

# **Calcul Numeric**

**Cursul 3**

**2017**

*Anca Ignat*

## Surse de erori în calculule numerice

### 1. Erori în datele de intrare:

- măsurători afectate de erori sistematice sau perturbații temporare,
- erori de rotunjire:  $1/3$  ,  $\pi$  ,  $1/7, \dots$

### 2. Erori de rotunjire în timpul calculelor:

- datorate capacității limitate de memorare a datelor, operațiile nu sunt efectuate exact.

### 3. Erori de discretizare:

- limita unui șir , suma unei serii , funcții neliniare approximate de funcții liniare, aproximarea derivatei unei funcții

### 4. Simplificări în modelul matematic

- idealizări , ignorarea unor parametri.

### 5. Erori umane și erori ale bibliotecilor folosite.

## Eroare absolută , eroare relativă

$a$  – valoarea exactă,

$\tilde{a}$  – valoarea aproximativă.

***Eroare absolută :***  $a - \tilde{a}$  sau  $|a - \tilde{a}|$  sau  $\|a - \tilde{a}\|$

$$a = \tilde{a} \pm \Delta_a, |a - \tilde{a}| \leq \Delta_a$$

***Eroare relativă:***  $a \neq 0$   $\frac{a - \tilde{a}}{a}$  sau  $\frac{|a - \tilde{a}|}{|a|}$  sau  $\frac{\|a - \tilde{a}\|}{\|a\|}$

$$\frac{|a - \tilde{a}|}{|a|} \leq \delta_a \quad (\delta_a \text{ se exprimă de regulă în } \%).$$

În aproximările  $1\text{kg} \pm 5\text{g}$ ,  $50\text{g} \pm 5\text{g}$  erorile absolute sunt egale dar pentru prima cantitate eroarea relativă este 0,5% iar pentru a doua eroarea relativă este 10%.

$$a_1 = \tilde{a}_1 \pm \Delta_{a_1}, a_2 = \tilde{a}_2 \pm \Delta_{a_2},$$

$$a_1 \pm a_2 = (\tilde{a}_1 \pm \tilde{a}_2) \pm (\Delta_{a_1} \pm \Delta_{a_2})$$

$$\Delta_{a_1+a_2} \leq \Delta_{a_1} + \Delta_{a_2}.$$

$a_1$  cu eroare relativă  $\delta_{a_1}$  și  $a_2$  cu eroare relativă  $\delta_{a_2}$  :

$$a = a_1 * a_2 \text{ sau } \frac{a_1}{a_2} \text{ rezultă } \delta_a = \delta_{a_1} + \delta_{a_2}.$$

## Condiționare $\leftrightarrow$ stabilitate

Condiționarea unei probleme caracterizează sensibilitatea soluției în raport cu perturbarea datelor de intrare, în ipoteza unor calcule exacte (independent de algoritmul folosit pentru rezolvarea problemei).

Fie  $\mathbf{x}$  datele exacte de intrare,  $\tilde{\mathbf{x}}$  o aproximație cunoscută a acestora,  $\mathbf{P}(\mathbf{x})$  soluția exactă a problemei și  $\mathbf{P}(\tilde{\mathbf{x}})$  soluția problemei cu  $\tilde{\mathbf{x}}$  ca date de intrare. Se presupune că s-au făcut calcule exacte la obținerea soluțiilor  $\mathbf{P}(\mathbf{x})$  și  $\mathbf{P}(\tilde{\mathbf{x}})$ .

O problemă se consideră a fi *prost condiționată* dacă  $P(\mathbf{x})$  și  $P(\tilde{\mathbf{x}})$  diferă mult chiar dacă eroarea relativă  $\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|}$  este mică.

Condiționarea numerică a unei probleme este exprimată prin amplificarea erorii relative:

$$k(\mathbf{x}) = \frac{\frac{\|P(\mathbf{x}) - P(\tilde{\mathbf{x}})\|}{\|P(\mathbf{x})\|}}{\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|}} \quad \text{pentru } \mathbf{x} \neq \mathbf{0} \text{ și } P(\mathbf{x}) \neq \mathbf{0}$$

O valoare mică pentru  $k(\mathbf{x})$  caracterizează o problemă bine-condiționată.

Condiționarea este o proprietate locală (se evaluează pentru diverse date de intrare  $\mathbf{x}$ ). O problemă este bine-condiționată dacă este bine-condiționată în orice punct.

Pentru rezolvarea unei probleme  $P$ , calculatorul execută un algoritm  $\tilde{P}$ . Deoarece se folosesc numere în virgulă mobilă, calculele sunt afectate de erori:

$$P(\mathbf{x}) \neq \tilde{P}(\mathbf{x})$$



*Stabilitatea numerică* exprimă mărimea erorilor numerice introduse de algoritm, în ipoteza unor date de intrare exacte,

$$\|P(x) - \tilde{P}(x)\| \text{ sau } \frac{\|P(x) - \tilde{P}(x)\|}{\|P(x)\|}.$$

O eroare relativă de ordinul erorii de rotunjire caracterizează un *algoritm numeric stabil*.

Un **algoritm numeric stabil** aplicat unei **probleme bine condiționate** conduce la **rezultate cu precizie foarte bună**.

