

# **Calcul Numeric**

**Cursul 1**

**2022**

*Anca Ignat*

## **Echipa**

**Andreea Arusoaie**

**Ștefan Bălăucă**

**Sebastian Ciobanu**

**Andrei Luca**

**Valentin Roșca**

**Augustus Tabarcea (en)**

**Anca Ignat**

**[ancai@info.uaic.ro](mailto:ancai@info.uaic.ro) , [ancai\\_fii@yahoo.ro](mailto:ancai_fii@yahoo.ro)**

**[numericalcalculus2019@gmail.com](mailto:numericalcalculus2019@gmail.com) – teme de laborator**

**<http://profs.info.uaic.ro/~ancai/CN/>**

**Detalii despre modul de desfășurare online: Discord**

**Consultații:**

- **prin e-mail la adresele de mai sus sau**
- **Online, Zoom – Luni 16:30-18**

## **Regulament - 2022**

### **Laborator**

- 8 teme
- cei care prezintă temele până la termenul limită precizat la fiecare temă, punctajul maxim ce poate fi obținut este punctajul afișat pentru fiecare temă
- cei care prezintă temele după termenul limită precizat punctarea se va face din 50% din punctajul temei prezentate

## Examen

- teză scrisă de 1 oră, cu 3 sau 4 exerciții din materia predată,  
*"cu cursurile pe masă"*
- teza scrisă este notată între 1 și 10
- testul scris va avea loc:
  - în săptămâna 10 sau 11 - din primele 9 cursuri
  - 16 mai - 22 mai 2022 - pentru cei care nu obțin punctaj de promovare sau pentru mărirea notei - din toate cursurile

## **Calculul punctajului / notei final(e)**

**Punctaj final = punctaj laborator + 45\*nota examen**

**Promovează** disciplina acei studenți care au:

- **nota la examen  $\geq 3$**

**și**

- **punctaj final  $\geq 410$  pt**

**Nota finală** se calculează din punctajul final aplicând ”**curba lui Gauss**”.

## **Desfășurarea semestrului**

**Săptămânile 1-7 și 9-13 – școală conform orarului**

**Săptămâna a 8-a (prima săptămână de evaluare) – liberă**

**Săptămânile 10, 11, 12 - test scris**

## Bibliografie

1. Elemente de informatică și calcul numeric, vol. 1 - C.Ignat, C.Ilioi, T.Jucan - Ed. Univ. 'Al. I. Cuza' Iași, 1987
2. Matrix Computations - G.H. Golub, C.F. van Loan - John Hopkins Univ. Press, 2012
3. Numerical Analysis – R.L. Burden, J.D. Faires – Brooks/Cole, Thomson Learning (10-th edition, 2015)
4. Calcul numeric în C - T. A. Beu - Ed. Albastră, Cluj, 2004
5. Numerical analysis with algorithms and programming, Santanu Saha Ray, CRC Press, 2016.
6. Numerical Recipes 3rd Edition: The Art of Scientific Computing, W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, Cambridge University Press, NY, USA, 2007 (<http://numerical.recipes/>)
7. Linear Algebra, Ideas and Applications - R.C. Penney, 4-th ed., Wiley, 2016
8. Numerical Optimization – J. Nocedal, S.J. Wright, Springer-Verlag, New York, 1999



## Capitolele cursului

- ❖ Rezolvarea sistemelor liniare ( $Ax=b$ )
- ❖ Optimizare numerică ( $\min \{ F(x) ; x \in R^n \}$ )
- ❖ Valori și vectori proprii ( $Au = \lambda u$ )
- ❖ Interpolare numerică
- ❖ Ecuații neliniare ( $f(x)=0$ )

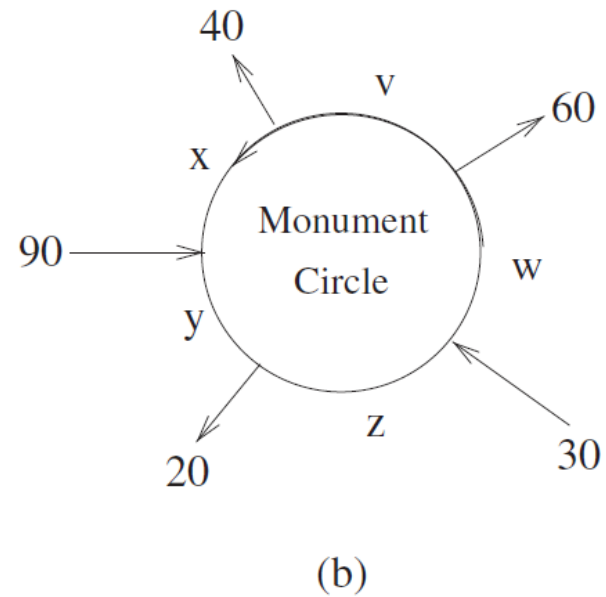
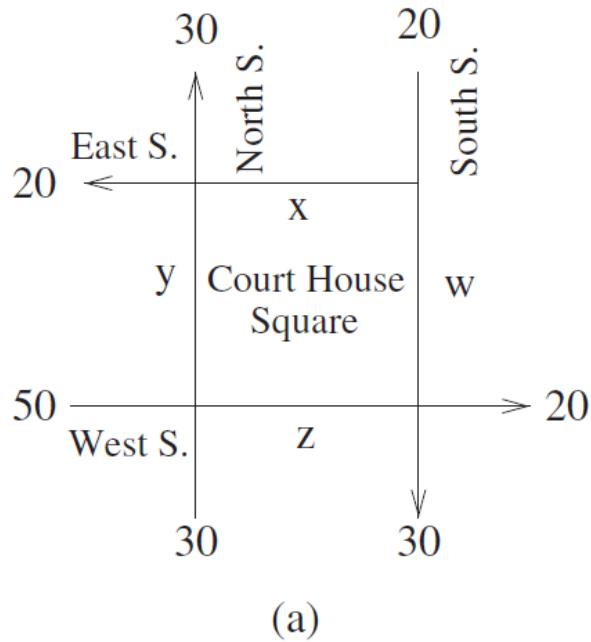
*“The world cannot be understood without numbers. But the world cannot be understood with numbers alone.”*

Hans Rosling, Anna Rosling Ronnlund, Ola Rosling

Factfulness. Zece motive pentru care interpretăm greșit lumea și de ce lucrurile stau mai bine decât crezi

“The world cannot be understood without numbers, nor through numbers alone.”

# Traffic Flow



**FIGURE 1.24** Two traffic patterns.

(R.C. Penney – Linear Algebra, Ideas and Applications, 4-th ed., Wiley, 2016)

- one way streets;
- the numbers represent the average number of cars per minute that enter or leave a given street at 3:30pm;
- $x, y, z, w, \dots$  - average number of cars per minute on a certain street
- no. of cars entering = no. of cars leaving

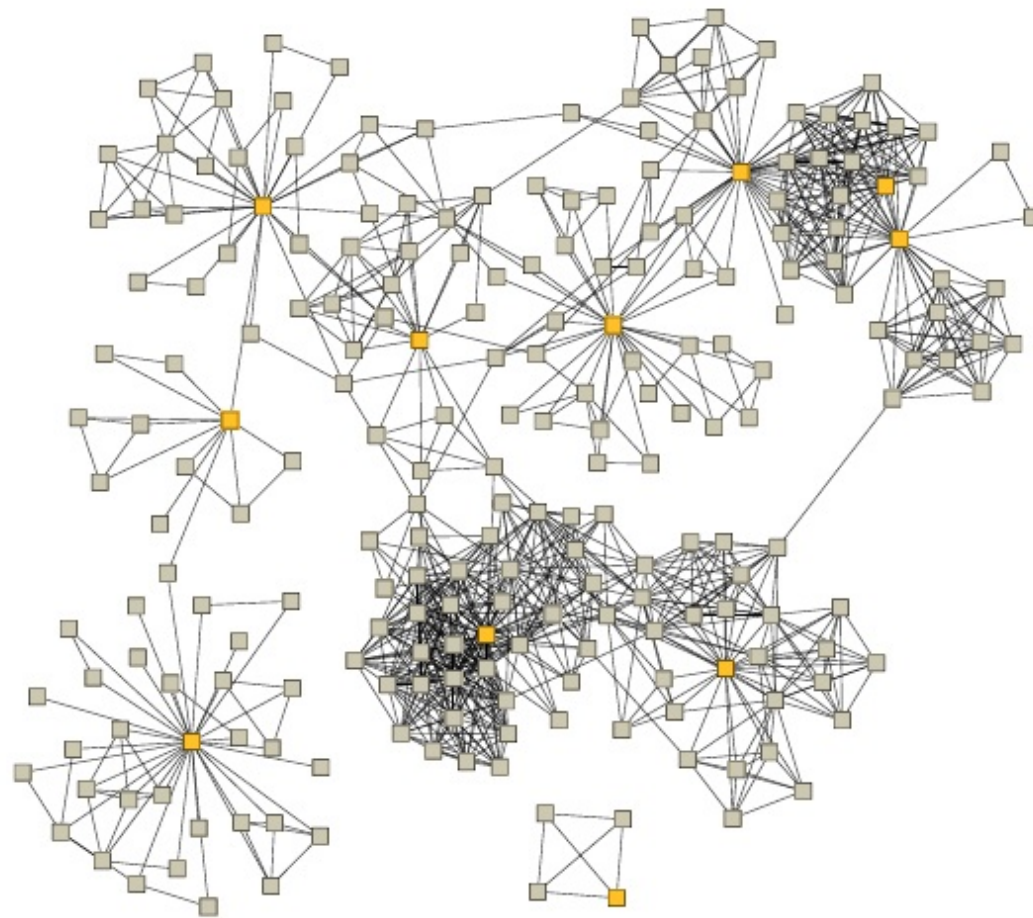
$$x + y = 50$$

$$y + z = 80$$

$$z + w = 50$$

$$x + w = 20$$

$$x + y + z + w = 100$$



## Centralitate în rețelele sociale

- noțiune introdusă de A. Bavelas în 1948 studiind comunicarea între oameni

$(V, E)$  – graful care modelează rețeaua,  $A$  – matricea de adiacență asociată,  $A = (a_{ij})_{i,j=1}^N$ ,  $N=/?/$

- care sunt cele mai ‘importante’ noduri din rețea?

1. Centralitate de grad (*'degree centrality'*) – se bazează pe noțiunea de grad/grade asociate nodurilor în grafuri

Se numește **drum geodesic** între două vârfuri orice drum de lungime minimă (număr minim de muchii) dintre cele 2 vârfuri.

2. Centralitate de apropiere (*'closeness centrality'*)

Centralitatea de apropiere a unui nod este suma lungimilor drumurilor geodesice de la nodul respectiv la toate celelalte noduri.



### 3. Centralitate de interrelație (*'betweenness centrality'*)

$$b(v) = \sum_{\substack{s \neq v \neq t, \\ s, t \in V}} \frac{n_{st}(v)}{n_{st}}$$

unde  $n_{st}$  este numărul total de drumuri geodesice între nodurile  $s$  și  $t$ , iar  $n_{st}(v)$  este numărul de drumuri geodesice care trec prin nodul  $v$ .

- măsoară controlul pe care îl deține nodul  $v$  în circulația informațiilor în rețea

#### 4. Centralitate de vector propriu (*'eigenvector centrality'*)

- se ține cont de faptul că nu toate muchiile (conexiunile) sunt la fel de importante (ca în cazul centralității de grad)
- conexiunile către persoane influente vor 'împrumuta' importanță mai mare decât conexiunie către persoanele mai puțin influente

$x(i)$  = centralitatea de vector propriu a nodului  $v_i$

$$x(i) = \frac{1}{\lambda} \sum_{j \in \Gamma(v_i)} x(j) = \frac{1}{\lambda} \sum_{j=1}^N a_{ij} x(j), \quad \lambda > 0 \text{ o constanta}$$

$$\mathbf{x} = (x(1), x(2), \dots, x(N))^T$$

$$\mathbf{x} = \frac{1}{\lambda} A\mathbf{x} \Leftrightarrow A\mathbf{x} = \lambda \mathbf{x}$$

$\lambda_A > 0$  - valoarea proprie Perron (cea mai mare valoare proprie),  $\mathbf{x}$  - vectorul propriu asociat

**Compresia imaginilor digitale și descompunerea după valori singulare**  
o imagine digitală  $\leftrightarrow$  matrice de pixeli  $A$  cu  $m$  linii și  $n$  coloane

$$A = \left( a_{ij} \right)_{\substack{i=1,\dots,m \\ j=1,\dots,n}}, \quad a_{ij} \in \mathbb{R} \text{ sau } \mathbb{R}^3$$

$$(a_{ij} \in \{0,1,\dots,255\} \text{ sau } a_{ij} \in \{0,1,\dots,255\}^3 \text{ sau } a_{ij} \in [0,1]^{(3)})$$

Memorarea lui  $A$ :  $m \cdot n \cdot \text{mem}(\text{int/double})(\cdot 3)$  bytes

Descompunerea după valori singulare (**SVD**) a unei matrici

$$A = U S V^T, \quad U \in \mathbb{R}^{m \times m}, \quad S \in \mathbb{R}^{m \times n}, \quad V \in \mathbb{R}^{n \times n}$$

$$U = [u_1, u_2, \dots, u_m], \quad V = [v_1, v_2, \dots, v_n] - \text{matrici ortogonale}$$

$$\left(u_i, u_j\right)_{\mathbb{R}^m} = \delta_{ij} = \begin{cases} 1 & \text{daca } i = j \\ 0 & \text{daca } i \neq j \end{cases}, \quad \left(v_i, v_j\right)_{\mathbb{R}^n} = \delta_{ij} \quad \forall i, j$$

$$S = \begin{pmatrix} \sigma_1 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sigma_2 & & \mathbf{0} & \mathbf{0} \\ & & \ddots & & \\ \mathbf{0} & \mathbf{0} & & \sigma_r & \mathbf{0} \\ & & & & \ddots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > \mathbf{0}$ ,  $r \leq \min\{m, n\}$  - valorile singulare ale matr.  $A$

$$A = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \cdots + \sigma_r u_r v_r^T$$

$$A \approx A_k = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \cdots + \sigma_k u_k v_k^T$$

memorarea lui  $A_k$  necesită  $k(m+n+1) \cdot \text{mem}(\text{double})$

$$m=1265, n=538 \quad , \quad r_c = \frac{mn}{k(m+n+1)},$$

$$k=50 \quad r_c=7.54; \quad k=100 \quad r_c=3.77; \quad k=200 \quad r_c=1.88;$$

$$k=300 \quad r_c=1.25; \quad k=400 \quad r_c=0.94; \quad k=538 \quad r_c=0.70$$

## Vectori și matrici

Fie  $x_i, y_i, \lambda \in \mathbb{R}$ . Se definesc vectorii  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  și operațiile de adunare și înmulțire cu scalari astfel:

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{x} + \mathbf{y} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{pmatrix}, \lambda \mathbf{x} = \begin{pmatrix} \lambda x_1 \\ \lambda x_2 \\ \vdots \\ \lambda x_n \end{pmatrix} .$$

Fie vectorul  $z \in \mathbb{C}^n$ :

$$z = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} \text{ cu } z_1, z_2, \dots, z_n \in \mathbb{C}.$$

Pentru  $z \in \mathbb{C}$  utilizăm notațiile:

$$z = a + ib, \operatorname{Re} z = a, \operatorname{Im} z = b,$$

$$\bar{z} = a - ib - \text{conjugatul numărului } z$$

$$|z| = \sqrt{a^2 + b^2} - \text{modulul numărului complex } z$$



Notăm cu  $\mathbb{R}^{m \times n} / \mathbb{C}^{m \times n}$  spațiul matricilor cu elemente reale / complexe cu  $m$  linii și  $n$  coloane

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}, \quad a_{ij} \in \mathbb{R} / \mathbb{C}, i = 1, 2, \dots, m, j = 1, 2, \dots, n$$

## Definiție

$X$  se numește *spațiu vectorial* (*spațiu liniar*)

$$+ : X \times X \rightarrow X \text{ și } \cdot : K \times X \rightarrow X, \quad (K = \mathbb{R})$$

astfel încât  $(X, +)$  este un grup comutativ :

$$a + b = b + a, \quad \forall a, b \in X - \text{comutativitate,}$$

$$(a + b) + c = a + (b + c), \quad \forall a, b, c \in X - \text{asociativitate,}$$

$$\exists 0 \in X \text{ a. î. } a + 0 = 0 + a = a, \quad \forall a \in X - \text{element neutru,}$$

$$\forall a \in X, \exists -a \in X \text{ a. î. } a + (-a) = (-a) + a = 0 - \text{element opus.}$$

iar pentru operația de înmulțire cu scalari au loc relațiile:

$$\lambda(a+b) = \lambda a + \lambda b, \forall \lambda \in K, \forall a, b \in X,$$

$$(\lambda + \mu) a = \lambda a + \mu a, \forall \lambda, \mu \in K, \forall a \in X,$$

$$\lambda(\mu a) = (\lambda \mu) a, \forall \lambda, \mu \in K, \forall a \in X,$$

$$\exists 1 \in K \text{ astfel încât } 1 \cdot a = a, \forall a \in X.$$

## Definiție

Fie  $X$  un spațiu liniar. Spunem că vectorii  $x_1, x_2, \dots, x_p \in X$  sunt *liniar independenți* dacă:

$$\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_p x_p = \mathbf{0} \Rightarrow \alpha_1 = \alpha_2 = \dots = \alpha_p = 0, \alpha_i \in K$$

*Spațiul* vectorial  $X$  este *finit dimensional* dacă există  $p$  vectori liniar independenți în  $X$ ,  $x_1, x_2, \dots, x_p \in X$ , și orice mulțime de  $q$  elemente din  $X$  cu  $q > p$  este liniar dependentă. În acest caz dimensiunea spațiului  $X$  este  $p$  ( $\dim X = p$ ).

Fie spațiul vectorial  $X$  finit dimensional cu  $\dim X = p$ . Orice sistem de  $p$  vectori liniar independenți din  $X$  se numește *bază* a spațiului  $X$ .

Fie  $x_1, x_2, \dots, x_p \in X$  o bază pentru spațiul  $X$ . Atunci pentru  $\forall x \in X, \exists$  unice constantele  $\alpha_1, \alpha_2, \dots, \alpha_p \in K$  astfel încât

$$x = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_p x_p = \sum_{i=1}^p \alpha_i x_i .$$

$\mathbb{R}^n$  este un spațiu vectorial finit dimensional,  $\dim \mathbb{R}^n = n$  cu baza canonică:

$$e_1 = \begin{pmatrix} \mathbf{1} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, e_2 = \begin{pmatrix} \mathbf{0} \\ \mathbf{1} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \dots, e_k = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{1} \\ \vdots \\ \mathbf{0} \end{pmatrix} \text{ - poziția } k, \dots, e_n = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{1} \end{pmatrix}$$

## Calcul matricial

Fie matricea  $A \in \mathbb{R}^{m \times n}$ :

$$A = \begin{pmatrix} \mathbf{a}_{11} & \cdots & \mathbf{a}_{1n} \\ \vdots & \ddots & \vdots \\ \mathbf{a}_{m1} & \cdots & \mathbf{a}_{mn} \end{pmatrix}, \quad A = (\mathbf{a}_{ij})_{i=1..m, j=1..n}$$

Se definește *matricea transpusă*:

$$A^T = \begin{pmatrix} \mathbf{a}_{11} & \cdots & \mathbf{a}_{m1} \\ \vdots & \ddots & \vdots \\ \mathbf{a}_{1n} & \cdots & \mathbf{a}_{mn} \end{pmatrix}, \quad A^T = (\mathbf{a}_{ji})_{i=1..m, j=1..n} \in \mathbb{R}^{n \times m}$$

Pentru matricea:

$$A \in \mathbb{C}^{m \times n}, A = \left( a_{ij} \right)_{\substack{i=1 \dots m \\ j=1 \dots n}}$$

se definește *matricea adjunctă*  $A^H$ :

$$A^H = \overline{A^T} = \left( \overline{a_{ji}} \right)_{\substack{j=1 \dots n \\ i=1 \dots m}}$$

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}, \quad A^H = \begin{pmatrix} \overline{a_{11}} & \cdots & \overline{a_{m1}} \\ \vdots & \ddots & \vdots \\ \overline{a_{1n}} & \cdots & \overline{a_{mn}} \end{pmatrix}$$

Pentru  $A \in \mathbb{R}^{m \times n}$  matricea adjunctă coincide cu transpusa,

$$A^H = A^T.$$



Fie vectorul  $\mathbf{x} \in \mathbb{R}^n$ , acesta este considerat vector coloană,

$\mathbf{x} \in \mathbb{R}^{n \times 1}$ :

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \Rightarrow \mathbf{x}^T = (x_1 \ x_2 \ \cdots \ x_n)$$



**O primăvară frumoasă!**

# **Calcul Numeric**

**Cursul 2**

**2022**

*Anca Ignat*

Dacă facem înmulțirea matricială  $Ae_j$  obținem coloana  $j$  a matricei  $A$ :

$$Ae_j = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{1}_{\text{poziția } j} \\ \vdots \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{pmatrix}$$

$Ae_j$  este coloana  $j$  a matricei  $A$ ,  $j=1, \dots, n$  ;

$e_i^T A$  este linia  $i$  a matricei  $A$ ,  $i=1, \dots, m$ .

Fie vectorii  $\mathbf{x}$ ,  $\mathbf{y}$ , cu ajutorul lor definim produsele scalare în  $\mathbb{C}^n$  și  $\mathbb{R}^n$ :

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{C}^n, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{C}^n$$

$$(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i \overline{y_i} = \mathbf{y}^H \mathbf{x} = (\overline{y_1} \quad \overline{y_2} \quad \cdots \quad \overline{y_n}) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

$$(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i y_i = \mathbf{y}^T \mathbf{x} = (y_1 \ y_2 \ \cdots \ y_n) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

## Proprietățile matricei $A^H$

1.  $(A + B)^H = A^H + B^H$

2.  $(A^H)^H = A$

3.  $(AB)^H = B^H A^H$

4.  $(A^{-1})^H = (A^H)^{-1}$

## Proprietăți ale matricei $A^T$

$$(A + B)^T = A^T + B^T$$

$$(A^T)^T = A$$

$$(AB)^T = B^T A^T$$

$$(A^{-1})^T = (A^T)^{-1}$$

## Propoziție

Fie  $A \in \mathbb{C}^{m \times n}$ ,  $x \in \mathbb{C}^n$ ,  $y \in \mathbb{C}^m$  atunci:

$$(Ax, y)_{\mathbb{C}^m} = (x, A^H y)_{\mathbb{C}^n}.$$

Pentru cazul real avem:

$$A \in \mathbb{R}^{m \times n}, x \in \mathbb{R}^n, y \in \mathbb{R}^m \Rightarrow (Ax, y)_{\mathbb{R}^m} = (x, A^T y)_{\mathbb{R}^n}$$

Demonstrație

$$\begin{aligned} (Ax, y) &= y^H (Ax) = y^H A x = y^H (A^H)^H x = \\ &= (A^H y)^H x = (x, A^H y). \end{aligned}$$



## Tipuri de matrice

### Definiții

O matrice  $A \in \mathbb{R}^{n \times n}$  se numește *simetrică* dacă  $A = A^T$ .

O matrice  $A \in \mathbb{C}^{n \times n}$  se numește *autoadjunctă* dacă  $A = A^H$ .

O matrice  $A \in \mathbb{C}^{n \times n}$  se numește *unitară* dacă  $A^H A = A A^H = I_n$ .

O matrice  $A \in \mathbb{C}^{n \times n}$  se numește *ortogonală* dacă

$$A^T A = A A^T = I_n.$$

O matrice  $A \in \mathbb{C}^{n \times n}$ ,  $A = (a_{ij})$  se numește matrice *triunghiulară inferior* (sau *inferior triunghiulară*) dacă

$$a_{ij} = 0 \text{ pentru } j > i$$

$$A = \begin{pmatrix} a_{11} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ a_{21} & a_{22} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ a_{31} & a_{32} & a_{33} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & & & & & \\ a_{(n-1)1} & a_{(n-1)2} & a_{(n-1)3} & \cdots & a_{(n-1)(n-1)} & \mathbf{0} \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{n(n-1)} & a_{nn} \end{pmatrix}$$

O matrice  $A \in \mathbb{C}^{n \times n}$ ,  $A = (a_{ij})$  se numește matrice *triunghiulară superior* (sau *superior triunghiulară*) dacă

$$a_{ij} = 0 \text{ pentru } j < i$$

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1(n-1)} & a_{1n} \\ 0 & a_{22} & a_{23} & \cdots & a_{2(n-1)} & a_{2n} \\ 0 & 0 & a_{33} & \cdots & a_{3(n-1)} & a_{3n} \\ \vdots & & & & & \\ 0 & 0 & 0 & \cdots & a_{(n-1)(n-1)} & a_{(n-1)n} \\ 0 & 0 & 0 & \cdots & 0 & a_{nn} \end{pmatrix}$$

Notăm cu  $I_n$  matricea unitate:

$$I_n \in \mathbb{R}^{n \times n}, I_n = \begin{pmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{1} \end{pmatrix}$$

Matrice diagonală  $D = \text{diag}[d_1, d_2, \dots, d_n]$

$$D \in \mathbb{R}^{n \times n}, D = \begin{pmatrix} d_1 & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & d_2 & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & d_{n-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & d_n \end{pmatrix}$$

## Norme

### Definiție

Fie  $X$  un spațiu vectorial real. Se numește *normă* aplicația:

$$\| \cdot \| : X \rightarrow \mathbb{R}_+$$

care îndeplinește condițiile:

- (1)  $\|x\| \geq 0$ ;  $\|x\| = 0 \Leftrightarrow x = 0$ ;
- (2)  $\|x + y\| \leq \|x\| + \|y\|, \forall x, y \in X$ ;
- (3)  $\|\lambda x\| = |\lambda| \|x\|, \forall x \in X, \forall \lambda \in \mathbb{R}$ .

Vom numi *norme vectoriale* normele definite pe spațiile  $X = \mathbb{C}^n$  sau  $\mathbb{R}^n$ .

## Exemple

Fie spațiile vectoriale  $\mathbb{C}^n$  sau  $\mathbb{R}^n$ . Pe aceste spații următoarele aplicații sunt norme vectoriale:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |\mathbf{x}_i| - \text{city-block (Manhattan)}$$

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n |\mathbf{x}_i|^2} - \text{euclidiană}$$

$$\|\mathbf{x}\|_\infty = \max\{|\mathbf{x}_i|, i = 1..n\} - \text{Cebyshev (a tablei de șah)}$$

Dacă  $\|\cdot\|_v$  este o normă vectorială și  $P \in \mathbb{R}^{n \times n}$  este o matrice nesingulară ( $\det P \neq 0$ ) atunci aplicația:

$$\|\cdot\|_P : \mathbb{R}^n \rightarrow \mathbb{R}, \quad \|x\|_P = \|Px\|_v$$

este de asemenea o normă vectorială.



## Definiție

Se numește *produs scalar* în spațiul vectorial  $X$  aplicația:

$$(\cdot, \cdot): X \times X \rightarrow K$$

care satisface condițiile :

$$(a) \quad (x, x) \geq 0, \forall x \in X, \quad (x, x) = 0 \Leftrightarrow x = 0;$$

$$(b) \quad (x, y) = \overline{(y, x)}, \forall x, y \in X,$$

$$(c) \quad (\lambda x, y) = \lambda (x, y), \forall x, y \in X, \forall \lambda \in K,$$

$$(d) \quad (x + y, z) = (x, z) + (y, z), \forall x, y, z \in X.$$

*Inegalitatea lui Cauchy-Buniakovski-Schwarz:*

$$|(\mathbf{x}, \mathbf{y})| \leq \sqrt{(\mathbf{x}, \mathbf{x})} \sqrt{(\mathbf{y}, \mathbf{y})} \quad \forall \mathbf{x}, \mathbf{y} \in X$$

Într-un spațiu vectorial dotat cu produs scalar se poate induce o normă numită euclidiană:

$$\|\mathbf{x}\|_2 = |\mathbf{x}| := \sqrt{(\mathbf{x}, \mathbf{x})}.$$

Reamintim definiția produselor scalare pe  $\mathbb{C}^n$  și pe  $\mathbb{R}^n$  introduse anterior:

$$(\mathbf{x}, \mathbf{y})_{\mathbb{C}^n} = \sum_{i=1}^n x_i \overline{y_i} \quad , \quad (\mathbf{x}, \mathbf{y})_{\mathbb{R}^n} = \sum_{i=1}^n x_i y_i$$

Obținem norma euclidiană (valabilă în spațiile  $\mathbb{C}^n$  și  $\mathbb{R}^n$ ):

$$\|\mathbf{x}\|_2 = |\mathbf{x}| = \sqrt{\sum_{i=1}^n |x_i|^2} .$$

## Norme matriceale

### Definiție

Aplicația  $\|\cdot\|: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  se numește *normă matriceală* dacă:

$$(1) \|A\| \geq 0 \quad \forall A \in \mathbb{R}^{n \times n} \quad ; \quad \|A\| = 0 \Leftrightarrow A = \mathbf{0}.$$

$$(2) \|\alpha A\| = |\alpha| \|A\| \quad , \quad \forall \alpha \in \mathbb{R} \quad , \quad \forall A \in \mathbb{R}^{n \times n} .$$

$$(3) \|A + B\| \leq \|A\| + \|B\| \quad , \quad \forall A, B \in \mathbb{R}^{n \times n} .$$

$$(4) \|A * B\| \leq \|A\| \cdot \|B\| \quad , \quad \forall A, B \in \mathbb{R}^{n \times n} .$$

## Exemple

Norma Frobenius definită de relația  $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}$  este o normă matriceală.

Aplicația  $\|A\|_{\max} = \max\{|a_{ij}|; i = 1, \dots, n, j = 1, \dots, n\}$  NU este o normă matriceală.

Pentru  $n = 2$  fie:

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}, B = A^T = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}$$

$$A * B = I_2, \|A\|_{\max} = \|B\|_{\max} = \frac{1}{\sqrt{2}}$$

$$\|A * B\|_{\max} = 1 > \|A\|_{\max} \cdot \|B\|_{\max} = \frac{1}{2}.$$

## Norme matriceale naturale

-  $\|\cdot\|_v : \mathbb{R}^n \rightarrow \mathbb{R}_+$  o normă vectorială  $\rightarrow \|\cdot\|_i : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}_+$  *normă matriceală naturală* sau *indusă*.

$$\|A\|_i = \max\left\{ \frac{\|Ax\|_v}{\|x\|_v} ; x \in \mathbb{R}^n, x \neq \mathbf{0} \right\}$$

Definiții echivalente :

$$\begin{aligned} \|A\|_i &= \max\{ \|Ax\|_v ; x \in \mathbb{R}^n, \|x\|_v \leq 1 \} \\ &= \max\{ \|Ax\|_v ; x \in \mathbb{R}^n, \|x\|_v = 1 \} \end{aligned}$$

$\|A\|_i$  se numește *normă matriceală naturală* sau *normă indusă* de norma vectorială  $\|\cdot\|_v$

Avem următoarea relație:

$$\|Ax\|_v \leq \|A\|_i \|x\|_v, \forall A \in \mathbb{R}^{n \times n}, \forall x \in \mathbb{R}^n.$$

Norma Frobenius  $\|\cdot\|_F$  nu este o normă naturală.

$$\|I_n\|_i = \max\left\{ \frac{\|I_n x\|_v}{\|x\|_v}; x \neq \mathbf{0} \right\} = 1, \quad \forall \|\cdot\|_i,$$

$$\|I_n\|_F = (1 + 1 + \cdots + 1)^{1/2} = \sqrt{n} \neq 1 \text{ pentru } n \geq 2.$$



Pentru  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$  norma matriceală indusă este:

$$\|A\|_1 = \max\left\{\sum_{i=1}^n |a_{ij}|; j = 1, 2, \dots, n\right\}$$

Pentru  $\|\mathbf{x}\|_\infty = \max\{|x_i|; i = 1, \dots, n\}$  norma matriceală indusă este:

$$\|A\|_\infty = \max\left\{\sum_{j=1}^n |a_{ij}|; i = 1, 2, \dots, n\right\}.$$

-  $\|\cdot\|_{\mathbf{v}}$  și  $\|\cdot\|_{\mathbf{v},\mathbf{P}}$  - norme vectoriale  $\rightarrow$   $\|\cdot\|_{\mathbf{i}}$  și respectiv  $\|\cdot\|_{\mathbf{i},\mathbf{P}}$   
 norme matriciale induse

$$\|\cdot\|_{\mathbf{v}} \rightarrow \|\cdot\|_{\mathbf{i}}$$

$$\downarrow \quad ?$$

$$\|\cdot\|_{\mathbf{v},\mathbf{P}} \rightarrow \|\cdot\|_{\mathbf{i},\mathbf{P}}$$

$$\|\mathbf{x}\|_{\mathbf{v},\mathbf{P}} = \|\mathbf{P}\mathbf{x}\|_{\mathbf{v}} \rightarrow \|A\|_{\mathbf{i},\mathbf{P}} = \|PAP^{-1}\|_{\mathbf{i}}$$

## Valori și vectori proprii

### Definiții

Fie  $A \in \mathbb{R}^{n \times n}$ . Se numește *valoare proprie (autovaloare)* a matricei  $A$  un număr complex  $\lambda \in \mathbb{C}$  pentru care există un vector nenul  $x \in \mathbb{C}^n$ ,  $x \neq \mathbf{0}$  a.î.:

$$Ax = \lambda x.$$

Vectorul  $x$  se numește *vector propriu (autovector)* asociat val. proprii  $\lambda$ .

$$Ax = \lambda x \Leftrightarrow (\lambda I_n - A)x = \mathbf{0}, x \neq \mathbf{0} \Leftrightarrow \det(\lambda I_n - A) = 0$$

→ Matricea  $\lambda I_n - A$  este singulară.

Polinomul:

$$p_A(\lambda) = \det(\lambda I_n - A) = \lambda^n - a_1 \lambda^{n-1} - a_2 \lambda^{n-2} - \dots - a_{n-1} \lambda - a_n$$

se numește *polinom caracteristic* asociat matricei  $A$ .

→ **grad**  $p_A = n$  → are  $n$  rădăcini care sunt valorile proprii ale matricei  $A$ .

Se numește *rază spectrală* a matricei  $A$ :

$$\rho(A) = \max\{|\lambda_i|, i = 1, \dots, n, \lambda_i - \text{valorile proprii ale matricei } A\}$$

$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2} \quad \text{norma indusă este}$$

$$\|A\|_2 = |A| = \sqrt{\rho(A^T A)} \quad \text{se numește } \mathbf{norma\ spectrală}.$$

## Propoziție

Fie  $\|\cdot\|$  o normă matriceală naturală. Atunci:

$$\rho(A) \leq \|A\|, \forall A \in \mathbb{R}^{n \times n}.$$

## Numere în format binar

În 1985 IEEE a publicat un raport numit Binary Floating Point Arithmetic Standard 754-1985 și o actualizare în 2008 IEEE 754-2008 care furnizează standarde pentru numere în virgulă mobilă binare și decimale, formate de interschimbare a tipului de date, algoritmi de rotunjire aritmetică, tratarea excepțiilor. Aceste standarde sunt respectate de toți fabricanții de calculatoare care folosesc arhitectura în virgulă mobilă.



Cel mai mic număr pozitiv care poate fi reprezentat este cu  $s = 0, c = 1, f = 0$  adică

$$z = 2^{-1022} (1 + 0) \approx 0.22251 \times 10^{-307}$$

iar cel mai mare este pentru  $s = 0, c = 2046, f = 1 - 2^{-52}$

$$Z = 2^{1023} (2 - 2^{-52}) \approx 0.17977 \times 10^{309}.$$

Numerele care apar în calcule și sunt mai mici decât  $z$  sunt setate în general la 0 (*underflow*) iar cele mai mari decât  $Z$  duc, de obicei, la oprirea calculelor (*overflow*).



Se observă că numărul 0 are două reprezentări:  
 $s = 0, c = 1, f = 0$  și  $s = 1, c = 1, f = 0$ .

### Reprezentarea zecimală

$$\pm 0.d_1 d_2 \dots d_k \times 10^n \quad 1 \leq d_1 \leq 9, \quad 0 \leq d_i \leq 9, \quad i = 2, \dots, k \quad -$$

reprezentarea zecimală folosind  $k$  cifre. Orice număr real  $y$ :

$$y = 0.d_1 d_2 \dots d_k d_{k+1} d_{k+2} \dots \times 10^n$$

poate fi reprezentat folosind  $k$  cifre printr-o simplă trunchiere

$$fl(y) = 0.d_1 d_2 \dots d_k \times 10^n .$$

O altă metodă de a obține o reprezentare cu  $k$  cifre este prin rotunjire:

$$fl(y) = 0.\delta_1\delta_2\dots\delta_k \times 10^n$$

Dacă  $d_{k+1} \geq 5$  se adaugă  $1$  la  $d_k$  pentru a obține  $fl(y)$  (*round up*), altfel se face trunchierea la  $k$  cifre (*round down*).

