

# Metode numerice (notițe curs)

– UPB 2020

## I. Introducere

Determinarea practică a soluțiilor matematice al unor probleme științifice și ingineresti necesită, în cele mai multe cazuri, construirea unui algoritm pentru găsirea lor, adaptat posibilităților de lucru ale sistemului de calcul numeric.

Astfel, **algoritm** reprezintă totalitatea raționamentelor matematice și logice, efectuate pas cu pas pentru rezolvarea problemei, deci o rețetă generală, precis formulată, care conduce la rezultate concrete într-un interval finit de timp.

Formularea matematică (modelul) a unei probleme fizice sau ingineresti conține inevitabil aproximări datorate:

- înțelegerii incomplete a fenomenelor naturale sau a naturii aleatoare ale acestora;
- descrierii matematice exclusive a trăsăturilor esențiale ale procesului fizic, prin neglijarea deliberată a detaliilor legate de efectele secundare

Erorilor asociate acestor aproximări li se adaugă cele rezultate din prelucrarea datelor cu ajutorul unui sistem de calcul numeric. Acesta poate realiza numai un număr limitat de operații aritmetice simple (adunare, scădere, împărțire, înmulțire, ridicare la putere) cu numere finite raționale. Operații importante ca derivarea, integrarea, evaluarea seriilor infinite nu pot fi implementate direct pe un calculator numeric. Acesta este compus din registre de calcul și memorii cu capacitate limitată, încât este imposibil să se reprezinte cantitățile infinite mici sau infinite mari, ori un interval finit de numere reale.

Un algoritm adaptat posibilităților calculatorului numeric, utilizând numai operații aritmetice și unele operații logice se spune că utilizează **Metode Numerice**.

Lucrarea tratează, în continuare, prin această prismă, principalele capitole matematice utilizate în inginerie.

## II. Erori

**Eroarea absolută,  $e_x$ ,** este definită ca fiind diferența dintre valoarea adevărată (exactă),  $x$  și valoarea aproximativă,  $\bar{x}$  (determinată prin calcul sau prin măsurare):

$$e_x = x - \bar{x} \quad (\text{II.1})$$

Dintre cele 3 mărimi numai valoarea aproximativă este cunoscută de obicei, uneori cunoscându-se și marginea erorii (modulul maxim probabil).

**Eroarea relativă** este eroarea absolută împărțită la valoarea aproximativă (prin convenție), pentru că aceasta se cunoaște în general.

$$\varepsilon_x = e_x / \bar{x} \quad (\text{II.2})$$

### II.1 Erori inițiale

Erorile inițiale sunt erori din datele inițiale, cauzate de incertitudinea măsurărilor, din greșeli (grosolane) sau din reprezentarea unei valori transcendente sau iraționale printr-un număr finit de cifre.

Exemple:

- a) o măsurătoare fizică nu poate fi exactă; instrumentul folosit transferă precizia sa mărimii măsurate; în cazul dimensiunilor liniare, un șubler indică cu certitudine 0,1 mm, un micrometru sau un comparator 0,01 mm, iar un traductor capacitiv ajunge la o rezoluție nanometrică. Numărul de cifre semnificative corecte variază deci în funcție de precizia măsurătorii, recomandându-se limitelor toleranței mărimii.
- b) multe numere nu pot fi reprezentate printr-un număr finit de cifre. Dacă exemplele lui  $\pi$  sau  $\sqrt{2}$  sunt evidente, o simplă fracție poate să nu aibă o reprezentare exactă, în funcție de baza de numerație folosită. În sistemul zecimal, numărul  $1/10$  este 0,1, în timp ce în sistemul binar el este 0,0001100110011....., deci o reprezentare infinită imposibilă în calculator. Deci, dacă se face suma a 10 numere reprezentând aproximații binare ale lui 0,1, nu se va obține 1,0.

### II.2 Erori de trunchiere

O procedură numerică poate introduce ea însăși erori. De exemplu, pentru calculul sinusului unui unghi (în radiani) se poate utiliza seria Maclaurin:

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

Aceasta este infinită și nu se pot utiliza toți termenii seriei în calcul. Limitarea expresiei la un număr finit de termeni conduce la o eroare reprezentată de termenii omiși, numită eroare de trunchiere.

Deși nu se poate calcula precis, eroarea de trunchiere se poate estima prin metode specifice fiecărui algoritm, limitele acestora fiind puse în evidență la fiecare problemă analizată.

## II.3 Erori de calcul

### II.3.1. Calculul în virgulă mobilă

La reprezentarea numerelor în dispozitive cu capacitate limitată de memorare și prelucrare este specifică tăierea cifrelor ne semnificative de la sfârșitul acestora, care poate apărea atât la introducerea numărului, cât și în urma efectuării unor operații aritmetice. În unele cazuri particulare, eroarea rezultată prin tăiere poate avea valori importante.

Rezolvarea limitării tehnice a reprezentării numerelor reale în calculator o constituie utilizarea numerelor în virgulă mobilă. Acestea sunt compuse dintr-o fracție numită mantisă și un întreg numit exponent sau caracteristică.

$$N = m \cdot b^c \quad (\text{II.3})$$

Unde:  $m$  = mantisă,  $c$  = exponent,  $b$  = baza de numerație.