Un algoritm  $\tilde{P}$  destinat rezolvării problemei  $P$  este numeric stabil dacă este îndeplinită una din condițiile:

1.  $\tilde{P}(\mathbf{x}) \approx P(\mathbf{x})$  pentru orice intrare  $\mathbf{x}$ ;

2. există  $\tilde{\mathbf{x}}$  apropiat de  $\mathbf{x}$ , astfel ca  $\tilde{P}(\mathbf{x}) \approx P(\tilde{\mathbf{x}})$

$\mathbf{x}$  = datele exacte,

$P(\mathbf{x})$  = soluția exactă folosind date exacte,

$\tilde{P}(\mathbf{x})$  = soluția „*calculată*” folosind algoritmul  $\tilde{P}$  cu date  
exacte de intrare

# Rezolvarea sistemelor liniare

## Istoric

- 1900 î.Hr., Babilon - apar primele probleme legate de ecuații liniare simultane
- 300 î.Hr. Babilon - tăbliță cu următoarea problemă: *”Avem două câmpuri de arie totală 1800 ha. Producția la hectar pe primul câmp este de  $\frac{2}{3}$  bușel (=36,3l) iar pe al doilea este de  $\frac{1}{2}$  bușel. Dacă producția totală este de 1100 bușeli, să se determine aria fiecărui teren în parte.*

- 200-100 î.Hr. China – *9 capitole despre arta matematică* – metodă de rezolvare foarte asemănătoare eliminării Gauss  
(„Avem 3 tipuri de grâu. Știm că 3 baloturi din primul tip, 2 baloturi din al doilea tip și 1 balot din al treilea tip cântăresc 39 măsuri. Deasemenea, 2 baloturi din primul tip, 3 baloturi din al doilea tip și 1 balot din al treilea tip cântăresc 34 măsuri și 1 balot din primul tip, 2 baloturi din al doilea tip și 3 baloturi din al treilea tip cântăresc 26 măsuri. Câte măsuri cântărește un balot din fiecare tip de grâu”)

- 1545, Cardan – în *Ars Magna*, propune o regulă (*regula de modo*) pentru rezolvarea unui sistem de 2 ecuații cu 2 necunoscute (seamănă cu regula lui Cramer)
- 1683, Seki Kowa, Japonia - ideea de „*determinant*”- „*Method of solving the dissimulated problems*”. Calculează ceea ce astăzi cunoaștem sub numele de determinant, determinanții matricilor  $2 \times 2$ ,  $3 \times 3$ ,  $4 \times 4$ ,  $5 \times 5$  în legătură cu rezolvarea unor ecuații dar nu a sistemelor de ecuații.

- 1683, Leibniz într-o scrisoare către l'Hôpital explică faptul că sistemul de ecuații:

$$10 + 11x + 12y = 0$$

$$20 + 21x + 22y = 0$$

$$30 + 31x + 32y = 0$$

are soluție deoarece :

$$10 \cdot 21 \cdot 32 + 11 \cdot 22 \cdot 30 + 12 \cdot 20 \cdot 31 = 10 \cdot 22 \cdot 31 + 11 \cdot 20 \cdot 32 + 12 \cdot 21 \cdot 30$$

(condiția ca determinantul matricii coeficienților este 0).

Leibniz era convins că o notație matematică bună este cheia progresului și experimentează mai mult de 50 de moduri diferite de a scrie coeficienții unui sistem de ecuații. Leibniz folosește termenul de „*rezultant*” în loc de determinant și a demonstrat regula lui Cramer pentru „rezultanți”. Știa că orice determinant poate fi dezvoltat în raport cu o coloană – operația se numește azi dezvoltarea Laplace.

- 1750, Cramer prezintă o formulă bazată pe determinanți pentru rezolvarea unui sistem de ecuații liniare – *regula lui Cramer* – „*Introduction in the analysis of algebraic curves*” (dă o regulă generală pentru sisteme  $n \times n$ :  
*„One finds the value of each unknown by forming  $n$  fractions of which the common denominator has as many terms as there are permutations of  $n$  things”*
- 1764 Bezout, 1771 Vandermonde, 1772 Laplace – reguli de calcul al determinanților



- 1773 Lagrange – prima utilizare implicită a matricilor în legătură cu formele biliniare ce apar la optimizarea unei funcții reale de 2 sau mai multe variabile (dorea să caracterizeze punctele de maxim și minim a funcțiilor de mai multe variabile)

- 1800-1801, Gauss introduce noțiunea de „*determinant*” (determină proprietățile forme pătratică) – *Disquisitiones arithmeticae*(1801); descrie operațiile de înmulțire matricială și inversă a unei matrici în contextul tabloului coeficienților unei forme pătratică. Gauss dezvoltă *eliminarea Gaussiană* pe când studia orbita asteroidului Pallas de unde obține un sistem liniar cu 6 ecuații cu 6 necunoscute.

- 1812, Cauchy folosește termenul de „*determinant*” în sensul cunoscut azi.
- 1826, Cauchy găsește valorile proprii și deduce rezultate legate de diagonalizarea unei matrici. Introduce noțiunea de matrici asemenea și demonstrează ca acestea au aceeași ecuație caracteristică. Demonstrează că orice matrice reală simetrică este diagonalizabilă.

- 1850, Sylvester introduce pentru prima data termenul de *matrice* (din latină, „uter” – un loc unde ceva se formează sau este produs, „*an oblong arrangement of terms*”)
- 1855, Cayley – algebră matricială, prima definiție abstarctă a unei matrici. Studiază transformările liniare și compunerea lor ceea ce îl conduce la operațiile cu matrici (adunare, înmulțire, înmulțirea cu un scalar, inversa)

- 1858, Cayley în *Memoriu asupra teoriei matricilor* : „*Sunt multe lucruri de spus despre această teorie a matricilor și, după părerea mea, această teorie ar trebui să preceadă teoria determinanților*”
- Jordan (1870 – Treatise on substitutions and algebraic equations – forma canonică Jordan), Frobenius (1878 – On linear substitutions and bilinear forms, rangul unei matrici), Peano

- 1890, Weierstrass – On determinant theory, definiția axiomatică a determinantului
- 1925, Heisenberg reinventează algebra matricială pentru mecanica cuantică
- 1947, vonNeuman & Goldstine introduc numerele de condiționare atunci când analizează erorile de rotunjire
- 1948, Turing introduce descompunerea  $LU$  a unei matrici
- 1958, Wilkinson dezvoltă factorizarea  $QR$

....