Un număr  $r^*$  aproximează numărul  $r$  cu  $t$  cifre exacte dacă  $t$  este cel mai mare întreg nenegativ pentru care:

$$\frac{|r - r^*|}{|r|} \leq 5 \times 10^{-t} .$$

În cazul trunchierii avem

$$\left| \frac{y - fl(y)}{y} \right| \leq 10^{-k+1}$$

iar când se face rotunjirea:

$$\left| \frac{y - fl(y)}{y} \right| \leq 0.5 \times 10^{-k+1}.$$

### **Operațiile elementare**

$$x +_c y = fl(fl(x) + fl(y))$$

$$x -_c y = fl(fl(x) - fl(y))$$

$$x \times_c y = fl(fl(x) \times fl(y))$$

$$x \div_c y = fl(fl(x) \div fl(y))$$

## Surse de erori în calculele 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 modelare

- erori de discretizare:

limita unui șir , suma unei serii , funcții neliniare  
aproximate de funcții liniare, aproximarea  
derivatei unei funcții ...

- simplificări în modelul matematic

idealizări , ignorarea unor parametri...

### 3. Erori în timpul calculelor:

- de rotunjire datorate capacității limitate de memorare a datelor, operațiile nu sunt efectuate exact.
- erori ale bibliotecilor folosite (,bug').

### 4. Erori umane (date, algoritmi, înțelegerea problemei)

## 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  $\delta_{a_2}$  :

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



# **Calcul Numeric**

**Cursul 3**

**2022**

*Anca Ignat*

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

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

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

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

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.

Eroarea relativă în datele de ieșire  $\approx$

Număr de condiționare  $\times$  Eroarea relativă în datele de intrare

Se consideră polinomul Wilkinson:

$$w(x) = (x - 1)(x - 2) \cdots (x - 20) = x^{20} - 210x^{19} + P_{18}(x)$$

Dacă se schimbă coeficientul **210** al lui  $x^{19}$  cu

$$-210 - 2^{-23} = -210.0000001192$$

soluțiile (cu 5 zecimale exacte) noului polinom sunt:

**1.00000 2.00000, 3.00000, 4.00000, 5.00000, 6.00001, 6.99970, 8.00727,  
8.91725, 20.84691, 10.09527 ± i0.64350 11.79363 ± i1.65233,  
13.99236 ± i2.51883, 16.73074 ± i2.81262, 19.50244 ± i1.94033**

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

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

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

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

$x$  = datele exacte,

$P(x)$  = soluția exactă folosind date exacte,

$\tilde{P}(x)$  = 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. De asemenea, 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 matricelor  $2 \times 2$ ,  $3 \times 3$ ,  $4 \times 4$ ,  $5 \times 5$  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 :

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

(condiția ca determinantul matricei 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 matricelor 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 formei pătraticе) – *Disquisitiones arithmeticae*(1801); descrie operațiile de înmulțire matricială și inversă a unei matrice în contextul tabloului coeficienților unei forme pătraticе. 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 matrice. Introduce noțiunea de matrice asemenea și demonstrează ca acestea au aceeași ecuație caracteristică. Demonstrează că orice matrice reală simetrică este diagonalizabilă.
- 1850, Sylvester introduce pentru prima dată 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 abstractă a unei matrice. Studiază transformările liniare și compunerea lor ceea ce îl conduce la operațiile cu matrice (adunare, înmulțire, înmulțirea cu un scalar, inversa)
- 1858, Cayley în *Memoriu asupra teoriei matricelor* : „Sunt multe lucruri de spus despre această teorie a matricelor și, după părerea mea, această teorie ar trebui să preceadă teoria determinantilor”

- Jordan (1870 – Treatise on substitutions and algebraic equations – forma canonică Jordan), Frobenius (1878 – On linear substitutions and bilinear forms, rangul unei matrici)
- 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 matrice
- 1958, Wilkinson dezvoltă factorizarea  $QR$

....

## Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$  pentru care există o normă matricială naturală astfel ca  $\|A\| < \mathbf{1}$ . Atunci există matricele  $(I_n \pm A)^{-1}$  și avem evaluările:

$$\frac{\mathbf{1}}{\mathbf{1} + \|A\|} \leq \|(I \pm A)^{-1}\| \leq \frac{\mathbf{1}}{\mathbf{1} - \|A\|}.$$

## 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.

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

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

### Propoziție

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

nesingulară și avem:



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

Demonstrație. Avem:

$$\|\Delta A\| < \frac{\mathbf{1}}{\|A^{-1}\|} \Rightarrow \|A^{-1} \Delta A\| \leq \|A^{-1}\| \cdot \|\Delta A\| < \mathbf{1} \stackrel{\text{Prop}}{\Rightarrow} \exists \left( I + A^{-1} \Delta A \right)^{-1}$$

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

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

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

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 matricei  $A$ .

### Propoziție

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  $\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_n}$ .

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

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

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

$A$  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**):

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

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



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

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

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

$$\det A = 1$$

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

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

$$\begin{cases} 400x - 201y = 200 \\ -800x + 401y = -200 \end{cases}, \quad \begin{cases} 401x - 201y = 200 \\ -800x + 401y = -200 \end{cases}$$

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

$$k_2(A) = 2503 \qquad \qquad \qquad k_2(A) = 1002000$$

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

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

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

## Matricea Hilbert

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

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

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  $\mathbf{b} \in \mathbb{R}^n$ . Rezolvarea sistemului de ecuații liniare  $A\mathbf{x}=\mathbf{b}$  se poate face folosind *regula lui Cramer*:

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

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

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  $x$ . Unul din cele mai folosiți algoritmi este **algoritmul de eliminare Gauss** :

$Ax=b \Leftrightarrow \tilde{A}x = \tilde{b}$  cu  $\tilde{A}$  matrice superior triunghiulară

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

n = 10

Eliminarea Gauss

04.03.2022 21:02:25

Timp in milisecunde (eliminare Gauss) 2

04.03.2022 21:02:25

Eroarea  $|Ax-b| = 1,00442975964863E-14$

Eroarea  $|x-x_{exact}| = 7,1581622508791E-14$

Solutia exacta : 1,2, ...

Regula lui Cramer

04.03.2022 21:02:25

Timp in milisecunde (Cramer) 4364

04.03.2022 21:02:29

Eroarea  $|Ax-b| = 2,85393991926588E-11$

Eroarea  $|x-x_{exact}| = 3,99554995866884E-11$

n =



```
C:\D\cn info3 2021-2022 sem ii\curs\curs 03\Cramer.exe
n = 11
Eliminarea Gauss
05.03.2022 16:14:13
Timp in milisecunde (eliminare Gauss) 0
05.03.2022 16:14:13
Eroarea |Ax-b| = 1,56226172104547E-14
Eroarea |x-x_exact| = 4,84508249897898E-13
Solutia exacta : 1,2, ...
Regula lui Cramer
05.03.2022 16:14:13
Timp in milisecunde (Cramer) 50199
05.03.2022 16:15:04
Eroarea |Ax-b| = 9,72812910066041E-10
Eroarea |x-x_exact| = 7,05987026507931E-10
n =
```

n = 12

Eliminarea Gauss

04.03.2022 21:52:28

Timpe in miliseconde (eliminarea Gauss) 0

04.03.2022 21:52:28

Eroarea  $|Ax-b| = 2,01688903424503E-14$

Eroarea  $|x-x_{exact}| = 1,34706499390866E-13$

Solutia exacta : 1,2, ...

Regula lui Cramer

04.03.2022 21:52:28

Timpe in miliseconde (Cramer) 921179

04.03.2022 22:07:53

Eroarea  $|Ax-b| = 1,03068652311345E-09$

Eroarea  $|x-x_{exact}| = 5,65280022570973E-10$

n =

n = 13

Eliminarea Gauss

05.03.2022 11:42:13

Timp in milisecunde (eliminare Gauss) 2

05.03.2022 11:42:13

Eroarea  $|Ax-b| = 2,36547860027221E-14$

Eroarea  $|x-x_{exact}| = 6,02527783109332E-14$

Solutia exacta : 1,2, ...

Regula lui Cramer

05.03.2022 11:42:13

Timp in milisecunde (Cramer) 12236765

05.03.2022 15:06:09

Eroarea  $|Ax-b| = 3,44101391273652E-09$

Eroarea  $|x-x_{exact}| = 9,95924703928746E-10$

n =

C:\D\cn info3 2021-2022 sem ii\curs\curs 03\Cramer.exe

n = 8000

Eliminarea Gauss

05.03.2022 16:17:20

  Timp in milisecunde (eliminare Gauss) 2285277

05.03.2022 16:55:25

Eroarea  $|Ax-b| = 1,45671149537004E-05$

Eroarea  $|x-x_{exact}| = 0,000268972468485747$

Solutia exacta : 1,2, ...

Regula lui Cramer

05.03.2022 16:55:26

## Metoda substituției

Fie sistemul liniar  $Ax = b$  unde matricea sistemului  $A$  este triunghiulară. Pentru a găsi soluția unică a sistemului, trebuie ca matricea să fie nesingulară. Determinantul matricelor 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 :

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

$\vdots$

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

$\vdots$

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

$$a_{nn}x_n = b_n$$

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  $i$ :

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

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

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

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

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

# **Calcul Numeric**

**Cursul 4**

**2022**

*Anca Ignat*

## 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:

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



## *Pas 2*

se construiește sistemul

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

⋮

## *Pasul $r$*

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

⋮

## *Pasul $n-1$ :*

se obține sistemul

$A^{(n-1)} \mathbf{x} = \mathbf{b}^{(n-1)} \sim \mathbf{Ax} = \mathbf{b}$ , 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 ecuații și/sau două coloane în matricea  $A$ .

## Pasul 1

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

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

Fie ecuația  $i$ , cu  $i=1, \dots, n$

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

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:

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

⋮

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

⋮

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

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

Se presupune  $a_{22}^{(1)} \neq \mathbf{0}$  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 :

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

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

Sistemul are următoarea formă:



$$\left\{ \begin{array}{l}
 a_{11}^{(r-1)} x_1 + \cdots + a_{1r}^{(r-1)} x_r + \cdots + a_{1n}^{(r-1)} x_n = b_1^{(r-1)} \\
 \cdot \\
 \cdot \\
 \cdot \\
 \cdot \\
 \cdot \\
 a_{rr}^{(r-1)} x_r + \cdots + a_{rn}^{(r-1)} x_n = b_r^{(r-1)} \\
 a_{r+1r}^{(r-1)} x_r + \cdots + a_{r+1n}^{(r-1)} x_n = b_{r+1}^{(r-1)} \\
 \vdots \\
 a_{ir}^{(r-1)} x_r + \cdots + a_{in}^{(r-1)} x_n = b_i^{(r-1)} \\
 \vdots \\
 a_{nr}^{(r-1)} x_r + \cdots + a_{nn}^{(r-1)} x_n = b_n^{(r-1)}
 \end{array} \right.$$

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 (stabilitate + $a_{rr}^{(r-1)} \neq 0$ ?)

### 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 ecuațiile  $i_0$  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  $i_0$ :

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

și se inter-schimbă ecuațiile  $i_0, r$  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 inter-schimbă ecuațiile  $i_0, r$  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

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

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

Deoarece operațiile efectuate (cele de interschimbare de ecuații 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:



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

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

$\det A = \pm \det A^{(r-1)} = 0 \Rightarrow A$  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. De asemenea, 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:

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

## **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:

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

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

## Descompuneri $LU$

$$A \in \mathbb{R}^{n \times n}, \quad A=L \cdot U,$$

$L$  inferior triunghiulară și  $U$  superior triunghiulară

$$L, U \in \mathbb{R}^{n \times n}$$

$$Ax = b \Leftrightarrow L U x = b \Leftrightarrow \begin{cases} Ly = b \Rightarrow \text{soluția } y^* \\ Ux = y^* \Rightarrow \text{soluția } x^* \end{cases}, \quad x^* = A^{-1}b$$



Fie minorul principal principal al matricei  $A$ :

$$A_p = \det \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & & & \\ a_{p1} & a_{p2} & \cdots & a_{pp} \end{bmatrix} \in \mathbb{R}^{p \times p}, \quad p = 1, \dots, n$$

## Teoremă (descompunere $LU$ )

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice reală pătratică de dimensiune  $n$  astfel încât  $\det A_p \neq 0, \forall p = 1, \dots, n$ . Atunci există o unică matrice inferior triunghiulară  $L = (l_{ij})_{i,j=1,\dots,n}$  și o unică matrice superior triunghiulară  $U = (u_{ij})_{i,j=1,\dots,n}$  cu  $u_{ii} = 1, i = 1, \dots, n$  astfel încât

$$A = L \cdot U \quad (1)$$

Demonstrație. *Existența*: demonstrația se face prin inducție după  $n$  dimensiunea matricii  $A$ .

## *Algoritmul Crout de calcul al descompunerii LU*

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice reală pătratică de dimensiune  $n$  care satisface ipotezele teoremei de mai sus. Algoritmul de calcul al matricelor  $L$  și  $U$  are  $n$  etape. La fiecare pas se determină simultan:

- câte o coloană din matricea  $L$  și
- câte o linie din matricea  $U$ .

Descriem în continuare, un pas oarecare.

## **Pasul $p$ ( $p = 1, 2, \dots, n$ )**

Se determină elementele coloanei  $p$  ale matricii  $L$ :

$$l_{ip}, i = p, \dots, n$$

și elementele liniei  $p$  ale matricii  $U$ ,

$$u_{pp} = 1, u_{pi}, i = p + 1, \dots, n,$$

$$(u_{pi} = l_{ip} = 0, i = 1, \dots, p-1).$$

$$\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ l_{pp} \\ l_{p+1p} \\ \vdots \\ l_{np} \end{pmatrix}$$

col.  $p$  a matr.  $L$

$$\left( \mathbf{0} \quad \mathbf{0} \quad \dots \quad \mathbf{1} \quad u_{pp+1} \quad \dots \quad u_{pn} \right) \text{ lin. } p \text{ a matr. } U$$

Se cunosc de la pașii anteriori:

- elementele primelor  $p-1$  coloane din  $L$   
(elemente  $l_{ik}$  cu  $k = 1, \dots, p-1, \forall j$ ).
- elementele primelor  $p-1$  linii din  $U$   
(elemente  $u_{kj}$  cu  $k = 1, \dots, p-1, \forall j$ )

Calculul elementelor *coloanei*  $p$  din matricea  $L$ ,  $l_{ip}$   $i = p, \dots, n$  se face folosind elementul  $a_{ip}$  și  $(LU)_{ip}$ . Avem:

$$\begin{aligned}
a_{ip} &= (LU)_{ip} = \sum_{k=1}^n l_{ik} u_{kp} \quad (u_{kp} = 0, k = p+1, \dots, n) = \sum_{k=1}^p l_{ik} u_{kp} = \\
&= \sum_{k=1}^{p-1} l_{ik} u_{kp} + l_{ip} u_{pp} = \sum_{k=1}^{p-1} l_{ik} u_{kp} + l_{ip} \quad (u_{pp} = 1)
\end{aligned}$$

Pentru  $i = p, \dots, n$  avem:

$$l_{ip} = a_{ip} - \sum_{k=1}^{p-1} l_{ik} u_{kp}, \quad i = p, \dots, n \quad (2)$$

( $u_{pp} = 1$ ,  $l_{ik}$ ,  $u_{kp}$   $k = 1, \dots, p-1$  sunt elemente de pe coloane din  $L$  și linii din  $U$  calculate la pașii anteriori)

Calculul elementelor *liniei p* din matricea  $U$ :

$$u_{pi}, i = p + 1, \dots, n \quad (u_{pi} = 0, i = 1, \dots, p - 1, u_{pp} = 1)$$

se face analog:

$$\begin{aligned} a_{pi} &= (LU)_{pi} = \sum_{k=1}^n l_{pk} u_{ki} \quad (l_{pk} = 0, k = p + 1, \dots, n) = \sum_{k=1}^p l_{pk} u_{ki} = \\ &= \sum_{k=1}^{p-1} l_{pk} u_{ki} + l_{pp} u_{pi} \end{aligned}$$

Dacă  $l_{pp} \neq 0$  putem calcula elementele nenule ale coloanei  $p$  din matricea  $U$  astfel:



$$u_{pi} = \frac{(a_{pi} - \sum_{k=1}^{p-1} l_{pk} u_{ki})}{l_{pp}}, \quad i = p+1, \dots, n \quad (3)$$

(elementele  $l_{pk}$ ,  $u_{ki}$   $k = 1, \dots, p-1$  sunt calculate anterior pasului  $p$ )

Dac  $l_{pp} = 0$ , calculele se opresc, descompunerea  $LU$  nu poate fi calculată - matricea  $A$  are un minor  $A_p$  cu determinantul  $0$ .

**Unicitatea:** Demonstrație prin reducere la absurd.

Facem observația că inversa unei matrici nesingulare triunghiulară inferior (superior) este o matrice de același tip.

Presupunem că

$$A = L \cdot U = L_1 \cdot U_1 \quad (4)$$

Din ipoteza  $A$  nesingulară rezultă existența inverselor matricelor  $L, L_1, U, U_1$ . Înmulțind egalitatea (4) la stânga cu  $L^{-1}$  și cu  $U_1^{-1}$  la dreapta obținem

$$U U_1^{-1} = L^{-1} L_1.$$

Matricea  $U U_1^{-1}$  este superior triunghiulară cu elementele diagonale egale cu  $I$  iar matricea  $L^{-1} L_1$  este inferior triunghiulară. Rezultă că:

$$U U_1^{-1} = L^{-1} L_1 = I_n, \quad \text{deci} \quad L=L_1, \quad U=U_1.$$

Numărul de operații efectuate:

$A$  (adunări, scăderi):

$$\sum_{p=1}^n [(n-p+1)(p-1) + (n-p)(p-1)] = \frac{n(n-1)(2n-1)}{6} = \frac{1}{3}n^3 + O(n^2)$$

$M$  (înmulțiri, împărțiri):

$$\sum_{p=1}^n [(n-p+1)(p-1) + (n-p)p] = \frac{(n-1)n(n+1)}{3} = \frac{1}{3}n^3 + O(n^2)$$

## Descompunerea Cholesky

O matrice  $A \in \mathbb{R}^{n \times n}$  se numește *pozitiv definită* dacă:

$$(Ax, x)_{\mathbb{R}^n} > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0$$

Notăție:  $A > 0$

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice simetrică ( $A=A^T$ ) și pozitiv definită.

*Descompunerea Cholesky* pentru matricea  $A$  este de forma:

$$A = LL^T, \quad L \text{ matrice inferior triunghiulară}$$

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} = \mathbf{L}\mathbf{L}^T = \begin{pmatrix} l_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ l_{21} & l_{22} & \cdots & \mathbf{0} \\ \vdots & & & \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & \cdots & l_{n1} \\ \mathbf{0} & l_{22} & \cdots & l_{n2} \\ \vdots & & & \\ \mathbf{0} & \mathbf{0} & \cdots & l_{nn} \end{pmatrix}$$

Matricea  $\mathbf{L}$  se calculează în  $n$  pași, coloană după coloană.

## Pas $r$ ( $r=1, \dots, n$ )

Se calculează elementele coloanei  $r$  a matricii  $L$ :  
întâi elementul diagonal  $l_{rr}$  apoi  
celelalte elemente  $l_{ir}$  ( $i=r+1, \dots, n$ )

Coloana  $r$  a matricii  $L$ :

$$\left( \mathbf{0} \quad \dots \quad \mathbf{0} \quad l_{rr} \quad l_{r+1r} \quad \dots \quad l_{ir} \quad \dots \quad l_{nr} \right)^T$$

- se cunosc elementele primelor  $(r-1)$  coloane ale matricii  $L$

Calcul  $l_{rr}$ :

$$a_{rr} = (LL^T)_{rr} = (l_{r1} \quad \dots \quad l_{rr-1} \quad l_{rr-1} \quad \mathbf{0} \quad \dots \quad \mathbf{0}) \begin{pmatrix} l_{r1} \\ \vdots \\ l_{rr-1} \\ l_{rr} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix} =$$

$$= l_{r1}^2 + l_{r2}^2 + \dots + l_{rr-1}^2 + l_{rr}^2 \quad \Rightarrow \quad l_{rr} = \pm \sqrt{a_{rr} - \sum_{k=1}^{r-1} l_{rk}^2}$$

Calcul  $l_{ir}$  ( $i=r+1, \dots, n$ ):

$$\begin{aligned}
\mathbf{a}_{ir} &= (\mathbf{L}\mathbf{L}^T)_{ir} = (l_{i1} \quad \cdots \quad l_{ir-1} \quad l_{ir} \quad \cdots \quad l_{ii} \quad \cdots \quad \mathbf{0}) \begin{pmatrix} l_{r1} \\ \vdots \\ l_{rr-1} \\ l_{rr} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix} = \\
&= l_{i1}l_{r1} + l_{i2}l_{r2} + \cdots + l_{ir-1}l_{rr-1} + l_{ir}l_{rr} \quad \Rightarrow \quad l_{ir} = \frac{\left( \mathbf{a}_{ir} - \sum_{k=1}^{r-1} l_{ik}l_{rk} \right)}{l_{rr}}
\end{aligned}$$



# **Calcul Numeric**

**Cursul 5**

**2022**

*Anca Ignat*

## **Algoritmul de eliminare Gauss fără schimbare de ecuații ⇔ descompunere LU**

Presupunem că la fiecare pas al algoritmului de eliminare Gauss pivotul este nenul ( $a_{rr}^{(r-1)} \neq 0$ ), deci nu e nevoie de schimbare de ecuații.

Algoritmul se poate scrie astfel:

**for  $r = 1, \dots, n - 1$**

**for  $i = r + 1, \dots, n$**

- $f = - \frac{a_{ir}}{a_{rr}};$

//  $E_i = E_i + f * E_r$

- **for  $j = r + 1, \dots, n$**

$$a_{ij} = a_{ij} + f * a_{rj};$$

- $a_{ir} = 0;$

- $b_i = b_i + f * b_r;$

Considerăm vectorul și matricea:

$$\mathbf{t}^{(r)} = \begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ t_{r+1}^{(r)} \\ \vdots \\ t_n^{(r)} \end{pmatrix} \in \mathbb{R}^n, \quad \mathbf{T}_r := \mathbf{I}_n + \mathbf{t}^{(r)} \mathbf{e}_r^T \in \mathbb{R}^{n \times n}$$

$$\mathbf{t}^{(r)} \mathbf{e}_r^T = \begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ t_{r+1}^{(r)} \\ \vdots \\ t_n^{(r)} \end{pmatrix} (\mathbf{0} \ \dots \ \mathbf{1} \ \mathbf{0} \ \dots \ \mathbf{0}) = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \dots & \overset{\text{col } r}{\mathbf{0}} & \dots & \mathbf{0} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & t_{r+1}^{(r)} & \dots & \mathbf{0} - \text{lin}(r+1) \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \dots & t_n^{(r)} & \dots & \mathbf{0} \end{pmatrix}$$

Matricea  $T_r$  este matrice triunghiulară inferior cu  $\mathbf{1}$  pe diagonala principală:

$$T_r = \begin{pmatrix} \mathbf{1} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{1} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & t_{r+1}^{(r)} & \dots & \mathbf{0} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \dots & t_n^{(r)} & \dots & \mathbf{1} \end{pmatrix}$$

Inversa matricei  $T_r$  este

$$T_r^{-1} = I_n - t^{(r)} e_r^T .$$

$$\begin{aligned} T_r T_r^{-1} &= (I_n + t^{(r)} e_r^T)(I_n - t^{(r)} e_r^T) = I_n - t^{(r)} e_r^T + t^{(r)} e_r^T - t^{(r)} e_r^T t^{(r)} e_r^T = \\ &= I_n - t^{(r)} t_r^{(r)} e_r^T = I_n \quad (0 = t_r^{(r)} = e_r^T t^{(r)}) \end{aligned}$$

Dacă  $A$  este o matrice oarecare, vrem să vedem cum se poate construi matricea  $B = T_r A$  fără a face înmulțire matricială. Vom studia legătura între liniile matricelor  $A$  și  $B$ .

$$\begin{aligned} \mathbf{e}_i^T \mathbf{B} &= \mathbf{e}_i^T (\mathbf{T}_r \mathbf{A}) = \mathbf{e}_i^T (\mathbf{I}_n + t^{(r)} \mathbf{e}_r^T) \mathbf{A} = \mathbf{e}_i^T \mathbf{A} + \mathbf{e}_i^T t^{(r)} \mathbf{e}_r^T \mathbf{A} = \\ &= \mathbf{e}_i^T \mathbf{A} + t_i^{(r)} (\mathbf{e}_r^T \mathbf{A}) \end{aligned}$$

Linia  $i$  a noii matrice  $\mathbf{B}$  se obține din linia  $i$  a matricei  $\mathbf{A}$  la care se adaugă linia  $r$  a matricei  $\mathbf{A}$  înmulțită cu factorul  $t_i^{(r)}$ .

$$\mathbf{e}_i^T \mathbf{B} = \begin{cases} \mathbf{e}_i^T \mathbf{A} & i = 1, \dots, r \quad (t_i^{(r)} = 0) \\ \mathbf{e}_i^T \mathbf{A} + t_i^{(r)} (\mathbf{e}_r^T \mathbf{A}) & i = r + 1, \dots, n \end{cases} .$$



Operația  $T_r A$  descrie Pasul  $r$  al algoritmului de eliminare Gauss dacă:

$$t_{r+1}^{(r)} = -\frac{a_{r+1r}^{(r)}}{a_{rr}^{(r)}}, \dots, t_i^{(r)} = -\frac{a_{ir}^{(r)}}{a_{rr}^{(r)}}, \dots, t_n^{(r)} = -\frac{a_{nr}^{(r)}}{a_{rr}^{(r)}}.$$

Algoritmul de eliminare Gauss fără schimbare de ecuații poate fi descris astfel:

$$T_{n-1} \cdots T_2 T_1 A = U \quad \text{cu} \quad T_r = I_n + t^{(r)} e_r^T,$$

$$t^{(r)} = \left( \mathbf{0} \ \cdots \ \mathbf{0} \ \left(-\frac{a_{r+1r}^{(r)}}{a_{rr}^{(r)}}\right) \ \cdots \ \left(-\frac{a_{ir}^{(r)}}{a_{rr}^{(r)}}\right) \ \cdots \ \left(-\frac{a_{nr}^{(r)}}{a_{rr}^{(r)}}\right) \right)^T.$$

Avem:

$$A = T_1^{-1}T_2^{-1} \cdots T_{n-1}^{-1}U = LU \quad , \quad L := T_1^{-1}T_2^{-1} \cdots T_{n-1}^{-1}$$

$$\begin{aligned} T_1^{-1}T_2^{-1} &= (I_n - t^{(1)}e_1^T)(I_n - t^{(2)}e_2^T) = I_n - t^{(1)}e_1^T - t^{(2)}e_2^T + t^{(1)}e_1^T t^{(2)}e_2^T = \\ &= I_n - t^{(1)}e_1^T - t^{(2)}e_2^T + t^{(1)}t_1^{(2)}e_2^T = I_n - t^{(1)}e_1^T - t^{(2)}e_2^T \quad (t_1^{(2)} = 0) \end{aligned}$$

Prin inducție se arată că:

$$L = T_1^{-1}T_2^{-1} \cdots T_{n-1}^{-1} = I_n - t^{(1)}e_1^T - t^{(2)}e_2^T - \cdots - t^{(n-1)}e_{n-1}^T$$

$$L = \begin{pmatrix} \mathbf{1} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \frac{a_{21}}{a_{11}} & \mathbf{1} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \frac{a_{31}}{a_{11}} & \frac{a_{32}^{(1)}}{a_{22}^{(1)}} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{a_{r1}}{a_{11}} & \frac{a_{r2}^{(1)}}{a_{22}^{(1)}} & \dots & \mathbf{1} & \dots & \mathbf{0} \\ \frac{a_{r+11}}{a_{11}} & \frac{a_{r+12}^{(1)}}{a_{22}^{(1)}} & \dots & \frac{a_{r+1r}^{(r-1)}}{a_{rr}^{(r-1)}} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{a_{n1}}{a_{11}} & \frac{a_{n2}^{(1)}}{a_{22}^{(1)}} & \dots & \frac{a_{nr}^{(r-1)}}{a_{rr}^{(r-1)}} & \dots & \mathbf{1} \end{pmatrix}$$

## Descompuneri $QR$

### Definiție

Se numește *matrice ortogonală*, o matrice  $Q \in \mathbb{R}^{n \times n}$  care satisface relația:

$$Q^T Q = Q Q^T = I_n \quad (Q^{-1} = Q^T) .$$

Matricele ortogonale au următoarele proprietăți:

- Dacă  $Q$  este matrice ortogonală atunci și matricea transpusă  $Q^T$  este ortogonală.

$$Q^T Q = Q^T (Q^T)^T = Q Q^T = (Q^T)^T Q^T = I_n$$

- Dacă  $Q_1$  și  $Q_2$  sunt matrice ortogonale atunci  $Q_1Q_2$  este tot matrice ortogonală.

$$(Q_1Q_2)^T (Q_1Q_2) = Q_2^T Q_1^T Q_1Q_2 = Q_2^T I Q_2 = I_n$$

$$(Q_1Q_2)(Q_1Q_2)^T = Q_1Q_2Q_2^T Q_1^T = Q_1^T I Q_1 = I_n$$

- Dacă  $Q \in \mathbb{R}^{n \times n}$  este matrice ortogonală și  $x \in \mathbb{R}^n$  atunci  $\|Qx\|_2 = \|x\|_2$ .

$$\|Qx\|_2^2 = (Qx, Qx) = (x, Q^T Qx) = (x, x) = \|x\|_2^2, \quad \|\cdot\|_2 \geq 0 \Rightarrow$$

$$\|Qx\|_2 = \|x\|_2$$

Fie  $A$  o matrice reală pătratică de dimensiune  $n$ . Pp. că avem:

$$A = QR$$

unde  $Q$  este o matrice ortogonală iar  $R$  este o matrice superior triunghiulară.

$$Ax = b \Leftrightarrow QRx = b \Leftrightarrow Q^T QRx = Q^T b \Leftrightarrow Rx = Q^T b$$

## Algoritmul lui Householder

*Matrice de reflexie* -  $P \in \mathbb{R}^{n \times n}$  de forma:

$$P = I_n - 2vv^T, \quad v \in \mathbb{R}^n, \quad \|v\|_2 = \sqrt{\sum_{j=1}^n |v_j|^2} = 1$$

$$vv^T = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} (v_1, v_2, \dots, v_n) = \begin{pmatrix} v_1^2 & v_1 v_2 & \dots & v_1 v_n \\ v_2 v_1 & v_2^2 & \dots & v_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ v_n v_1 & v_n v_2 & \dots & v_n^2 \end{pmatrix}.$$

Matricile de reflexie sunt :

simetrice -  $P = P^T$  și

ortogonale -  $PP^T = P^T P = P^2 = I_n$ .

$$P^T = (I_n - 2vv^T)^T = I_n - 2(vv^T)^T = I_n - 2(v^T)^T v^T = I_n - 2vv^T = P$$



$$\begin{aligned}
P^2 &= (I_n - 2vv^T)(I_n - 2vv^T) = I_n - 2vv^T - 2vv^T + 4(vv^T)(vv^T) = \\
&= I_n - 4vv^T + 4v(v^T v)v^T = I_n - 4vv^T + 4v\|v\|_2^2 v^T = \\
&= I_n - 4vv^T + 4vv^T = I_n \quad (\|v\|_2 = 1)
\end{aligned}$$

$$\begin{aligned}
n = 2, \quad v &= \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - 2 \begin{pmatrix} 1 \\ 0 \end{pmatrix} (1 \ 0) = \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}, \quad x \in \mathbb{R}^2, \quad y = Px \\
x &= \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -x_1 \\ x_2 \end{pmatrix}
\end{aligned}$$

Vectorul  $y=Px$  este reflectatul vectorului  $x$  în raport cu axa  $Ox_2$ .

Algoritmul care folosește matricele de reflexie pentru a obține o descompunere  $QR$  pentru o matrice  $A \in \mathbb{R}^{n \times n}$  a fost descris de Alston S. Householder în articolul "*Unitary triangularization of a nonsymmetric matrix*" apărut în Journal of the Assoc. of Computing Machinery 5 (1958), 339-342.

Transformarea matricei  $A$  într-una superior triunghiulară se face în  $(n-1)$  pași, la fiecare pas folosindu-se o matrice de reflexie.

**Pas 1:**  $A^{(1)} = P_1 A$  (matricea  $P_1$  se alege astfel încât col.  $1$  să fie transformată în formă superior triunghiulară)

**Pas 2:**  $A^{(2)} = P_2 A^{(1)} = P_2 (P_1 A)$  ( $P_2$  transformă col.  $2$  în formă superior triunghiulară, fără să schimbe col.  $1$ )

**Pas  $r$ :**  $A^{(r)} = P_r A^{(r-1)} = P_r (P_{r-1} \dots P_1 A)$  (se transformă col  $r$  în formă superior triunghiulară fără să schimbe primele  $(r-1)$  coloane)

Descompunerea  $QR$  construită cu algoritmul Householder este următoarea:

$$P_{n-1} \cdots P_r \cdots P_2 P_1 A = \tilde{Q}A = R$$

unde

$$\tilde{Q} = P_{n-1} \cdots P_r \cdots P_2 P_1$$

$\tilde{Q}$  este matrice ortogonală ca produs de matrice ortogonale.

$$\tilde{Q}A = R \Leftrightarrow \tilde{Q}^T \tilde{Q}A = \tilde{Q}^T R \Leftrightarrow A = \tilde{Q}^T R = QR$$

$$Q = \tilde{Q}^T = (P_{n-1} \cdots P_r \cdots P_2 P_1)^T = P_1 P_2 \cdots P_r \cdots P_{n-1}$$

## Pasul $r$

La intrarea în pasul  $r$  matricea  $A$  are forma:

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1r} & \cdots & a_{1n} \\ \mathbf{0} & a_{22} & \cdots & a_{2r} & \cdots & a_{2n} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & a_{rr} & \cdots & a_{rn} \\ \mathbf{0} & \mathbf{0} & \cdots & a_{r+1r} & \cdots & a_{r+1n} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & a_{ir} & \cdots & a_{in} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & a_{nr} & \cdots & a_{nn} \end{pmatrix}$$

Pasul  $r$  constă în:

$$A := P_r A$$

$$P_r = I_n - 2v^r (v^r)^T, \quad v^r \in R^n, \quad \|v^r\|_2 = 1$$

unde vectorul  $v^r$  se alege astfel ca matricea  $A$  să aibă și coloana  $r$  în formă superior triunghiulară:

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1r} & \cdots & a_{1n} \\ \mathbf{0} & a_{22} & \cdots & a_{2r} & \cdots & a_{2n} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & a_{rr} & \cdots & a_{rn} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & a_{r+1n} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & a_{in} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \cdots & a_{nn} \end{pmatrix}$$

## Calculul matricei $P_r$

Pentru simplitate vom nota  $P_r = P$ ,  $v^r = v$ .

$$Ae_r = \begin{pmatrix} a_{1r} \\ a_{2r} \\ \vdots \\ a_{r-1r} \\ a_{rr} \\ a_{r+1r} \\ \vdots \\ a_{ir} \\ \vdots \\ a_{nr} \end{pmatrix} \rightarrow (PA)e_r = \bar{A}e_r = \begin{pmatrix} \bar{a}_{1r} = a_{1r} \\ \bar{a}_{2r} = a_{2r} \\ \vdots \\ \bar{a}_{r-1r} = a_{r-1r} \\ \bar{a}_{rr} = k \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}$$



Aplicând proprietatea matricelor ortogonale :

$$Q \in \mathbb{R}^{n \times n}, \mathbf{x} \in \mathbb{R}^n, \|Q\mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

pentru matricea  $Q=P$  și vectorul  $\mathbf{x}=A\mathbf{e}_r$  avem:

$$\|P A \mathbf{e}_r\|_2^2 = a_{1r}^2 + a_{2r}^2 + \dots + a_{r-1r}^2 + k^2 =$$

$$\|A \mathbf{e}_r\|_2^2 = a_{1r}^2 + a_{2r}^2 + \dots + a_{r-1r}^2 + a_{rr}^2 + a_{r+1r}^2 + \dots + a_{ir}^2 + \dots + a_{nr}^2$$

Din relația de mai sus rezultă:

$$k^2 = \sigma = a_{rr}^2 + a_{r+1r}^2 + \dots + a_{ir}^2 + \dots + a_{nr}^2 = \sum_{i=r}^n a_{ir}^2 \Rightarrow k = \pm \sqrt{\sigma}$$

Determinarea vectorului  $\mathbf{v}$  ce definește matricea  $\mathbf{P}$ :

$$\begin{aligned}(\mathbf{PA})\mathbf{e}_r &= (\mathbf{I}_n - 2\mathbf{v}\mathbf{v}^T)(\mathbf{A}\mathbf{e}_r) = \mathbf{A}\mathbf{e}_r - 2(\mathbf{v}\mathbf{v}^T)(\mathbf{A}\mathbf{e}_r) = \\ &= \mathbf{A}\mathbf{e}_r - 2\mathbf{v}(\mathbf{v}^T(\mathbf{A}\mathbf{e}_r)) = \mathbf{A}\mathbf{e}_r - (2\alpha)\mathbf{v} = \mathbf{A}\mathbf{e}_r - \mathbf{u}\end{aligned}$$

unde cu  $\alpha$  și  $\mathbf{u}$  am notat:

$$\alpha := \mathbf{v}^T (A\mathbf{e}_r) = ((A\mathbf{e}_r), \mathbf{v})_{\mathbb{R}^n} = \left( \begin{pmatrix} \mathbf{a}_{1r} \\ \mathbf{a}_{2r} \\ \vdots \\ \mathbf{a}_{ir} \\ \vdots \\ \mathbf{a}_{nr} \end{pmatrix}, \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_i \\ \vdots \\ \mathbf{v}_n \end{pmatrix} \right)_{\mathbb{R}^n} =$$

$$= \mathbf{a}_{1r} \mathbf{v}_1 + \mathbf{a}_{2r} \mathbf{v}_2 + \cdots + \mathbf{a}_{ir} \mathbf{v}_i + \cdots + \mathbf{a}_{nr} \mathbf{v}_n$$

$$u := (2\alpha) v = Ae_r - (PA)e_r = \begin{pmatrix} a_{1r} \\ a_{2r} \\ \vdots \\ a_{r-1r} \\ a_{rr} \\ a_{r+1r} \\ \vdots \\ a_{ir} \\ \vdots \\ a_{nr} \end{pmatrix} - \begin{pmatrix} a_{1r} \\ a_{2r} \\ \vdots \\ a_{r-1r} \\ k \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ a_{rr} - k \\ a_{r+1r} \\ \vdots \\ a_{ir} \\ \vdots \\ a_{nr} \end{pmatrix}$$

Cu aceste notații, matricea  $P$  devine:

$$P = I_n - 2\left(\frac{1}{2\alpha}\right)u \left(\frac{1}{2\alpha}\right)u^T = I_n - \left(\frac{1}{2\alpha^2}\right)uu^T = I_n - \frac{1}{\beta}uu^T$$

$$\beta := 2\alpha^2$$

Pentru a cunoaște matricea  $P$  trebuie să mai determinăm constanta  $\beta$ . Din condiția:

$$\|v\|_2^2 = 1 \Rightarrow \left\|\frac{1}{2\alpha}u\right\|_2^2 = 1 \Rightarrow \frac{1}{4\alpha^2}\|u\|_2^2 = 1 \Rightarrow 2\beta = \|u\|_2^2$$

$$\begin{aligned}
\|u\|_2^2 &= \left\| \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ a_{rr} - k \\ a_{r+1r} \\ \vdots \\ a_{ir} \\ \vdots \\ a_{nr} \end{pmatrix} \right\|_2^2 = (a_{rr} - k)^2 + a_{r+1r}^2 + \cdots + a_{ir}^2 + \cdots + a_{nr}^2 = \\
&= a_{rr}^2 + a_{r+1r}^2 + \cdots + a_{ir}^2 + \cdots + a_{nr}^2 - 2ka_{rr} + k^2 = \\
&= \sigma - 2ka_{rr} + \sigma = 2(\sigma - ka_{rr})
\end{aligned}$$

de unde obținem:

$$\beta = \sigma - ka_{rr}$$

Vom alege semnul constantei  $k$  astfel încât  $\beta$  să fie cât mai mare posibil deoarece constanta  $\beta$  apare în operația de împărțire. Avem:

$$\beta \text{ "mare"} \rightarrow \beta = \sigma - ka_{rr} \geq \sigma \quad (\sigma \geq 0) \rightarrow$$
$$\text{semn } k = -\text{semn } a_{rr}$$

Ce înseamnă  $\beta = 0$  ?

$$\beta = \frac{1}{2} \|u\|_2^2 = 0 \rightarrow \|u\|_2^2 = 0 \rightarrow u = 0 \rightarrow$$

$$\rightarrow a_{rr} = k, a_{r+1r} = 0, \dots, a_{ir} = 0, \dots, a_{nr} = 0$$

Cum  $a_{rr}=k$  și  $\text{semn } k = -\text{semn } a_{rr}$  obținem:

$$a_{ir} = 0, \quad \forall i = r, \dots, n$$

adică avem coloana  $r$  deja în formă superior triunghiulară, se poate trece la pasul următor. În acest caz matricea  $A$  este singulară.



