.91 E+07
7	9.85 E+08

Ca exemplu, calculând inversa matricii pentru  $n = 4$ , prin eliminare Gauss, cu elementele lui  $\mathbf{H}_4$  reprezentate în simplă precizie, se obține:

$$\hat{\mathbf{H}}_4^{-1} = \begin{bmatrix} 15.99979 & -119.9980 & 239.9954 & -139.9971 \\ -119.9979 & 1199.981 & -2699.957 & 1679.974 \\ 239.9954 & -2699.958 & 6479.907 & -4199.943 \\ -139.9972 & 1679.974 & -4199.943 & 2799.965 \end{bmatrix}$$

Inversa calculată analitic  $\mathbf{H}_4^{-1}$ , are elemente întregi:  $\alpha_{11} = 16, \dots, \alpha_{44} = 2800$ .

Calculând numărul de condiție al matricii  $\mathbf{H}_4$ , cu 1-norma, se obține:

$$\text{Cond}(\mathbf{H}_4)_1 = \|\mathbf{H}_4\|_1 \cdot \|\mathbf{H}_4^{-1}\|_1 = (25/12) \cdot 13260 = 28375$$

## 2 Perturbare în $\mathbf{A}$ și $\mathbf{b}$

Să presupunem că atât matricea  $\mathbf{A}$  a sistemului, cât și termenul liber  $\mathbf{b}$ , suferă mici schimbări  $\delta\mathbf{A}$ , respectiv,  $\delta\mathbf{b}$ , iar soluția devine  $\mathbf{x} + \delta\mathbf{x}$ :

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$$

Avem următoarea

### **Teoremă**

Presupunem  $\mathbf{A}$  nesingulară și fie perturbarea  $\delta\mathbf{A}$  satisfăcând condiția

$$\|\delta\mathbf{A}\| < 1 / \|\mathbf{A}^{-1}\|.$$

Atunci:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\text{Cond}(\mathbf{A})}{1 - \text{Cond}(\mathbf{A}) \cdot \|\delta\mathbf{A}\| / \|\mathbf{A}\|} \left( \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \right) \quad (12)$$

■

Pentru demonstrație – v. Atkinson (1978), Ralston & Rabinowitz (1978).

În particular, dacă  $\delta\mathbf{b} = 0$ , atunci efectul unei perturbări în  $\mathbf{A}$  este dat de:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\text{Cond}(\mathbf{A})}{1 - \text{Cond}(\mathbf{A}) \cdot \|\delta\mathbf{A}\| / \|\mathbf{A}\|} \cdot \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \quad (13)$$

■

În ceea ce urmează vom da rezultate privind efectul erorilor de rotunjire asupra soluției calculate prin eliminarea Gauss.

Rezultatul (13) va fi utilizat în estimarea *a priori* a erorii în eliminarea Gauss.