## **Evaluarea erorii în rezolvarea sistemelor liniare**

(condiționarea sistemelor liniare)

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$ ,  $x \in \mathbb{R}^n$  și sist. de ec. liniare:

$$Ax = b$$

$A$  nesingulară  $\Leftrightarrow \det A \neq 0 \Rightarrow \exists$  sol. sist.  $x = A^{-1}b$

Pentru erorile în datele de intrare facem notațiile:

- $\Delta A \in \mathbb{R}^{n \times n}$  eroarea absolută pentru  $A$ ;
- $\Delta b \in \mathbb{R}^n$  eroarea absolută pentru  $b$ ;

În realitate se rezolvă sistemul:

$$(A + \Delta A) \tilde{x} = b + \Delta b$$

soluția fiind  $\tilde{x}$ :

$$\tilde{x} = x + \Delta x$$

În mod natural se ridică următoarele probleme :

1. Dacă  $A$  este matrice nesingulară,  $\Delta A = ?$  a.î.  $A + \Delta A$  să fie nesingulară ?

2. Pp.  $A$  și  $A + \Delta A$  nesingulare care sunt relațiile între

$$\frac{\|\Delta A\|}{\|A\|}, \frac{\|\Delta b\|}{\|b\|} \text{ și } \frac{\|\Delta x\|}{\|x\|} ?$$



1. Pp.  $A$  nesingulară.

$$A + \Delta A = A(I_n + A^{-1}\Delta A) \rightarrow$$

$$A + \Delta A \text{ nesingulară} \Leftrightarrow (I_n + A^{-1}\Delta A) \text{ nesingulară}$$

### Propoziția 5

Fie  $A$  nesingulară și  $\|\Delta A\| < \frac{1}{\|A^{-1}\|}$ . Atunci  $I + A^{-1}\Delta A$  este

nesingulară și avem:

$$\left\| (I_n + A^{-1}\Delta A)^{-1} \right\| \leq \frac{1}{1 - \|A^{-1}\| \cdot \|\Delta A\|}$$

Demonstrație. Avem:

$$\|\Delta A\| < \frac{1}{\|A^{-1}\|} \Rightarrow \|A^{-1}\Delta A\| \leq \|A^{-1}\| \cdot \|\Delta A\| < 1 \stackrel{\text{Pr.4}}{\Rightarrow} \exists (I + A^{-1}\Delta A)^{-1}$$

$$\left\| (I + A^{-1}\Delta A)^{-1} \right\| \leq \frac{1}{1 - \|A^{-1}\Delta A\|} \leq \frac{1}{1 - \|A^{-1}\| \cdot \|\Delta A\|} .$$

Pp. că  $A$  este nesingulară și  $\|\Delta A\| < \frac{1}{\|A^{-1}\|}$ .

$$\begin{aligned}
(A + \Delta A)(x + \Delta x) &= b + \Delta b \Rightarrow (A + \Delta A)\Delta x + Ax + (\Delta A)x = b + \Delta b \Rightarrow \\
A(I + A^{-1}\Delta A)\Delta x &= \Delta b - (\Delta A)x \Rightarrow \Delta x = (I + A^{-1}\Delta A)^{-1} A^{-1}[\Delta b - (\Delta A)x] \Rightarrow \\
\|\Delta x\| &\leq \left\| (I + A^{-1}\Delta A)^{-1} \right\| \|A^{-1}\| (\|\Delta b\| + \|\Delta A\| \|x\|) \Rightarrow \\
\frac{\|\Delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|} \left( \frac{\|\Delta b\|}{\|x\|} + \|\Delta A\| \right) \tag{1}
\end{aligned}$$

Din  $Ax = b$  obținem  $\|b\| \leq \|A\| \|x\| \Rightarrow \frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$  și ținând seamă

de acest rezultat, din (1) deducem:

$$\frac{\|\Delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1}\| \|\Delta A\|} \left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right).$$

$k(A) = \|A^{-1}\| \|A\|$  *numărul de condiționare* al matricii  $A$ .

### Propoziția 6

Dacă matricea  $A$  este nesingulară și  $\|\Delta A\| < \frac{1}{\|A^{-1}\|}$  atunci:

$$\frac{\|\Delta x\|}{\|x\|} \leq \frac{k(A)}{1 - k(A) \cdot \frac{\|\Delta A\|}{\|A\|}} \left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right).$$

Din  $I_n = A A^{-1}$  rezultă  $1 = \|I_n\| \leq \|A\| \|A^{-1}\| = k(A)$ .

$k(A) \geq 1$ ,  $\forall A$  dar dep. de norma matricială naturală utilizată.

O matrice  $A$  pentru care numărul de condiționare este mare se numește matrice *prost condiționată* ( $k(A)$ , mare').

$Ax=b$  cu  $k(A)$  mare  $\rightarrow \frac{\|\Delta x\|}{\|x\|}$  poate fi mare chiar dacă erorile

relative  $\frac{\|\Delta b\|}{\|b\|}$  și  $\frac{\|\Delta A\|}{\|A\|}$  sunt mici.