Ne interesează cum se efectuează operația  $A=PrA$  fără a face înmulțire matricială. Vom pune în evidență schimbările în raport cu coloanele.

$$\begin{aligned}
 (PA)e_j &= \text{noua col. } j \text{ a matr. } A = \left(I_n - \frac{1}{\beta}uu^T\right)(Ae_j) = \\
 &= Ae_j - \frac{1}{\beta}(uu^T)(Ae_j) = Ae_j - \frac{1}{\beta}u(u^T(Ae_j)) = \\
 &= Ae_j - \frac{\gamma_j}{\beta}u
 \end{aligned}$$

$$\gamma_j := \mathbf{u}^T (A \mathbf{e}_j) = \left( \begin{array}{c} \mathbf{a}_{1j} \\ \mathbf{a}_{2j} \\ \vdots \\ \mathbf{a}_{rj} \\ \mathbf{a}_{r+1j} \\ \vdots \\ \mathbf{a}_{nj} \end{array} \right), \left( \begin{array}{c} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{a}_{rr} - k \\ \mathbf{a}_{r+1r} \\ \vdots \\ \mathbf{a}_{nr} \end{array} \right)_{\mathbb{R}^n} =$$

$$= \mathbf{a}_{rj}(\mathbf{a}_{rr} - k) + \cdots + \mathbf{a}_{ij}\mathbf{a}_{ir} + \cdots + \mathbf{a}_{nj}\mathbf{a}_{nr} = \sum_{i=1}^n u_i \mathbf{a}_{ij} = \sum_{i=r}^n u_i \mathbf{a}_{ij}$$

$$(\mathbf{u}_i = \mathbf{0}, i = 1, \dots, r-1, \mathbf{u}_r = \mathbf{a}_{rr} - k, \mathbf{u}_i = \mathbf{a}_{ir}, i = r+1, \dots, n)$$

Noua coloană  $j$  se obține din vechea coloană  $j$  din care scădem vectorul  $u$  înmulțit cu constanta  $\gamma_j$ . Ne interesează ca primelor  $(r-1)$  coloane să nu li se schimbe forma superior triunghiulară deja obținută.

Pentru  $j=1, \dots, (r-1)$  avem:

$$\gamma_j := u^T (Ae_j) = \left( \begin{array}{c} a_{1j} \\ a_{2j} \\ \vdots \\ a_{jj} \\ a_{j+1j} = \mathbf{0} \\ \vdots \\ a_{rj} = \mathbf{0} \\ \vdots \\ a_{nj} = \mathbf{0} \end{array} \right), \left( \begin{array}{c} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ a_{rr} - k \\ a_{r+1r} \\ \vdots \\ a_{nr} \end{array} \right)_{\mathbb{R}^n} =$$

$$= a_{1j} \mathbf{0} + \dots + a_{jj} \mathbf{0} + \dots + \mathbf{0}(a_{rr} - k) + \dots + \mathbf{0}a_{ir} + \dots + \mathbf{0}a_{nr} = \mathbf{0}$$

Din faptul că  $\gamma_j = 0, j = 1, \dots, r-1$  rezultă că primele  $(r-1)$  coloane ale matricei  $A$  nu se schimbă când facem operația  $A = P_r A$ , rămân în formă superior triunghiulară.

Algoritmul de trecere de la matricea  $A$  la matricea  $P_r A$  este următorul:

$$(P_r A)e_j = \begin{cases} Ae_j & \text{pentru } j = 1, \dots, r-1 \\ (a_{1r}, a_{2r}, \dots, a_{r-1r}, k, 0, \dots, 0)^T & \text{pentru } j = r \\ Ae_j - \frac{\gamma_j}{\beta} u & \text{pentru } j = r+1, \dots, n \end{cases}$$

$$\gamma_j = (Ae_j, u)_{\mathbb{R}^n} = \sum_{i=r}^n u_i a_{ij}$$

$$u_i = \mathbf{0}, \quad i = 1, \dots, r-1, \quad u_r = a_{rr} - k, \quad u_i = a_{ir}, \quad i = r+1, \dots, n$$

Operația de transformare a vectorului termenilor liberi

$$b := P_r b:$$

$$P_r b = \left( I_n - \frac{1}{\beta} (uu^T) \right) b = b - \frac{1}{\beta} (uu^T) b = b - \frac{1}{\beta} u (u^T b) = b - \frac{\gamma}{\beta} u$$

$$\gamma = u^T b = (b, u)_{\mathbb{R}^n} = \sum_{i=r}^n u_i b_i$$

## *Algoritmul lui Householder*

$$\tilde{Q} = I_n;$$

for  $r = 1, \dots, n - 1$

$$\left\{ \begin{array}{l} \text{calculează matricea } P_r \text{ (constanta } \beta \text{ și vectorul } u) \\ A = P_r * A; \\ b = P_r * b; \\ \tilde{Q} = P_r * \tilde{Q}; \end{array} \right.$$

La sfârșitul acestui algoritm, în matricea  $A$  vom avea matricea superior triunghiulară  $R$ , în vectorul  $b$  vom avea  $Q^T b^{\text{init}}$ ,  $b^{\text{init}}$  este vectorul inițial al termenilor liberi, iar matricea  $\tilde{Q}$  va conține matricea  $Q^T$  din factorizarea  $QR$  a matricei  $A$ .

## *Algoritmul Householder detaliat*

$$\tilde{Q} = I_n;$$

**for**  $r = 1, \dots, n - 1$

*//* construcția matricii  $P_r$  – constanta  $\beta$  și vectorul  $u$

- $\sigma = \sum_{i=r}^n a_{ir}^2;$

- **if** ( $\sigma \leq \varepsilon$ ) **break**; *//*  $r = r + 1 \leftrightarrow P_r = I_n$  ( $A$  singulară)

- $k = \sqrt{\sigma};$

- **if** ( $a_{rr} > 0$ )  $k = -k;$

- $\beta = \sigma - k a_{rr};$

- $u_r = a_{rr} - k; u_i = a_{ir}, i = r + 1, \dots, n;$



//  $A = P_r * A$

// transformarea coloanelor  $j = r + 1, \dots, n$

• for  $j = r + 1, \dots, n$

$$* \gamma = (\gamma_j / \beta) = (Ae_j, u) / \beta = \left( \sum_{i=r}^n u_i a_{ij} \right) / \beta;$$

\* for  $i = r, \dots, n$

$$a_{ij} = a_{ij} - \gamma * u_i;$$

// transformarea coloanei  $r$  a matricii  $A$

•  $a_{rr} = k$ ;  $a_{ir} = 0$ ,  $i = r + 1, \dots, n$ ;

//  $b = P_r * b$

$$• \gamma = (\gamma / \beta) = (b, u) / \beta = \left( \sum_{i=r}^n u_i b_i \right) / \beta;$$

• for  $i = r, \dots, n$   $b_i = b_i - \gamma * u_i$ ;

$$// \tilde{Q} = P_r * \tilde{Q}$$

• for  $j = 1, \dots, n$

$$* \gamma = (\tilde{Q}e_j, u) / \beta = \left( \sum_{i=r}^n u_i \tilde{q}_{ij} \right) / \beta;$$

\* for  $i = r, \dots, n$

$$\tilde{q}_{ij} = \tilde{q}_{ij} - \gamma * u_i;$$

Numărul de operații efectuate:

**A** (adunări, scăderi):

$$(n-1) + 2n(n-1) + \frac{n(n-1)(2n-1)}{3} = \frac{(n-1)(n+1)(2n+3)}{3} = \\ = \frac{2}{3}n^3 + O(n^2)$$

**M** (înmulțiri, împărțiri):

$$4(n-1) + 3n(n-1) + \frac{n(n-1)(2n-1)}{6} = \frac{2}{3}n^3 + O(n^2)$$

**R** (radicali):  $(n-1)$

## Algoritmul lui Givens

Fie  $A$  o matrice reală pătratică de dimensiune  $n$ . Pp. că avem:

$$A = QR$$

unde  $Q$  este o matrice ortogonală iar  $R$  este o matrice superior triunghiulară.

$$Ax = b \Leftrightarrow QRx = b \Leftrightarrow Q^T QRx = Q^T b \Leftrightarrow Rx = Q^T b$$

În cazul algoritmului Givens, pentru a aduce sistemul  $Ax=b$  la forma  $Rx = Q^T b$  se folosesc matricele de rotație. O matrice de rotație  $R_{pq}(\theta) = (r_{ij})_{i,j=1,n}$  are următoarea formă :

$$\mathbf{R}_{pq}(\boldsymbol{\theta}) = \begin{pmatrix}
& & & \mathbf{p} & & \mathbf{q} & & \\
\mathbf{1} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\
\mathbf{0} & \mathbf{1} & \dots & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\
\vdots & & & & & & & \\
\mathbf{0} & \mathbf{0} & \dots & \mathbf{c} & \dots & \mathbf{s} & \dots & \mathbf{0} \\
\vdots & & & & & \mathbf{1} & & \\
\mathbf{0} & \mathbf{0} & \dots & (-\mathbf{s}) & \dots & \mathbf{c} & \dots & \mathbf{0} \\
\vdots & & & & & & & \\
\mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{1}
\end{pmatrix} \begin{matrix} \\ \\ \\ \mathbf{p} \\ \\ \mathbf{q} \\ \\ \end{matrix}$$

$$r_{ij} = \begin{cases} \mathbf{1} & \text{pentru } i = j, \quad i \neq p, i \neq q \\ \mathbf{c} & \text{pentru } i = j, \quad i = p, i = q \\ \mathbf{s} & \text{pentru } i = p, j = q \\ -\mathbf{s} & \text{pentru } i = q, j = p \\ \mathbf{0} & \text{în rest} \end{cases}$$

unde  $p, q \in \{1, \dots, n\}$  iar  $c$  și  $s$  sunt două numere reale care satisfac relația  $c^2 + s^2 = 1$ . Constantele  $c$  și  $s$  pot fi alese astfel încât  $c = \cos \theta$ ,  $s = \sin \theta$ . Se arată ușor, folosind relația  $c^2 + s^2 = 1$ , că matricea  $R_{pq}(\theta)$  este ortogonală:

$$R_{pq}(\theta) R_{pq}^T(\theta) = R_{pq}^T(\theta) R_{pq}(\theta) = I_n$$

Calculul matricei:

$$\mathbf{B} = \mathbf{R}_{pq}(\theta)\mathbf{A},$$

$\mathbf{B}$  se obține din  $\mathbf{A}$  modificând doar liniile  $p$  și  $q$ .

Fie

$$\mathbf{A}_i = \mathbf{e}_i^T \mathbf{A}, \text{ - linia } i \text{ a matricei } \mathbf{A}$$

$$\mathbf{B}_i = \mathbf{e}_i^T \mathbf{B} \text{ - linia } i \text{ a matricei } \mathbf{B}.$$

Liniile matricei  $\mathbf{B}$ :

$$\mathbf{B}_i = \mathbf{A}_i, \quad i = 1, \dots, n, \quad i \neq p, i \neq q$$

$$\mathbf{B}_p = c\mathbf{A}_p + s\mathbf{A}_q$$

$$\mathbf{B}_q = -s\mathbf{A}_p + c\mathbf{A}_q$$

$$b_{pj} = c a_{pj} + s a_{qj} , j = 1, \dots, n$$

$$b_{qj} = -s a_{pj} + c a_{qj} , j = 1, \dots, n$$

$$b_{ij} = a_{ij} \quad \text{în rest}$$

Calculul matricei :

$$D = A R_{pq}^T(\theta),$$

$D$  se obține din  $A$  modificând doar coloanele  $p$  și  $q$ .

Notăm  $Ae_j$ ,  $De_j$  – coloana  $j$  a matricei  $A$  și respectiv  $D$ .



Coloanele matricei  $D$ :

$$D_j = A_j, \quad j = 1, \dots, n, \quad j \neq p, j \neq q$$

$$De_p = cAe_p + sAe_q$$

$$De_q = -sAe_p + cAe_q$$

$$d_{ip} = c a_{ip} + s a_{iq}, \quad i = 1, \dots, n$$

$$d_{iq} = -s a_{ip} + c a_{iq}, \quad i = 1, \dots, n$$

$$d_{ij} = a_{ij} \quad \text{în rest}$$

Algoritmul lui Givens se desfășoară în  $(n-1)$  pași - la pasul  $r$  se transformă coloana  $r$  a matricei  $A$  în formă superior triunghiulară fără a modifica primele  $(r-1)$  coloane.

## Pasul 1

*Intrare:* matricea  $A$ , vectorul  $\mathbf{b}$

*Ieșire:* matricea  $A^{(1)}$  (cu prima coloană în formă superior triunghiulară),  $\mathbf{b}^{(1)}$

Se efectuează următoarele operații de înmulțire cu matrice de rotație:

$$\mathbf{R}_{1n}(\theta_{1n}) \cdots \mathbf{R}_{13}(\theta_{13}) \mathbf{R}_{12}(\theta_{12}) A = A^{(1)}$$

$$\mathbf{R}_{1n}(\theta_{1n}) \cdots \mathbf{R}_{13}(\theta_{13}) \mathbf{R}_{12}(\theta_{12}) \mathbf{b} = \mathbf{b}^{(1)}$$

Unghiurile  $\theta_{li}$  (constantele  $c_{li}$  și  $s_{li}$ ) se aleg astfel ca elementul de pe poziția  $(i, 1)$  din matricea rezultat să devină  $0$ .

## Pasul $r$

*Intrare:* matricea  $A^{(r-1)}$  (are primele  $(r-1)$  coloane în formă superior triunghiulară),  $\mathbf{b}^{(r-1)}$

*Ieșire:* matricea  $A^{(r)}$  (cu primele  $r$  coloane în formă superior triunghiulară),  $\mathbf{b}^{(r)}$

La acest pas matricea  $A^{(r-1)}$  și vectorul  $\mathbf{b}^{(r-1)}$  se transformă astfel:

$$\mathbf{R}_{rn}(\theta_{rn}) \cdots \mathbf{R}_{ri}(\theta_{ri}) \cdots \mathbf{R}_{rr+1}(\theta_{rr+1}) A^{(r-1)} = A^{(r)}$$

$$\mathbf{R}_{rn}(\theta_{rn}) \cdots \mathbf{R}_{ri}(\theta_{ri}) \cdots \mathbf{R}_{rr+1}(\theta_{rr+1}) \mathbf{b}^{(r-1)} = \mathbf{b}^{(r)}$$

unde elementele  $c = c_{ri}$  și  $s = s_{ri}$  din matricele de rotație se aleg astfel ca după înmulțirea cu  $R_{ri}(\theta_{ri}), i = r + 1, \dots, n$  elementul  $(i, r)$  să devină  $0$ .

Considerăm operația  $B = R_{ri}(\theta_{ri})A$ , unde  $\theta_{ri}$  se alege astfel ca  $b_{ir} = 0$ :

$$b_{rj} = c a_{rj} + s a_{ij} ,$$

$$b_{ij} = -s a_{rj} + c a_{ij} , j = 1, \dots, n$$

$$(b_{kl} = a_{kl} \quad \text{în rest})$$

$$(j = r) \quad b_{ir} = -s a_{rr} + c a_{ir}$$

Cea mai simplă alegere pentru  $c$  și  $s$  astfel ca să obținem  $b_{ir} = \mathbf{0}$  este:

$$c = f a_{rr} \quad , \quad s = f a_{ir} \quad ,$$

$$f \text{ ales astfel ca } c^2 + s^2 = 1 \quad \Rightarrow \quad f = \frac{1}{\sqrt{a_{rr}^2 + a_{ir}^2}}$$

$$c = \frac{a_{rr}}{\sqrt{a_{rr}^2 + a_{ir}^2}} \quad , \quad s = \frac{a_{ir}}{\sqrt{a_{rr}^2 + a_{ir}^2}}$$

$$\sqrt{a_{rr}^2 + a_{ir}^2} = 0 \quad \Rightarrow \quad a_{rr} = a_{ir} = 0 \quad \Rightarrow \quad c = 1 \quad , \quad s = 0$$

$$(R_{ir}(\theta) = I_n)$$

Deci elementul de pe poziția  $(i,r)$  este deja nul. Nu avem ce schimba în matricea  $A$ .

Operația  $B = R_{ri}(\theta_{ri})A$ , nu afectează forma superior triunghiulară a primelor  $(r-1)$  coloane. În matricea  $B$  aceste coloane vor continua să fie în formă superior triunghiulară.

$$b_{rj} = c a_{rj} + s a_{ij} = 0, \quad b_{ij} = -s a_{rj} + c a_{ij} = 0, \quad j = 1, \dots, r-1$$

**deoarece  $a_{rj} = a_{ij} = 0$**

Înmulțirea  $B = R_{ri}(\theta_{ri})A$  nu schimbă decât liniile  $r$  și  $i$  ale matricei  $B$ . În concluzie, operația  $B = R_{ri}(\theta_{ri})A$  nu schimbă elementele nule deja obținute, ci doar face ca elementul de pe poziția  $(i,r)$  să devină  $0$ .

Algoritmul lui Givens poate fi descris astfel:

$$\mathbf{R}_{n-1n}(\theta_{n-1n}) \cdots \mathbf{R}_{rn}(\theta_{rn}) \cdots \mathbf{R}_{r+1}(\theta_{r+1}) \cdots \mathbf{R}_{1n}(\theta_{1n}) \cdots \mathbf{R}_{12}(\theta_{12}) \mathbf{A} = \mathbf{R}$$

$$\mathbf{R}_{n-1n}(\theta_{n-1n}) \cdots \mathbf{R}_{rn}(\theta_{rn}) \cdots \mathbf{R}_{r+1}(\theta_{r+1}) \cdots \mathbf{R}_{1n}(\theta_{1n}) \cdots \mathbf{R}_{12}(\theta_{12}) \mathbf{b} = \bar{\mathbf{b}} = \mathbf{Q}^T \mathbf{b}$$

Notăm cu  $\tilde{\mathbf{Q}}$  următoarea matrice:

$$\tilde{\mathbf{Q}} = \mathbf{R}_{n-1n}(\theta_{n-1n}) \cdots \mathbf{R}_{rn}(\theta_{rn}) \cdots \mathbf{R}_{r+1}(\theta_{r+1}) \cdots \mathbf{R}_{1n}(\theta_{1n}) \cdots \mathbf{R}_{12}(\theta_{12})$$

Matricea  $\tilde{Q}$  este matrice ortogonală ca produs de matrice ortogonale. Descompunerea  $QR$  a matricei  $A$  este următoarea:

$$\tilde{Q} A = R \quad (\tilde{Q}^{-1} = \tilde{Q}^T) \Rightarrow A = \tilde{Q}^T R = QR$$

$$\begin{aligned} Q = \tilde{Q}^T &= \left( R_{n-1n}(\theta_{n-1n}) \cdots R_{rn}(\theta_{rn}) \cdots R_{r+1}(\theta_{r+1}) \cdots R_{1n}(\theta_{1n}) \cdots R_{12}(\theta_{12}) \right)^T \\ &= R_{12}^T(\theta_{12}) \cdots R_{1n}^T(\theta_{1n}) \cdots R_{r+1}^T(\theta_{r+1}) \cdots R_{rn}^T(\theta_{rn}) \cdots R_{n-1n}^T(\theta_{n-1n}) \end{aligned}$$



Pe scurt, *algoritmul lui Givens* este următorul:

$$\tilde{Q} = I_n;$$

**for**  $r = 1, \dots, n - 1$

**for**  $i = r + 1, \dots, n$

$$\left\{ \begin{array}{l} A = R_{ri}(\theta_{ri}) * A; \\ b = R_{ri}(\theta_{ri}) * b; \\ \tilde{Q} = R_{ri}(\theta_{ri}) * \tilde{Q}; \end{array} \right.$$

$$\tilde{Q} = I_n;$$

for  $r = 1, \dots, n - 1$

for  $i = r + 1, \dots, n$

// construcția matricii  $R_{ri}(\theta_{ri})$  – constantele  $c$  și  $s$

- $f = \sqrt{a_{rr}^2 + a_{ir}^2};$

- if ( $f \leq \varepsilon$ ) {  $c = 1; s = 0;$  } //  $R_{ri}(\theta_{ri}) = I$

- else {  $c = a_{rr} / f; s = a_{ir} / f; }$

//  $A = R_{ri}(\theta_{ri}) * A$

- for  $j = r + 1, \dots, n$

$$\begin{cases} a_{rj} = c * a_{rj} + s * a_{ij}; \\ a_{ij} = -s * a_{rj}^{vechi} + c * a_{ij}; \end{cases}$$

- $\mathbf{a}_{ir} = \mathbf{0}$ ;  $\mathbf{a}_{rr} = \mathbf{f}$  ;

//  $\mathbf{b} = \mathbf{R}_{ri}(\boldsymbol{\theta}_{ri}) * \mathbf{b}$

- $\mathbf{b}_r = \mathbf{c} * \mathbf{b}_r + \mathbf{s} * \mathbf{b}_i$ ;

- $\mathbf{b}_i = -\mathbf{s} * \mathbf{b}_r^{\text{vechi}} + \mathbf{c} * \mathbf{b}_i$ ;

//  $\tilde{\mathbf{Q}} = \mathbf{R}_{ri}(\boldsymbol{\theta}_{ri}) * \tilde{\mathbf{Q}}$

- for  $j = 1, \dots, n$

$$\begin{cases} \tilde{\mathbf{q}}_{rj} = \mathbf{c} * \tilde{\mathbf{q}}_{rj} + \mathbf{s} * \tilde{\mathbf{q}}_{ij}; \\ \tilde{\mathbf{q}}_{ij} = -\mathbf{s} * \tilde{\mathbf{q}}_{rj}^{\text{vechi}} + \mathbf{c} * \tilde{\mathbf{q}}_{ij}; \end{cases}$$

La sfârșitul acestui algoritm, în matricea  $A$  vom avea matricea superior triunghiulară  $R$ , în vectorul  $b$  vom avea  $\tilde{Q}b^{\text{init}}$  ( $b^{\text{init}}$  - vectorul termenilor liberi inițial), iar matricea  $\tilde{Q}$  va conține matricea  $Q^T$  din factorizarea  $QR$  a matricei  $A$ .

Numărând operațiile efectuate (exceptând calculul matricei  $\tilde{Q}$ ) obținem:

$$\frac{n(n-1)}{2} \text{ radicali}$$

$$\frac{n(n-1)(4n+7)}{6} = \frac{2}{3}n^3 + O(n^2) - \text{adunări/scăderi}$$

$$\frac{2n(n-1)(2n+5)}{3} = \frac{4}{3}n^3 + O(n^2) \text{ înmulțiri/împărțiri}$$

## ***QR* – algoritmul Gram Schmidt**

$$A \in \mathbb{R}^{n \times n} \text{ cu } \det A \neq 0$$

$$A = QR$$

$$\begin{bmatrix} a^1 & a^2 & \cdots & a^n \end{bmatrix} = \begin{bmatrix} q^1 & q^2 & \cdots & q^n \end{bmatrix} \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ \mathbf{0} & r_{22} & \cdots & r_{2n} \\ \vdots & & & \\ \mathbf{0} & \mathbf{0} & \cdots & r_{nn} \end{pmatrix} \quad (1)$$

$a^j = Ae_j$  – coloana  $j$  a matricei  $A$

$q^j = Qe_j$  – coloana  $j$  a matricei  $Q$

Relația (1) poate fi rescrisă astfel:

$$\left\{ \begin{array}{l} r_{11}q^1 = a^1 \\ r_{12}q^1 + r_{22}q^2 = a^2 \\ \vdots \\ r_{1p}q^1 + \cdots + r_{jp}q^j + \cdots + r_{pp}q^p = a^p \\ \vdots \\ r_{1n}q^1 + \cdots + r_{jn}q^j + \cdots + r_{nn}q^n = a^n \end{array} \right. \quad (2)$$

Algoritmul de calcul al descompunerii  $QR$  cu metoda Gram-Schmidt se desfășoară în  $n$  pași, la fiecare pas calculându-se:

- coloana  $p$  din matricea  $R$
- coloana  $p$  din matricea  $Q$

Avem:

$$\det A \neq 0, A = QR, Q - \text{ortogonală} \Rightarrow \det R \neq 0 (r_{ii} \neq 0 \forall i)$$

### ***Pasul 1***

Se folosește prima ecuație a sistemului (2)

$$r_{11}q^1 = a^1$$

Se face produsul scalar al acestei ecuații cu vectorii  $q^1$  și  $a^1$ . Se folosește proprietatea coloanelor matricelor ortogonale:

$$(q^i, q^j)_{\mathbb{R}^n} = \begin{cases} \|q^i\|_2^2 = 1 & \text{pentru } i = j \\ \mathbf{0} & \text{pentru } i \neq j \end{cases}$$

$$q^1 = \frac{1}{r_{11}}a^1 \quad \text{și} \quad \|q^1\|_2 = 1$$

$$r_{11} = \pm \|a^1\|_2, \quad q^1 = \frac{1}{r_{11}}a^1 \quad (r_{11} \neq 0 \text{ deoarece } \det A \neq 0)$$



### *Pasul p*

Se folosește ecuația  $p$  sistemului (2):

$$r_{1p}q^1 + \cdots + r_{jp}q^j + \cdots + r_{pp}q^p = a^p$$

La acest pas se cunosc deja coloanele  $q^1, q^2, \dots, q^{p-1}$ . Se face, pe rând, produsul scalar al acestei ecuații cu vectorii:

$$q^1, q^2, \dots, q^{p-1}$$

$$q^p \text{ și } a^p.$$

$$\left( \sum_{k=1}^p r_{kp} q^k, q^j \right)_{\mathbb{R}^n} = (a^p, q^j)_{\mathbb{R}^n} \quad j = 1, \dots, p-1$$

$$\left( \sum_{k=1}^p r_{kp} \mathbf{q}^k, \mathbf{q}^j \right)_{\mathbb{R}^n} = \sum_{\substack{k=1 \\ k \neq j}}^p r_{kp} \left( \mathbf{q}^k, \mathbf{q}^j \right)_{\mathbb{R}^n} + r_{jp} \left( \mathbf{q}^j, \mathbf{q}^j \right)_{\mathbb{R}^n} = r_{jp}$$

$$r_{jp} = \left( \mathbf{a}^p, \mathbf{q}^j \right)_{\mathbb{R}^n} \quad j = 1, \dots, p-1$$

Avem:

$$\mathbf{q}^p = \frac{1}{r_{pp}} \left( \mathbf{a}^p - r_{1p} \mathbf{q}^1 - \dots - r_{p-1p} \mathbf{q}^{p-1} \right)$$

$$\| \mathbf{q}^p \|_2^2 = 1 = \frac{1}{r_{pp}^2} \left\| \left( \mathbf{a}^p - r_{1p} \mathbf{q}^1 - \dots - r_{p-1p} \mathbf{q}^{p-1} \right) \right\|_2^2$$

Obținem:

$$\mathbf{r}_{pp} = \pm \left\| \mathbf{a}^p - \mathbf{r}_{1p} \mathbf{q}^1 - \dots - \mathbf{r}_{p-1p} \mathbf{q}^{p-1} \right\|_2$$

$$\mathbf{q}^p = \frac{\mathbf{1}}{\mathbf{r}_{pp}} \left( \mathbf{a}^p - \mathbf{r}_{1p} \mathbf{q}^1 - \dots - \mathbf{r}_{p-1p} \mathbf{q}^{p-1} \right) = \frac{\mathbf{1}}{\mathbf{r}_{pp}} \left( \mathbf{a}^p - \sum_{j=1}^{p-1} \mathbf{r}_{jp} \mathbf{q}^j \right)$$

## Algoritmul Gram Schmidt modificat

*for*  $i = 1, n$

$$\mathbf{v}^i = \mathbf{a}^i;$$

*for*  $i = 1, n$

$$r_{ii} = \|\mathbf{v}^i\|_2;$$

$$\mathbf{q}^i = (1/r_{ii}) \mathbf{v}^i;$$

*for*  $j = (i + 1), n$

$$r_{ij} = (\mathbf{q}^i, \mathbf{v}^j);$$

$$\mathbf{v}^j = \mathbf{v}^j - r_{ij}\mathbf{q}^i;$$

$$\mathbf{M}: \frac{(n^3 + 3n^2)}{2}$$

$$\mathbf{A}: \frac{(n^3 + n^2 - 2)}{2}$$

# **Calcul Numeric**

**Cursul 6**

**2022**

*Anca Ignat*

## Metode iterative pentru rezolvarea sistemelor de ecuații liniare

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n \quad (1)$$

- se presupune cunoscut că  $A$  este nesingulară,  $\det A \neq 0$ ;
- soluția exactă a sistemului (1) se notează cu  $x^*$ :

$$x^* = A^{-1}b \quad (2)$$

- $n$  - dimensiunea sistemului este "*mare*";
- $A$  este matrice rară - cu "*puține*" elemente  $a_{ij} \neq 0$ ;

- pentru a aproxima soluția  $\mathbf{x}^*$  matricea  $A$  nu se schimbă (transformă) ci doar se folosesc elementele nenule ale matricei ;
- se construiește un șir de vectori  $\{\mathbf{x}^{(k)}\} \subseteq \mathbb{R}^n$ , șir care în anumite cazuri, converge la  $\mathbf{x}^*$  :

$$\mathbf{x}^{(k)} \rightarrow \mathbf{x}^* \quad \text{pentru } k \rightarrow \infty$$

## O schemă generală de deducere a unei metode iterative

Fie descompunerea:

$$A = B - C, \quad B, C \in R^{n \times n}, \quad B \text{ "ușor" inversabilă} \quad (3)$$

Ce înseamnă  $B$  "ușor" inversabilă ? Sistemul liniar, având ca matrice a sistemului matricea  $B$ :

$$Bx = f$$

se rezolvă 'ușor' (adică repede) – ca în cazul sistemelor cu matrice diagonale sau triunghiulare, de exemplu.



$$Ax^* = b \Leftrightarrow Bx^* - Cx^* = b \Leftrightarrow$$

$$Bx^* = Cx^* + b \Leftrightarrow x^* = B^{-1}Cx^* + B^{-1}b = Mx^* + d$$

unde

$$M := B^{-1}C \in \mathbb{R}^{n \times n}, \quad d := B^{-1}b \in \mathbb{R}^n \quad (4)$$

Șirul  $\{x^{(k)}\}$  se construiește astfel:

$$x^{(k+1)} := Mx^{(k)} + d, \quad k = 0, 1, 2, \dots \quad x^{(0)} \in \mathbb{R}^n \text{ ales arbitrar} \quad (5)$$

Vectorul  $x^{(k+1)}$  poate fi privit și ca soluția sistemului liniar:

$$Bx = f \quad \text{cu } f := Cx^{(k)} + b \quad (6)$$

Cunoscând vectorul  $\mathbf{x}^{(k)}$ , următorul element din șir,  $\mathbf{x}^{(k+1)}$ , se poate construi fie utilizând relația (5) (dacă putem construi matricea  $\mathbf{M}$  explicit), fie rezolvând sistemul liniar (6).

Matricea  $\mathbf{M}$  poartă numele de *matricea iterației* iar vectorul  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  se numește *iterația inițială*.

Ne punem problema convergenței șirului  $\mathbf{x}^{(k)}$ :

$$\mathbf{x}^{(k)} \rightarrow \mathbf{x}^* \quad , \quad k \rightarrow \infty$$

Se știe că această convergență nu are loc pentru orice matrice  $\mathbf{B}$ . Avem următorul rezultat general de convergență.

