

Calcul Numeric

Cursul 2

2019

Anca Ignat

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, $\dim \mathbb{R}^n = \mathbf{n}$ 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 $A \in \mathbb{R}^{m \times n}$:

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

Se definește *matricea transpusă*:

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

Pentru matricea:

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

se definește *matricea adjunctă* A^H :

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

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

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,

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

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

Dacă facem înmulțirea matricială Ae_j obținem coloana j a matricii 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 matricii A , $j=1,\dots,n$;

$e_i^T A$ este linia i a matricii 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 :

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

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

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

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

Proprietățile matricii 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 matricii 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 matrici

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 \mathbf{I}_n matricea unitate:

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

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

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

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:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|;$$

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

$$\|\mathbf{x}\|_\infty = \max\{|x_i|, i = 1..n\}.$$

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

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

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:

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

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

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

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

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

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

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

Norme matriciale

Definiție

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

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

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

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

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

Exemple

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

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

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 matriciale naturale

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

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

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 matricială 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.

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

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

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

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

Pentru $\|x\|_\infty = \max\{|x_i|; i = 1, \dots, n\}$ norma matricială 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 $\|\cdot\|_{\mathbf{i},\mathbf{P}}$
normele matriciale induse

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

Valori și vectori proprii

Definiții

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

$$Ax = \lambda x.$$

Vectorul x se numește *vector propriu (autovector)* asociat val. proprii λ .

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

→ 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 matricii A .

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

Se numește *rază spectrală* a matricii A :

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

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

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

Propoziția 1

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

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

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

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

Propoziția 2

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 matricială naturală pentru care $\|A\| < 1$ atunci:

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

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

Propoziția 3

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

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

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

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

Propoziția 4

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

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

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.

O reprezentare binară pe 64 de biți a unui număr real se face în felul următor: primul bit este bitul de semn, următorii 11 biți reprezintă exponentul c iar următorii 52 de biți conțin informații despre partea fracționară, f , numită și mantisă

$$(-1)^s 2^{c-1023} (1+f) \text{ .}$$

$$0\ 10000000011\ 10111001000100000000000000000000000000000000 = \\ 27.56640625$$

[27.566406249999982236431605997495353221893310546875,
27.5664062500000017763568394002504646778106689453125).

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

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

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 intreg 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 calculule numerice

1. Erori în datele de intrare:

- măsurători afectate de erori sistematice sau perturbații temporare,
- erori de rotunjire: $1/3$, π , $1/7$,...

2. Erori de rotunjire în timpul calculelor:

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

3. Erori de discretizare:

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

4. Simplificări în modelul matematic

- idealizări , ignorarea unor parametri.

5. Erori umane și erori ale bibliotecilor folosite.

Eroare absolută , eroare relativă

a – valoarea exactă,

\tilde{a} – valoarea aproximativă.

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

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

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

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

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

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

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

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

a_1 cu eroare relativă δ_{a_1} și a_2 cu eroare relativă δ_{a_2} :

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

Condiționare \leftrightarrow stabilitate

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

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

O problemă se consideră a fi *prost condiționată* dacă $P(x)$ 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:

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

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

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

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 ± 0.64350i 11.79363 ± 1.65233i
13.99236 ± 2.51883i 16.73074 ± 2.81262i 19.50244 ± 1.94033i

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