Fie  $A$  o matrice simetrică  $A = A^T$ , nesară. Utilizând norma matricială subordonată normei vectoriale euclidiene:

$$\|A\|_2 = \sqrt{\rho(A^T A)} = \sqrt{\rho(A^2)}$$

$$k(A) = \|A\|_2 \cdot \|A^{-1}\|_2$$

Matricea simetrică  $A$  are valorile proprii reale  $\lambda_1, \lambda_2, \dots, \lambda_n$ ,

$A^2$  are valorile proprii  $\lambda_1^2, \lambda_2^2, \dots, \lambda_n^2$

$A^{-1}$  are valorile proprii  $\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_n}$ .

$$|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_n| \Rightarrow \rho(A) = |\lambda_n| \quad \text{și} \quad \rho(A^{-1}) = \frac{1}{|\lambda_1|}$$

$$A = A^T \rightarrow \|A\|_2 = \rho(A) = |\lambda_n|, \quad \|A^{-1}\|_2 = \rho(A^{-1}) = \frac{1}{|\lambda_1|},$$

$$k_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \frac{|\lambda_n|}{|\lambda_1|} \text{ număr de condiționare spectral.}$$

$$A \text{ matrice ortogonală} \rightarrow k_2(A) = 1$$

$$A^T A = A \cdot A^T = I_n \Rightarrow A^{-1} = A^T$$

$$\|A\|_2 = \sqrt{\rho(A^T A)} = \sqrt{\rho(I)} = 1 = \|A^T\|_2$$

$$k_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \|A\|_2 \cdot \|A^T\|_2 = 1,$$

Matrice aproape singulară dar cu număr de condiționare mic

$$A = \text{diag} [1, 0.1, 0.1, \dots, 0.1] \in \mathbb{R}^{100 \times 100} \Rightarrow \det A = 1 \cdot (0.1)^{99} = 10^{-99}$$

$$\|A\|_2 = 1, \|A^{-1}\|_2 = 10 \Rightarrow k_2(A) = \|A\|_2 \|A^{-1}\|_2 = 10$$



Matrice foarte prost condiționată cu det. nenul (**det**  $A=1$ )

$$A = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 1 & 2 & \dots & 0 \\ \vdots & & & & \\ 0 & 0 & 0 & \dots & 2 \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix},$$

$$A^{-1} = \begin{pmatrix} 1 & -2 & 4 & \dots & (-2)^{i-1} & \dots & (-2)^{n-1} \\ 0 & 1 & -2 & \dots & (-2)^{i-2} & \dots & (-2)^{n-2} \\ \vdots & & & & & & \\ 0 & 0 & 0 & \dots & 0 & \dots & 1 \end{pmatrix}$$

$$\|A\|_{\infty} = \|A\|_1 = 3 \quad ,$$

$$\|A^{-1}\|_{\infty} = \|A^{-1}\|_1 = 1 + 2 + \cdots + 2^{n-1} = 2^n - 1$$

$$n = 100 \Rightarrow k(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty} = \|A\|_1 \|A^{-1}\|_1 = 3 \cdot (2^{100} - 1)$$

$$\det A = 1$$

$$\begin{cases} x + y = 2 \\ x + 1.001y = 2 \end{cases}, \quad \begin{cases} x + y = 2 \\ x + 1.001y = 2.001 \end{cases} \quad k(A) = 4002$$

$$x = 2, y = 0 \qquad \qquad \qquad x = 1, y = 1$$

$$\begin{cases} 400x - 201y = 200 \\ -800x + 401y = -200 \end{cases}, \quad \begin{cases} 401x - 201y = 200 \\ -800x + 401y = -200 \end{cases} \quad k(A) = 4002$$

$$x = -100, y = -200 \qquad \qquad \qquad x = 40000, y = 79800$$

$$\begin{cases} 1.2969x + 0.8648y = 0.8642 \\ 0.2161x + 0.1441y = 0.1440 \end{cases} \quad x = 2, \quad y = -2 \quad k_2(A) = 249730000$$

$$\bar{x} = 0.9911, \quad \bar{y} = -0.4870,$$

$$r = b - Az = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix} - \begin{pmatrix} 1.2969 & 0.8648 \\ 0.1441 & 0.1441 \end{pmatrix} \begin{pmatrix} 0.9911 \\ -0.4870 \end{pmatrix} = \begin{pmatrix} 10^{-8} \\ -10^{-8} \end{pmatrix}$$

## Matricea Hilbert

$$H = (h_{ij})_{i,j=1}^n \quad , \quad h_{ij} = \frac{1}{i+j-1} = \int_0^1 x^{i+j-2} dx$$

$$k_2(H_n) \approx \frac{(\sqrt{2} + 1)^{4(n+1)}}{2^{\frac{15}{4}} \sqrt{\pi n}} \sim e^{3.5n}$$

n	$k_2(H_n)$	n	$k_2(H_n)$
1	1	7	$4.753 \cdot 10^8$
2	19.281	8	$1.526 \cdot 10^{10}$
3	$5.241 \cdot 10^2$	9	$4.932 \cdot 10^{11}$
4	$1.551 \cdot 10^4$	10	$1.602 \cdot 10^{13}$
5	$4.766 \cdot 10^5$	11	$5.220 \cdot 10^{14}$
6	$1.495 \cdot 10^7$	12	$1.678 \cdot 10^{16}$

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

## Metode numerice de rezolvarea sistemelor liniare

Fie matricea nesingulară  $A \in \mathbb{R}^{n \times n}$  și  $b \in \mathbb{R}^n$ . Rezolvarea sistemului de ecuații liniare  $Ax=b$  se poate face folosind *regula lui Cramer*:

$$x_i = \frac{\det A_i(b)}{\det A}, i = 1, \dots, n,$$

în care  $A_i(b)$  se obține din matricea  $A$  prin înlocuirea coloanei  $i$  cu vectorul  $b$ .

Algoritmul dat de regula lui Cramer este foarte costisitor din punct de vedere al resurselor și instabil numeric.