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $\{A^k\}$  un șir de matrici.

$$A^k \rightarrow \mathbf{0}_{n \times n}, k \rightarrow \infty \Leftrightarrow \|A^k\| \rightarrow 0, k \rightarrow \infty.$$

### Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$ . Atunci:

$$A^k \rightarrow \mathbf{0}, k \rightarrow \infty \Leftrightarrow \rho(A) < 1.$$

Dacă există o normă matriceală naturală pentru care  $\|A\| < 1$  atunci:

$$A^k \rightarrow \mathbf{0} \text{ pentru } k \rightarrow \infty.$$

$$(n = 1 \rightarrow a \in \mathbb{R}, a^k \rightarrow 0 \text{ pentru } k \rightarrow \infty \Leftrightarrow |a| < 1.)$$

## Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$ . Seria  $\sum_{k=0}^{\infty} A^k$  converge dacă și numai dacă raza spectrală a matricei  $A$  este subunitară:

$$\sum_{k=0}^{\infty} A^k = S \Leftrightarrow \rho(A) < 1.$$

Dacă există o normă a matricei  $A$  astfel încât  $\|A\| < 1$  atunci seria converge. În cazul convergenței avem :

$$\sum_{k=0}^{\infty} A^k = S = (I - A)^{-1}.$$

## Teorema de convergență

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice nesingulară și  $B, C \in \mathbb{R}^{n \times n}$ ,  $\det B \neq 0$ , astfel ca  $A=B-C$ . Fie  $x^{(0)} \in \mathbb{R}^n$  un vector oarecare și  $\{x^{(k)}\}$  șirul de vectori dat de relația (5) cu  $M$  și  $d$  dați de (4). Atunci:

$$x^{(k)} \rightarrow x^*, k \rightarrow \infty, \forall x^{(0)} \Leftrightarrow \rho(M) < 1 \quad (7)$$

unde  $\rho(M) = \max\{|\lambda|; \lambda - \text{valoare proprie a matricii } M\}$  este raza spectrală a matricii  $M$ . Dacă există o normă matricială naturală astfel ca  $\|M\| < 1$  atunci șirul  $\{x^{(k)}\}$  converge la soluția  $x^*$  a sistemului (1).

$$\|M\| < 1 \Rightarrow x^{(k)} \rightarrow x^*, k \rightarrow \infty, \forall x^{(0)}. \quad (8)$$

Demonstrație: Scăzând relațiile (5) și  $\mathbf{x}^* = M\mathbf{x}^* + \mathbf{d}$  obținem:

$$\mathbf{x}^{(k+1)} - \mathbf{x}^* = M(\mathbf{x}^{(k)} - \mathbf{x}^*), k = 0, 1, 2, \dots$$

Avem:

$$\mathbf{x}^{(p)} - \mathbf{x}^* = M(\mathbf{x}^{(p-1)} - \mathbf{x}^*) = M^2(\mathbf{x}^{(p-2)} - \mathbf{x}^*) = \dots = M^p(\mathbf{x}^{(0)} - \mathbf{x}^*)$$

$$\mathbf{x}^{(p)} - \mathbf{x}^* = M^p(\mathbf{x}^{(0)} - \mathbf{x}^*), \forall p$$

Prin urmare:

$$\mathbf{x}^{(p)} \rightarrow \mathbf{x}^*, p \rightarrow \infty \Leftrightarrow M^p \rightarrow \mathbf{0}, p \rightarrow \infty$$

$$M^p \rightarrow \mathbf{0}, p \rightarrow \infty \Leftrightarrow \rho(M) < 1$$

Dacă:

$$\|M\| < 1 \Rightarrow M^p \rightarrow \mathbf{0}, p \rightarrow \infty \Rightarrow \mathbf{x}^{(p)} \rightarrow \mathbf{x}^*, p \rightarrow \infty \quad \forall \mathbf{x}^{(0)}$$

## Evaluarea erorii absolute $\| \mathbf{x}^{(k)} - \mathbf{x}^* \|$

Presupunem  $\|M\| < 1$  (șirul  $\{\mathbf{x}^{(k)}\}$  converge la  $\mathbf{x}^*$ ).

Avem din (5):

$$\mathbf{x}^{(l+1)} = M\mathbf{x}^{(l)} + d$$

$$\mathbf{x}^{(l)} = M\mathbf{x}^{(l-1)} + d$$

$$\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)} = M(\mathbf{x}^{(l)} - \mathbf{x}^{(l-1)}) \quad \forall l$$

Pentru orice  $k, j$ , folosind relațiile de mai sus, avem:

$$\mathbf{x}^{(k+j+1)} - \mathbf{x}^{(k+j)} = M(\mathbf{x}^{(k+j)} - \mathbf{x}^{(k+j-1)}) = \dots = M^j(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) \quad \forall k, j$$

Aplicând succesiv relația precedentă obținem:

$$\begin{aligned} \mathbf{x}^{(k+p)} - \mathbf{x}^{(k)} &= \mathbf{x}^{(k+p)} - \mathbf{x}^{(k+p-1)} + \mathbf{x}^{(k+p-1)} - \mathbf{x}^{(k+p-2)} + \dots + \\ &\quad + \mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)} + \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = \end{aligned}$$

$$= \sum_{j=0}^{p-1} (\mathbf{x}^{(k+j+1)} - \mathbf{x}^{(k+j)})$$

$$\mathbf{x}^{(k+p)} - \mathbf{x}^{(k)} = \sum_{j=0}^{p-1} (\mathbf{x}^{(k+j+1)} - \mathbf{x}^{(k+j)}) = \left( \sum_{j=0}^{p-1} M^j \right) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$



Făcând  $p \rightarrow \infty$  obținem:

$$\mathbf{x}^* - \mathbf{x}^{(k)} = \left( \sum_{j=0}^{\infty} \mathbf{M}^j \right) \mathbf{M} (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

$$\|\mathbf{M}\| < 1 \Rightarrow \sum_{j=0}^{\infty} \mathbf{M}^j = (\mathbf{I}_n - \mathbf{M})^{-1}$$

Mai avem și evaluarea:

$$\|\mathbf{M}\| < 1 \Rightarrow \frac{1}{1 + \|\mathbf{M}\|} \leq \|(\mathbf{I}_n - \mathbf{M})^{-1}\| \leq \frac{1}{1 - \|\mathbf{M}\|}$$

Prin urmare:

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{\|\mathbf{M}\|}{1 - \|\mathbf{M}\|} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|$$

Această relație ne spune că din punct de vedere practic putem opri algoritmul atunci când diferența dintre două iterații succesive devine suficient de mică, acest lucru asigurând apropierea de soluție.

În continuare vedea cum se memorează matricele rare și vom particulariza matricea ***B***.

## Memorarea matricelor rare

- se memorează doar valorile nenule și suficiente informații despre indici astfel ca să se poată reconstitui complet matricea

Pp. că matricea  $A$  are  $NN$  elemente nenule.

### *Memorare comprimată pe linii*

Se folosesc 3 vectori:

*valori* – vector de numere reale de dimensiune  $NN$

*ind\_col* – vector de indici de dimensiune  $NN$

*inceput\_linii* – vector de întregi de dimensiune  $n+1$

În vectorul *valori* se memorează elementele nenule ale matricii *A* în ordinea liniilor iar în vectorul *ind\_col* se memorează indicii de coloană ai elementelor din *valori*. În vectorul *inceput\_linii* se stochează indicele/poziția în vectorul *valori / ind\_col* al/a primului element de pe linia *i* memorat în vectorii *valori / ind\_col*.

$$- \textit{inceput\_linii}(n+1) = NN+1$$

$$- \textit{inceput\_linii}(i+1) - \textit{inceput\_linii}(i) =$$

numărul de elemente nenule de pe linia *i*,  $i=1,n$

$$A = \begin{pmatrix} 102.5 & 0.0 & 2.5 & 0.0 & 0.0 \\ 3.5 & 104.88 & 1.05 & 0.0 & 0.33 \\ 0.0 & 0.0 & 100.0 & 0.0 & 0.0 \\ 0.0 & 1.3 & 0.0 & 101.3 & 0.0 \\ 0.73 & 0.0 & 0.0 & 1.5 & 102.23 \end{pmatrix}$$

$n=5, NN=12$

$$\mathit{valori} = ( 102.5, 2.5, 0.33, 1.05, 104.88, 3.5, 100.0, 101.3, 1.3, 1.5, 0.73, 102.23 )$$

$$\mathit{ind\_col} = ( 1, 3, 5, 3, 2, 1, 3, 4, 2, 4, 1, 5 )$$

$$\mathit{inceput\_linii} = ( 1, 3, 7, 8, 10, 13 )$$

Dacă se știe că matricea are maxim  $n\_max$  elemente nenule pe fiecare linie se pot folosi 2 matrice pentru memorarea rară:

*valori* – matrice de numere reale de dimensiune  $n \times n\_max$

*ind\_col* – matrice de indici de dimensiune  $n \times n\_max$

În matricea *valori* se memorează pe linia  $i$  elementele nenule de pe linia  $i$  a matricei  $A$  iar în matricea *ind\_col* se memorează indicii de coloană ai elementelor corespunzătoare din matricea *valori*.

$$A = \begin{pmatrix} 102.5 & 0.0 & 2.5 & 0.0 & 0.0 \\ 0.0 & 104.88 & 1.05 & 0.0 & 0.33 \\ 0.0 & 0.0 & 100.0 & 0.0 & 0.0 \\ 0.0 & 1.3 & 0.0 & 101.3 & 0.0 \\ 0.73 & 0.0 & 0.0 & 1.5 & 102.23 \end{pmatrix}$$

$$valori = \begin{pmatrix} 102.5 & 2.5 & 0 \\ 104.88 & 1.05 & 0.33 \\ 100.0 & 0 & 0 \\ 101.3 & 1.3 & 0 \\ 102.23 & 1.5 & 0.73 \end{pmatrix} \quad ind\_col = \begin{pmatrix} 1 & 3 & 0 \\ 2 & 3 & 5 \\ 3 & 0 & 0 \\ 4 & 2 & 0 \\ 5 & 4 & 1 \end{pmatrix}$$

Diagonalele matricei  $A$ :

$$d_0 : (a_{11}, a_{22}, \dots, a_{nn})$$

$$d_1 : (a_{12}, a_{23}, \dots, a_{n-1n})$$

$$d_{-1} : (a_{21}, a_{32}, \dots, a_{nn-1})$$

$$d_2 : (a_{13}, a_{24}, \dots, a_{n-2n})$$

$$d_{-2} : (a_{31}, a_{42}, \dots, a_{nn-2})$$

$\vdots$

Pentru matricele care au elementele nenule plasate pe câteva din diagonalele matricei  $A$  ( $n_d$  diagonale cu elemente nenule) se pot folosi pentru memorare o matrice și un vector:



*diag* – matrice cu numere reale de dimensiune  $n \times n\_d$

*diag\_no* – vector de întregi de dimensiune  $n\_d$

În matricea *diag* se memorează pe coloane diagonalele cu elemente nenule iar în *diag\_no* este specificat numărul diagonalei care e memorat în coloana *j* a matricei *diag*.

$$diag(i, j) = a_{i+diag\_no(j)}$$

$$A = \begin{pmatrix} 20.5 & 2.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 40.5 & 3.0 & 0.0 & 0.0 \\ 1.0 & 0.0 & 100.0 & 0.0 & 0.0 \\ 0.0 & 2.3 & 0.0 & 101.5 & 4.0 \\ 0.0 & 0.0 & 3.0 & 0.0 & 102.5 \end{pmatrix}$$

$$diag = \begin{pmatrix} * & 20.5 & 2.0 \\ * & 40.5 & 3.0 \\ 1.0 & 100.0 & 0.0 \\ 2.3 & 101.5 & 4.0 \\ 3.0 & 102.5 & * \end{pmatrix} \quad diag\_no = (-2, 0, 1)$$

Alte tipuri de memorări rare:

[http://netlib.org/linalg/html\\_templates/node90.html](http://netlib.org/linalg/html_templates/node90.html)

## Metoda Jacobi pentru rezolvarea sistemelor liniare

Fie sistemul:

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n$$

cu

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

Alegem:

$$B = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}] = \begin{pmatrix} a_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & a_{22} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & a_{nn} \end{pmatrix}$$

$$\det B = a_{11}a_{22} \cdots a_{nn} \neq 0$$

$$B^{-1} = \text{diag}\left[\frac{1}{a_{11}}, \frac{1}{a_{22}}, \dots, \frac{1}{a_{nn}}\right] = \begin{pmatrix} \frac{1}{a_{11}} & 0 & \cdots & 0 \\ 0 & \frac{1}{a_{22}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{a_{nn}} \end{pmatrix}$$

Matricea  $C$  este:

$$C = B - A = \begin{pmatrix} \mathbf{0} & -a_{12} & -a_{13} & \cdots & -a_{1n} \\ -a_{21} & \mathbf{0} & -a_{23} & \cdots & -a_{2n} \\ -a_{31} & -a_{32} & \mathbf{0} & \cdots & -a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -a_{n1} & -a_{n2} & -a_{n3} & \cdots & \mathbf{0} \end{pmatrix}$$

$$C = (c_{ij}) \in \mathbb{R}^{n \times n} \quad c_{ij} = \begin{cases} -a_{ij} & \text{dacă } i \neq j \\ \mathbf{0} & \text{dacă } i = j \end{cases}$$

Matricea iterației se poate calcula și are forma:

$$M := B^{-1}C = \begin{pmatrix} \mathbf{0} & -\frac{a_{12}}{a_{11}} & -\frac{a_{13}}{a_{11}} & \dots & -\frac{a_{1n}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & \mathbf{0} & -\frac{a_{23}}{a_{22}} & \dots & -\frac{a_{2n}}{a_{22}} \\ -\frac{a_{31}}{a_{33}} & -\frac{a_{32}}{a_{33}} & \mathbf{0} & \dots & -\frac{a_{3n}}{a_{33}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -\frac{a_{n1}}{a_{nn}} & -\frac{a_{n2}}{a_{nn}} & -\frac{a_{n3}}{a_{nn}} & \dots & \mathbf{0} \end{pmatrix}$$

$$M = (m_{ij}) \in \mathbb{R}^{n \times n} \quad m_{ij} = \begin{cases} -\left(\frac{a_{ij}}{a_{ii}}\right) & \text{dacă } i \neq j \\ \mathbf{0} & \text{dacă } i = j \end{cases}$$

Construim vectorul  $\mathbf{g}$ :

$$\mathbf{g} := M\mathbf{x}^{(k)} \in \mathbb{R}^n, \quad M\mathbf{x}^{(k)} = (\mathbf{g}_i)_{i=1}^n$$

Componentele vectorului  $\mathbf{g}$  sunt:

$$\mathbf{g}_i = \sum_{j=1}^n m_{ij} x_j^{(k)} = -\sum_{\substack{j=1 \\ j \neq i}}^n \frac{a_{ij}}{a_{ii}} x_j^{(k)} = -\left(\sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)}\right) / a_{ii}, \quad i = 1, \dots, n$$

Vectorul  $d$  este:

$$d = B^{-1}b = (d_i)_{i=1}^n \in \mathbb{R}^n, \quad d_i = \frac{b_i}{a_{ii}}, i = 1, \dots, n$$

Șirul  $\{x^{(k)}\} \subseteq \mathbb{R}^n$  se construiește folosind formula:

$$x^{(k+1)} = Mx^{(k)} + d \quad \leftrightarrow \quad x_i^{(k+1)} = g_i + d_i, i = 1, \dots, n$$

$$x_i^{(k+1)} = \frac{\left( b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right)}{a_{ii}}, \quad i = 1, \dots, n$$



$$x_i^{(k+1)} = \frac{(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)})}{a_{ii}}, \quad i = 1, \dots, n \quad (1)$$

Formula (1) descrie *metoda lui Jacobi* de aproximare a soluției unui sistem liniar.

## Condiții suficiente de convergență

### Propoziția 1

$$\|M\| < 1 \quad \Rightarrow \quad x^{(k)} \rightarrow x^*, k \rightarrow \infty.$$

Demonstrație. Fie  $x^*$  soluția sistemului  $Ax=b$ . Din relația  $A=B-C$  rezultă  $Bx^* = Cx^* + b$  sau  $x^* = Mx^* + d$ .

Procesul iterativ  $x^{(k+1)} = Mx^{(k)} + d$  conduce la relația:

$$\|x^* - x^{(k+1)}\| = \|M(x^* - x^{(k)})\| \leq \|M\| \|x^* - x^{(k)}\| \leq \dots \leq \|M\|^{k+1} \|x^* - x^{(0)}\|$$

În continuare vom aplica această propoziție pentru diverse norme.

• Din  $\|M\|_F = \left(\sum_{i=1}^n \sum_{j=1}^n m_{ij}^2\right)^{\frac{1}{2}} < 1$  deducem:

$$\sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \left(\frac{a_{ij}}{a_{ii}}\right)^2 < 1 \quad \Rightarrow \quad x^{(k)} \rightarrow x^*, k \rightarrow \infty$$

- Din  $\|M\|_1 = \max\{\sum_{i=1}^n |m_{ij}|; j = 1, \dots, n\} < 1$  deducem:

$$\sum_{\substack{i=1 \\ i \neq j}}^n \left( \frac{|a_{ij}|}{|a_{ii}|} \right) < 1 \quad \forall j = 1, \dots, n \Rightarrow x^{(k)} \rightarrow x^*, k \rightarrow \infty$$

- (*Criteriul dominanței diagonalei pe linii*)

$$\text{Din } \|M\|_\infty = \max\{\sum_{j=1}^n |m_{ij}|; i = 1, \dots, n\} < 1 \text{ deducem:}$$

$$\sum_{\substack{j=1 \\ j \neq i}}^n \left( \frac{|a_{ij}|}{|a_{ii}|} \right) < 1 \quad \forall i = 1, \dots, n \Rightarrow x^{(k)} \rightarrow x^*, k \rightarrow \infty$$

$$\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| < |a_{ii}| \quad \forall i = 1, \dots, n \Rightarrow \lim_{k \rightarrow \infty} x^{(k)} = x^*$$

- (*Criteriul dominanței diagonalei pe coloane*)

$$\sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| < |a_{jj}| \quad \forall j = 1, \dots, n \Rightarrow \|M\|_1 < 1 \Rightarrow \lim_{k \rightarrow \infty} x^{(k)} = x^*$$

# Metoda Gauss-Seidel pentru rezolvarea sistemelor liniare

Considerăm din nou sistemul liniar:

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n$$

cu

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

Putem deduce metoda Gauss-Seidel din metoda lui Jacobi astfel:

$$x_i^{(k+1)} = (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)}) / a_{ii}, i = 1, \dots, n - \text{Jacobi}$$

↕

$$x_i^{(k+1)} = (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)}) / a_{ii}, i = 1, \dots, n - \text{Gauss-Seidel}$$

Când calculăm  $x_i^{(k+1)}$  cunoaștem deja  $x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}$  și putem folosi aceste valori în prima sumă.

Deducerea metodei Gauss-Seidel din schema generală se face luând:

$$B = \begin{pmatrix} a_{11} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ a_{21} & a_{22} & \mathbf{0} & \cdots & \mathbf{0} \\ a_{31} & a_{32} & a_{33} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{pmatrix}$$

$$B = (b_{ij}) \in \mathbb{R}^{n \times n} \quad b_{ij} = \begin{cases} a_{ij} & \text{dacă } j \leq i \\ \mathbf{0} & \text{dacă } j > i \end{cases}$$

Matricea  $B$  este nesingulară ( $a_{ii} \neq \mathbf{0}, \forall i$ ):

$$\det B = a_{11} a_{22} \cdots a_{nn} \neq \mathbf{0}$$



Matricea  $C$  este:

$$C = B - A = \begin{pmatrix} \mathbf{0} & -a_{12} & -a_{13} & \cdots & -a_{1n} \\ \mathbf{0} & \mathbf{0} & -a_{23} & \cdots & -a_{2n} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -a_{n-1n} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix}$$

$$C = (c_{ij}) \in \mathbb{R}^{n \times n} \quad c_{ij} = \begin{cases} -a_{ij} & \text{dacă } i < j \\ \mathbf{0} & \text{dacă } i \geq j \end{cases}$$

În cazul metodei Gauss-Seidel, vectorul  $\mathbf{x}^{(k+1)}$  se obține din  $\mathbf{x}^{(k)}$  rezolvând sistemul inferior triunghiular (7) din schema generală:

$$\mathbf{B}\mathbf{x} = \mathbf{C}\mathbf{x}^{(k)} + \mathbf{b} = \mathbf{f} \quad (1)$$

Soluția sistemului (1) este dată de formula:

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

Vectorul  $f$  este:

$$f_i = (Cx^{(k)})_i + b_i, \quad i = 1, 2, \dots, n \quad (3)$$

$$(Cx^{(k)})_i = \sum_{j=1}^n c_{ij} x_j^{(k)} = - \sum_{j=i+1}^n a_{ij} x_j^{(k)}, \quad \forall i = 1, \dots, n \quad (4)$$

Folosind formula de rezolvare a sistemelor inferior triunghiulare (2), din relațiile (3) și (4) avem:

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

## Condiții suficiente de convergență pentru metoda Gauss-Seidel

### Propoziția 1

Dacă matricea  $A$  este astfel încât:

$$\sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \left( \frac{a_{ij}}{a_{ii}} \right)^2 < 1$$

atunci are loc convergența șirului construit cu metoda Gauss-Seidel la soluția sistemului  $Ax=b$ :

$$x^{(k)} \rightarrow x^*, k \rightarrow \infty \quad \forall x^{(0)} \in \mathbb{R}^n$$

## Propoziția 2 (*Criteriul dominanței diagonalei pe linii*)

Dacă matricea  $A$  este astfel încât:

$$\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| < |a_{ii}| \quad \forall i = 1, \dots, n$$

atunci:

$$x^{(k)} \rightarrow x^*, k \rightarrow \infty \quad \forall x^{(0)} \in \mathbb{R}^n$$

### Propoziția 3 (*Criteriul dominanței diagonalei pe coloane*)

Dacă matricea  $A$  este astfel încât:

$$\sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| < |a_{jj}| \quad \forall j = 1, \dots, n$$

atunci metoda Gauss-Seidel converge:

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*, \quad \forall \mathbf{x}^{(0)}$$

# **Calcul Numeric**

**Cursul 7**

**2022**

*Anca Ignat*

## **Metode iterative pentru matrice simetrice și pozitiv definite**

Considerăm cazul sistemelor liniare cu matricea sistemului simetrică și pozitiv definită:

$$A = A^T \text{ – matrice simetrică – } a_{ij} = a_{ji} \quad \forall i, j = 1, 2, \dots, n$$



$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{pmatrix} = A^T = \begin{pmatrix} a_{11} & a_{21} & a_{31} & \cdots & a_{n1} \\ a_{12} & a_{22} & a_{32} & \cdots & a_{n2} \\ a_{13} & a_{23} & a_{33} & \cdots & a_{n3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & a_{3n} & \cdots & a_{nn} \end{pmatrix}$$

$$A = A^T \quad \Rightarrow \quad A = L + D + L^T$$

$$D = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}] = \begin{pmatrix} a_{11} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & a_{22} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & a_{nn} \end{pmatrix}$$

$$L = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ a_{21} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ a_{31} & a_{32} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & \mathbf{0} \end{pmatrix} \quad L^T = \begin{pmatrix} \mathbf{0} & a_{12} & a_{13} & \cdots & a_{1n} \\ \mathbf{0} & \mathbf{0} & a_{23} & \cdots & a_{2n} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & a_{n-1n} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix}$$

## Definiții

Matricea  $A \in \mathbb{R}^{n \times n}$  se numește *pozitiv semidefinită* ( $A \geq 0$ ):

$$(Ax, x)_{\mathbb{R}^n} \geq 0 \quad \forall x \in \mathbb{R}^n.$$

Matricea  $A$  se numește *pozitiv definită* ( $A > 0$ ) dacă:

$$(Ax, x)_{\mathbb{R}^n} > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0.$$

## Propoziție

Dacă matricea  $A \in \mathbb{R}^{n \times n}$  este pozitiv definită atunci matricea  $A$  este nesingulară.

Demonstrație: Presupunem prin reducere la absurd că matricea  $A$  este pozitiv definită și singulară. Atunci, sistemul de ecuații liniare  $Ax=0$  are pe lângă soluția banală  $x=0$  și o soluție  $x^0 \neq 0$ . Avem:

$$x^0 \neq 0 \Rightarrow 0 < (Ax^0, x^0) = (0, x^0) = 0 \text{ contradicție!}$$

$$A > 0 \Rightarrow a_{ii} = (Ae_i, e_i) > 0 \quad \forall i = 1, \dots, n$$

## Lemă

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice simetrică și  $B \in \mathbb{R}^{n \times n}$  o matrice nesară astfel încât matricea  $P = B + B^T - A$  este pozitiv definită. Fie matricea  $M = I_n - B^{-1}A$ . Atunci raza spectrală a matricei  $M$  este strict subunitară dacă și numai dacă matricea  $A$  este pozitiv definită:

$$\rho(M) < 1 \Leftrightarrow A > 0$$

## Teoremă

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice simetrică, nesară, cu diagonala pozitivă,  $a_{ii} > 0$ ,  $\forall i = 1, \dots, n$  și  $b \in \mathbb{R}^n$  vectorul termenilor liberi. Atunci metoda lui Gauss-Seidel generează șiruri convergente la soluția  $x^* = A^{-1}b$ ,  $\forall x^{(0)}$  dacă și numai dacă  $A$  este pozitiv definită .

Demonstrație: Din teorema de convergență avem:

$$\mathbf{x}^{(k)} \rightarrow \mathbf{x}^* , \quad k \rightarrow \infty \quad \Leftrightarrow \quad \rho(\mathbf{M}) < 1$$

Dacă matricea  $\mathbf{A}$  se scrie sub forma:

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

matricele  $\mathbf{B}$  și  $\mathbf{C}$  sunt date de:

$$\mathbf{B} = \mathbf{L} + \mathbf{D} \quad , \quad \mathbf{C} = \mathbf{B} - \mathbf{A} = -\mathbf{L}^T$$

Matricea iterației  $\mathbf{M}$  este:

$$\mathbf{M} = \mathbf{B}^{-1}\mathbf{C} = \mathbf{B}^{-1}(\mathbf{B} - \mathbf{A}) = \mathbf{I}_n - \mathbf{B}^{-1}\mathbf{A}$$

Încercăm să aplicăm **Lema** de mai sus. Pentru aceasta verificăm dacă matricea  $P$  este pozitiv definită:

$$P = B + B^T - A = L + D + (L + D)^T - L - D - L^T = D$$

$$(Px, x)_{\mathbb{R}^n} = (Dx, x)_{\mathbb{R}^n} = ((a_{ii}x_i)_i, (x_i)_i)_{\mathbb{R}^n} = \sum_{i=1}^n a_{ii}x_i^2$$

$$a_{ii} > 0 \quad \forall i \Rightarrow (Px, x)_{\mathbb{R}^n} > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0 \Rightarrow P > 0$$

Putem aplica **Lema** de unde deducem convergența șirului construit cu metoda Gauss-Seidel doar în cazul în care matricea  $A$  este pozitiv definită:

$$x^{(k)} \rightarrow x^*, \quad k \rightarrow \infty \Leftrightarrow \rho(M) < 1 \Leftrightarrow A \text{ pozitiv definită}$$

## Metodele relaxării

Fie  $A \in \mathbb{R}^{n \times n}$  o matrice reală pătratică de dimensiune  $n$ , simetrică,  $A=A^T$  și pozitiv definită,  $A > \mathbf{0}$  și  $\mathbf{b} \in \mathbb{R}^n$  un vector real. Considerăm sistemul de ecuații liniare:

$$Ax = b$$

Deoarece matricea  $A$  este pozitiv definită sistemul de mai sus are soluție unică,  $\mathbf{x}^* = A^{-1}\mathbf{b}$ . Vom considera funcția  $f : \mathbb{R}^n \rightarrow \mathbb{R}_+$ :

$$f(y) = \left( A(\mathbf{x}^* - y), \mathbf{x}^* - y \right)_{\mathbb{R}^n}, \quad y \in \mathbb{R}^n$$



Din faptul că matricea  $A$  este pozitiv definită avem:

$$f(\mathbf{y}) \geq \mathbf{0}, \forall \mathbf{y} \in \mathbb{R}^n \quad \text{și} \quad f(\mathbf{y}) > \mathbf{0} = f(\mathbf{x}^*), \forall \mathbf{y} \neq \mathbf{x}^*$$

Prin urmare  $\mathbf{x}^*$  este și unica soluție a problemei de minimizare:

$$\min \{ f(\mathbf{y}); \mathbf{y} \in \mathbb{R}^n \} = \mathbf{0} = f(\mathbf{x}^*)$$

Vom căuta soluția sistemului  $A\mathbf{x}=\mathbf{b}$ ,  $\mathbf{x}^* = A^{-1}\mathbf{b}$  ca fiind soluția problemei de minimizare de mai sus folosind o metodă de tip relaxare de forma:

$$\mathbf{y}^{(0)} \in \mathbb{R}^n - \text{dat}, \quad \mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} + c_k \mathbf{e}_l, \quad l = l_k, \quad k = 0, 1, \dots$$

$$y_j^{(k+1)} = y_j^{(k)}, \quad \forall j \neq l, \quad y_l^{(k+1)} = y_l^{(k)} + c_k$$

Constanta  $c_k$  se determină astfel încât  $f(\mathbf{y}^{(k+1)}) < f(\mathbf{y}^{(k)})$  în speranța că șirul  $\mathbf{y}^{(k)}$  astfel construit converge la  $\mathbf{x}^*$ . Notăm cu

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{y}^{(k)} \quad \text{vectorul reziduu.}$$

Avem:

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{y}^{(k)} = \mathbf{A}\mathbf{x}^* - \mathbf{A}\mathbf{y}^{(k)} = \mathbf{A}(\mathbf{x}^* - \mathbf{y}^{(k)})$$

$$f(\mathbf{y}^{(k+1)}) = f(\mathbf{y}^{(k)}) - 2c_k r_l^{(k)} + c_k^2 a_{ll}$$

Pentru ca  $f(\mathbf{y}^{(k+1)}) < f(\mathbf{y}^{(k)})$  este necesar și suficient să alegem  $\mathbf{c}_k$  astfel ca:

$$\mathbf{c}_k^2 \mathbf{a}_{ll} - 2 \mathbf{c}_k \mathbf{r}_l^{(k)} < 0 \quad \Leftrightarrow \quad (\mathbf{a}_{ll} > 0) \mathbf{c}_k \in \left( \mathbf{0}, 2 \frac{\mathbf{r}_l^{(k)}}{\mathbf{a}_{ll}} \right) \text{ sau } \left( 2 \frac{\mathbf{r}_l^{(k)}}{\mathbf{a}_{ll}}, \mathbf{0} \right)$$

$$\mathbf{c}_k = \omega_k \frac{\mathbf{r}_l^{(k)}}{\mathbf{a}_{ll}}, \quad \text{cu } \omega_k \in (0, 2)$$

Metoda de relaxare obținută este următoarea:

$$\mathbf{y}^{(0)} \in \mathbb{R}^n - \text{dat}, \quad \mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} + \omega_k \frac{\mathbf{r}_l^{(k)}}{\mathbf{a}_{ll}} \mathbf{e}_l, \quad k = 0, 1, \dots, \quad \omega_k \in (0, 2)$$

Pentru a aproxima  $\mathbf{x}^*$  se deduce o clasă de metode numite *metodele relaxării succesive*. Aceste metode se obțin aplicând metodele de relaxare de mai sus. Vom considera:

$$\omega_k = \omega, \forall k$$

Vom construi un șir  $\{\mathbf{x}^{(k)}\} \subseteq \mathbb{R}^n$  astfel:

$\mathbf{x}^{(0)} = \mathbf{y}^{(0)}$  un vector din  $\mathbb{R}^n$  dat

$$l = 1 \quad \mathbf{y}^{(1)} = \mathbf{y}^{(0)} + \omega \frac{\mathbf{r}_1^{(0)}}{a_{11}} \mathbf{e}_1$$

$$l = 2 \quad \mathbf{y}^{(2)} = \mathbf{y}^{(1)} + \omega \frac{\mathbf{r}_2^{(1)}}{a_{22}} \mathbf{e}_2$$

⋮

$$l = n \quad \mathbf{y}^{(n)} = \mathbf{y}^{(n-1)} + \omega \frac{\mathbf{r}_n^{(n-1)}}{a_{nn}} \mathbf{e}_n$$

$$\mathbf{x}^{(1)} = \mathbf{y}^{(n)}$$

Trecerea de la iterația  $k$  la iterația următoare se face astfel:

$$\mathbf{x}^{(k)} = \mathbf{y}^{(kn)}$$

$$l = 1 \quad \mathbf{y}^{(kn+1)} = \mathbf{y}^{(kn)} + \omega \frac{\mathbf{r}_1^{(kn)}}{a_{11}} \mathbf{e}_1$$

$$l = 2 \quad \mathbf{y}^{(kn+2)} = \mathbf{y}^{(kn+1)} + \omega \frac{\mathbf{r}_2^{(kn+1)}}{a_{22}} \mathbf{e}_2$$

⋮

$$l = n \quad \mathbf{y}^{(kn+n)} = \mathbf{y}^{(kn+n-1)} + \omega \frac{\mathbf{r}_n^{(kn+n-1)}}{a_{nn}} \mathbf{e}_n$$

$$\mathbf{x}^{(k+1)} = \mathbf{y}^{((k+1)n)}, \quad k = 0, 1, 2, \dots$$

Acum putem scrie dependența vectorului  $\mathbf{x}^{(k+1)}$  de  $\mathbf{x}^{(k)}$ :

$\mathbf{x}^{(0)} \in \mathbb{R}^n$ ,  $\omega \in (0, 2)$  date,

$$x_i^{(k+1)} = x_i^{(k)} + \frac{\omega}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i}^n a_{ij} x_j^{(k)} \right), \quad i = 1, 2, \dots, n,$$

$$x_i^{(k+1)} = (1 - \omega) x_i^{(k)} + \frac{\omega}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right), \quad i = 1, 2, \dots, n,$$

$$k = 0, 1, 2, \dots$$

Metodele de mai sus poartă numele de *metodele relaxării succesive*. Pentru  $\omega = 1$  obținem metoda Gauss-Seidel.

- $0 < \omega < 1$  metodele se numesc de *sub-relaxare* și pot fi folosite în cazul când metoda Gauss-Seidel diverge.
- $1 < \omega < 2$  metodele se numesc de *supra-relaxare* și pot fi folosite pentru accelerarea convergenței în cazul când metoda Gauss-Seidel converge.

Rearanjând formulele de mai sus avem:

$$\sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} + \frac{a_{ii}}{\omega} x_i^{(k+1)} = \left( B x^{(k+1)} \right)_i = \frac{(1-\omega)}{\omega} a_{ii} x_i^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} + b_i =$$

$$= \left( C x^{(k)} \right)_i + (b)_i$$



Matricea  $A$  fiind simetrică, poate fi scrisă sub forma:

$$A = L + D + L^T \quad \text{cu} \quad L = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & 0 & \mathbf{0} \\ a_{21} & \mathbf{0} & \cdots & 0 & \mathbf{0} \\ \vdots & & & & \\ a_{n1} & a_{n2} & \cdots & a_{nn-1} & \mathbf{0} \end{pmatrix},$$

$$D = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}]$$

Cu aceste notații, matricile  $B$  și  $C$  de mai sus pot fi scrise astfel:

$$B = L + \frac{1}{\omega} D \quad , \quad C = \frac{1-\omega}{\omega} D - L^T$$

Vom verifica dacă metodele relaxării succesive se înscriu în clasa generală de metode iterative, adică vom verifica dacă  $A=B-C$  :

$$B - C = L + \frac{1}{\omega} D - \frac{1-\omega}{\omega} D + L^T = L + D + L^T = A$$

Convergența șirului  $x^{(k)}$  la soluția  $x^* = A^{-1}b$  ?

## Teoremă

Fie o matrice  $A \in \mathbb{R}^{n \times n}$ , simetrică,  $A=A^T$  cu  $\det A \neq 0$ ,  $a_{ii} > 0$ ,  $\forall i = 1, \dots, n$ ,  $b \in \mathbb{R}^n$  un vector real și  $\omega \in (0, 2)$ . Atunci șirul  $x^{(k)}$  construit cu o metoda de relaxare succesivă converge la soluția  $x^*$  a sistemului liniar  $Ax=b$  oricare ar fi iterația inițială  $x^{(0)}$  dacă și numai dacă matricea  $A$  este pozitiv definită.

$$x^{(k)} \rightarrow x^*, k \rightarrow \infty, \forall x^{(0)} \Leftrightarrow (Ax, x) > 0, \forall x \in \mathbb{R}^n, x \neq 0$$

Demonstrație: Vom verifica dacă raza spectrală a matricei iterației este subunitară folosind **Lema**. Avem:

$$M = B^{-1}C = B^{-1}(B - A) = I_n - B^{-1}A$$

$$B = L + \frac{1}{\omega}D, \quad \det B = \frac{1}{\omega^n} a_{11}a_{22} \cdots a_{nn} \neq 0 \quad (a_{ii} > 0, \forall i)$$

Matricea  $A$  este simetrică iar matricea  $B$  este nesingulară. Pentru a fi îndeplinite ipotezele Lemei trebuie să verificăm că matricea  $P$  este pozitiv definită:

$$P = B + B^T - A = L + \frac{1}{\omega}D + L^T + \frac{1}{\omega}D - L - D - L^T = \frac{2 - \omega}{\omega}D$$

$$\begin{aligned} (Px, x) &= \frac{(2-\omega)}{\omega} \sum_{i=1}^n a_{ii} x_i^2 > 0, \quad x \neq 0 \quad (a_{ii} > 0, \forall i) \Leftrightarrow \\ &\Leftrightarrow \frac{(2-\omega)}{\omega} > 0 \quad \Leftrightarrow \omega \in (0, 2) \end{aligned}$$

Toate ipotezele lemei sunt îndeplinite, prin urmare avem convergența dorită.

## Valori și vectori proprii (eigenvalues, eigenvectors)

### Definiție

Fie  $A \in \mathbb{R}^{n \times n}$ . Numărul complex  $\lambda \in \mathbb{C}$  se numește *valoare proprie* a matricei  $A$  dacă există un vector  $u \in \mathbb{C}^n$ ,  $u \neq \mathbf{0}$  astfel ca:

$$Au = \lambda u$$

Vectorul  $u$  se numește *vector propriu* asociat valorii proprii  $\lambda$ .