Din aceste motive s-au căutat alte metode de aproximare a soluției  $\mathbf{x}$ . Unul din cele mai folosiți algoritmi este **algoritmul de eliminare Gauss** :

$A\mathbf{x}=\mathbf{b} \iff \tilde{A}\mathbf{x}=\tilde{\mathbf{b}}$  cu  $\tilde{A}$  matrice superior triunghiulară

$$\mathbf{x} = A^{-1}\mathbf{b} = \tilde{A}^{-1}\tilde{\mathbf{b}} \quad (\text{notăm } A\mathbf{x} = \mathbf{b} \sim \tilde{A}\mathbf{x} = \tilde{\mathbf{b}})$$



## **Metoda substituției**

Fie sistemul liniar  $A\mathbf{x} = \mathbf{b}$  unde matricea sistemului  $A$  este triunghiulară. Pentru a găsi soluția unică a sistemului, trebuie ca matricea să fie nesingulară. Determinantul matricilor triunghiulare este dat de formula:

$$\det A = a_{11}a_{22} \cdots a_{nn}$$

$$\det A \neq 0, a_{ii} \neq 0 \quad \forall i = 1, 2, \dots, n$$

Vom considera întâi cazul când matricea  $A$  este inferior triunghiulară. Sistemul are forma:

$$\begin{aligned} a_{11}x_1 &= b_1 \\ a_{21}x_1 + a_{22}x_2 &= b_2 \\ \vdots & \\ a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{ii}x_i &= b_i \\ \vdots & \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{ni}x_i + \cdots + a_{nn}x_n &= b_n \end{aligned} \tag{1}$$

Necunoscutele  $x_1, x_2, \dots, x_n$ , se deduc folosind ecuațiile sistemului de la prima către ultima.

Din prima ecuație se deduce  $x_1$ :

$$x_1 = \frac{b_1}{a_{11}} \quad (2)$$

Din a doua ecuație , utilizând valoarea  $x_1$  din (2), obținem  $x_2$ :

$$x_2 = \frac{b_2 - a_{21}x_1}{a_{22}}$$

Când ajungem la ecuația  $i$ :

$$a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{ii-1}x_{i-1} + a_{ii}x_i = b_i$$

folosind variabilele  $x_1, x_2, \dots, x_{i-1}$  calculate anterior, avem:

$$x_i = \frac{b_i - a_{i1}x_1 - \cdots - a_{ii-1}x_{i-1}}{a_{ii}}$$

Din ultima ecuație se deduce  $x_n$  astfel:

$$x_n = \frac{b_n - a_{n1}x_1 - a_{n2}x_2 - \cdots - a_{nn-1}x_{n-1}}{a_{nn}}$$

Algoritmul de calcul al soluției sistemelor (1) cu matrice inferior triunghiulară este următorul:

$$x_i = \frac{b_i - \sum_{j=1}^{i-1} a_{ij}x_j}{a_{ii}}, i = 1, 2, \dots, n-1, n$$

Acest algoritm se numește *metoda substituției directe*.

Vom considera, în continuare sistemul (1) cu matrice superior triunghiulară :

$$a_{11}x_1 + \cdots + a_{1i}x_i + \cdots + a_{1n-1}x_{n-1} + a_{1n}x_n = b_1$$

$$\ddots$$

$$a_{ii}x_i + \cdots + a_{in-1}x_{n-1} + a_{in}x_n = b_i$$

$$\ddots$$

$$a_{n-1n-1}x_{n-1} + a_{n-1n}x_n = b_{n-1}$$

$$a_{nn}x_n = b_n$$

Necunoscutele  $x_1, x_2, \dots, x_n$  se deduc pe rând, folosind ecuațiile sistemului, de la ultima către prima.

Din ultima ecuație găsim  $x_n$ :

$$x_n = \frac{b_n}{a_{nn}}$$

Folosind valoarea lui  $x_n$  dedusă mai sus, din penultima ecuație obținem:

$$x_{n-1} = \frac{b_{n-1} - a_{n-1n}x_n}{a_{n-1n-1}}$$

Când ajungem la ecuația  $i$ :

$$a_{ii}x_i + a_{i,i+1}x_{i+1} + \cdots + a_{in}x_n = b_i$$

se cunosc deja  $x_{i+1}, x_{i+2}, \dots, x_n$  de unde ducem:

$$x_i = \frac{b_i - a_{i,i+1}x_{i+1} - \cdots - a_{in}x_n}{a_{ii}}$$

Din prima ecuație găsim valoarea lui  $x_1$ :

$$x_1 = \frac{b_1 - a_{12}x_2 - \cdots - a_{1n}x_n}{a_{11}}$$



Procedeul descris mai sus se numește de *metoda substituției inverse* pentru rezolvarea sistemelor liniare cu matrice superior triunghiulară:

$$x_i = \frac{b_i - \sum_{j=i+1}^n a_{ij}x_j}{a_{ii}}, i = n, n-1, \dots, 2, 1.$$

**M** - numărul de operații \*, / (înmulțiri/împărțiri) efectuate

**A** - numărul operațiilor ± (adunări/scăderi) efectuate.

Atunci pentru calculul componentei  $x_i$  se efectuează  $M=n-i+1$ ,  $A=n-i$  și în total:

$$M = \sum_{i=n}^1 (n-i+1) = \sum_{k=1}^n k = \frac{n(n+1)}{2},$$

$$A = \sum_{i=n}^1 (n-i) = \sum_{k=1}^{n-1} k = \frac{n(n-1)}{2}$$

Efortul de calcul pentru metoda substituției directe este

$$M = \frac{n(n+1)}{2} \quad A = \frac{n(n-1)}{2}.$$

## Algoritmul de eliminare Gauss

Algoritmul se realizează în  $n-1$  pași prin transformarea sistemului dat într-un sistem echivalent cu matrice triunghiulară superior.