Pentru existența vectorului  $u \neq \mathbf{0}$  este necesar și suficient ca matricea  $(\lambda I_n - A)$  să fie singulară, adică  $\det(\lambda I_n - A) = 0$ .

Polinomul de grad  $n$ :

$$p_A(\lambda) = \det(\lambda I_n - A)$$

se numește *polinom caracteristic* al matricei  $A$ .

### Propoziția 1

Fie rădăcinile polinomului caracteristic  $\lambda_1, \lambda_2, \dots, \lambda_n$  distincte,  $\lambda_i \neq \lambda_j$  pentru  $1 \leq i < j \leq n$  și  $u_1, u_2, \dots, u_n$  vectorii proprii corespunzători. Atunci  $u_1, u_2, \dots, u_n$  sunt liniar independenți. (demonstrația se face prin inducție)

## Propoziția 2

Fie valorile proprii  $\lambda_i$  ale matricei  $A \in \mathbb{R}^{n \times n}$  distincte. Atunci există o matrice nesingulară  $T$  astfel ca:

$$T^{-1}AT = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_n].$$

Demonstrație. Fie  $u_1, u_2, \dots, u_n$  vectorii proprii ai matricei  $A$ . Considerăm matricea  $T$  ale cărei coloane sunt vectorii proprii  $u_i$ ,  $T = [u_1 \ u_2 \ \dots \ u_n]$ . Deoarece vectorii proprii sunt liniar independenți conform propoziției 1 rezultă că matricea  $T$  este nesingulară. Vom avea:

$$AT = [Au_1 \ Au_2 \ \dots \ Au_n] = [\lambda_1 u_1 \ \lambda_2 u_2 \ \dots \ \lambda_n u_n] = T \cdot \text{diag}[\lambda_1 \ \lambda_2 \ \dots \ \lambda_n]$$

Înmulțind la stânga cu  $T^{-1}$  obținem concluzia propoziției 2.



## Definiție

**Matricile**  $A$  și  $B$  sunt *asemenea* (notație  $A \sim B$ ) dacă și numai dacă există o matrice nesingulară  $T$  ( $\det T \neq 0$ ) astfel ca:

$$A = T B T^{-1}$$

## Propoziția 3

$$A \sim B \Rightarrow p_A(\lambda) = p_B(\lambda).$$

Demonstrație.

$$\begin{aligned} p_A(\lambda) &= \det(\lambda I_n - A) = \det(\lambda I_n - T B T^{-1}) = \det(\lambda T T^{-1} - T B T^{-1}) \\ &= \det(T(\lambda I_n - B)T^{-1}) = \det(T) \det(\lambda I_n - B) \det(T^{-1}) = p_B(\lambda) \end{aligned}$$

Propoziția 3 ne spune că matricile asemenea au același polinom caracteristic și aceleași valori proprii.

## Teorema lui Gershgorin

Fie  $A \in \mathbb{R}^{n \times n}$  și  $\lambda \in \mathbb{C}$  o valoare proprie oarecare a matricei  $A$ .  
Atunci:

$$\exists i_0 \in \{1, 2, \dots, n\} \text{ astfel încât } |\lambda - a_{i_0 i_0}| \leq r_{i_0}, \quad r_{i_0} = \sum_{\substack{j=1 \\ j \neq i_0}}^n |a_{i_0 j}|.$$

(Valoarea proprie  $\lambda$  se află în cercul din planul complex de centru  $a_{i_0 i_0}$  și rază  $r_{i_0}$ .)

Demonstrație. Fie  $\lambda \in \mathbb{C}$  o valoare proprie a matricei  $A$  și  $u \neq 0$  un vector propriu asociat valorii proprii  $\lambda$ ,  $Au = \lambda u$ .

Avem:

$$\lambda u_i = a_{ii}u_i + \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}u_j \Leftrightarrow (\lambda - a_{ii})u_i = \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}u_j, i = 1, \dots, n.$$

Fie  $i_0$  astfel ca  $|u_{i_0}| = \|u\|_\infty = \max\{|u_k|; k = 1, \dots, n\} > 0$  ( $u \neq 0$ ).

Vom avea:

$$|\lambda - a_{i_0 i_0}| = \left| \sum_{\substack{j=1 \\ j \neq i_0}}^n a_{i_0 j} \frac{u_j}{u_{i_0}} \right| \leq \sum_{\substack{j=1 \\ j \leq i_0}}^n |a_{i_0 j}| \frac{|u_j|}{|u_{i_0}|} \leq r_{i_0}, \text{ ținând seama că } \frac{|u_j|}{|u_{i_0}|} \leq 1.$$

**Observație.** Presupunem că matricea  $A$  are  $n$  vectori proprii liniar independenți  $u^1, u^2, \dots, u^n$  asociați valorilor proprii  $\lambda_1, \lambda_2, \dots, \lambda_n$ . Fie  $U = [u^1 \ u^2 \ \dots \ u^n]$ . Datorită independenței vectorilor  $u^k$  rezultă că matricea  $U$  este nesingulară și avem:

$$\Lambda = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_n] \quad U^{-1}AU = \Lambda.$$

Considerăm matricea perturbată:

$$A(\varepsilon) = A + \varepsilon B.$$

$$U^{-1}A(\varepsilon)U = \Lambda + \varepsilon U^{-1}BU = \Lambda + \varepsilon C.$$

$A(\varepsilon) \sim U^{-1}A(\varepsilon)U \Rightarrow$  au aceleași valori proprii  $\lambda_i(\varepsilon)$

$$|\lambda(\varepsilon) - \lambda_i - \varepsilon c_{ii}| \leq \varepsilon \sum_{\substack{j=1 \\ j \neq i}}^n |c_{ij}| \Rightarrow |\lambda(\varepsilon) - \lambda_i| = \mathcal{O}(\varepsilon).$$

## Metoda puterii pentru matrice simetrice

### Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ . Atunci toate valorile proprii ale matricei  $A$  sunt numere reale.

Demonstrație. Fie  $\lambda \in \mathbb{C}$  și  $u \in \mathbb{C}^n$ ,  $u \neq \mathbf{0}$ ,  $Au = \lambda u$ .

Considerăm produsul scalar:

$$(Au, u)_{\mathbb{C}^n} = \lambda (u, u)_{\mathbb{C}^n} = \lambda \|u\|_2^2.$$

$$(Au, u)_{\mathbb{C}^n} = (u, A^T u)_{\mathbb{C}^n} = (u, Au)_{\mathbb{C}^n} = \overline{(Au, u)_{\mathbb{C}^n}} \Rightarrow (Au, u)_{\mathbb{C}^n} \in \mathbb{R}$$

$$\lambda = \frac{(Au, u)_{\mathbb{C}^n}}{\|u\|_2^2} \in \mathbb{R}.$$

### Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ . Atunci există o bază ortonormată de vectori proprii ai matricei  $A$ ,  $\{u^1, u^2, \dots, u^n\}$  :

$$(u^i, u^j)_{\mathbb{R}^n} = \delta_{ij} = \begin{cases} \mathbf{1} & \text{dacă } i = j \\ \mathbf{0} & \text{dacă } i \neq j \end{cases}$$

Echivalent, putem scrie ca există vectori proprii  $\{u^1, u^2, \dots, u^n\}$  asociați valorilor proprii reale  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  astfel ca:

$$AU = U\Lambda \Leftrightarrow U^T AU = \Lambda$$

cu  $\Lambda = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_n]$  și  $U = [u^1 \ u^2 \ \dots \ u^n]$  matrice ortogonală.

## Definiție

Se numește *coeficient Rayleigh* al vectorului  $\mathbf{u} \in \mathbb{R}^n$  pentru matricea  $A$  următoarea mărime scalară:

$$r(\mathbf{u}) = \frac{\mathbf{u}^T A \mathbf{u}}{\mathbf{u}^T \mathbf{u}} = \frac{(A\mathbf{u}, \mathbf{u})_{\mathbb{R}^n}}{(\mathbf{u}, \mathbf{u})_{\mathbb{R}^n}} = \frac{(A\mathbf{u}, \mathbf{u})_{\mathbb{R}^n}}{\|\mathbf{u}\|_2^2}$$

Se verifică ușor că dacă  $\mathbf{u} \in \mathbb{R}^n$  este vector propriu al matricei  $A$  asociat valorii proprii  $\lambda$  atunci  $r(\mathbf{u}) = \lambda$ .

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ . Matricea are valori proprii reale  $\lambda_1, \lambda_2, \dots, \lambda_n$ . Presupunem în plus că:

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n| \geq 0$$

Metoda puterii este un algoritm care aproximează valoarea proprie de modul maxim  $\lambda_1$  și un vector propriu asociat.

Se pornește de la un vector nenul de normă euclidiană 1,  $\mathbf{u}^{(0)} \in \mathbb{R}^n$ ,  $\|\mathbf{u}^{(0)}\|_2 = 1$  și se construiește următorul șir de vectori de normă euclidiană 1:

$$\mathbf{u}^{(0)}, \mathbf{u}^{(1)} = \frac{1}{\|\mathbf{A}\mathbf{u}^{(0)}\|_2} \mathbf{A}\mathbf{u}^{(0)}, \mathbf{u}^{(2)} = \frac{1}{\|\mathbf{A}\mathbf{u}^{(1)}\|_2} \mathbf{A}\mathbf{u}^{(1)}, \dots,$$

$$\mathbf{u}^{(k)} = \frac{1}{\|\mathbf{A}\mathbf{u}^{(k-1)}\|_2} \mathbf{A}\mathbf{u}^{(k-1)}, \dots$$

În anumite condiții acest șir converge la un vector propriu asociat valorii proprii  $\lambda_1$ , iar coeficienții Rayleigh corespunzători converg către  $\lambda_1$ .



## Teoremă

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$  o matrice simetrică pentru care valorile proprii îndeplinesc condiția:

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n| \geq 0.$$

Dacă  $\mathbf{u}^{(0)} \in \mathbb{R}^n$ ,  $\|\mathbf{u}^{(0)}\|_2 = 1$ ,  $(\mathbf{u}^{(0)}, \mathbf{u}^1)_{\mathbb{R}^n} \neq 0$  ( $\mathbf{u}^1$  vector propriu asociat lui  $\lambda_1$ ) atunci:

$$\mathbf{u}^{(k)} = \frac{1}{\|A^k \mathbf{u}^{(0)}\|_2} A^k \mathbf{u}^{(0)} \rightarrow \mathbf{u}^1 \text{ (vector propriu asociat lui } \lambda_1 \text{)}$$

$$r(\mathbf{u}^{(k)}) \rightarrow \lambda_1$$

Demonstrație. Fie  $\{u^1, u^2, \dots, u^n\}$  vectori proprii asociați valorilor proprii  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  care formează o bază ortonormată în  $\mathbb{R}^n$ . Avem:

$$u^{(0)} = a_1 u^1 + a_2 u^2 + \dots + a_n u^n, \quad a_i \in \mathbb{R}$$

Deoarece  $(u^{(0)}, u^1)_{\mathbb{R}^n} \neq 0$  rezultă că  $a_1 \neq 0$ . Din construcția șirului  $u^{(k)}$  deducem că există o constantă  $c_k$  astfel ca:

$$\begin{aligned}
\mathbf{u}^{(k)} &= \mathbf{c}_k \mathbf{A}^k \mathbf{u}^{(0)} = \\
&= \mathbf{c}_k \mathbf{A}^k (a_1 \mathbf{u}^1 + a_2 \mathbf{u}^2 + \cdots + a_n \mathbf{u}^n) = \\
&= \mathbf{c}_k (a_1 \lambda_1^k \mathbf{u}^1 + a_2 \lambda_2^k \mathbf{u}^2 + \cdots + a_n \lambda_n^k \mathbf{u}^n) = \\
&= \mathbf{c}_k \lambda_1^k \left[ a_1 \mathbf{u}^1 + a_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k \mathbf{u}^2 + \cdots + a_n \left( \frac{\lambda_n}{\lambda_1} \right)^k \mathbf{u}^n \right]
\end{aligned}$$

Din această ultimă relație, din faptul că  $\lambda_1$  este valoare proprie dominantă și  $a_1 \neq 0$  deducem că pentru  $k$  suficient de mare vectorul  $\mathbf{u}^{(k)}$  se aliniază după vectorul propriu  $\mathbf{u}^1$ :

$$\mathbf{u}^{(k)} \approx \mathbf{c}_k \lambda_1^k a_1 \mathbf{u}^1$$

## Metoda puterii

$$\mathbf{u}^{(0)} \in \mathbb{R}^n, \|\mathbf{u}^{(0)}\|_2 = \mathbf{1};$$

$$k = \mathbf{0};$$

*do*

$$\circ k ++;$$

$$\circ \mathbf{w} = A\mathbf{u}^{(k-1)};$$

$$\circ \mathbf{u}^{(k)} = \frac{\mathbf{1}}{\|\mathbf{w}\|_2} \mathbf{w};$$

$$\circ \lambda_k = r(\mathbf{u}^{(k)}) = \left( A\mathbf{u}^{(k)}, \mathbf{u}^{(k)} \right)_{\mathbb{R}^n};$$

$$\mathbf{while} (\|A\mathbf{u}^{(k)} - \lambda_k \mathbf{u}^{(k)}\| > \varepsilon \text{ \textit{și} } k \leq k_{\max});$$

## Metoda iterației inverse

Considerăm o matrice simetrică  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$  și  $\mu \in \mathbb{R}$  un număr real care nu este valoare proprie a matricei  $A$ . Vom folosi metoda puterii pentru a aproxima valoarea proprie a matricei  $A$  care este cea mai apropiată de  $\mu$  și un vector propriu asociat.

$$\mu \neq \text{valoare proprie} \rightarrow \det(A - \mu I_n) \neq 0 \rightarrow \exists (A - \mu I_n)^{-1}$$

Fie  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  valorile proprii reale ale matricei  $A$ .

Valorile proprii ale matricei  $(A - \mu I_n)^{-1}$  sunt:

$$\left\{ \frac{1}{(\lambda_1 - \mu)}, \frac{1}{(\lambda_2 - \mu)}, \dots, \frac{1}{(\lambda_n - \mu)} \right\}$$

Matricele  $A$  și  $(A - \mu I_n)^{-1}$  au aceiași vectori proprii. Să presupunem că  $\lambda_I$  este valoarea proprie cea mai apropiată de  $\mu$  (și singura). Atunci:

$$\frac{1}{|\lambda_I - \mu|} > \frac{1}{|\lambda_j - \mu|} \quad \forall j \neq I$$

Această relație sugerează ideea aplicării metodei puterii matricei  $(A - \mu I_n)^{-1}$  pentru a aproxima valoarea proprie  $(\lambda_I - \mu)^{-1}$  și a unui vector propriu asociat. Algoritmul duce la aproximarea valorii proprii cea mai apropiată de  $\mu$ ,  $\lambda_I$  și a unui vector propriu asociat acestei autovalori,  $u^I$ .

## Metoda iterației inverse

$$\mathbf{u}^{(0)} \in \mathbb{R}^n, \|\mathbf{u}^{(0)}\|_2 = \mathbf{1};$$

$$k = \mathbf{0};$$

*do*

- $k++$ ;

- Se rezolvă sistemul  $(A - \mu I_n)\mathbf{w} = \mathbf{u}^{(k-1)}$ ;

- $\mathbf{u}^{(k)} = \frac{\mathbf{1}}{\|\mathbf{w}\|_2} \mathbf{w}$ ;

- $\lambda_k = r(\mathbf{u}^{(k)}) = \left( A\mathbf{u}^{(k)}, \mathbf{u}^{(k)} \right)_{\mathbb{R}^n}$ ;

*while* ( $\|A\mathbf{u}^{(k)} - \lambda_k \mathbf{u}^{(k)}\| > \varepsilon$  și  $k \leq k_{\max}$ );



# **Calcul Numeric**

**Cursul 8**

**2022**

*Anca Ignat*

## Forma superioară Hessenberg

Spunem că o matrice  $H \in \mathbb{R}^{n \times n}$  este în *formă superioară Hessenberg* dacă:

$$h_{ij} = 0, i = 1, \dots, n, j = 1, \dots, i - 2$$

O matrice în formă Hessenberg arată astfel:

$$H = \begin{pmatrix} h_{11} & h_{12} & h_{13} & \cdots & h_{1n-1} & h_{1n} \\ h_{21} & h_{22} & h_{23} & \cdots & h_{2n-1} & h_{2n} \\ \mathbf{0} & h_{32} & h_{33} & \cdots & h_{3n-1} & h_{3n} \\ \mathbf{0} & \mathbf{0} & h_{43} & \cdots & h_{4n-1} & h_{4n} \\ \vdots & & & & & \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & h_{nn-1} & h_{nn} \end{pmatrix}$$

Ne interesează un algoritm care să transforme o matrice pătratică  $A$  oarecare într-o matrice Hessenberg superioară  $H$  care să aibă aceleași valori proprii:

$A \rightarrow H$  a.î.  $H \sim A$ ,  $H = \tilde{P}A \tilde{P}^{-1}$ ,  $\tilde{P}$  matrice nesingulară  
Algoritmul este o adaptare a algoritmului lui Housholder și se desfășoară în  $(n-2)$  pași, folosind matricile de reflexie pentru a transforma matricea.

### **Pas 1**

se efectuează operațiile  $A = P_1 A P_1$  (matricea  $P_1$  se alege astfel încât coloana 1 să fie transformată în formă superior Hessenberg)

## Pas 2

$$A = P_2 A P_2 = P_2 (P_1 A^{init}) P_2$$

( $P_2$  transformă coloana 2 în formă superior Hessenberg fără să schimbe coloana 1)

## Pas $r$

$$A = P_r A P_r = P_r (P_{r-1} \cdots P_1 A^{init} P_1 \cdots P_{r-1}) P_r$$

(se transformă coloana  $r$  în formă superior Hessenberg fără să schimbe primele  $(r-1)$  coloane)

## Pasul $r$ ( $r=1,2,\dots,n-2$ )

La intrarea în pasul  $r$  matricea  $A$  are primele  $(r-1)$  coloane în formă superior Hessenberg. La ieșirea din pasul  $r$  matricea  $A$  va avea primele  $r$  coloane în formă superior Hessenberg:

$$A_{ies} = P_r A_{intr} P_r , \quad A_{ies} \sim A_{intr}$$

$$P_r = I_n - 2v^r (v^r)^T , \quad v^r \in R^n , \quad \|v^r\|_2 = 1$$

Vectorul  $v^r$  se alege astfel ca matricea  $A_{ies}$  să aibă coloana  $r$  în formă superior Hessenberg și să nu schimbe primele  $(r-1)$  coloane ale matricii  $A_{intr}$ .

## Calculul matricii $P_r$

$$P = I_n - \frac{1}{\beta} uu^T$$

$$\beta = \sigma - k \cdot a_{r+1r}$$

$$k^2 = \sigma = a_{r+1r}^2 + \cdots + a_{ir}^2 + \cdots + a_{nr}^2 = \sum_{i=r+1}^n a_{ir}^2 \Rightarrow k = \pm \sqrt{\sigma}$$

$$\text{semn } k = -\text{semn } a_{r+1r}$$

$$u := \begin{pmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ a_{r+1r} - k \\ \vdots \\ a_{ir} \\ \vdots \\ a_{nr} \end{pmatrix}$$

$$\beta = \mathbf{0} \rightarrow r = r + 1 \quad (P = I_n)$$



Algoritmul de trecere de la matricea  $A$  la matricea  $P_r A$  este următorul:

$$(P_r A)e_j = \begin{cases} Ae_j & \text{pentru } j = 1, \dots, r-1 \\ (a_{1r}, a_{2r}, \dots, a_{rr}, k, \mathbf{0}, \dots, \mathbf{0})^T & \text{pentru } j = r \\ Ae_j - \frac{\gamma_j}{\beta} u & \text{pentru } j = r+1, \dots, n \end{cases}$$

$$\gamma_j = (Ae_j, u)_{\mathbb{R}^n} = \sum_{i=r+1}^n u_i a_{ij}$$

$$u_i = \mathbf{0}, \quad i = 1, \dots, r, \quad u_{r+1} = a_{r+1r} - k, \quad u_i = a_{ir}, \quad i = r+2, \dots, n$$

Vom descrie în continuare cum se efectuează operația  $A := AP_r$  fără a face înmulțire matricială (matricea  $A$  este cea obținută mai sus având primele  $r$  coloane în formă superior Hessenberg).

Vom arata că această operație nu schimbă forma superior Hessenberg obținută. Vom pune în evidență transformările liniilor matricii  $A$ . Pentru  $i=1, \dots, n$  avem:

$$\begin{aligned} e_i^T (AP) &= \text{noua linie } i \text{ a matricii } AP = (e_i^T A) \left( I_n - \frac{1}{\beta} uu^T \right) = \\ &= e_i^T A - \frac{1}{\beta} (e_i^T A)u u^T = e_i^T A - \frac{\gamma_i}{\beta} u^T \end{aligned}$$

unde

$$\gamma_i = (e_i^T A)u = a_{ir+1}u_{r+1} + \cdots + a_{in}u_n$$

Elementele liniei  $i$  se schimbă astfel:

$$a_{ij} = a_{ij} - \frac{\gamma_i}{\beta} u_j, \quad j = r+1, \dots, n, \quad i = 1, \dots, n$$

Operația  $A := AP_r$  nu modifică primele  $r$  coloane ale matricii  $A$ , ele rămânând în formă superior Hessenberg.

## *Algoritmul de obținere a formei superior Hessenberg*

for  $r = 1, \dots, n - 2$

// construcția matricii  $P_r$  – constanta  $\beta$  și vectorul  $u$

- $\sigma = \sum_{i=r+1}^n a_{ir}^2;$

- if (  $\sigma \leq \varepsilon$  ) break ; //  $r = r + 1 \leftrightarrow P_r = I_n$

- $k = \sqrt{\sigma};$

- if (  $a_{r+1r} > 0$  )  $k = -k;$

- $\beta = \sigma - k a_{r+1r};$

- $u_{r+1} = a_{r+1r} - k; u_i = a_{ir}, i = r + 2, \dots, n;$

//  $A = P_r * A$

// transformarea coloanelor  $j = r + 1, \dots, n$

- for  $j = r + 1, \dots, n$

$$\left\{ \begin{array}{l} \gamma = (\gamma_j / \beta) = (Ae_j, u) / \beta = \left( \sum_{i=r+1}^n u_i a_{ij} \right) / \beta; \\ \text{for } i = r + 1, \dots, n \\ a_{ij} = a_{ij} - \gamma * u_i; \end{array} \right.$$

// transformarea coloanei  $r$  a matricii  $A$

- $a_{r+1r} = k$ ;  $a_{ir} = 0$ ,  $i = r + 2, \dots, n$ ;

$$// \mathbf{A} = \mathbf{A} * \mathbf{P}_r$$

// transformarea liniilor  $i = \mathbf{1}, \dots, n$

- for  $i = \mathbf{1}, \dots, n$

$$\left\{ \begin{array}{l} \boldsymbol{\gamma} = (\gamma_i / \boldsymbol{\beta}) = ((\mathbf{e}_i^T \mathbf{A})\mathbf{u}) / \boldsymbol{\beta} = \left( \sum_{j=r+1}^n u_j \mathbf{a}_{ij} \right) / \boldsymbol{\beta}; \\ \text{for } j = r + \mathbf{1}, \dots, n \\ \mathbf{a}_{ij} = \mathbf{a}_{ij} - \boldsymbol{\gamma} * u_j; \end{array} \right.$$

## Algoritmul *QR* de aproximare a valorilor proprii ale unei matrice oarecare

Prezentăm în continuare cel mai folosit algoritm de aproximare a valorilor proprii pentru matrice pătratice oarecare.

Spunem că o matrice  $S \in \mathbb{R}^{n \times n}$  este în *formă Schur reală* dacă matricea  $S$  este în formă superior Hessenberg și în plus este bloc-diagonală:

$$S = \begin{pmatrix} S_{11} & S_{12} & \cdots & S_{1p} \\ \mathbf{0} & S_{22} & \cdots & S_{2p} \\ \vdots & & & \\ \mathbf{0} & \mathbf{0} & \cdots & S_{pp} \end{pmatrix}$$

blocurile  $S_{ii}$  sunt astfel ca:

- $S_{ii} \in \mathbb{R}$  - este valoare proprie reală
- $S_{ii} \in \mathbb{R}^{2 \times 2}$  - este bloc corespunzător valorilor proprii complexe

Valorile proprii corespunzătoare blocului

$S_{ii} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in \mathbb{R}^{2 \times 2}$  sunt rădăcinile ecuației:



$$\begin{vmatrix} \lambda - a & -b \\ -c & \lambda - d \end{vmatrix} = (\lambda - a)(\lambda - d) - bc = \lambda^2 - (a + d)\lambda + ad - bc = 0$$

Se presupune că această ecuație de gradul 2 are rădăcini complexe.

**Algoritmul QR** de aproximare a valorilor proprii construiește un șir de matrici  $A^{(k)} \in \mathbb{R}^{n \times n}$ , matrici asemenea cu matricea  $A$ ,  $A^{(k)} \sim A, \forall k$ , șir care converge la o matrice în formă Schur reală,  $A^{(k)} \rightarrow S, k \rightarrow \infty$ . Matricea limită  $S$  este asemenea cu matricea  $A$ , valorile proprii ale matricii  $S$  fiind ușor de calculat. Șirul  $A^{(k)}$  se construiește astfel:

$$A^{(0)} := A, \quad A^{(0)} = Q_0 R_0 \text{ (descomp. } QR \text{ calc. pentru matricea } A^{(0)})$$

$$A^{(1)} := R_0 Q_0, \quad A^{(1)} = Q_1 R_1 \text{ (descomp. } QR \text{ calc. pentru matricea } A^{(1)})$$

$$A^{(2)} := R_1 Q_1$$

⋮

$$A^{(k)} = Q_k R_k \text{ (descomp. } QR \text{ calc. pentru matricea } A^{(k)}),$$

$$A^{(k+1)} := R_k Q_k, \quad k = 0, 1, 2, \dots$$

Matricile  $Q_k$  sunt matrici ortogonale ( $Q_k^{-1} = Q_k^T$ ) iar matricile  $R_k$  sunt superior triunghiulare.

Matricile  $A^{(k)}$  și  $A^{(k+1)}$  sunt asemenea:

$$Q_k^T * | A^{(k)} = Q_k R_k \Rightarrow R_k = Q_k^T A^{(k)}$$

$$A^{(k+1)} = R_k Q_k = Q_k^T A^{(k)} Q_k \Rightarrow A^{(k+1)} \sim A^{(k)}, \forall k$$

Matricile șirului construit sunt toate asemenea prin urmare au aceleași valori proprii anume cele ale matricii inițiale  $A = A^{(0)}$ :

$$A = A^{(0)} \sim A^{(1)} \sim \dots \sim A^{(k)} \sim \dots \sim S$$

Dacă matricea  $A^{(k)}$  este în formă superioară Hessenberg, atunci descompunerea  $QR$  realizată cu algoritmul lui Givens se simplifică. Reamintim algoritmul lui Givens:

$$R_{n-1n}(\theta_{n-1n}) \cdots R_{pn}(\theta_{pn}) \cdots R_{pp+1}(\theta_{pp+1}) \cdots R_{1n}(\theta_{1n}) \cdots R_{12}(\theta_{12}) A = R$$

Dacă matricea  $A$  este în formă Hessenberg în algoritmul lui Givens, din cele  $\frac{n(n-1)}{2}$  înmulțiri cu matrici de rotație rămân doar  $(n-1)$ :

$$R_{n-1n}(\theta_{n-1n}) \cdots R_{pp+1}(\theta_{pp+1}) \cdots R_{23}(\theta_{23}) R_{12}(\theta_{12}) A = R.$$

Problema care se pune este dacă pornind cu o matrice în formă Hessenberg, toate matricile șirului rămân în formă Hessenberg:

$A^{(k)}$  (în formă Hessenberg) =  $H = QR$  (cu Givens)  $\Rightarrow ?$

$A^{(k+1)} = \bar{H} = RQ = Q^T A^{(k)} Q = Q^T H Q$  – este tot în formă Hessenberg ?

Avem:

$$\bar{H} = Q^T H Q = R R_{12}^T(\theta_{12}) \cdots R_{rr+1}^T(\theta_{rr+1}) \cdots R_{n-1n}^T(\theta_{n-1n})$$

Notăm cu:

$$\bar{R} = R R_{12}^T(\theta_{12})$$

pentru care avem:

$$\begin{cases} \bar{r}_{i1} = cr_{i1} + sr_{i2}, \forall i \\ \bar{r}_{i2} = -sr_{i1} + cr_{i2}, \forall i \end{cases} + \begin{cases} r_{i1} = 0, i = 2, \dots, n \\ r_{i2} = 0, i = 3, \dots, n \end{cases} \Rightarrow \begin{cases} \bar{r}_{i1} = 0, i = 3, \dots, n \\ \bar{r}_{i2} = 0, i = 3, \dots, n \end{cases}$$

deci coloana  $\mathbf{1}$  se transformă în formă Hessenberg iar coloana  $\mathbf{2}$  rămâne în formă superior triunghiulară.

La pasul  $p$  avem:

$$\left( \mathbf{R} \mathbf{R}_{12}^T(\theta_{12}) \cdots \mathbf{R}_{p-1p}^T(\theta_{p-1p}) \right) \mathbf{R}_{pp+1}^T(\theta_{pp+1}) = \tilde{\mathbf{R}} \mathbf{R}_{pp+1}^T(\theta_{pp+1}) = \bar{\mathbf{R}},$$

$$\tilde{\mathbf{R}} = \mathbf{R} \mathbf{R}_{12}^T(\theta_{12}) \cdots \mathbf{R}_{p-1p}^T(\theta_{p-1p})$$

matricea  $\tilde{\mathbf{R}}$  are primele  $(p-1)$  coloane în formă Hessenberg iar restul coloanelor sunt în formă superior triunghiulară. Vom arata că la acest pas matricea  $\bar{\mathbf{R}}$  va avea primele  $p$  coloane în formă Hessenberg iar restul coloanelor în formă superior triunghiulară. Operația  $\bar{\mathbf{R}} := \tilde{\mathbf{R}} \mathbf{R}_{pp+1}^T(\theta_{pp+1})$  presupune doar schimbarea elementelor coloanelor  $p$  și  $p+1$ :

$$\begin{cases} \bar{\mathbf{r}}_{ip} = c\tilde{\mathbf{r}}_{ip} + s\tilde{\mathbf{r}}_{ip+1}, \forall i \\ \bar{\mathbf{r}}_{ip+1} = -s\tilde{\mathbf{r}}_{ip} + c\tilde{\mathbf{r}}_{ip+1}, \forall i \end{cases} + \begin{cases} \tilde{\mathbf{r}}_{ip} = \mathbf{0}, i = p+1, \dots, n \\ \tilde{\mathbf{r}}_{ip+1} = \mathbf{0}, i = p+2, \dots, n \end{cases} \Rightarrow$$

$$\begin{cases} \bar{\mathbf{r}}_{ip} = \mathbf{0}, i = p+2, \dots, n \\ \bar{\mathbf{r}}_{ip+1} = \mathbf{0}, i = p+2, \dots, n \end{cases}$$

Observăm din relația de mai sus că în matricea  $\bar{\mathbf{R}}$  coloana  $p$  are formă Hessenberg iar coloana  $p+1$  rămâne în formă superior triunghiulară (celelalte elemente din matrice nu se modifică).

Prin urmare după pasul  $n-1$  matricea  $\bar{\mathbf{H}} = \mathbf{A}^{(k+1)}$  este în formă superioară Hessenberg. Algoritmul  $QR$  de aproximare a valorilor proprii folosind descompunerea Givens păstrează forma Hessenberg.



## *Algoritmul QR pentru valori proprii*

// se aduce matricea  $A$  la forma Hessenberg

- $A = \bar{Q} A \bar{Q}^T$ ;

- $k = 0$ ;

- while (  $A \neq$  forma Schur reală )

- $A = QR$ ; // se calculează cu algoritmul Givens
  - $A = RQ$  sau  $Q^T A Q$ ;
  - $k = k + 1$ ;

În practică se presupune că matricea  $A$  este în **formă Hessenberg neredusă**, adică:

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

Dacă matricea nu este în formă neredusă, problema se decuplează:

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}, \quad p = n-1 \text{ sau } n-2$$

$p \quad n-p$

## Algoritmului $QR$ cu deplasare (“*shift*”) simplă

Algoritmul cu deplasare simplă este următorul:

- $A = \bar{Q} A \bar{Q}^T$ ; // aducerea la forma Hessenberg neredusă
- $k = 0$ ;
- while (  $A \neq$  forma Schur reală )
  - $A - d_k I_n = QR$ ; // se calc. cu alg. Givens
  - $A := RQ + d_k I_n$ ;
  - $k = k + 1$ ;

$d_k \in \mathbb{R}$  sunt constantele de deplasare.

Dacă  $A - d I_n = QR$  ( $A^{(k)}$ ) și  $\bar{A} = RQ + d I_n$  ( $A^{(k+1)}$ ), se pune problema dacă cele două matrice sunt asemenea ( $A \sim \bar{A}$ ) (șirul de matrice construit cu pasul  $QR$  cu deplasare simplă au aceleași valori proprii).

$$\bar{A} = Q^T QRQ + d Q^T Q = Q^T (QR + d I_n) Q = Q^T A Q \Rightarrow \bar{A} \sim A$$

Varianta cu deplasare se efectuează pentru a accelera convergența algoritmului. Dacă  $\lambda_1, \lambda_2, \dots, \lambda_n$  sunt valorile proprii ale matricii  $A$  ordonate astfel ca:

$$|\lambda_1 - d| \geq |\lambda_2 - d| \geq \dots \geq |\lambda_n - d|$$

Rapiditatea cu care  $\mathbf{a}_{p+1p}^{(k)} \rightarrow \mathbf{0}, k \rightarrow \infty$  este dată de rata de convergența a expresiei  $\left| \frac{\lambda_{p+1} - d}{\lambda_p - d} \right|^k$ . Dacă se alege  $d \approx \lambda_n$  convergența  $\mathbf{a}_{n-1n}^{(k)} \rightarrow \mathbf{0}$  este rapidă. Avem următorul rezultat:

### **Teoremă**

Fie  $d$  o valoare proprie a unei matrice Hessenberg nereduse  $H$ . Dacă  $\bar{H} = RQ + d I_n$ , cu  $H - d I_n = QR$  descompunerea  $QR$  a matricei  $H - d I_n = QR$ . Atunci:

$$\bar{h}_{nn-1} = \mathbf{0} \ , \ \bar{h}_{nn} = d$$

Algoritmul **QR** cu deplasare simplă găsește valoarea proprie  $d$  într-un singur pas. Euristic s-a constatat că la fiecare pas, cea mai bună aproximare a unei valori proprii este  $a_{nn}^{(k)}$ .

$$d_k = a_{nn}^{(k)}$$

### *Algoritmul QR cu deplasare simplă*

- $A = \bar{Q} A \bar{Q}^T$ ; // aducerea la forma Hessenberg neredusă
- $k = 0$ ;
- while (  $A \neq$  forma Schur reală )
  - $A - a_{nn} I_n = QR$ ; // se calc. cu algoritmul Givens
  - $A := RQ + a_{nn} I_n$ ;
  - $k = k + 1$ ;

## Algoritmului $QR$ cu deplasare (“*shift*”) dublă

În cazul când valorile proprii  $a_1, a_2$  corespunzătoare blocului:

$$G = \begin{bmatrix} a_{pp} & a_{pn} \\ a_{np} & a_{nn} \end{bmatrix}, \quad p = n - 1$$

sunt complexe,  $a_1, a_2 \in \mathbb{C}$ , abordarea cu deplasare simplă nu mai asigură accelerarea convergenței. Avem:

$$\begin{aligned} \det(\lambda I_2 - G) &= (\lambda - a_1)(\lambda - a_2) = (\lambda - a_{pp})(\lambda - a_{nn}) - a_{pn}a_{np} = \\ &= \lambda^2 - (a_1 + a_2)\lambda + a_1a_2 = \lambda^2 - (a_{pp} + a_{nn})\lambda + a_{pp}a_{nn} - a_{pn}a_{np} \end{aligned}$$

$$a_1 + a_2 = a_{pp} + a_{nn} = \text{trace}(G) \quad , \quad a_1a_2 = a_{pp}a_{nn} - a_{pn}a_{np} = \det G$$

Algoritmul **QR** cu deplasare dublă constă în trecerea de la matricea  $A = A^{(k)}$  la matricea  $A_2 = A^{(k+1)}$  realizând doi pași cu deplasare simplă :

$$A \rightarrow A_1 \text{ (deplasare simplă } a_1), A_1 \rightarrow A_2 \text{ (deplasare simplă } a_2)$$



$$A - a_1 I_n = Q_1 R_1$$

$$A_1 = R_1 Q_1 + a_1 I_n$$

$$A_1 - a_2 I_n = Q_2 R_2$$

$$A_2 = R_2 Q_2 + a_2 I_n$$

Fie matricea :

$$\begin{aligned} M &:= (Q_1 Q_2)(R_2 R_1) = Q_1 (Q_2 R_2) R_1 = Q_1 (A_1 - a_2 I_n) R_1 = \\ &= Q_1 (Q_1^T A Q_1 - a_2 I_n) R_1 = Q_1 Q_1^T A Q_1 R_1 - a_2 Q_1 R_1 = \\ &= (A - a_2 I_n) Q_1 R_1 = (A - a_2 I_n) (A - a_1 I_n) \\ M &= (Q_1 Q_2)(R_2 R_1) = (A - a_2 I_n) (A - a_1 I_n) = \\ &= A^2 - (a_1 + a_2)A + a_1 a_2 I_n \end{aligned}$$

Avem următoarele relații de asemănare:

$$A \sim A_1 = Q_1^T A Q_1 \sim A_2 = Q_2^T A_1 Q_2 = Q_2^T Q_1^T A Q_1 Q_2 = (Q_1 Q_2)^T A (Q_1 Q_2)$$

$$A_2 = (Q_1 Q_2)^T A (Q_1 Q_2) = Q^T A Q \quad , \quad Q := Q_1 Q_2$$

Matricea  $Q$  care asigură trecerea de la matricea  $A$  la matricea  $A_2$  este matricea ortogonală din descompunerea  $QR$  a matricii  $M = (A - a_2 I_n)(A - a_1 I_n)$ . Pasul  $QR$  cu deplasare dublă se face urmând etapele:

1) se calculează matricea  $M = A^2 - s A + qI_n$  cu

$$s = a_1 + a_2 = a_{pp} + a_{nn} \quad , \quad q = a_1 a_2 = a_{pp} a_{nn} - a_{pn} a_{np} ;$$

2) se calculează descompunerea  $QR$  a matricii  $M$ ;

3)  $A_2 := Q^T A Q$ .