### *Pas 1*

la acest pas se obține sistemul:

$A^{(1)}x = b^{(1)} \sim Ax = b$ , unde  $A^{(1)}$  are prima coloană în formă superior triunghiulară.

## ***Pas 2***

se construiește sistemul

$A^{(2)}\mathbf{x} = \mathbf{b}^{(2)} \sim A\mathbf{x} = \mathbf{b}$ , unde  $A^{(2)}$  are primele două coloane în formă superior triunghiulară.

⋮

## ***Pasul $r$***

se obține sistemul  $A^{(r)}\mathbf{x} = \mathbf{b}^{(r)} \sim A\mathbf{x} = \mathbf{b}$ , unde  $A^{(r)}$  are primele  $r$  coloane în formă superior triunghiulară.

⋮

## ***Pasul $n-1$ :***

se obține sistemul

$A^{(n-1)}\mathbf{x} = \mathbf{b}^{(n-1)} \sim A\mathbf{x} = \mathbf{b}$ , unde  $A^{(n-1)}$  are primele  $n-1$  coloane în formă superior triunghiulară.

Dacă la un anumit pas matricea  $A^{(r)}$  nu poate fi construită aceasta ne va arăta că matricea  $A$  este singulară.

În realizarea acestor pași se utilizează următoarele operații elementare:

- înmulțirea unei ecuații cu un factor și adunarea la altă ecuație;
- interschimbarea a două linii și/sau două coloane în matricea  $A$ .

## **Pasul 1**

Intrare : sistemul  $A\mathbf{x}=\mathbf{b}$

Ieșire : sistemul  $A^{(1)}\mathbf{x} = \mathbf{b}^{(1)} \sim A\mathbf{x} = \mathbf{b}$ , matr  $A^{(1)}$  are prima coloană în formă superior triunghiulară.

Fie ecuația  $i$ , cu  $i=1,\dots,n$

$$E_i : \quad a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n = b_i.$$

Presupunem  $a_{11} \neq 0$ . Operațiile efectuate au ca obiectiv anularea coeficienților lui  $x_1$  din ecuațiile de la 2 la  $n$  și sunt descrise în continuare:

$$\mathbf{E}_1 * \begin{pmatrix} -\mathbf{a}_{21} \\ \mathbf{a}_{11} \end{pmatrix} + \mathbf{E}_2 = \mathbf{E}_2^{(1)} \quad \Rightarrow \quad \mathbf{a}_{21}^{(1)} = \mathbf{0}$$

$$\vdots$$

$$\mathbf{E}_1 * \begin{pmatrix} -\mathbf{a}_{i1} \\ \mathbf{a}_{11} \end{pmatrix} + \mathbf{E}_i = \mathbf{E}_i^{(1)} \quad \Rightarrow \quad \mathbf{a}_{i1}^{(1)} = \mathbf{0}$$

$$\vdots$$

$$\mathbf{E}_1 * \begin{pmatrix} -\mathbf{a}_{n1} \\ \mathbf{a}_{11} \end{pmatrix} + \mathbf{E}_n = \mathbf{E}_n^{(1)} \quad \Rightarrow \quad \mathbf{a}_{n1}^{(1)} = \mathbf{0}$$

Sistemul obținut prin aceste operații are forma:

$$\left\{ \begin{array}{l} a_{11}^{(1)} x_1 + a_{12}^{(1)} x_2 + \dots + a_{1n}^{(1)} x_n = b_1^{(1)} \\ a_{22}^{(1)} x_2 + \dots + a_{2n}^{(1)} x_n = b_2^{(1)} \\ \vdots \\ a_{i2}^{(1)} x_2 + \dots + a_{in}^{(1)} x_n = b_i^{(1)} \\ \vdots \\ a_{n2}^{(1)} x_2 + \dots + a_{nn}^{(1)} x_n = b_n^{(1)} \end{array} \right.$$



## Pas 2

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

Ieșire :  $A^{(2)}\mathbf{x} = \mathbf{b}^{(2)} \sim A\mathbf{x} = \mathbf{b}$ ,  $A^{(2)}$  are primele două coloane în formă superior triunghiulară.

Se presupune  $a_{22}^{(1)} \neq 0$  și se urmărește anularea elementelor  $a_{32}^{(2)}, a_{42}^{(2)}, \dots, a_{n2}^{(2)}$  (transformarea coloanei 2 în formă superior triunghiulară). Operațiile efectuate asupra ecuațiilor  $E_i^{(1)}, i = 3, \dots, n$  sunt următoarele :

$$\left\{ \begin{array}{l} E_2^{(1)} * \left( -\frac{a_{32}^{(1)}}{a_{22}^{(1)}} \right) + E_3^{(1)} = E_3^{(2)} \Rightarrow a_{32}^{(2)} = \mathbf{0}; \\ \vdots \\ E_2^{(1)} * \left( -\frac{a_{i2}^{(1)}}{a_{22}^{(1)}} \right) + E_i^{(1)} = E_i^{(2)} \Rightarrow a_{i2}^{(2)} = \mathbf{0}; \\ \vdots \\ E_2^{(1)} * \left( -\frac{a_{n2}^{(1)}}{a_{22}^{(1)}} \right) + E_n^{(1)} = E_n^{(2)} \Rightarrow a_{n2}^{(2)} = \mathbf{0}; \end{array} \right.$$

Se observă că nu se schimbă forma superior triunghiulară a primei coloane.

### **Pas $r$**

Intrare :  $A^{(r-1)}\mathbf{x} = \mathbf{b}^{(r-1)}$

Ieșire :  $A^{(r)}\mathbf{x} = \mathbf{b}^{(r)} \sim A\mathbf{x} = \mathbf{b}$ ,  $A^{(r)}$  are primele  $r$  coloane în formă superior triunghiulară.