## Vectori proprii

Considerăm două matrici asemenea  $A$  și  $B$ :

$$A \sim B \Leftrightarrow A = PBP^{-1}, \quad P \text{ matrice nesingulară}$$

Știm că cele două matrici au același polinom caracteristic,  $p_A(\lambda) \equiv p_B(\lambda)$ , deci au aceleași valori proprii. Ne interesează care este legătura între vectorii proprii asociați aceleiași valori proprii. Fie  $u$  vector propriu asociat valorii proprii  $\lambda$  pentru matricea  $A$  și  $w$  vector propriu asociat valorii proprii  $\lambda$  pentru matricea  $B$ . Care este relația între  $u$  și  $w$ ?

$$Au = \lambda u, Bw = \lambda w, A = PBP^{-1} \Rightarrow PBP^{-1}u = \lambda u \Rightarrow$$

$$BP^{-1}u = \lambda P^{-1}u \Rightarrow w = P^{-1}u, u = Pw$$

Dacă se aplică algoritmul **QR** unei matrici simetrice, forma Schur reală la care se ajunge este o matrice diagonală:

$$S = A = \mathbf{diag}[\lambda_1, \lambda_2, \dots, \lambda_n]$$

Legătura dintre matricea simetrică inițială  $A$  și matricea diagonală este de forma:

$$S = A = \mathbf{diag}[\lambda_1, \lambda_2, \dots, \lambda_n] = U^T A U$$

unde  $U$  este o matrice ortogonală, coloanele matricii  $U$  fiind vectori proprii asociați valorilor proprii reale  $\lambda_1, \lambda_2, \dots, \lambda_n$ . Matricea  $U$  se poate calcula astfel:

## *Algoritmul QR pentru matrici simetrice (valori +vectori proprii)*

// se aduce matricea  $A$  la forma Hessenberg

- $A = \bar{Q} A \bar{Q}^T$  ;

- $U = \bar{Q}^T$  ;

- $k = 0$ ;

- while (  $A \neq$  matrice diagonală )

- $A = QR$ ; // se calculează cu algoritmul Givens
  - $A = RQ$  sau  $Q^T A Q$ ;
  - $U = UQ$ ;
  - $k = k + 1$ ;

## Descompunerea după valori singulare (Singular Value Decomposition)

### Teoremă

Fie  $A \in \mathbb{R}^{m \times n}$ . Atunci există o matrice ortogonală pătratică de dimensiune  $m$ ,  $U \in \mathbb{R}^{m \times m}$ , o matrice ortogonală pătratică de dimensiune  $n$ ,  $V \in \mathbb{R}^{n \times n}$  și constantele pozitive:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0, \quad r \leq \min\{m, n\} \text{ a.î.}$$
$$A = U \Sigma V^T, \quad \Sigma \in \mathbb{R}^{m \times n}, \quad \Sigma = \begin{bmatrix} D & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{bmatrix}, \quad (1)$$
$$D \in \mathbb{R}^{r \times r}, \quad D = \text{diag}[\sigma_1, \dots, \sigma_r]$$

Constanta  $r$  este chiar rangul matricii  $A$ ,  $r = \mathbf{rang}(A)$ .

Constantele  $\sigma_1, \sigma_2, \dots, \sigma_r$  poartă numele de *valori singulare* ale matricii  $A$ .

Folosind relația (1) avem:

$$A^T = (U\Sigma V^T)^T = V\Sigma^T U^T,$$

$$AA^T = U\Sigma V^T V\Sigma^T U^T = U\Sigma\Sigma^T U^T = U\Lambda_m U^T,$$

$$\Lambda_m = \Sigma\Sigma^T = \begin{bmatrix} D^2 & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (m-r)} \end{bmatrix} \in \mathbb{R}^{m \times m}$$



$$A^T A = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T = V \Lambda_n V^T ,$$

$$\Lambda_n = \Sigma^T \Sigma = \begin{bmatrix} D^2 & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (n-r)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Ținând cont de ortogonalitatea matricilor  $U$  și  $V$ , putem rescrie relațiile de mai sus astfel:

$$(AA^T)U = U \Lambda_m , \quad \Lambda_m = \text{diag}[\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, \mathbf{0}, \dots, \mathbf{0}] \in \mathbb{R}^{m \times m}$$

$$(A^T A)V = V \Lambda_n , \quad \Lambda_n = \text{diag}[\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, \mathbf{0}, \dots, \mathbf{0}] \in \mathbb{R}^{n \times n}$$

Din aceste relații deducem că  $\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2$  sunt valorile proprii strict pozitive ale matricilor  $AA^T$  și/sau  $A^T A$  iar matricile  $U$  și  $V$  sunt matrici ale căror coloane sunt vectorii proprii asociați (cei ce formează baze ortonormate). Matricile  $AA^T$  și  $A^T A$  sunt matrici simetrice:

$$\left( AA^T \right)^T = \left( A^T \right)^T A^T = AA^T \quad , \quad \left( A^T A \right)^T = A^T \left( A^T \right)^T = A^T A$$

și au toate valorile proprii nenegative:

$$(AA^T)u = \lambda u \Rightarrow (AA^T u, u) = (\lambda u, u) \Rightarrow$$

$$\lambda = \frac{(A^T u, A^T u)}{(u, u)} = \frac{\|A^T u\|_2^2}{\|u\|_2^2} \geq 0$$

Putem folosi descompunerea după valori singulare pentru a defini pseudo-inversa unei matrici oarecare  $A \in \mathbb{R}^{m \times n}$  ( $n \neq m$ ).

$$A = U\Sigma V^T, \quad A^{-1} =_? (U\Sigma V^T)^{-1} = (V^T)^{-1} \Sigma^{-1} U^{-1} = V \Sigma^{-1} U^T$$

Rămâne de definit matricea  $\Sigma^{-1}$ . Urmând acest raționament se definește *pseudoinversa Moore-Penrose* a unei matrici

$A \in \mathbb{R}^{m \times n}$  astfel:

$$A^I = V \Sigma^I U^T, \quad A^I \in \mathbb{R}^{n \times m}, \quad \Sigma^I = \begin{bmatrix} D^{-1} & 0_{r \times (m-r)} \\ \mathbf{0}_{(n-r) \times r} & 0_{(n-r) \times (m-r)} \end{bmatrix} \in \mathbb{R}^{n \times m},$$

$$D^{-1} \in \mathbb{R}^{r \times r}, \quad D^{-1} = \mathbf{diag} \left[ \frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r} \right].$$

Pseudoinversa definită mai sus satisface următoarele proprietăți:

$$\begin{aligned} (A^I)^I &= A, \quad \forall A \in \mathbb{R}^{m \times n} & ; & \quad (A^T)^I = (A^I)^T, \quad \forall A \in \mathbb{R}^{m \times n} \\ AA^I A &= A & , & \quad A^I A A^I = A^I \end{aligned}$$

Există o proprietate care nu mai este satisfăcută de pseudoinversă deși este respectată de inversa clasică:

$$\exists A, B \text{ a.î. } (AB)^I \neq B^I A^I.$$

Descompunerea după valori singulare poate fi utilizată și pentru rezolvarea sistemelor liniare cu matrici oarecare ( $m \neq n$ )

$$Ax = b \quad , \quad A \in \mathbb{R}^{m \times n} \quad , \quad b \in \mathbb{R}^m \quad , \quad x := A^I b \in \mathbb{R}^n.$$

## Problema celor mai mici pătrate

$$A \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m, \quad Ax = b, \quad x \in \mathbb{R}^n \quad (1)$$

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

$$a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2$$

⋮

$$a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n = b_i$$

⋮

$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m$$

Sistemul are soluții clasice dacă:

$$\mathbf{rang} A = \mathbf{rang} [A / b]$$

- $m < n$  - o infinitate de soluții
- $m \geq n$ 
  - dacă  $\mathbf{rang} A = \mathbf{rang} [A / b]$  soluții clasice
  - dacă  $\mathbf{rang} A \neq \mathbf{rang} [A / b]$  soluții în sensul celor mai mici pătrate

Vectorul reziduu:

$$r(x) = b - Ax \in \mathbb{R}^m$$

Vectorul  $x \in \mathbb{R}^n$  se numește *soluție în sensul celor mai mici pătrate* pentru sistemul (1) dacă este soluția următoarei probleme de optimizare:

$$\min\{\|r(x)\|_2^2 = \|b - Ax\|_2^2; x \in \mathbb{R}^n\} \quad (LSP)$$

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \in \mathbb{R}^{3 \times 2}, \quad b = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \quad m = 3, n = 2$$

$$\text{rang } A = 2 \neq \text{rang } [A / b] = 3$$



Sistemul:

$$\begin{aligned}x_1 + 4x_2 &= 0 \\2x_1 + 5x_2 &= 0 \\3x_1 + 6x_2 &= 1\end{aligned}\tag{2}$$

nu are soluție clasică (nu există  $x_1, x_2$  care să satisfacă toate cele 3 ecuații simultan). Vectorul reziduu are forma:

$$r(x) = b - Ax = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -x_1 - 4x_2 \\ -2x_1 - 5x_2 \\ 1 - 3x_1 - 6x_2 \end{pmatrix}$$

Soluția în sensul celor mai mici pătrate a acestui sistem este definită ca soluția problemei de optimizare:

$$\min\{(-x_1 - 4x_2)^2 + (-2x_1 - 5x_2)^2 + (1 - 3x_1 - 6x_2)^2; x_1, x_2 \in \mathbb{R}\}$$

$$\min\{1 - 6x_1 - 12x_2 + 64x_1x_2 + 14x_1^2 + 77x_2^2; x_1, x_2 \in \mathbb{R}\}$$

Această problemă de minimizare are soluția:

$$x_1 = \frac{13}{18}, \quad x_2 = -\frac{2}{9}, \quad \|r(x)\|_2^2 = \frac{1}{6}$$

și este soluția în sensul celor mai mici pătrate a sistemului (2).

$$\text{range}(A) = \{ y \in \mathbb{R}^m; y = a_1 A^1 + a_2 A^2 + \dots + a_n A^n, a_i \in \mathbb{R}, i = 1, n \}$$

$$A = [A^1 \ A^2 \ \dots \ A^n], \ A^i \in \mathbb{R}^m \text{ sunt coloanele matricii } A$$

## Teoremă

Fie  $A \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ),  $b \in \mathbb{R}^m$ . Un vector  $x \in \mathbb{R}^n$  minimizează norma euclidiană a vectorului reziduu  $\|r(x)\|_2 = \|b - Ax\|_2$ , rezolvând problema (*LSP*), dacă și numai dacă:

$$r(x) \perp \text{range}(A) \quad \Leftrightarrow \quad A^T r(x) = \mathbf{0}$$

sau echivalent

$$A^T Ax = A^T b \quad (3)$$

Sistemul (3) poartă numele de sistemul de *ecuații normale*.

Este un sistem pătratic de dimensiune  $n$ , matricea sistemului  $A^T A \in \mathbb{R}^{n \times n}$  este simetrică. Sistemul de ecuații normale (3) este nesingular dacă și numai dacă  $\text{rang} A = n$ , în acest caz soluția  $\mathbf{x}$  a sistemului (3) este unică.

$$\det A^T A \neq 0 \Leftrightarrow \text{rang} A = n$$

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}, \quad A^T A = \begin{bmatrix} 14 & 32 \\ 32 & 77 \end{bmatrix}, \quad A^T \mathbf{b} = \begin{pmatrix} 3 \\ 6 \end{pmatrix}$$

$$\begin{aligned} 14x_1 + 32x_2 &= 3 \\ 32x_1 + 77x_2 &= 6 \end{aligned} \Rightarrow x_1 = \frac{13}{18}, \quad x_2 = -\frac{2}{9}$$

## Pseudo-inversa matricii $A$

Presupunem că  $A$  are  $\text{rang } A = n$ . Atunci pseudo-inversa poate fi definită ca:

$$A^+ = (A^T A)^{-1} A^T \in \mathbb{R}^{n \times m} \quad (A^+ = A^I \quad ?)$$

$$A^+ = \begin{bmatrix} 14 & 32 \\ 32 & 77 \end{bmatrix}^{-1} * \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

## Rezolvarea sistemului de ecuații normale

1) Folosind factorizarea Cholesky (descompunere  $LU$ ) pentru matrici simetrice:

$$A^T A = LL^T, \quad L \in \mathbb{R}^{n \times n} \text{ matrice inferior triunghiulară}$$

- Se calculează matricea  $A^T A$  și vectorul  $A^T b$ ;
- Se calculează factorizarea Cholesky a matricii  $A^T A = LL^T$ ;
- Se rezolvă sistemul inferior triunghiular  $Ly = A^T b$  pentru  $y$ ;
- Se rezolvă sistemul superior triunghiular  $L^T x = y$  pentru  $x$ ;

2) Se calculează descompunerea  $QR$  (cu algoritmul lui

Householder adaptat) pentru matricea  $A$ :

$$A = QR, \quad Q \in \mathbb{R}^{m \times m} \text{ matrice ortogonală, } R \in \mathbb{R}^{m \times n},$$

$$R = \begin{bmatrix} \bar{R} \in \mathbb{R}^{n \times n} \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix}, \quad \bar{R} - \text{matrice superior triunghiulară}$$

- Se calculează factorizarea  $QR$  modificată a matricii  $A$ ;
- Se calculează vectorul  $Q^T b$ ;
- Se rezolvă sistemul sup. triunghiular  $\bar{R}x = (Q^T b)_{i=1,n}$ ;



3) Se folosește desc. după valori singulare a matricii  $A$

$$A = U\Sigma V^T, \quad \Sigma \in \mathbb{R}^{m \times n}, \quad U \in \mathbb{R}^{m \times m}, \quad V \in \mathbb{R}^{n \times n}$$

- Se calculează SVD pentru matricea  $A=U\Sigma V^T$ ;
- Se calculează vectorul  $U^T b$ ;
- Se rezolvă sistemul diagonal  $\Sigma w = U^T b$  pentru  $w$ ;
- Soluția este  $x=Vw$ ;

1), 2) sau 3) ?  $\rightarrow$  se recomandă 2)

# **Calcul Numeric**

**Cursul 9**

**2022**

*Anca Ignat*

## Interpolare numerică

Presupunem că despre o funcție  $f$  cunoaștem doar valorile într-un număr finit de puncte. Pornind de la aceste date, dorim să aproximăm funcția  $f$  într-un alt punct.

$x$	$x_0$	$x_1$	$x_2$	$\dots$	$x_{n-1}$	$x_n$
$f$	$y_0$	$y_1$	$y_2$	$\dots$	$y_{n-1}$	$y_n$

În tabelul de mai sus  $f(x_i) = y_i$ ,  $i=0,1,\dots,n$  și  $x_i \neq x_j$ ,  $i \neq j$ .

Dat un punct  $x \neq x_i$ ,  $i=0,1,\dots,n$  dorim să aproximăm  $f(x)$  cunoscând cele  $(n+1)$  perechi  $(x_i, y_i)$ ,  $i=0,\dots,n$ . Punctele  $x_i$  se numesc *noduri de interpolare*.

## Polinomul de interpolare Lagrange

Notăm cu  $\Pi_n$  mulțimea polinoamelor de grad cel mult  $n$ . Dimensiunea acestui spațiu este  $n+1$ , baza uzuală fiind dată de polinoamele  $1, x, x^2, \dots, x^n$ . Vom considera o altă bază în acest spațiu. Se consideră polinoamele  $p_i$ :

$$p_i \in \Pi_n \text{ astfel ca } p_i(x_j) = \begin{cases} 0 & \text{pentru } j \neq i \\ 1 & \text{pentru } j = i \end{cases}, j = 0, \dots, n, i = 0, \dots, n$$

Din relația  $p_i(x_j) = 0, \forall j \neq i$  și faptul că  $p_i$  este de grad  $n$  rezultă că  $x_0, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$  sunt cele  $n$  rădăcini ale polinomului  $p_i$ .

Avem:

$$p_i(x) = c_i(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n),$$
$$c_i \in \mathbb{R}, i = 0, \dots, n$$

Constanta  $c_i$  se determină din relația  $p_i(x_i) = 1$ :

$$p_i(x_i) = 1 = c_i(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n) \Rightarrow$$
$$c_i = \frac{1}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}$$

Polinoamele  $p_i$  au forma:

$$p_i(x) = \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)} = \prod_{\substack{j=0 \\ j \neq i}}^n \left( \frac{x - x_j}{x_i - x_j} \right)$$
$$i = 0, \dots, n$$

### Propoziție

Polinoamele  $p_0, p_1, \dots, p_n$  formează o bază în  $\Pi_n$ .

Demonstrație: Vom arăta că cele  $n+1$  polinoame sunt liniar independente:

$$q(x) = a_0 p_0(x) + a_1 p_1(x) + \cdots + a_n p_n(x) = \mathbf{0}, \forall x$$

$$\Rightarrow a_0 = a_1 = \cdots = a_n = \mathbf{0}$$

Vom face pe rând  $x = x_0, x = x_1, \dots, x = x_n$  în polinomul  $q$ :

$$\begin{aligned} x = x_0 \quad q(x_0) &= a_0 p_0(x_0) + a_1 p_1(x_0) + \cdots + a_n p_n(x_0) = \\ &= a_0 \mathbf{1} + a_1 \mathbf{0} + \cdots + a_n \mathbf{0} = a_0 = \mathbf{0} \Rightarrow a_0 = \mathbf{0} \end{aligned}$$

$$x = x_1 \quad q(x_1) = 0 \Rightarrow a_1 = \mathbf{0}$$

$\vdots$

$$\begin{aligned} x = x_k \quad q(x_k) &= a_0 p_0(x_k) + \cdots + a_k p_k(x_k) + \cdots + a_n p_n(x_k) = \\ &= a_0 \mathbf{0} + \cdots + a_k \mathbf{1} + \cdots + a_n \mathbf{0} = a_k = \mathbf{0} \Rightarrow a_k = \mathbf{0} \end{aligned}$$

$\vdots$

$$x = x_n \quad q(x_n) = 0 \Rightarrow a_n = \mathbf{0}$$

Toate constantele  $a_i$  sunt nule deci polinoamele  $\{p_i; i = 0, \dots, n\}$  formează o bază în  $\Pi_n$ .

Pentru a aproxima funcția  $f$  pornind de la tabelul de mai sus, vom construi un polinom  $l_n \in \Pi_n$  a.î. să satisfacă **condițiile de interpolare**:

$$l_n \in \Pi_n, \quad l_n(x_i) = y_i, \quad \forall i = 0, \dots, n \quad (1)$$

Odată construit acest polinom, vom aproxima  $f(x)$  prin  $l_n(x)$ ,  $f(x) \approx l_n(x)$

Vom scrie polinomul  $l_n$  în raport cu noua bază  $\{p_i; i = 0, \dots, n\}$ , deci există constantele reale  $a_0, a_1, \dots, a_n$  astfel ca:



$$l_n(x) = \sum_{i=0}^n a_i p_i(x)$$

Constantele  $a_k$  se determină astfel:

$$\begin{aligned} y_k = l_n(x_k) &= a_0 p_0(x_k) + \cdots + a_k p_k(x_k) + \cdots + a_n p_n(x_k) = \\ &= a_0 \mathbf{0} + \cdots + a_k \mathbf{1} + \cdots + a_n \mathbf{0} = a_k \Rightarrow a_k = y_k \end{aligned}$$

Prin urmare un polinom de grad  $n$  care îndeplinesc condițiile de interpolare (1) este:

$$\begin{aligned} l_n(x) &= \sum_{i=0}^n y_i p_i(x) = \sum_{i=0}^n y_i \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)} \\ &= \sum_{i=0}^n y_i \prod_{\substack{j=0 \\ j \neq i}}^n \left( \frac{x - x_j}{x_i - x_j} \right) \end{aligned} \quad (2)$$

Polinomul din formula (2) se numește *polinomul de interpolare Lagrange*.

### **Propoziție**

Polinomul  $l_n$  dat de formula (2) este unicul polinom de grad  $n$  care îndeplinește condițiile de interpolare (1).

Demonstrație: Presupunem că mai există un polinom  $q \in \Pi_n$  care îndeplinește condițiile (1):

$$q \in \Pi_n \quad , \quad q(x_i) = y_i \quad , \quad \forall i = 0, \dots, n$$

Fie polinomul  $p(x) = l_n(x) - q(x) \in \Pi_n$ .

$$p(x_k) = l_n(x_k) - q(x_k) = y_k - y_k = 0, \forall k = 0, \dots, n$$

Polinomul  $p$  are ca rădăcini toate nodurile de interpolare. Polinomul  $p$  este polinom de grad cel mult  $n$  și are  $(n+1)$  rădăcini distincte ( $x_i \neq x_j, \forall i \neq j$ ). Acest polinom nu poate fi decât polinomul identic nul:

$$p(x) = l_n(x) - q(x) \equiv 0 \quad \forall x, \quad l_n(x) = q(x) \quad \forall x$$

Polinomul  $l_n$  este unicul care satisface (2).

Fie  $w_{n+1}$  polinomul de grad  $(n+1)$  care are ca rădăcini nodurile de interpolare:

$$w_{n+1}(x) = (x - x_0)(x - x_1) \cdots (x - x_n) \in \Pi_{n+1}$$

Fie  $a = \min\{x_0, x_1, \dots, x_n\}$ ,  $b = \max\{x_0, x_1, \dots, x_n\}$ .

## Teorema restului (eroarea la interpolarea Lagrange)

Fie  $f \in C^{n+1}[a, b]$  și  $\bar{x} \in [a, b]$ ,  $\bar{x} \neq x_i$ ,  $\forall i = 0, \dots, n$ .

Atunci există un punct  $y \in [a, b]$ ,  $y = y(x_0, x_1, \dots, x_n, \bar{x})$  (punctul  $y$  depinde de nodurile de interpolare  $x_i$  și de punctul  $\bar{x}$ ) astfel că eroarea la interpolarea numerică este dată de:

$$f(\bar{x}) - l_n(\bar{x}) = \frac{f^{(n+1)}(y)}{(n+1)!} w_{n+1}(\bar{x}) \quad (3)$$

Demonstrație: Considerăm funcția  $F$ :

$$F(x) := f(x) - l_n(x) - cw_{n+1}(x)$$

Constanta reală  $c$  este aleasă astfel ca  $F(\bar{x}) = 0$  adică:

$$c = \frac{f(\bar{x}) - l_n(\bar{x})}{w_{n+1}(\bar{x})}, \quad (x \neq x_i \quad \forall i) \Rightarrow w_{n+1}(\bar{x}) \neq 0 \quad (4)$$

Funcția  $f$  fiind de clasă  $C^{n+1}$  pe intervalul  $[a,b]$  rezultă că și funcția  $F$  este din  $C^{n+1}[a,b]$ . Avem:

$$F(x_i) = f(x_i) - l_n(x_i) - cw_{n+1}(x_i) = y_i - y_i - c \cdot 0 = 0, \quad \forall i = 0, \dots, n$$

Funcția  $F$  are  $(n+2)$  zerouri,  $x_0, x_1, \dots, x_n, \bar{x}$ . Aplicând succesiv Teorema lui Rolle rezultă că  $F'$  are  $(n+1)$  zerouri,  $F''$  are  $n$  zerouri, ...,  $F^{(n+1)}$  are  $1$  zero în intervalul  $[a,b]$ . Vom nota această rădăcină a lui  $F^{(n+1)}$  cu  $y$ . Punctul  $y$  depinde de zerourile inițiale  $x_0, x_1, \dots, x_n, \bar{x}$  și:

$$y = y(x_0, x_1, \dots, x_n, \bar{x}) \in [a, b] \quad \text{a.î.} \quad F^{(n+1)}(y) = \mathbf{0}. \quad (5)$$

Derivata de ordinul  $(n+1)$  a funcției  $F$  se calculează astfel:

$$\begin{aligned} F^{(n+1)}(x) &= f^{(n+1)}(x) - l_n^{(n+1)}(x) - c w_{n+1}^{(n+1)}(x) = \\ &= f^{(n+1)}(x) - \mathbf{0} - c(n+1)! = f^{(n+1)}(x) - c(n+1)! \end{aligned} \quad (6)$$

(derivata de ordin  $(n+1)$  a polinomului de grad  $n$   $l_n$  este  $\mathbf{0}$ ).

Din relațiile (4), (5) și (6) rezultă că:

$$c = \frac{f^{(n+1)}(y)}{(n+1)!} = \frac{f(\bar{x}) - l_n(\bar{x})}{w_{n+1}(\bar{x})} \Rightarrow f(\bar{x}) - l_n(\bar{x}) = \frac{f^{(n+1)}(y)}{(n+1)!} w_{n+1}(\bar{x})$$

## Forma Newton a polinomului de interpolare Lagrange

Fie  $l_k(x, x_0, x_1, \dots, x_k, f)$  polinomul de interpolare Lagrange pentru funcția  $f$  pe sistemul de noduri distincte  $\{x_0, x_1, \dots, x_k\}$ .

### Propoziție

Fie  $l_{k-1}(x, x_0, x_1, \dots, x_{k-1}, f)$ ,  $l_{k-1}(x, x_1, x_2, \dots, x_k, f) \in \Pi_{k-1}$  polinoamele de interpolare Lagrange pentru funcția  $f$  pe sistemele de noduri  $\{x_0, x_1, \dots, x_{k-1}\}$  și respectiv  $\{x_1, x_2, \dots, x_k\}$ .

Atunci:

$$\begin{aligned} l_k(x, x_0, x_1, \dots, x_k, f) &= \\ &= \frac{(x - x_0)l_{k-1}(x, x_1, x_2, \dots, x_k, f) - (x - x_k)l_{k-1}(x, x_0, x_1, \dots, x_{k-1}, f)}{x_k - x_0} \end{aligned} \tag{1}$$



Demonstrație: Exercițiu.

Considerăm următoarele probleme de interpolare pentru  $f$ :

$$\{(x_0, y_0), (x_1, y_1), \dots, (x_{k-1}, y_{k-1})\} \rightarrow l_{k-1}(x, x_0, x_1, \dots, x_{k-1}, f)$$

$$\{(x_0, y_0), (x_1, y_1), \dots, (x_k, y_k)\} \rightarrow l_k(x, x_0, x_1, \dots, x_k, f)$$

Ne interesează să găsim o formulă de trecere rapidă de la polinomul de interpolare pe  $k$  noduri la cel care are un nod în plus. Deoarece polinomul de grad cel mult  $k$ :

$$q(x) = l_k(x, x_0, x_1, \dots, x_k, f) - l_{k-1}(x, x_0, x_1, \dots, x_{k-1}, f) \in \Pi_k$$

are ca rădăcini punctele  $x_0, x_1, \dots, x_{k-1}$  ( $q(x_i) = y_i - y_i = 0$ ,  $i=0, \dots, k-1$ ) avem relația:

$$l_k(x, x_0, x_1, \dots, x_k, f) = l_{k-1}(x, x_0, x_1, \dots, x_{k-1}, f) + A \prod_{j=0}^{k-1} (x - x_j) \quad (2)$$

în care  $A$  este dat de relația:

$$A = \frac{l_k(x_k, x_0, x_1, \dots, x_k, f) - l_{k-1}(x_k, x_0, x_1, \dots, x_{k-1}, f)}{\prod_{j=0}^{k-1} (x_k - x_j)} \quad (3)$$

$$\begin{aligned}
A &= \frac{y_k}{\prod_{j=0}^{k-1} (x_k - x_j)} - \frac{\sum_{i=0}^{k-1} y_i \prod_{\substack{j=0 \\ j \neq i}}^{k-1} \left( \frac{x_k - x_j}{x_i - x_j} \right)}{\prod_{j=0}^{k-1} (x_k - x_j)} = \\
&= \frac{y_k}{\prod_{j=0}^{k-1} (x_k - x_j)} - \sum_{i=0}^{k-1} \frac{y_i}{(x_k - x_i) \prod_{\substack{j=0 \\ j \neq i}}^{k-1} (x_i - x_j)}
\end{aligned}$$

$$A = \sum_{i=0}^k \frac{y_i}{\prod_{\substack{j=0 \\ j \neq i}}^k (x_i - x_j)} \quad (4)$$

$$[x_k]_f = y_k = f(x_k) ,$$

$$[x_0, x_1, \dots, x_k]_f = \frac{[x_1, x_2, \dots, x_k]_f - [x_0, x_1, \dots, x_{k-1}]_f}{x_k - x_0}$$

numită *diferență divizată de ordin  $k$  a funcției  $f$  pe nodurile  $\{x_0, x_1, \dots, x_k\}$ .*

## Propoziție

$$\left[ x_0, x_1, \dots, x_k \right]_f = \sum_{i=0}^k \frac{y_i}{\prod_{\substack{j=0 \\ j \neq i}}^k (x_i - x_j)} = \sum_{i=0}^k \frac{y_i}{(w_{n+1}(x_k))'}$$

pentru orice sistem de noduri  $\{x_0, x_1, \dots, x_k\}$  și orice  $k$ .

Demonstrație: Se face prin inducție. Pentru  $k=1$  avem:

$$\left[ x_0, x_1 \right]_f = \frac{y_0}{x_0 - x_1} + \frac{y_1}{x_1 - x_0} = \frac{\left[ x_1 \right]_f - \left[ x_0 \right]_f}{x_1 - x_0}$$

Presupunem că relația (8) este valabilă pentru orice  $k$  și pentru orice sistem de noduri  $\{x_0, x_1, \dots, x_k\}$ . Pentru  $k+1$  folosim relația de recurență și apoi aplicăm ipoteza inductivă:

$$\begin{aligned}
[x_0, x_1, \dots, x_{k+1}]_f &= \frac{[x_1, x_1, \dots, x_{k+1}]_f - [x_0, x_2, \dots, x_k]_f}{x_{k+1} - x_0} = \\
&= \frac{1}{x_{k+1} - x_0} \left( \sum_{i=1}^{k+1} \frac{y_i}{\prod_{\substack{j=1 \\ j \neq i}}^{k+1} (x_i - x_j)} - \sum_{i=0}^k \frac{y_i}{\prod_{\substack{j=0 \\ j \neq i}}^k (x_i - x_j)} \right) = \\
&= \frac{1}{x_{k+1} - x_0} \left\{ -\frac{y_0}{\prod_{\substack{j=0 \\ j \neq 0}}^k (x_0 - x_j)} + \frac{y_{k+1}}{\prod_{\substack{j=1 \\ j \neq k+1}}^{k+1} (x_{k+1} - x_j)} + \right. \\
&\quad \left. + \sum_{i=1}^k \left[ \frac{y_i}{\prod_{\substack{j=1 \\ j \neq i}}^k (x_i - x_j)} \left( \frac{1}{x_i - x_{k+1}} - \frac{1}{x_i - x_0} \right) \right] \right\} =
\end{aligned}$$

$$\begin{aligned}
&= \frac{y_0}{\prod_{\substack{j=0 \\ j \neq 0}}^{k+1} (x_0 - x_j)} + \frac{y_{k+1}}{\prod_{\substack{j=0 \\ j \neq k+1}}^{k+1} (x_{k+1} - x_j)} + \sum_{i=1}^k \frac{y_i}{\prod_{\substack{j=0 \\ j \neq i}}^{k+1} (x_i - x_j)} = \\
&= \sum_{i=0}^{k+1} \frac{y_i}{\prod_{\substack{j=0 \\ j \neq i}}^{k+1} (x_i - x_j)}
\end{aligned}$$

Inducția este completă.

Din definiție se observă că diferența divizată  $[x_0, x_1, \dots, x_k]_f$  nu depinde de ordinea nodurilor  $\{x_0, x_1, \dots, x_k\}$ .

Vom nota în continuare cu  $l_k(x)$  polinomul de interpolare Lagrange pe nodurile  $\{x_0, x_1, \dots, x_k\}$  pentru funcția  $f$ . Avem:

$$\begin{aligned} l_n(x) &= l_0(x) + [l_1(x) - l_0(x)] + \dots + [l_k(x) - l_{k-1}(x)] + \dots + [l_n(x) - l_{n-1}(x)] = \\ &= y_0 + [x_0, x_1]_f (x - x_0) + \dots + [x_0, x_1, \dots, x_k]_f (x - x_0) \cdots (x - x_{k-1}) + \dots \\ &\quad + [x_0, x_1, \dots, x_n]_f (x - x_0) \cdots (x - x_{n-1}) \end{aligned}$$



Am obținut *forma Newton* a polinomului de interpolare Lagrange:

$$l_n(x) = y_0 + [x_0, x_1]_f (x - x_0) + [x_0, x_1, x_2]_f (x - x_0)(x - x_1) + \cdots + [x_0, x_1, \dots, x_n]_f (x - x_0) \cdots (x - x_{n-1})$$

## Schema lui Aitken de calcul a diferențelor divizate

Ne propunem să calculăm diferențele divizate

$$[x_0, x_1]_f, [x_0, x_1, x_2]_f, \dots, [x_0, x_1, \dots, x_n]_f$$

necesare construirii polinomului de interpolare Lagrange în forma Newton. Procedeu folosește definiția recursivă a diferențelor divizate și se desfășoară în  $n$  pași. La pasul  $l$  se calculează numai diferențe divizate de ordinul  $l$ :

$$[x_0, x_1]_f, [x_1, x_2]_f, \dots, [x_{n-1}, x_n]_f.$$

În general, la pasul  $k$  se calculează diferențe divizate de ordin  $k$ :

$$\left[ x_0, x_1, \dots, x_k \right]_f, \left[ x_1, x_2, \dots, x_{k+1} \right]_f, \dots, \left[ x_{n-k}, x_{n-k+1}, \dots, x_n \right]_f.$$

La pasul  $n$  se calculează o singură diferență divizată de ordin  $n$  și anume  $\left[ x_0, x_1, \dots, x_n \right]_f$ .