Sistemul are forma următoare:

$$\left\{ \begin{array}{l} \mathbf{a}_{11}^{(r-1)} \mathbf{x}_1 + \cdots + \mathbf{a}_{1r}^{(r-1)} \mathbf{x}_r + \cdots + \mathbf{a}_{1n}^{(r-1)} \mathbf{x}_n = \mathbf{b}_1^{(r-1)} \\ \quad \cdot \quad \cdot \quad \cdot \\ \quad \quad \cdot \quad \cdot \quad \cdot \\ \mathbf{a}_{rr}^{(r-1)} \mathbf{x}_r + \cdots + \mathbf{a}_{rn}^{(r-1)} \mathbf{x}_n = \mathbf{b}_r^{(r-1)} \\ \mathbf{a}_{r+1r}^{(r-1)} \mathbf{x}_r + \cdots + \mathbf{a}_{r+1n}^{(r-1)} \mathbf{x}_n = \mathbf{b}_{r+1}^{(r-1)} \\ \quad \vdots \\ \mathbf{a}_{ir}^{(r-1)} \mathbf{x}_r + \cdots + \mathbf{a}_{in}^{(r-1)} \mathbf{x}_n = \mathbf{b}_i^{(r-1)} \\ \quad \vdots \\ \mathbf{a}_{nr}^{(r-1)} \mathbf{x}_r + \cdots + \mathbf{a}_{nn}^{(r-1)} \mathbf{x}_n = \mathbf{b}_n^{(r-1)} \end{array} \right.$$

Presupunem  $\mathbf{a}_{rr}^{(r-1)} \neq \mathbf{0}$ .

Vom urmări anularea elementelor  $\mathbf{a}_{r+1r}^{(r)}$ ,  $\mathbf{a}_{r+2r}^{(r)}$ ,  $\dots$ ,  $\mathbf{a}_{nr}^{(r)}$ .

$$\left\{ \begin{array}{l} \mathbf{E}_r^{(r-1)} * \left( -\frac{\mathbf{a}_{r+1r}^{(r-1)}}{\mathbf{a}_{rr}^{(r-1)}} \right) + \mathbf{E}_{r+1}^{(r-1)} = \mathbf{E}_{r+1}^{(r)} \quad \Rightarrow \quad \mathbf{a}_{r+1r}^{(r)} = \mathbf{0}; \\ \vdots \\ \mathbf{E}_r^{(r-1)} * \left( -\frac{\mathbf{a}_{ir}^{(r-1)}}{\mathbf{a}_{rr}^{(r-1)}} \right) + \mathbf{E}_i^{(r-1)} = \mathbf{E}_i^{(r)} \quad \Rightarrow \quad \mathbf{a}_{ir}^{(r)} = \mathbf{0}; \\ \vdots \\ \mathbf{E}_r^{(r-1)} * \left( -\frac{\mathbf{a}_{nr}^{(r-1)}}{\mathbf{a}_{rr}^{(r-1)}} \right) + \mathbf{E}_n^{(r-1)} = \mathbf{E}_n^{(r)} \quad \Rightarrow \quad \mathbf{a}_{nr}^{(r)} = \mathbf{0}; \end{array} \right.$$

Se observă că nu se schimbă forma superior triunghiulară a primelor  $r-1$  coloane.

La fiecare pas s-a făcut ipoteza  $a_{rr}^{(r-1)} \neq 0$ . Elementul  $a_{rr}^{(r-1)}$  poartă numele de *pivot*. În cazul în care elementul pivot este nul se pot aplica următoarele strategii, numite de *pivotare*:

## Pivotare ( $a_{rr}^{(r-1)} \neq 0$ ?)

### 1<sup>0</sup> *Fără pivotare*

Se caută primul indice  $i_0 \in \{r, r+1, \dots, n\}$  astfel încât  $a_{i_0 r}^{(r-1)} \neq 0$ . Se interschimbă liniile  $i_0$  și  $r$ .

Să observăm că în procesul de calcul la pasul  $r$  intervine factorul  $\frac{1}{a_{rr}^{(r-1)}}$  astfel că valori mici ale lui  $|a_{rr}^{(r-1)}|$  conduc la

amplificarea erorilor de calcul. Pentru a asigura stabilitatea numerică a procesului de calcul este de dorit ca  $|a_{rr}^{(r-1)}|$  să fie ‘mare’.

## *2<sup>0</sup> Pivotare parțială*

Se determină indicele  $i_0$ :

$$\left| a_{i_0 r}^{(r-1)} \right| = \max \left\{ \left| a_{ir}^{(r-1)} \right| ; i = r, \dots, n \right\}$$

și se interschimbă liniile  $i_0, r$  dacă  $i_0 \neq r$ .

## *3<sup>0</sup> Pivotare totală*

Se determină indicii  $i_0$  și  $j_0$ :

$$\left| a_{i_0 j_0}^{(r-1)} \right| = \max \left\{ \left| a_{ij}^{(r-1)} \right| ; i = r, \dots, n, j = r, \dots, n \right\}$$

și se interschimbă liniile  $i_0, r$  dacă  $i_0 \neq r$  și coloanele  $j_0, r$  dacă  $j_0 \neq r$



Schimbarea coloanelor implică schimbarea ordinii variabilelor astfel încât în final va trebui refăcută ordinea inițială a variabilelor.

Dacă după pivotare elementul pivot rămâne nul,  $a_{rr}^{(r-1)} = 0$ , atunci putem deduce că  $A^{(r-1)}$  este singulară.

În adevăr, dacă în procesul de pivotare parțială  $a_{rr}^{(r-1)} = 0$ , atunci

$$\begin{aligned}
 A^{(r-1)} &= \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ & \ddots & & \\ & & a_{rr} = 0 & \cdots a_{rn} \\ & & 0 & \\ & & 0 & \\ & & 0 & \cdots a_{nn} \end{bmatrix} \Rightarrow \\
 \det A^{(r-1)} &= a_{11}^{(r-1)} a_{22}^{(r-1)} \cdots a_{r-1, r-1}^{(r-1)} \det \begin{bmatrix} 0 & \cdots & a_{rn} \\ 0 & & \\ \vdots & & \\ 0 & \cdots & a_{nn} \end{bmatrix} = 0
 \end{aligned}$$

Deoarece operațiile efectuate (cele de inetrschimbare de linii și/sau coloane) nu au schimbat decât semnul determinantului avem:

$$\det A = \pm \det A^{(r-1)} = 0 \Rightarrow \det A = 0$$

prin urmare matricea  $A$  inițială este singulară.

Și în cazul procesului de pivotare totală dacă  $a_{rr}^{(r-1)} = 0$ , atunci:

$$A^{(r-1)} = \begin{bmatrix} a_{11} & a_{12} & \cdots & \cdots & a_{1n} \\ & \ddots & & & \\ & & a_{rr} = 0 \cdots 0 & & \\ & & 0 & & \\ & & 0 & & \\ & & 0 \cdots \cdots \cdots 0 & & \end{bmatrix} \Rightarrow$$

$$\det A^{(r-1)} = a_{11}^{(r-1)} a_{22}^{(r-1)} \cdots a_{r-1, r-1}^{(r-1)} \det \begin{bmatrix} 0 \cdots \cdots 0 \\ 0 \cdots \cdots 0 \\ \vdots \quad \quad \quad \vdots \\ 0 \cdots \cdots 0 \end{bmatrix} = 0$$

$\det A = \pm \det A^{(r-1)} = 0 \Rightarrow A$  este matrice singulară.

```

 $r = 1;$ 
pivotare( $r$ );
while ( $r \leq n - 1$  și  $|a_{rr}| > \varepsilon$ )
    // Pas  $r$ 
    * for  $i = r + 1, \dots, n$ 
        •  $f = -\frac{a_{ir}}{a_{rr}};$ 
        • for  $j = r + 1, \dots, n$ 
             $a_{ij} = a_{ij} + f * a_{rj};$ 
        •  $a_{ir} = 0;$ 
        •  $b_i = b_i + f * b_r;$ 
    *  $r = r + 1;$ 
    * pivotare( $r$ );
if ( $|a_{rr}| \leq \varepsilon$ ) 'MATRICE SINGULARA'
else {  $A \leftarrow A^{(n-1)}$  ,  $b \leftarrow b^{(n-1)}$ 
    se rezolvă sistemul triunghiular superior  $Ax = b$ }

```

**Numărul de operații** efectuate la pasul  $r$  și în total este:

$$(n-r)[1M + (n-r)A + (n-r)M + 1A + 1M] \Rightarrow$$

$$\mathbf{M}: \quad \sum_{r=1}^{n-1} (n-r)^2 + 2 \sum_{r=1}^{n-1} (n-r) = \frac{(n-1)n(2n+5)}{6},$$

$$\mathbf{A}: \quad \sum_{r=1}^{n-1} (n-r)^2 + \sum_{r=1}^{n-1} (n-r) = \frac{(n-1)n(n+1)}{3},$$

$$\mathbf{M}: \quad \frac{n^3}{3} + \mathcal{O}(n^2) \quad ; \quad \mathbf{A}: \quad \frac{n^3}{3} + \mathcal{O}(n^2)$$

## **Eliminarea „chinezească”**

200-100 î.Cr. China – *9 capitole despre arta matematică* – metodă de rezolvare foarte asemănătoare eliminării Gauss

*„Avem 3 tipuri de grâu. Știm că 3 baloturi din primul tip, 2 baloturi din al doilea tip și 1 balot din al treilea tip cântăresc 39 măsuri. Deasemenea, 2 baloturi din primul tip, 3 baloturi din al doilea tip și 1 balot din al treilea tip cântăresc 34 măsuri și 1 balot din primul tip, 2 baloturi din al doilea tip și 3 baloturi din al treilea tip cântăresc 26 măsuri. Câte măsuri cântărește un balot din fiecare tip de grâu”*

Notăția actuală:

$$3b_1 + 2b_2 + b_3 = 39$$

$$2b_1 + 3b_2 + b_3 = 34$$

$$b_1 + 2b_2 + 3b_3 = 26$$

Notăția *chinezească*

<b>1</b>	<b>2</b>	<b>3</b>
<b>2</b>	<b>3</b>	<b>2</b>
<b>3</b>	<b>1</b>	<b>1</b>
<b>26</b>	<b>34</b>	<b>39</b>



## **Pasul 1**

Se înmulțește coloana a doua cu 3 și se scade din ea coloana a treia atât timp cât este posibil.

Se înmulțește prima coloană cu 3 și se scade din ea coloana a treia atât timp cât este posibil.

Se ajunge la forma:

$$\begin{array}{ccc} 0 & 0 & 3 \\ 4 & 5 & 2 \\ 8 & 1 & 1 \\ 39 & 24 & 39 \end{array}$$

## **Pasul 2**

Se înmulțește prima coloană cu 5 și se scade din ea coloana a doua atât timp cât este posibil.

Se ajunge la forma:

$$\begin{array}{ccc} 0 & 0 & 3 \\ 0 & 5 & 2 \\ 36 & 1 & 1 \\ 99 & 24 & 39 \end{array}$$

Pentru rezolvare se folosește metoda substituției inverse pe sistemul obținut mai sus.