	<b>Pas 1</b>	$\dots$	<b>Pas <math>k</math></b>	$\dots$	<b>Pas <math>n</math></b>
$x_0$	$y_0$				
$x_1$	$y_1$		$[x_0, x_1]_f$		
$x_2$	$y_2$		$[x_1, x_2]_f$		
$\vdots$					
$x_k$	$y_k$		$[x_{k-1}, x_k]_f$		$[x_0, x_1, \dots, x_k]_f$
$\vdots$				$\vdots$	$\ddots$
$x_{n-1}$	$y_{n-1}$		$[x_{n-2}, x_{n-1}]_f$		$[x_{n-k-1}, \dots, x_{n-1}]_f$
$x_n$	$y_n$		$[x_{n-1}, x_n]_f$		$[x_{n-k}, \dots, x_{n-1}]_f$
				$\dots$	$[x_0, x_1, \dots, x_n]_f$

Notăm  $dd[i,k] = [x_i, x_{i+1}, \dots, x_{i+k}]_f$  diferența divizată de ordin  $k$ , pe nodurile consecutive  $\{x_i, x_{i+1}, \dots, x_{i+k}\}$   $i=0, \dots, n-k$ ,  $k=1, \dots, n$ , cu  $dd[i,0]=y_i$ ,  $i=0, \dots, n$ . Schema lui Aitken se implementează astfel:

$$dd[i,0] = y_i, \quad i = 0, \dots, n;$$

for  $k = 1, \dots, n$

for  $i = 0, \dots, n - k$

$$dd[i,k] = \frac{dd[i+1,k-1] - dd[i,k-1]}{x_{i+k} - x_i}$$

Putem face aceleași calcule folosind un singur vector, de exemplu rescriind vectorul  $y$  astfel:

**for  $k = 1, \dots, n$**

**for  $i = n, \dots, k$**

$$y_i = \frac{y_i - y_{i-1}}{x_i - x_{i-k}}$$

La finalul acestei secvențe de program, vectorul  $y$  va conține elementele:

$$y_0, [x_0, x_1]_f, [x_0, x_1, x_2]_f, \dots, [x_0, x_1, \dots, x_n]_f$$

$(y_k = [x_0, x_1, \dots, x_k]_f, k=0, \dots, n).$

## Interpolare Newton pe noduri echidistante

Pp. că nodurile de interpolare sunt echidistante:

$$x_i = x_0 + i h \quad , \quad i = 0, 1, \dots, n$$

În relația de mai sus fie se dă  $h$  distanța între 2 noduri succesive, fie se precizează primul și ultimul nod,  $x_0$  și  $x_n$  iar  $h$  se calculează:

$$h = \frac{(x_n - x_0)}{n}$$
$$\left[ x_i, x_{i+1} \right]_f = \frac{f(x_{i+1}) - f(x_i)}{(x_{i+1} - x_i)} = \frac{y_{i+1} - y_i}{h}$$

Se introduce noțiunea de **diferență finită de ordinul 1**:

$$\Delta f(x) = f(x+h) - f(x)$$

Pornind de la această definiție se pot introduce și diferențe finite de ordin superior:

$$\begin{aligned}\Delta^2 f(x) &= \Delta(\Delta f(x)) = \Delta(f(x+h) - f(x)) = \\ &= f(x+2h) - 2f(x+h) + f(x)\end{aligned}$$

și în general se pot introduce recursiv **diferențele finite de ordin  $k$** :

$$\Delta^k f(x) = \Delta(\Delta^{k-1} f(x)) = \Delta^{k-1} f(x+h) - \Delta^{k-1} f(x).$$

Prin inducție după  $k$ , se poate deduce formula de calcul a diferențelor finite de ordin  $k$  folosind doar valorile funcției  $f$ :



$$\Delta^k f(x) = \sum_{i=0}^k (-1)^{k-i} C_k^i f(x + ih)$$

**Observație:** Dacă funcția  $f$  este polinom de grad  $m$  atunci  $\Delta f(x)$  este polinom de grad  $m-1$ ,  $\Delta^2 f(x)$  este polinom de grad  $m-2$ , ș.a.m.d. Prin urmare:

$\Delta^k f(x) \equiv 0$ , pentru  $k > m$ ,  $f$  – polinom de grad  $m$ .

Legătura între diferențele divizate și cele finite:

$$[x_i, x_{i+1}]_f = \frac{f(x_{i+1}) - f(x_i)}{(x_{i+1} - x_i)} = \frac{\Delta f(x_i)}{h}$$

$$[x_i, x_{i+1}, x_{i+2}]_f = \frac{[x_{i+1}, x_{i+2}]_f - [x_i, x_{i+1}]_f}{(x_{i+2} - x_i)} = \frac{\Delta^2 f(x_i)}{2h^2}$$

Prin inducție se poate arăta următoarea legătură între diferențele divizate de ordin  $k$  și cele finite:

$$\left[ x_i, x_{i+1}, \dots, x_{i+k} \right]_f = \frac{\Delta^k f(x_i)}{k! h^k}.$$

### Polinoame de interpolare pe noduri echidistante

$$\begin{aligned} l_n(x) = & y_0 + \left[ x_0, x_1 \right]_f (x - x_0) + \left[ x_0, x_1, x_2 \right]_f (x - x_0)(x - x_1) + \dots + \\ & + \left[ x_0, x_1, \dots, x_k \right]_f (x - x_0)(x - x_1) \dots (x - x_{k-1}) + \dots + \\ & + \left[ x_0, x_1, \dots, x_n \right]_f (x - x_0)(x - x_1) \dots (x - x_{n-1}) \end{aligned}$$

Consideră că punctul de interpolare este de forma:

$$\bar{x} = x_0 + t h$$

și înlocuim diferențele divizate cu diferențe finite în forma Newton a polinomului de interpolare:

$$\begin{aligned}(\bar{x} - x_0) \cdots (\bar{x} - x_{k-1}) &= (x_0 + t h - x_0) \cdots (x_0 + t h - x_0 - (k-1)h) = \\ &= h^k t(t-1) \cdots (t-k+1)\end{aligned}$$

$$\begin{aligned}
l_n(\bar{x}) = l_n(x_0 + th) = & y_0 + \Delta f(x_0)t + \Delta^2 f(x_0)\frac{t(t-1)}{2} + \dots + \\
& + \Delta^k f(x_0)\frac{t(t-1)\dots(t-k+1)}{k!} + \dots + \\
& + \Delta^n f(x_0)\frac{t(t-1)\dots(t-n+1)}{n!}
\end{aligned}$$

Această relație poartă numele de **formula lui Newton progresivă pe noduri echidistante**.

Considerăm polinomul de interpolare Lagrange pe nodurile în ordine inversă  $\{x_n, x_{n-1}, \dots, x_0\}$ :

$$l_n(x) = y_n + [x_n, x_{n-1}]_f (x - x_n) + [x_n, x_{n-1}, x_{n-2}]_f (x - x_n)(x - x_{n-1}) + \dots + [x_n, x_{n-1}, \dots, x_0]_f (x - x_n)(x - x_{n-1}) \dots (x - x_0)$$

Dacă punctul de interpolare este de forma:

$$\bar{x} = x_n + th$$

analog ca mai sus obține **formula lui Newton regresivă pe noduri echidistante:**

$$\begin{aligned}
l_n(\bar{x}) = l_n(x_n + th) = & y_n + \Delta f(x_{n-1})t + \Delta^2 f(x_{n-2})\frac{t(t+1)}{2} + \dots + \\
& + \Delta^k f(x_{n-k})\frac{t(t+1)\cdots(t+k-1)}{k!} + \dots + \\
& + \Delta^n f(x_0)\frac{t(t+1)\cdots(t+n-1)}{n!}
\end{aligned}$$

## Schema lui Aitken pentru diferențe finite

Pentru formulele lui Newton progresive/regresive, avem nevoie de calculul următoarelor diferențe finite:

$\Delta f(x_0), \Delta^2 f(x_0), \dots, \Delta^k f(x_0), \dots, \Delta^n f(x_0)$  - pentru formula progresivă

$\Delta f(x_{n-1}), \Delta^2 f(x_{n-2}), \dots, \Delta^k f(x_{n-k}), \dots, \Delta^n f(x_0)$  - pentru formula regresivă

Metoda de calcul a acestor diferențe finite este similară schemei Aitken pentru diferențele divizate.

	<b>Pas 1</b>	...	<b>Pas <math>k</math></b>	...	<b>Pas <math>n</math></b>
$y_0$					
$y_1$	$\Delta f(x_0)$				
$y_2$	$\Delta f(x_1)$				
$\vdots$					
$y_k$	$\Delta f(x_{k-1})$		$\Delta^k f(x_0)$		
$\vdots$			$\vdots$	$\ddots$	
$y_{n-1}$	$\Delta f(x_{n-2})$		$\Delta^k f(x_{n-k-1})$		
$y_n$	$\Delta f(x_{n-1})$		$\Delta^k f(x_{n-k})$	...	$\Delta^n f(x_0)$



## Funcții spline

Fie nodurile:

$$x_i \in [a, b], i = 0, 1, \dots, n,$$

cu

$$a = x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n = b$$

Se consideră funcția continuă polinomială pe porțiuni:

$$S(x) = P_i(x) \text{ pentru } x \in [x_i, x_{i+1}] \quad \forall i = 0, \dots, n-1$$

$$S(x) = \begin{cases} P_0(x), & x \in [x_0, x_1], \\ P_1(x), & x \in [x_1, x_2], \\ P_2(x), & x \in [x_2, x_3], \\ \vdots & \\ P_{n-2}(x), & x \in [x_{n-2}, x_{n-1}], \\ P_{n-1}(x), & x \in [x_{n-1}, x_n]. \end{cases}$$

$P_i(x)$ ,  $i=0, \dots, n$  sunt polinoame. O asemenea funcție poartă numele de *funcție spline*.

## Funcții spline liniare continue

### Definiție

Funcția  $S(x)$  definită mai sus se numește *funcție spline liniară continuă* dacă polinoamele  $P_i(x), i = 0, \dots, n-1$  sunt polinoame de gradul  $1$  și  $S(x) \in C[a, b]$ , adică:

$$\lim_{\substack{x \rightarrow x_i \\ x < x_i}} S(x) = \lim_{\substack{x \rightarrow x_i \\ x > x_i}} S(x), i = 1, \dots, n-1.$$

Fie funcția  $f : [a, b] \rightarrow \mathbb{R}$  pentru care se cunosc valorile:

$$y_i = f(x_i), i = 0, \dots, n.$$

Funcția spline liniară de interpolare  $S$  pentru funcția  $f$  îndeplinește condițiile de interpolare:

$$S(x_i) = y_i, i = 0, \dots, n.$$

Ținând seamă că polinoamele  $P_i(x)$  sunt polinoame de gradul  $1$  și  $S(x)$  este continuă vom avea condițiile:

$$\begin{cases} P_i(x_i) = y_i, \\ P_i(x_{i+1}) = y_{i+1}, & i = 0, \dots, n-1, \\ P_i(x) - \text{polinom de gradul } 1. \end{cases}$$

Din aceste condiții rezultă:

$$P_i(x) = \frac{x - x_i}{x_{i+1} - x_i} y_{i+1} + \frac{x_{i+1} - x}{x_{i+1} - x_i} y_i, i = 0, \dots, n-1$$

$$S(x_k) = P_{k-1}(x_k) = P_k(x_k) = y_k, k = 1, \dots, n-1,$$

$$S(x_0) = P_0(x_0) = y_0, S(x_n) = P_{n-1}(x_n) = y_n.$$

### Funcții spline cubice de clasă $C^2$

Se consideră sistemul de noduri distincte din intervalul  $[a, b]$ :

$$\Delta = \{a = x_0 < x_1 < \dots < x_{n-1} < x_n = b\}$$

Funcția  $S(x)$  asociată divizării  $\Delta$  care îndeplinește condițiile :

$$S(x) \in C^2[a, b],$$

**polinoamele  $P_i(x)$  au gradul 3,  $i = 0, \dots, n-1$ ,**

se numește *funcție spline cubică*.

Data fiind o funcție  $f : [a, b] \rightarrow \mathbb{R}$  cu valorile:

$$y_i = f(x_i), \quad i = 0, \dots, n,$$

se consideră funcția spline cubică  $S(x)$  de interpolare ce satisface

$$S(x_i) = y_i, \quad i = 0, \dots, n.$$

Pentru determinarea funcției spline cubice de interpolare observăm că polinoamele:

$$P_i(x) = \alpha_i x^3 + \beta_i x^2 + \gamma_i x + \delta_i, \quad x \in [x_i, x_{i+1}], \quad i = 0, \dots, n-1,$$

implică determinarea a celor  $4n$  necunoscute  $\{\alpha_i, \beta_i, \gamma_i, \delta_i; i = 0, \dots, n-1\}$  pentru care se impun:

$$\left\{ \begin{array}{l} n + 1 \text{ condiții din relațiile de interpolare } S(x_i) = y_i, i = 0, \dots, n, \\ 3(n - 1) \text{ condiții de continuitate pentru } S(x), S'(x) \text{ și } S''(x) \\ \quad \text{în nodurile } x_i, i = 1, \dots, n - 1, \end{array} \right.$$

în total  $4n-2$  condiții.

Se pot avea în vedere pentru adăugarea a două condiții suplimentare următoarele abordări :

- fixarea pantelor în extremitățile intervalului  $[a, b]$ . Se presupune că funcția  $f$  este derivabilă și se cunosc valorile  $f'(a), f'(b)$ . Se impun condițiile:

$$S'(x_0) = P'_0(x_0) = f'(a), S'(x_n) = P'_{n-1}(x_n) = f'(b);$$

- periodicitatea primelor două derivate:

$$f'(a) = f'(b) \quad (S'(x_0) = P'_0(x_0) = P'_{n-1}(x_n) = S'(x_n)),$$

$$f''(a) = f''(b) \quad (S''(x_0) = P''_0(x_0) = P''_{n-1}(x_n) = S''(x_n)) ;$$

- anularea derivatei secunde în capetele intervalului:

$$f''(a) = f''(b) = 0$$

$$(S''(x_0) = P''_0(x_0) = 0, S''(x_n) = P''_{n-1}(x_n) = 0).$$

Funcțiile spline care îndeplinesc aceste condiții se numesc *funcții spline cubice normale*.

- derivata de ordinul al treilea a funcției  $S$  este continuă în punctele  $x_1$  și  $x_{n-1}$ . Aceasta înseamnă că polinoamele  $P_0$ ,  $P_1$  respectiv  $P_{n-2}$ ,  $P_{n-1}$  coincid. Acest tip de funcție spline se numește „*not-a-knot*” și este utilizat în MATLAB.



Vom calcula în cele ce urmează funcția spline cubică în cazul în care cunoaștem suplimentar valorile celei de-a doua derivate a funcției  $f$  în capetele intervalului de interpolare:

$$a_0 = f''(a), \quad a_n = f''(b).$$

Recapitulând, vom avea următoarele condiții:

$$\left\{ \begin{array}{l} P_i(x_i) = y_i, i = 0, \dots, n-1, P_{n-1}(x_n) = y_n - \text{interpolare}, \\ P_{i-1}(x_i) = P_i(x_i), i = 1, \dots, n-1, - \text{continuitatea funcției } S, \\ P'_{i-1}(x_i) = P'_i(x_i), i = 1, \dots, n-1, - \text{continuitatea primei derivatei}, \\ P''_{i-1}(x_i) = P''_i(x_i), i = 1, \dots, n-1, - \text{continuitatea derivatei secunde}, \\ P''_0(x_0) = a_0 = f''(a), P''_{n-1}(x_n) = a_n = f''(b). \end{array} \right.$$

Vom nota:

$$S''(x_i) = a_i, i = 0, \dots, n.$$

Ținând seama de faptul că funcția  $S'' \in C[a, b]$  este o funcție liniară pe fiecare din intervalele  $[x_i, x_{i+1}]$  rezultă că:

$$S''(x) = \frac{x - x_i}{h_i} a_{i+1} + \frac{x_{i+1} - x}{h_i} a_i, \quad x \in [x_i, x_{i+1}], \quad \forall i = 0, \dots, n-1$$

$$h_i = x_{i+1} - x_i, \quad i = 0, \dots, n-1$$

iar din

$$S'(x) = \int S''(x) dx, \quad S(x) = \int S'(x) dx$$

rezultă:

$$S(x) = \frac{(x - x_i)^3}{6h_i} a_{i+1} + \frac{(x_{i+1} - x)^3}{6h_i} a_i + b_i x + c_i,$$

$$x \in [x_i, x_{i+1}], b_i, c_i \in \mathbb{R}, i = 0, \dots, n-1,$$

$$P_i(x) = \frac{(x - x_i)^3}{6h_i} a_{i+1} + \frac{(x_{i+1} - x)^3}{6h_i} a_i + b_i x + c_i ,$$

$$b_i, c_i \in \mathbb{R} , i = 0, \dots, n-1,$$

Vom calcula funcția spline pentru cazul:

$$S''(a) = a_0 = P_0''(x_0),$$

$$S''(b) = a_n = P_{n-1}''(x_n).$$

( $a_0$  și  $a_n$  sunt două constante cunoscute)

Impunând condițiile de interpolare și de continuitate vom obține:

$$P_i(x_i) = \frac{h_i^2}{6} a_i + b_i x_i + c_i = y_i ,$$

$$P_i(x_{i+1}) = \frac{h_i^2}{6} a_{i+1} + b_i x_{i+1} + c_i = y_{i+1} , \quad i = 0, \dots, n-1.$$

Din aceste relații calculăm  $b_i$  și  $c_i$  în funcție de  $a_i$ ,  $a_{i+1}$ ,  $y_i$ ,  $y_{i+1}$ :

$$b_i = \frac{y_{i+1} - y_i}{h_i} - \frac{h_i}{6} (a_{i+1} - a_i),$$

$$i = 0, \dots, n-1$$

$$c_i = \frac{x_{i+1} y_i - x_i y_{i+1}}{h_i} - \frac{h_i}{6} (x_{i+1} a_i - x_i a_{i+1}).$$

Din condiția de continuitate a primei derivate a funcției spline cubice  $P'_{i-1}(x_i) = P'_i(x_i)$ ,  $i = 1, \dots, n-1$ , ținând seama de:

$$P'_{i-1}(x) = \frac{(x - x_{i-1})^2}{2h_{i-1}} a_i - \frac{(x_i - x)^2}{2h_{i-1}} a_{i-1} + b_{i-1},$$

$$P'_i(x) = \frac{(x - x_i)^2}{2h_i} a_{i+1} - \frac{(x_{i+1} - x)^2}{2h_i} a_i + b_i,$$

rezultă, utilizând formulele pentru  $b_{i-1}$  și  $b_i$  deduse mai sus:

$$P'_{i-1}(x_i) = \frac{h_{i-1}}{2} a_i + \frac{y_i - y_{i-1}}{h_{i-1}} - \frac{h_{i-1}}{6} (a_i - a_{i-1}) =$$

$$P'_i(x_i) = -\frac{h_i}{2} a_i + \frac{y_{i+1} - y_i}{h_i} - \frac{h_i}{6} (a_{i+1} - a_i)$$

sau

$$(h_{i-1} + h_i) a_i + h_i a_{i+1} = 6 \left( \frac{y_{i+1} - y_i}{h_i} - \frac{y_i - y_{i-1}}{h_{i-1}} \right), \quad i = 1, \dots, n-1. \quad (1)$$

Pentru  $i=1$  și  $i=n$  din (1) avem:

$$2(h_0 + h_1)a_1 + h_1a_2 = 6 \left( \frac{y_2 - y_1}{h_1} - \frac{y_1 - y_0}{h_0} \right) - h_0a_0$$

$$h_{n-2}a_{n-2} + 2(h_{n-2} + h_{n-1})a_{n-1} = 6 \left( \frac{y_n - y_{n-1}}{h_{n-1}} - \frac{y_{n-1} - y_{n-2}}{h_{n-2}} \right) - h_{n-1}a_n$$

Sistemul liniar format din ecuațiile (1) cu necunoscutele  $\{a_1, a_2, \dots, a_{n-1}\}$  are forma:

$$Ha = f, \text{ cu } H \in \mathbb{R}^{n \times n}, f \in \mathbb{R}^n$$



$$2(h_0 + h_1)a_1 + h_1a_2 = 6\left(\frac{y_2 - y_1}{h_1} - \frac{y_1 - y_0}{h_0}\right) - h_0f''(a)$$

$$(h_{i-1} + h_i)a_i + h_ia_{i+1} = 6\left(\frac{y_{i+1} - y_i}{h_i} - \frac{y_i - y_{i-1}}{h_{i-1}}\right), \quad i = 2, \dots, n-1$$

$$h_{n-2}a_{n-2} + 2(h_{n-2} + h_{n-1})a_{n-1} = 6\left(\frac{y_n - y_{n-1}}{h_{n-1}} - \frac{y_{n-1} - y_{n-2}}{h_{n-2}}\right) - h_{n-1}f''(b)$$

$$H = \begin{bmatrix}
2(h_0 + h_1) & h_1 & \mathbf{0} & \mathbf{0} \cdots \mathbf{0} & \mathbf{0} & \mathbf{0} \\
h_1 & 2(h_1 + h_2) & h_2 & \mathbf{0} \cdots \mathbf{0} & \mathbf{0} & \mathbf{0} \\
\mathbf{0} & h_2 & 2(h_2 + h_3) & h_3 \cdots \mathbf{0} & \mathbf{0} & \mathbf{0} \\
\vdots & & & & & \\
\mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \cdots h_{n-3} & 2(h_{n-3} + h_{n-2}) & h_{n-2} \\
\mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \cdots \mathbf{0} & h_{n-2} & 2(h_{n-2} + h_{n-1})
\end{bmatrix}$$

$$f = \begin{bmatrix} 6 \left( \frac{y_2 - y_1}{h_1} - \frac{y_1 - y_0}{h_0} \right) - h_0 f''(a) \\ 6 \left( \frac{y_{i+1} - y_i}{h_i} - \frac{y_i - y_{i-1}}{h_{i-1}} \right) \quad i = 1, \dots, n-1 \\ 6 \left( \frac{y_n - y_{n-1}}{h_{n-1}} - \frac{y_{n-1} - y_{n-2}}{h_{n-2}} \right) - h_{n-1} f''(b) \end{bmatrix}$$

Matricea  $H$  are diagonala dominantă atât pe linii cât și pe coloane, este simetrică și pozitiv definită prin urmare putem utiliza metoda Gauss-Seidel sau o metodă de relaxare pentru rezolvarea sistemului  $Ha=f$ .

## Interpolare în sensul celor mai mici pătrate

$x$	$x_0$	$x_1$	$x_2$	$\dots$	$x_{n-1}$	$x_n$
$f$	$y_0$	$y_1$	$y_2$	$\dots$	$y_{n-1}$	$y_n$

$$f(x_i) = y_i, \quad i=0, \dots, n$$

$$f(x) \approx S_f(x; a_0, a_1, \dots, a_m)$$

$$S_f(x; a_0, a_1, \dots, a_m) = a_m x^m + a_{m-1} x^{m-1} + \dots + a_1 x + a_0$$

Coeficienții  $a_0, a_1, \dots, a_m$  se găsesc rezolvând problema de minimizare în sensul celor mai mici pătrate:

$$\min\left\{ \sum_{r=0}^n \left( S_f(x_r; a_0, a_1, \dots, a_m) - y_r \right)^2 ; a_0, a_1, \dots, a_m \in \mathbb{R} \right\} \quad (\text{LSP})$$

$$g : \mathbb{R}^{m+1} \rightarrow \mathbb{R}_+ ,$$

$$g(a_0, a_1, \dots, a_m) = \sum_{r=0}^n \left( S_f(x_r; a_0, a_1, \dots, a_m) - y_r \right)^2$$

$$g(a_0, a_1, \dots, a_m) = \sum_{r=0}^n \left( a_m x_r^m + \dots + a_k x_r^k + \dots + a_1 x_r + a_0 - y_r \right)^2$$

$$\frac{\partial g}{\partial a_k}(a_0, a_1, \dots, a_m) = 2 \sum_{r=0}^n \left( a_m x_r^m + \dots + a_k x_r^k + \dots + a_1 x_r + a_0 - y_r \right) x_r^k$$

Soluția problemei de minimizare a problemei (LSP) este obținută rezolvând sistemul de ecuații liniare, de dimensiune  $(m+1)$ :

$$\frac{\partial g}{\partial a_k}(a_0, a_1, \dots, a_m) = 0, \quad k = 0, 1, \dots, m$$

$$\sum_{r=0}^n (a_m x_r^m + \dots + a_k x_r^k + \dots + a_1 x_r + a_0) x_r^k = \sum_{r=0}^n y_r x_r^k, \quad k = 0, \dots, m$$

$$a_0 \sum_{r=0}^n x_r^k + a_1 \sum_{r=0}^n x_r^{k+1} + \dots + a_{m-1} \sum_{r=0}^n x_r^{k+m-1} + a_m \sum_{r=0}^n x_r^{k+m} = \sum_{r=0}^n y_r x_r^k, \quad k = 0, \dots, m$$

Constantele  $\{a_0, a_1, \dots, a_m\}$  sunt soluția sistemului liniar:

$$Ba = z ,$$

$$B \in \mathbb{R}^{(m+1) \times (m+1)} , B = (b_{kj})_{k,j=0}^m , z \in \mathbb{R}^{(m+1)} z = (z_k)_{k=0}^m$$

$$b_{kj} = \sum_{r=0}^n x_r^{k+j} , z_k = \sum_{r=0}^n y_r x_r^k , k, j = 0, \dots, m$$

# **Calcul Numeric**

**Cursul 10**

**2022**

*Anca Ignat*



## Rezolvarea ecuațiilor neliniare

Fie  $f : [a, b] \rightarrow \mathbb{R}$  o funcție continuă pe intervalul  $[a, b]$  astfel ca  $f(a)f(b) < 0$ . În aceste condiții, există  $x^* \in (a, b)$  astfel ca  $f(x^*) = 0$ . În cele ce urmează ne propunem să aproximăm soluția  $x^*$  a ecuației neliniare  $f(x) = 0$ .

### Metoda biseției (a înjumătățirii intervalului)

Pp.  $f(a)f(b) < 0!!!$  Pentru a aproxima soluția  $x^*$  căutată, vom construi un șir de intervale  $\{[a_k, b_k]; k \geq 0\}$  ce satsifac:

$$x^* \in [a_k, b_k]$$

$$[a_{k+1}, b_{k+1}] \subset [a_k, b_k]$$

$$b_{k+1} - a_{k+1} = \frac{b_k - a_k}{2}$$

Pentru primul interval vom considera:

$$a_0 = a, \quad b_0 = b, \quad k = 0.$$

Considerăm punctul  $c$  de mijloc al intervalului  $[a_k, b_k]$ :

$$c = \frac{a_k + b_k}{2}$$

Avem următoarele 3 variante:

1.  $f(c)=0$  - soluția căutată este  $x^*=c$ , algoritmul se oprește;
2.  $f(a_k)f(c) < 0 \rightarrow$  soluția se găsește în intervalul  $(a_k, c)$ ,  
continuăm procedeul cu intervalul  $[a_{k+1} = a_k, b_{k+1} = c]$ ;
3.  $f(b_k)f(c) < 0 \rightarrow$  soluția se găsește în intervalul  $x^* \in (c, b_k)$   
procedeul continuă cu intervalul  $[a_{k+1} = c, b_{k+1} = b_k]$ .

Dat  $\varepsilon > 0$  există un interval  $[a_{\bar{k}}, b_{\bar{k}}]$  astfel ca  $x^* \in (a_{\bar{k}}, b_{\bar{k}})$  și

$b_{\bar{k}} - a_{\bar{k}} < \varepsilon$  ( $\bar{k} > \log_2 \frac{b-a}{\varepsilon}$ ). Pentru  $\varepsilon$  suficient de mic atât

$a_{\bar{k}}$  cât și  $b_{\bar{k}}$  pot fi considerate aproximări ale soluției  $x^*$

( $a_{\bar{k}} \approx x^*$  prin lipsă iar  $b_{\bar{k}} \approx x^*$  prin adaos).

## Metoda tangentei (Newton-Raphson)

Vom presupune că funcția  $f \in C^1[a,b]$  este derivabilă pe  $[a,b]$  cu derivata continuă în acest interval și satisface relația  $f(a)f(b) < 0$ . Pentru a aproxima soluția  $x^*$  a ecuației  $f(x)=0$  vom construi un șir  $\{x_k\} \subseteq \mathbb{R}$  care să convergă la  $x^*$ ,  $x_k \rightarrow x^*$ , pentru  $k \rightarrow \infty$ . Primul element din șir,  $x_0$ , considerăm că este dat. Următorul element din șir se construiește ca fiind punctul de intersecție al tangentei la graficul funcției  $f$  în punctul  $(x_0, f(x_0))$  cu axa absciselor. Procedurul se repetă cu  $x_1$  pentru a-l obține pe  $x_2$ , ș.a.m.d.

$x_1 = \mathbf{Ox} \cap \text{tangenta la graficul funcției } f \text{ în punctul } (x_0, f(x_0))$

$x_2 = \mathbf{Ox} \cap \text{tangenta la graficul funcției } f \text{ în punctul } (x_1, f(x_1))$

$\vdots$

$x_{k+1} = \mathbf{Ox} \cap \text{tangenta la graficul fct. } f \text{ în pt. } (x_k, f(x_k)), k = 0, 1, 2,$

Ecuția tangentei la graficul funcției  $f$  într-un punct  $(a, f(a))$  este următoarea (pentru o funcție derivabilă):

$$y = f(a) + f'(a)(x - a)$$

Pentru a calcula  $\mathbf{x}_{k+1}$  din  $\mathbf{x}_k$  vom considera ecuația tangentei:

$$y = f(\mathbf{x}_k) + f'(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k)$$

unde luăm  $y = \mathbf{0}$ . Avem:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{f(\mathbf{x}_k)}{f'(\mathbf{x}_k)}, \quad k = 0, 1, 2, \dots, \quad \mathbf{x}_0 \text{ dat}$$

Formula de mai sus poate fi folosită doar dacă la fiecare pas  $f'(\mathbf{x}_k) \neq \mathbf{0}$ . Dacă la un pas avem  $f'(\mathbf{x}_k) = \mathbf{0}$  putem calcula câteva iterații  $\mathbf{x}_k$  ( $k \geq \bar{k}$ ) folosind  $f'(\mathbf{x}_{\bar{k}-1})$ .

## Teoremă de convergență

Fie  $f \in C^2[a, b]$ , cu  $f(a)f(b) < 0$ ,  $f'(x) \neq 0$  și  $f''(x) \neq 0$

$\forall x \in [a, b]$ . Dacă alegem

$$x_0 = \begin{cases} a & \text{pentru } f(a)f''(a) > 0 \\ b & \text{pentru } f(b)f''(b) > 0 \end{cases}$$

atunci șirul  $\{x_k; k \geq 0\}$  construit cu metoda tangentei este monoton, mărginit și convergent la unica soluție  $x^*$  a ecuației  $f(x)=0$ . Ordinul de convergență este mai mare decât 2.



## Metoda falsei poziții (a coardei)

Presupunem că funcția  $f$  este continuă,  $f \in C[a, b]$  și satisface relația  $f(a)f(b) < 0$ . Vom construi un șir  $\{x_k\} \subseteq \mathbb{R}$  care să convergă la soluția căutată  $x^*$ ,  $x_k \rightarrow x^*$ , pentru  $k \rightarrow \infty$ . Considerăm date primul element din șir,  $x_0$  și un alt punct  $\tilde{x}$ . Procedeu de construire a șirului este următorul:

$x_1 = \mathbf{Ox} \cap$  dreapta ce unește punctele  $(\tilde{x}, f(\tilde{x})), (x_0, f(x_0))$

$x_2 = \mathbf{Ox} \cap$  dreapta ce unește punctele  $(\tilde{x}, f(\tilde{x})), (x_1, f(x_1))$

$\vdots$

$x_{k+1} = \mathbf{Ox} \cap$  dreapta ce unește pt.  $(\tilde{x}, f(\tilde{x})), (x_k, f(x_k)), k = 0, 1, 2, \dots$

Ecuția dreptei ce trece prin punctele  $(a, f(a))$  cu  $(b, f(b))$  este:

$$\frac{y - f(a)}{f(a) - f(b)} = \frac{x - a}{a - b} .$$

Pentru a-l obține pe  $\mathbf{x}_{k+1}$  din  $\mathbf{x}_k$  avem:

$$\frac{\mathbf{y} - f(\mathbf{x}_k)}{f(\mathbf{x}_k) - f(\tilde{\mathbf{x}})} = \frac{\mathbf{x} - \mathbf{x}_k}{\mathbf{x}_k - \tilde{\mathbf{x}}} \quad \text{cu} \quad \mathbf{y} = \mathbf{0}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{f(\mathbf{x}_k)(\mathbf{x}_k - \tilde{\mathbf{x}})}{f(\mathbf{x}_k) - f(\tilde{\mathbf{x}})} = \frac{\tilde{\mathbf{x}}f(\mathbf{x}_k) - \mathbf{x}_k f(\tilde{\mathbf{x}})}{f(\mathbf{x}_k) - f(\tilde{\mathbf{x}})},$$

$k = 0, 1, 2, \dots, \quad \mathbf{x}_0, \tilde{\mathbf{x}} - \text{date}$

## Teoremă de convergență

Fie  $f \in C^2[a, b]$ , cu  $f(a)f(b) < 0$ ,  $f'(x) \neq 0$  și  $f''(x) \neq 0 \forall x \in [a, b]$ . Dacă alegem:

$$\begin{cases} \tilde{x} = a \text{ și } x_0 = b & \text{pentru } f(a)f''(a) > 0 \\ \tilde{x} = b \text{ și } x_0 = a & \text{pentru } f(b)f''(b) > 0 \end{cases}$$

atunci șirul  $\{x_k; k \geq 0\}$  construit cu metoda falsei poziții este monoton, mărginit deci convergent la unica soluție  $x^*$  a ecuației  $f(x)=0$ .

## Metoda secantei

Presupunem că funcția  $f$  este continuă,  $f \in C[a, b]$  și satisface relația  $f(a)f(b) < 0$ . Vom construi un șir  $\{x_k\} \subseteq \mathbb{R}$  care să convergă la soluția căutată  $x^*$ ,  $x_k \rightarrow x^*$ , pentru  $k \rightarrow \infty$ . Considerăm date primele două elemente din șir,  $x_0$  și  $x_1$ .

Procedeul de construire a șirului este următorul:

$$\begin{aligned}
x_2 &= \mathbf{Ox} \cap \text{dreapta ce unește punctele } (x_0, f(x_0)), (x_1, f(x_1)), \\
x_3 &= \mathbf{Ox} \cap \text{dreapta ce unește punctele } (x_1, f(x_1)), (x_2, f(x_2)), \\
&\vdots \\
x_{k+1} &= \mathbf{Ox} \cap \text{dreapta ce unește pct. } (x_{k-1}, f(x_{k-1})), (x_k, f(x_k)) , \\
&\qquad\qquad\qquad k = 1, 2, \dots
\end{aligned}$$

Obținem elementul  $x_{k+1}$  din  $x_k$  și  $x_{k-1}$  astfel:

$$\frac{y - f(x_k)}{f(x_k) - f(x_{k-1})} = \frac{x - x_k}{x_k - x_{k-1}} \quad \text{cu } y = 0$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}_{k-1})}{f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})} =$$

$$= \frac{\mathbf{x}_{k-1}f(\mathbf{x}_k) - \mathbf{x}_k f(\mathbf{x}_{k-1})}{f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})},$$

$k = 1, 2, \dots$ ,  $\mathbf{x}_0, \mathbf{x}_1$  dađı

## Teoremă de convergență

Fie  $x^*$  o soluție a ecuației  $f(x)=0$ . Pp. că  
 $f \in C^2[x^* - r, x^* + r]$ ,  $f'(x) \neq 0$  și  $f''(x) \neq 0$   
 $\forall x \in [x^* - r, x^* + r]$ . Atunci există

$0 < r_0 \leq r$  pentru care, dacă  $x_0, x_1 \in [x^* - r_0, x^* + r_0]$  atunci  
 $x_k \in [x^* - r_0, x^* + r_0], \forall k \geq 2$  și  $x_k \rightarrow x^*$ , pentru  $k \rightarrow \infty$ .

Ordinul de convergență este  $q = \frac{1 + \sqrt{5}}{2} \approx 1.61803$ .



## Metoda lui Laguerre

Fie polinomul:

$$p(x) = a_0(x - x_1)(x - x_2)\cdots(x - x_n) , a_0 \neq 0$$

Metoda lui Laguerre propune construirea unui șir de numere care să convergă la una din rădăcinile polinomului  $p$ .

Considerăm derivata polinomului  $p$ :

$$p'(x) = p(x) \left[ \frac{1}{x - x_1} + \frac{1}{x - x_2} + \dots + \frac{1}{x - x_n} \right]$$

Avem:

$$\ln | p(x) | = \ln | a_0 | + \ln | x - x_1 | + \ln | x - x_2 | + \dots + \ln | x - x_n |$$

$$\frac{d}{dx} \ln | p(x) | = \frac{1}{x - x_1} + \frac{1}{x - x_2} + \dots + \frac{1}{x - x_n} = \frac{p'(x)}{p(x)} = G(x)$$

$$\begin{aligned}
-\frac{d^2}{dx^2} \ln |p(x)| &= \frac{1}{(x-x_1)^2} + \frac{1}{(x-x_2)^2} + \dots + \frac{1}{(x-x_n)^2} = \\
&= \frac{[p'(x)]^2 - p(x)p''(x)}{[p(x)]^2} = H(x)
\end{aligned}$$

Fie  $x_1$  rădăcina pe care vrem s-o aproximăm și  $y_k$  valoarea aproximativă curentă. Notăm cu  $a = y_k - x_1$  și facem presupunerea că  $y_k$  se află la aceeași distanță de toate celelalte rădăcini, adică:

$$y_k - x_i = b \quad \forall i = 2, \dots, n.$$

Prin urmare avem:

$$G(y_k) = \frac{p'(y_k)}{p(y_k)} = \frac{1}{a} + \frac{n-1}{b}$$

$$H(y_k) = \frac{[p'(y_k)]^2 - p(y_k)p''(y_k)}{[p(y_k)]^2} = \frac{1}{a^2} + \frac{n-1}{b^2}$$

Rezolvăm acest sistem în raport cu  $a$  și obținem:

$$a = \frac{n}{\max \left[ G(y_k) \pm \sqrt{(n-1)(nH(y_k) - G^2(y_k))} \right]}$$

Semnul la numitor este ales astfel ca expresia să aibă magnitudine maximă. Dacă ținem cont de expresiile pentru  $G$  și  $H$  obținem pentru  $a$  următoarea formulă:

$$a = \frac{n p(y_k)}{\max \left[ p'(y_k) \pm \sqrt{(n-1)^2 [p'(y_k)]^2 - n(n-1) p(y_k) p''(y_k)} \right]}$$

Următorul element din șir va fi:

$$y_{k+1} = y_k - a = y_k - \frac{n}{\max \left[ G(y_k) \pm \sqrt{(n-1)(nH(y_k) - G^2(y_k))} \right]}$$

$$y_{k+1} = y_k - \frac{n p(y_k)}{\max \left[ p'(y_k) \pm \sqrt{(n-1)^2 [p'(y_k)]^2 - n(n-1)p(y_k)p''(y_k)} \right]}$$

Procedeul se oprește când  $a$  devine suficient de mic. Metoda lui Laguerre se poate apl. și ptr. aprox. răd. complexe și de asemenea pentru polinoame cu coeficienți complecși. Pentru rădăcini simple metoda lui Laguerre are ordinul de convergență 3.

## Sisteme de ecuații neliniare

Consideră sistemul neliniar:

$$\begin{cases} f_1(x_1, x_2, \dots, x_n) = 0 \\ f_2(x_1, x_2, \dots, x_n) = 0 \\ \vdots \\ f_n(x_1, x_2, \dots, x_n) = 0 \end{cases} \Leftrightarrow F(X) = \mathbf{0}, \quad F = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix}, \quad X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

Fie matricea jacobiană asociată funcției  $F$  (presupunem că funcțiile  $f_i$  sunt diferențiabile):

$$\nabla F(X) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & & & \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix} (x_1, x_2, \dots, x_n)$$

Pentru a găsi soluția  $X^*$  a sistemului de ecuații neliniare  $F(X) = 0$  se construiește un șir de vectori  $X^{(k)} \in \mathbb{R}^n$  astfel:



$X^{(0)}$  – dat

$$X^{(k+1)} = X^{(k)} - [\nabla F(X^{(k)})]^{-1} F(X^{(k)}) = X^{(k)} + \Delta^{(k)},$$

$$\Delta^{(k)} = -[\nabla F(X^{(k)})]^{-1} F(X^{(k)}) \in \mathbb{R}^n, \quad k = 0, 1, \dots$$

Vectorul de corecție  $\Delta^{(k)}$  poate fi calculat și ca soluție a sistemului liniar:

$$\nabla F(X^{(k)}) \Delta = -F(X^{(k)})$$

unde matricea sistemului este matricea jacobiană calculată în punctul  $X^{(k)}$  iar vectorul termenilor liberi este  $(-F(X^{(k)}))$ .

Metoda descrisă mai sus poartă numele de metoda Newton. Pentru  $n=1$  metoda Newton este chiar metoda tangentei descrisă anterior.

## Optimizare numerică

$$\min \{ f(x) ; x \in \mathbb{R}^n \} , f : \mathbb{R}^n \rightarrow \mathbb{R}$$

Un punct  $x^* \in \mathbb{R}^n$  se numește **punct de minim global** pentru funcția  $f$  dacă  $f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^n$ .

Un punct  $x^* \in \mathbb{R}^n$  se numește **punct de minim local** pentru funcția  $f$  dacă există o vecinătate  $V$  a punctului  $x^*$  pentru care  $f(x^*) \leq f(x) \quad \forall x \in V$ .

Un punct  $x^* \in \mathbb{R}^n$  se numește **punct de minim strict local** pentru funcția  $f$  dacă există o vecinătate  $V$  a punctului  $x^*$  pentru care  $f(x^*) < f(x) \quad \forall x \in V, x \neq x^*$ .

$$V = S(x^*, r) = \{ x \in \mathbb{R}^n ; \|x - x^*\| \leq r \}.$$

Dacă funcția  $f \in C^1(\mathbb{R}^n)$  se numește **gradient** al funcției  $f$  în punctul  $\mathbf{x}$  vectorul:

$$\nabla f(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)^T.$$

Dacă funcția  $f \in C^2(\mathbb{R}^n)$  se numește **matrice hessiană** a funcției  $f$  în punctul  $\mathbf{x}$  matricea:

$$\nabla^2 f(\mathbf{x}) = \left( \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) \right)_{i,j=1,\dots,n}$$

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{pmatrix}$$

## Teorema lui Taylor

Dacă funcția  $f \in C^1(\mathbb{R}^n)$  este continuu diferențiabilă atunci există  $t \in (0,1)$  astfel ca:

$$f(x + p) = f(x) + \nabla f(x + tp)^T p. \quad (1)$$

Dacă  $f \in C^2(\mathbb{R}^n)$  este de două ori continuu diferențiabilă atunci există  $t \in (0,1)$  astfel ca:

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + tp) p. \quad (2)$$

## Teorema 1 (condiții necesare de optim de ordinul întâi)

Fie  $f \in C^1(\mathbb{R}^n)$  și  $\mathbf{x}^* \in \mathbb{R}^n$  un punct de minim local pentru  $f$ .  
Atunci  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .

Un punct  $\mathbf{x}^*$  pentru care  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  se numește **punct staționar**. Teorema de mai sus spune că orice punct de minim local trebuie să fie punct staționar.

## **Teorema 2** (condiții necesare de optim de ordinul al doilea)

Fie  $f \in C^2(\mathbb{R}^n)$  și  $x^* \in \mathbb{R}^n$  un punct de minim local pentru  $f$ .  
Atunci  $\nabla f(x^*) = \mathbf{0}$  și matricea hessiană  $\nabla^2 f(x^*)$  este pozitiv semidefinită.

O matrice  $A \in \mathbb{R}^{n \times n}$  este **pozitiv semidefinită** dacă

$$(Ax, x)_{\mathbb{R}^n} \geq \mathbf{0} \quad \forall x \in \mathbb{R}^n$$

și este **pozitiv definită** dacă

$$(Ax, x)_{\mathbb{R}^n} > \mathbf{0} \quad \forall x \in \mathbb{R}^n, x \neq \mathbf{0}.$$



### **Teorema 3** (condiții suficiente de ordinul al doilea)

Fie  $f \in C^2(\mathbb{R}^n)$  și  $x^* \in \mathbb{R}^n$  un punct pentru care  $\nabla f(x^*) = \mathbf{0}$  și matricea hessiană  $\nabla^2 f(x^*)$  este pozitiv definită. Atunci  $x^*$  este punct de minim strict local pentru  $f$ .

O funcție  $f$  se numește **convexă** dacă:

$$f(ax_1 + (1-a)x_2) \leq af(x_1) + (1-a)f(x_2), \forall x_1, x_2 \in \mathbb{R}^n, \forall a \in [0,1].$$

## Teorema

Dacă  $f$  este funcție convexă orice punct de minim local este punct de minim global. Dacă în plus  $f \in C^1(\mathbb{R}^n)$  atunci orice punct staționar este punct de minim global.

## Metode de descreștere

Se numește **direcție de descreștere** a funcției  $f$  în punctul  $x$  un vector  $d \in \mathbb{R}^n$  pentru care

$$f(x + \alpha d) \leq f(x) \quad , \quad \forall \alpha \in [0, \bar{\alpha}).$$

## Teoremă

Fie  $f \in C^1(\mathbb{R}^n)$ . Un vector  $d$  este direcție de descreștere pentru  $f$  în  $x$  dacă și numai dacă:

$$\nabla f(x)^T d = (d, \nabla f(x))_{\mathbb{R}^n} < 0.$$

## *Algoritm de descreștere*

1. alege  $x \in \mathbb{R}^n$

2. do

- găsește o direcție de descreștere  $d$  a lui  $f$  în  $x$

- găsește  $\tilde{\alpha} > 0$  astfel ca

$$f(x + \tilde{\alpha}d) = \min\{f(x + \alpha d); \alpha \in (0, \bar{\alpha})\}$$

(ajustarea pasului – line search)

-  $x = x + \tilde{\alpha}d$

while (nu am găsit soluția)

## Ajustarea pasului (line search)

Trebuie rezolvată o problemă de minimizare unidimensională

$$\min\{g(a) ; a \in [b,c]\}$$

### Metoda Newton

Fie  $a_k$  aproximarea curentă a soluției problemei de minimizare de mai sus. Vom folosi dezvoltarea în serie Taylor a funcției  $g$

$$g(a) = g(a_k) + g'(a_k)(a - a_k) + \frac{1}{2}g''(a_k)(a - a_k)^2 + \frac{1}{3!}g'''(a_k)(a - a_k)^3$$

$$q(a) = g(a_k) + g'(a_k)(a - a_k) + \frac{1}{2}g''(a_k)(a - a_k)^2$$

Pentru  $a \approx a_k$  putem considera că funcția  $q$  aproximează funcția  $g$ ,  $q(a) \approx g(a)$  și

$$\min\{g(a) ; a \in [c, b]\} \approx \min\{q(a) ; a \in [c, b]\}.$$

Construim elementul  $a_{k+1}$  ca fiind soluția problemei:

$$q(a_{k+1}) = \min\{q(a) ; a \in [c, b]\}$$

$a_{k+1}$  este soluția ecuației

$$q'(a) = 0 \Leftrightarrow q'(a) = g'(a_k) + g''(a_k)(a - a_k) \Rightarrow a_{k+1} = a_k - \frac{g'(a_k)}{g''(a_k)}$$

## Metoda secantei

Dacă în relațiile de mai sus se face următoarea aproximare:

$$g''(a_k) \approx \frac{g'(a_k) - g'(a_{k-1})}{a_k - a_{k-1}}$$

obținem următoarea metodă de aproximare, numită și metoda secantei:

$$a_{k+1} = a_k - g'(a_k) \frac{a_k - a_{k-1}}{g'(a_k) - g'(a_{k-1})}$$

## Aproximare spline cubică

Vom folosi pentru a construi  $a_{k+1}$  pe  $a_{k-1}$ ,  $a_k$ ,  $g(a_{k-1})$ ,  $g(a_k)$ ,  $g'(a_{k-1})$ ,  $g'(a_k)$ . Aproximăm funcția  $g$  cu un polinom de grad 3

$$g(a) \approx q(a) = c_0 a^3 + c_1 a^2 + c_2 a + c_3$$

Funcția  $q$  (respectiv constantele  $c_i$ ) se calculează a.î.

$$q(a_{k-1}) = g(a_{k-1}) \quad , \quad q(a_k) = g(a_k)$$

$$q'(a_{k-1}) = g'(a_{k-1}) \quad , \quad q'(a_k) = g'(a_k)$$



$a_{k+1}$  este punctul de minim al funcției  $q$

$$a_{k+1} = a_k - (a_k - a_{k-1}) \left[ \frac{g'(a_k) + u_2 - u_1}{g'(a_k) - g'(a_{k-1}) + 2u_2} \right]$$

$$u_1 = g'(a_{k-1}) + g'(a_k) - 3 \frac{g(a_k) - g(a_{k-1})}{a_k - a_{k-1}}$$

$$u_2 = \sqrt{u_1^2 - g'(a_{k-1})g'(a_k)}$$

## Ajustarea inexactă a pasului

$$f(x + \tilde{\alpha}d) \cong \min\{f(x + \alpha d); \alpha \in (0, \bar{\alpha})\}$$

- nu se obține  $\tilde{\alpha}$  optimal
- pentru reducerea timpului de calcul optimizarea se oprește înainte de a ajunge la soluție, în funcție de anumite criterii/teste de oprire

## Regula lui Armijo

$$g(a) = f(x_k + a d_k) \quad , \quad \bar{g}(a) = g(0) + \varepsilon a g'(0) \quad , \quad \varepsilon \in (0,1)$$

$\bar{a}$  este acceptabil după regula lui Armijo dacă

$$(1) \quad g(\bar{a}) \leq \bar{g}(\bar{a})$$

$$(2) \quad g(\sigma \bar{a}) \geq \bar{g}(\sigma \bar{a})$$

**$k = 0;$**

**se alege  $a_0$**

***while* (  $g(a_k) > \bar{g}(a_k)$  )**

**$\cdot a_{k+1} = \frac{1}{\sigma} a_k$**

**$\cdot k = k + 1$**

**$\sigma = 2$  sau  $10$  ,  $\varepsilon = 0.2$**

## Testul Goldstein

Dat  $\varepsilon \in (0, 2)$ ,  $\bar{\alpha}$  este considerat acceptabil dacă:

$$g(\bar{\alpha}) \leq \bar{g}(\bar{\alpha}) \text{ si } g(\bar{\alpha}) > g(\mathbf{0}) + (1 - \varepsilon)g'(\mathbf{0})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \bar{\alpha}d_k$$

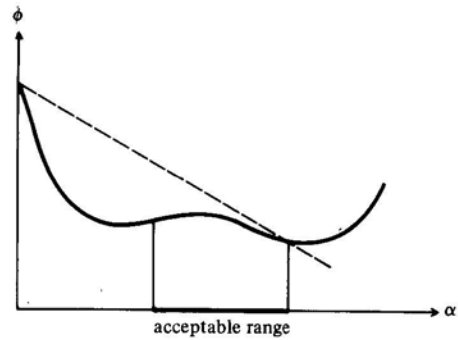
$\bar{\alpha}$  este acceptat dacă

$$\varepsilon \leq \frac{f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)}{\alpha \nabla f(\mathbf{x}_k)d_k} \leq 1 - \varepsilon$$

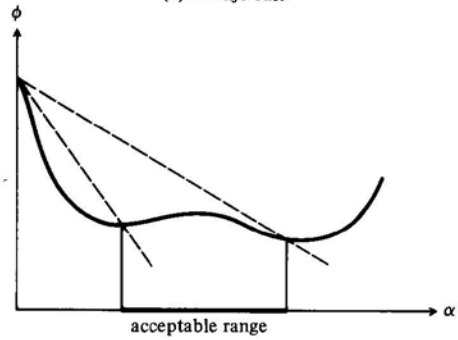
## Testul Wolfe

$$\varepsilon \in (0, 2)$$

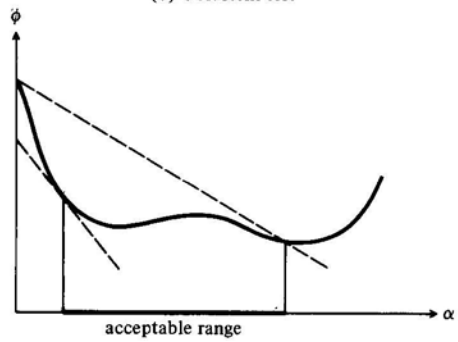
$$\bar{\alpha} : \quad g(\bar{\alpha}) \leq \bar{g}(\bar{\alpha}) = g(\mathbf{0}) + \varepsilon \bar{\alpha} g'(\mathbf{0})$$
$$g'(\bar{\alpha}) \geq (1 - \varepsilon)g'(\mathbf{0})$$



(a) Armijo rule



(b) Goldstein test



(c) Wolfe test

# **Calcul Numeric**

**Cursul 11**

**2022**

*Anca Ignat*



## Metoda pantei maxime

$$\min\{ f(x) ; x \in \mathbb{R}^n \} , f : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$x_0 \text{ dat} , x_{k+1} = x_k + \alpha_k d_k , k = 0, 1, \dots$$

$d_k$  direcție de descreștere a lui  $f$  în  $x_k$

$$\alpha_k > 0 \quad f(x_k + \alpha_k d_k) = \min\{ f(x_k + \alpha d_k) ; \alpha \in [0, \bar{\alpha}] \}$$

$d_k$  direcție de descreștere dacă :

$$\nabla f(x_k)^T d_k < 0.$$

Metoda pantei maxime:

$$d_k = -\nabla f(x_k)$$

Cazul pătratic:

$$f(x) = \frac{1}{2} x^T A x - x^T b = \frac{1}{2} (Ax, x)_{\mathbb{R}^n} - (b, x)_{\mathbb{R}^n}$$

$b \in \mathbb{R}^n$  ,  $A \in \mathbb{R}^{n \times n}$  simetrică și pozitiv definită

$A$  pozitiv definită  $\rightarrow \det A \neq 0$ ,  $f$  este strict convexă

$$\nabla f(x) = Ax - b \quad , \quad \nabla^2 f(x) = A$$

$f(x^*) = \min\{f(x); x \in \mathbb{R}^n\}$   $x^*$  punct unic de minim  $\Leftrightarrow$

$x^*$  soluția sistemului liniar  $Ax = b$  ,  $x^* = A^{-1}b$

$$g(x) = Ax - b$$

$$x_{k+1} = x_k - \alpha_k g_k , \quad g_k = Ax_k - b$$

$$\alpha_k = \operatorname{argmin}\{f(x_k - \alpha g_k); \alpha \in [0, \bar{\alpha}]\}$$

$$\begin{aligned}
f(\mathbf{x}_k - \alpha \mathbf{g}_k) &= \frac{1}{2}(\mathbf{x}_k - \alpha \mathbf{g}_k)^T A(\mathbf{x}_k - \alpha \mathbf{g}_k) - (\mathbf{x}_k - \alpha \mathbf{g}_k)^T \mathbf{b} = \\
&= \frac{1}{2}(\mathbf{g}_k^T A \mathbf{g}_k) \alpha^2 - (\mathbf{g}_k^T A \mathbf{x}_k - \mathbf{g}_k^T \mathbf{b}) \alpha + f(\mathbf{x}_k) = \\
&= \frac{1}{2}(\mathbf{g}_k^T A \mathbf{g}_k) \alpha^2 - (\mathbf{g}_k^T \mathbf{g}_k) \alpha + f(\mathbf{x}_k)
\end{aligned}$$

$f(\mathbf{x}_k - \alpha \mathbf{g}_k)$  ecuație de gr. 2 în  $\alpha$ , coef. lui  $\alpha^2$ ,  $\mathbf{g}_k^T A \mathbf{g}_k > 0$

$$\alpha_{\min} = \alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T A \mathbf{g}_k}$$

Metoda pantei maxime pentru funcționale pătratice:

$$\mathbf{x}_0 - \text{dat}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \left( \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k} \right) \mathbf{g}_k, \quad k = 0, 1, \dots$$

$$\mathbf{g}_k = \mathbf{A} \mathbf{x}_k - \mathbf{b}$$

$$E(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \mathbf{A}(\mathbf{x} - \mathbf{x}^*) = f(\mathbf{x}) + \frac{1}{2} \mathbf{x}^{*T} \mathbf{A} \mathbf{x}^*$$

$$\nabla E(\mathbf{x}) = \nabla f(\mathbf{x}) = \mathbf{g}(\mathbf{x})$$

Șirul construit cu metoda pantei maxime satisface:

$$E(\mathbf{x}_{k+1}) = \left[ 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T \mathbf{A} \mathbf{g}_k)(\mathbf{g}_k^T \mathbf{A}^{-1} \mathbf{g}_k)} \right] E(\mathbf{x}_k)$$

## Inegalitatea lui Kantorovich

$A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ ,  $A > \mathbf{0}$  pozitiv definită

$$\frac{(x^T x)^2}{(x^T A x)(x^T A^{-1} x)} \geq \frac{4cC}{(c + C)^2}$$

$c, C$  - cea mai mica și cea mai mare valoare proprie a lui  $A$ .

## Teoremă

Cazul pătratic: Pentru orice iterație inițială  $\mathbf{x}_0$ , șirul construit cu metoda pantei maxime converge la  $\mathbf{x}^*$  unicul punct de minim al funcției  $f$ . Avem:

$$E(\mathbf{x}_{k+1}) \leq \left( \frac{C - c}{C + c} \right)^2 E(\mathbf{x}_k)$$

Cazul general:  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f \in C^2(\mathbb{R}^n)$ ,  $f$  are un punct de minim local  $\mathbf{x}^*$ . Presupunem că  $F(\mathbf{x}^*) = \nabla^2 f(\mathbf{x}^*)$  are  $c > 0$  cea mai mică valoare proprie și  $C > 0$  cea mai mare valoare proprie. Dacă șirul  $\{\mathbf{x}_k\}$  construit cu metoda pantei maxime converge la  $\mathbf{x}^*$  (la fiecare pas  $A = \nabla^2 f(\mathbf{x}_k)$ ,  $\mathbf{b} = \nabla f(\mathbf{x}_k)$ ),  $\mathbf{x}_k \rightarrow \mathbf{x}^*$



atunci  $f(\mathbf{x}_k) \rightarrow f(\mathbf{x}^*)$  converge liniar cu o rată de convergență  $\leq \left(\frac{C-c}{C+c}\right)^2$ .

## Metoda Newton

$\mathbf{x}_k$  - punctul curent de aproximare

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) = g(\mathbf{y})$$

$g$  – funcțională pătratică,  $\mathbf{y} = \mathbf{x} - \mathbf{x}_k$ ,

$$g(y) = \frac{1}{2} y^T A y - y^T b + c$$

$$A = \nabla^2 f(x_k) \quad , \quad b = -\nabla f(x_k) \quad , \quad c = f(x_k)$$

$y^* = \operatorname{argmin}\{g(y); y \in \mathbb{R}^n\}$  este unica soluție a sistemului  
liniar  $Ay = b$  ,  $y^* = A^{-1}b = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$

Metoda Newton

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k) \quad , \quad k = 0, 1, \dots \quad , \quad x_0 \text{ - dat}$$

## Teoremă

Fie  $f \in C^3(\mathbb{R}^n)$  care are un punct de minim local  $\mathbf{x}^*$  astfel ca matricea  $\nabla^2 f(\mathbf{x}^*) > \mathbf{0}$  este pozitiv definită. Dacă punctul de început  $\mathbf{x}_0$  este suficient de aproape de  $\mathbf{x}^*$ ,  $\|\mathbf{x} - \mathbf{x}^*\| \leq r$ , atunci șirul  $\{\mathbf{x}_k\}$  generat cu metoda lui Newton converge la  $\mathbf{x}^*$  și ordinul de convergență este cel puțin 2.

## Metoda gradientilor conjugați (a direcțiilor conjugate)

$A \in \mathbb{R}^{n \times n}$  ,  $A = A^T$  ,  $A > \mathbf{0}$  pozitiv definită

$$\min\{f(x) = \frac{1}{2}x^T Ax - x^T b; x \in \mathbb{R}^n\}$$

### Definiție

Pentru o matrice  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$  simetrică, doi vectori  $d_1, d_2 \in \mathbb{R}^n$  se numesc **A-ortogonali** sau **conjugați în raport cu A** dacă:

$$d_1^T A d_2 = (A d_2, d_1)_{\mathbb{R}^n} = (d_2, A d_1)_{\mathbb{R}^n} = (A d_1, d_2)_{\mathbb{R}^n} = \mathbf{0}$$

$A = \mathbf{0}_{n \times n} \Rightarrow \forall d_1, d_2$  sunt  $A$  – ortogonali

$A = I_n \Rightarrow$  ortogonalitate clasică  $(d_1, d_2)_{\mathbb{R}^n} = \mathbf{0}$

Vectorii  $\{d_0, d_1, \dots, d_k\}$  se numesc  $A$ -ortogonali sau  $A$ -conjugăți dacă:

$$d_i^T A d_j = (A d_j, d_i)_{\mathbb{R}^n} = \mathbf{0} , \forall i \neq j , i, j = 0, 1, \dots, k$$

## Propoziție

Fie  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ ,  $A > \mathbf{0}$ , și  $\{d_0, d_1, \dots, d_k\}$  direcții  $A$ -conjugate,  $d_i \neq \mathbf{0}$ ,  $\forall i = 0, \dots, k$ . Atunci vectorii  $\{d_0, d_1, \dots, d_k\}$  sunt liniar independenți.

$$\left. \begin{array}{l} A \in \mathbb{R}^{n \times n}, A = A^T, A > \mathbf{0} \\ \{d_0, d_1, \dots, d_{n-1}\} \text{ - direcții } A \text{ - conjugate} \end{array} \right\} \Rightarrow \begin{array}{l} \{d_0, d_1, \dots, d_{n-1}\} \\ \text{bază în } \mathbb{R}^n \end{array}$$

$$x^* = \operatorname{argmin}\{f(x); x \in \mathbb{R}^n\} \Leftrightarrow x^* \text{ soluția sist. } Ax = b$$

$$x^* = \alpha_0 d_0 + \alpha_1 d_1 + \cdots + \alpha_{n-1} d_{n-1}$$

$$\left( \underbrace{Ax^*}_{=b}, d_i \right)_{\mathbb{R}^n} = d_i^T Ax^* = \alpha_i d_i^T Ad_i$$

$$\alpha_i = \frac{d_i^T b}{d_i^T Ad_i} = \frac{(b, d_i)_{\mathbb{R}^n}}{(Ad_i, d_i)_{\mathbb{R}^n}}$$

$$x^* = \sum_{i=0}^{n-1} \frac{(b, d_i)_{\mathbb{R}^n}}{(Ad_i, d_i)_{\mathbb{R}^n}} d_i$$

Considerăm procesul iterativ

$$\mathbf{x}_0 \in \mathbb{R}^n - \text{dat}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k, \quad k = 0, 1, 2, \dots, n-1$$

$$\alpha_k = -\frac{\mathbf{g}_k^T \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}, \quad \mathbf{g}_k = \mathbf{A} \mathbf{x}_k - \mathbf{b}$$

are proprietatea că  $\mathbf{x}_n = \mathbf{x}^*$ .

$$\mathbf{x}_k = \mathbf{x}_0 + \alpha_0 \mathbf{d}_0 + \alpha_1 \mathbf{d}_1 + \dots + \alpha_{k-1} \mathbf{d}_{k-1}$$

$$\mathbf{x}^* = \mathbf{x}_0 + \beta_0 \mathbf{d}_0 + \beta_1 \mathbf{d}_1 + \dots + \beta_{n-1} \mathbf{d}_{n-1}$$



$$\alpha_i = \frac{d_i^T A(x_k - x_0)}{d_i^T A d_i}, \quad \beta_i = \frac{d_i^T A(x^* - x_0)}{d_i^T A d_i}$$

$$d_k^T A(x_k - x_0) = \mathbf{0}$$

$$x^* - x_k = (\beta_0 - \alpha_0)d_0 + \cdots + (\beta_{k-1} - \alpha_{k-1})d_{k-1} + \beta_k d_k + \cdots + \beta_{n-1}d_{n-1}$$

**Corolar**

$$g_k^T d_i = \mathbf{0} \quad \forall i < k$$

## Algoritmul gradientilor conjugați

$$\mathbf{x}_0 \in \mathbb{R}^n, \mathbf{g}_0 = A\mathbf{x}_0 - \mathbf{b}$$

$$\mathbf{d}_0 = -\mathbf{g}_0 = \mathbf{b} - A\mathbf{x}_0$$

$$\alpha_k = -\frac{\mathbf{g}_k^T \mathbf{d}_k}{\mathbf{d}_k^T A \mathbf{d}_k}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

$$\mathbf{g}_{k+1} = A\mathbf{x}_{k+1} - \mathbf{b} \text{ sau } \mathbf{g}_{k+1} = \mathbf{g}_k + \alpha_k A \mathbf{d}_k$$

$$\beta_k = \frac{\mathbf{g}_{k+1}^T A \mathbf{d}_k}{\mathbf{d}_k^T A \mathbf{d}_k}$$

$$\mathbf{d}_{k+1} = -\mathbf{g}_{k+1} + \beta_k \mathbf{d}_k$$

$$y_1, y_2, \dots, y_p \in \mathbb{R}^n$$

$$\begin{aligned} \text{span}\{y_1, y_2, \dots, y_p\} &= \{y = a_1 y_1 + \dots + a_p y_p \in \mathbb{R}^n; a_i \in \mathbb{R}, i = 1, \dots, p\} \\ &= \text{subspațiul generat de vectorii } y_1, y_2, \dots, y_p \end{aligned}$$

### Teoremă

Presupunem că  $x_k \neq x^*$ . Avem următoarele relații:

$$(1) \quad g_k^T g_i = 0 \quad \forall i = 0, 1, \dots, k-1$$

$$(2) \quad \text{span}\{g_0, g_1, \dots, g_k\} = \text{span}\{g_0, Ag_0, \dots, A^k g_0\}$$

$$(3) \quad \text{span}\{d_0, d_1, \dots, d_k\} = \text{span}\{g_0, Ag_0, \dots, A^k g_0\}$$

$$(4) \quad d_k^T A d_i = 0 \quad \forall i = 0, 1, \dots, k-1$$

Șirul  $x_k \rightarrow x^*$  în cel mult  $n$  pași.

$\mathcal{K}(g_0, k) = \text{span}\{g_0, Ag_0, \dots, A^k g_0\}$  - se numește **subspațiu Krylov de grad  $k$**  pentru  $g_0$ .

## Forma practică a metodei gradientilor conjugați

$$\mathbf{x}_0 \in \mathbb{R}^n - \text{dat}, \quad k = 0$$

$$\mathbf{g}_0 = A\mathbf{x}_0 - \mathbf{b}, \quad \mathbf{d}_0 = -\mathbf{g}_0$$

*while* ( $\mathbf{g}_k \neq \mathbf{0}$ )

$$\left\{ \begin{array}{l} \alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{d}_k^T A \mathbf{d}_k} \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \\ \mathbf{g}_{k+1} = \mathbf{g}_k + \alpha_k A \mathbf{d}_k \\ \beta_k = \frac{\mathbf{g}_{k+1}^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k} \\ \mathbf{d}_{k+1} = -\mathbf{g}_{k+1} + \beta_k \mathbf{d}_k \\ \mathbf{k} = \mathbf{k} + \mathbf{1}; \end{array} \right.$$

## Teoremă

Dacă matricea  $A$  are doar  $r$  valori proprii distincte, algoritmul gradientilor conjugați calculează soluția  $\mathbf{x}^*$  în cel mult  $r$  iterații.

## Teoremă

Dacă  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  sunt valorile proprii ale matricei  $A$  atunci:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_A^2 \leq \left( \frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right)^2 \|\mathbf{x}_0 - \mathbf{x}^*\|_A^2$$

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_A^2 = 2E(\mathbf{x}_{k+1})$$

$$E(\mathbf{x}_{k+1}) \leq \left( \frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right)^2 E(\mathbf{x}_0).$$

## Metodele gradientilor conjugați neliniare

$$f(x) \approx f(x_k) + \nabla f(x_k)(x - x_k) + \frac{1}{2}(x - x_k)^T \nabla^2 f(x_k)(x - x_k)$$

$$g_k \leftrightarrow -\nabla f(x_k)$$

$$A \leftrightarrow \nabla^2 f(x_k)$$

Varianta pentru funcții oarecare a metodei gradientilor conjugați:

$$\mathbf{x}_0 \in \mathbb{R}^n - \text{dat}, \quad k = 0$$

$$\mathbf{g}_0 = \nabla f(\mathbf{x}_0), \quad \mathbf{d}_0 = -\mathbf{g}_0$$

*while* ( $\mathbf{g}_k \neq \mathbf{0}$ )

$$\left\{ \begin{array}{l} \mathbf{A} = [\nabla^2 f(\mathbf{x}_k)] \\ \alpha_k = -\frac{\mathbf{g}_k^T \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \\ \mathbf{g}_{k+1} = \nabla f(\mathbf{x}_{k+1}) \\ \beta_k = \frac{\mathbf{g}_{k+1}^T \mathbf{A} \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} \\ \mathbf{d}_{k+1} = -\mathbf{g}_{k+1} + \beta_k \mathbf{d}_k \\ k = k + 1; \end{array} \right.$$



## Metoda Fletcher-Reeves

$\alpha_k$  se calculează folosind metoda ajustării pasului

$$\beta_k = \frac{\mathbf{g}_{k+1}^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k}$$

$\mathbf{x}_0 \in \mathbb{R}^n$  – dat,

$\mathbf{g}_0 = \nabla f(\mathbf{x}_0)$ ,  $\mathbf{d}_0 = -\mathbf{g}_0$ ,  $k = 0$

***while* ( $\mathbf{g}_k \neq \mathbf{0}$ )**

$$\left\{ \begin{array}{l} \alpha_k = \min\{f(\mathbf{x}_k + \alpha \mathbf{d}_k); \alpha \in [0, \bar{\alpha}]\} \\ \text{(exact sau inexact cu testul lui Wolfe)} \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \\ \mathbf{g}_{k+1} = \nabla f(\mathbf{x}_{k+1}) \\ \beta_k^{FR} = \frac{\mathbf{g}_{k+1}^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{g}_k} \\ \mathbf{d}_{k+1} = -\mathbf{g}_{k+1} + \beta_k^{FR} \mathbf{d}_k \\ \mathbf{k} = \mathbf{k} + \mathbf{1}; \end{array} \right.$$

Se pune problema dacă  $d_k$  sunt direcții de descreștere?

$$d_k = -g_k + \beta_{k-1} d_{k-1}$$

$$g_k^T d_k = -g_k^T g_k + \beta_{k-1} g_k^T d_{k-1}$$

Dacă se folosește ajustarea pasului exactă:

$\alpha_{k-1}$  este punct de minim local pentru  $f$  pe direcția  $d_{k-1}$  prin urmare  $g_k^T d_{k-1} = 0$  ( $g_k = \nabla f(x_k)$ ).

$\Rightarrow g_k^T d_k = -g_k^T g_k = -\|g_k\|_2^2 < 0 \Rightarrow d_k$  direcție de descreștere

Dacă se folosește ajustarea pasului inexactă am putea avea  $\mathbf{g}_k^T \mathbf{d}_k > \mathbf{0}$  ( $\mathbf{d}_k$  direcție de creștere!!) dar folosind testul lui Wolfe deducem:

$$\begin{aligned} f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) &\leq f(\mathbf{x}_k) + \varepsilon \alpha_k \mathbf{g}_k^T \mathbf{d}_k \\ |\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k| &\leq (1 - \varepsilon) |\mathbf{g}_k^T \mathbf{d}_k| \\ \varepsilon \in \left(0, \frac{1}{2}\right) &\Rightarrow \mathbf{g}_k^T \mathbf{d}_k < \mathbf{0} \end{aligned}$$

## Metoda Polak-Ribière

- variantă a metodei Fletcher-Reeves

$$\beta_k^{PR} = \frac{\mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k)}{\mathbf{g}_k^T \mathbf{g}_k}$$

Dacă se face ajustarea pasului inexactă cu testul lui Wolfe nu putem deduce că  $\mathbf{d}_k$  sunt direcții de descreștere.

Se folosește  $\beta_k^+ = \max\{\beta_k^{PR}, 0\}$  și un test Wolfe adaptat pentru a obține  $\mathbf{d}_k$  direcții de descreștere.

## Varianta Hestenes-Stiefel

$$\beta_k^{HS} = \frac{\mathbf{g}_{k+1}^T (\mathbf{g}_{k+1} - \mathbf{g}_k)}{(\mathbf{g}_{k+1} - \mathbf{g}_k)^T \mathbf{d}_k}$$

## Precondiționare

Se consideră norma:

$$\| \mathbf{x} \|_A = \sqrt{(A\mathbf{x}, \mathbf{x})_{\mathbb{R}^n}}$$

Evaluarea erorii în metoda pantei maxime:

$$\| \mathbf{x}^{(k)} - \mathbf{x}^* \|_A \leq \left( \frac{k(A) - 1}{k(A) + 1} \right)^k \| \mathbf{x}^{(0)} - \mathbf{x}^* \|_A$$

$$k(A) = \frac{\lambda_n}{\lambda_1} - \text{numărul de condiționare spectral}$$

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \text{ valorile proprii ale matricii } A$$

Avem convergență rapidă dacă numărul de condiționare spectrală al matricei  $A$  este apropiat de 1 ( $k(A) \geq 1$  întotdeauna).

Ideea preconditionării este de a transforma sistemul  $Ax=b$  astfel încât să îmbunătățim proprietățile spectrale.

$$Ax = b \quad \Leftrightarrow \quad \tilde{A}x = \tilde{b} \quad , \quad \text{cu } k(\tilde{A}) \ll k(A)$$



## *Precondiționare*

$$Ax = b \rightarrow M^{-1}Ax = M^{-1}b \text{ ( la stânga)}$$

$$\rightarrow AM^{-1}y = b \text{ , } x = M^{-1}y \text{ ( la dreapta)}$$

$$\rightarrow M_1^{-1}AM_2^{-1}y = M_1^{-1}b \text{ , } x = M_2^{-1}y \text{ (split) , } M = M_1M_2$$

cu  $M$  matrice nesingulară ,  $M$  “ $\approx$ ”  $A$ . Matricea  $M$  sau  $M^{-1}$  poartă numele de *matrice de precondiționare*.

Cum trebuie să alegem matricea  $M$ ?

- sistemul preconditionat ( $\tilde{A}x = \tilde{b}$ ) să fie ușor de rezolvat (convergență rapidă)
- matricea de preconditionare să fie economic de construit și aplicat – ietrațiile să nu fie costisitor de construit

### *Matricea de preconditionare Jacobi*

$$M = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}]$$

### *Matrice de preconditionare SSOR*

$$M = (D + L) D^{-1} (D + L)^T \quad (A = L + D + L^T)$$

$$M(\omega) = \frac{1}{2 - \omega} \left( \frac{1}{\omega} D + L \right) \left( \frac{1}{\omega} D \right)^{-1} \left( \frac{1}{\omega} D + L \right)^T, \quad \omega \in (0, 2)$$

Pentru  $\omega$  - optimal, în anumite cazuri:

$$k(M(\omega_{opt})^{-1} A) = O(\sqrt{k(A)})$$

( $\omega_{opt}$  - foarte costisitor de calculat)