*Exemplu: În baza 10, numărul  $N=36,82$  se va rescrie  $N=0,3682 \cdot 10^2$ . El este normalizat pentru că prima cifră după virgulă a mantisei este diferită de zero, adică:*

$$\frac{1}{10} \leq |m| < 1 \quad (\text{II.4})$$

Procedeul de adunare a două numere reale în virgulă mobilă va fi ilustrat pe un calculator ipotetic, zecimal cu patru poziții afectate mantisei pentru numerele a)  $x_1=165,2$  ;  $x_2=0,21$  și respectiv b)  $x_3=26,31$  ;  $x_4=-19,76$ .

Etapele procedurii sunt:

- 1) normalizarea numerelor:  
a)  $x_1=0,1652 \cdot 10^3$  ;  $x_2=0,2100 \cdot 10^0$   
b)  $x_3=0,2631 \cdot 10^2$  ;  $x_4=-0,1976 \cdot 10^2$

- 2) compararea exponenților:  
a)  $c_1=3 > c_2=0$   
b)  $c_3=c_4=2$

- 3) de-normalizarea numărului cu exponent mai mic prin deplasarea mantisei spre dreapta virgulei, cu atâtea poziții cât reprezintă diferența exponenților:  
a)  $x_2=0,2100 \cdot 10^0 = 0,0002 \cdot 10^3$

Prin de-normalizare se pot pierde cifre semnificative ale numărului.

- 4) adunare propriu-zisă  
a)  $S_{12}=x_1+x_2=0,1652 \cdot 10^3+0,0002 \cdot 10^3=0,1654 \cdot 10^3$   
b)  $S_{34}=x_3+x_4=0,2632 \cdot 10^2+(-0,1976) \cdot 10^2=0,0655 \cdot 10^2$
- 5) re-normalizarea rezultatului (dacă este cazul)  
b)  $S_{34}=0,6550 \cdot 10^1$

În exemplul a, suma este afectată de tăierea efectuată la pasul 3, iar în exemplul b ultima cifră (0) nu mai are nici o semnificație, ea apărând în procesul deplasării. Aceste erori se propagă în următoarele operații, lor adăugându-li-se altele, specifice fiecărei etape. Tăierea ultimei cifre este evitată la unele sisteme de calcul printr-o poziție pentru așa-zisa cifră de gardă (în exemplul dat o a cincea cifră în mantisa calculatorului ipotetic), permițând salvarea unei cifre din mantisa deplasată spre dreapta.

Probleme deosebite la calculul în virgulă mobilă mai pot apare în următoarele cazuri:

- a) scăderea a două numere aproximativ egale cu exponent minim conduce la depășire inferioară în virgulă mobilă (*underflow*).

Exemplu: convenindu-se că exponentul minim reprezentabil în calculatorul ipotetic este

$$-615, \text{ atunci: } 0,1459 \cdot 10^{-615} - 0,1024 \cdot 10^{-615} = 0,0425 \cdot 10^{-615} = 0,4250 \cdot 10^{-616}.$$

Diverse sisteme de calcul rezolvă diferit această problemă, dar, de obicei, atribuie rezultatului valoarea zero în virgulă mobilă, ceea ce poate conduce la rezultate foarte ciudate cu un algoritm aparent corect.

- b) adunarea a două numere cu exponent maxim poate conduce la depășire superioară în virgulă mobilă (*overflow*).

Exemplu: Dacă exponentul maxi reprezentabil în calculatorul ipotetic este 615, atunci  $0,9999 \cdot 10^{615} + 0,9999 \cdot 10^{615} = 1,9998 \cdot 10^{615} = 0,1999 \cdot 10^{616}$

De obicei, depășirea exponentului maxim reprezentabil în calculator este semnalată prin întreruperea execuției programului.

La înmulțire/împărțire mantisele normalizate se înmulțesc/împart, iar exponenții se adună/scad, normalizarea rezultatului fiind simplă. Depășirea superioară poate apare atât la înmulțire, dar și la împărțire prin zero (când împărțitorul poate fi rezultatul unei depășiri inferioare).

În mod normal, calculatoarele numerice lucrează în baza 2 sau 16, încât reprezentarea numerelor în virgulă mobilă este:

$$x = m \cdot 2^c \quad \text{cu } \frac{1}{2} \leq |m| < 1, \text{ sau} \quad (\text{II.5})$$

$$x = m \cdot 16^c \quad \text{cu } \frac{1}{16} \leq |m| < 1. \quad (\text{II.6})$$

### II.3.2. Erori de rotunjire

Rezultatul oricăreia dintre cele 4 operații, într-o reprezentare a numerelor cu  $t$  cifre semnificative pentru mantisă, poate fi scris ca:

$$x = m_x \cdot 10^c + p_x \cdot 10^{c-t},$$

unde  $p_x$  este partea din mantisă asupra căreia se acționează la rotunjire, care nu este normalizată, deci  $0 \leq p_x < 1$ .

Unele compilatoare pentru limbaje de nivel înalt folosesc tăierea termenului  $p_x \cdot 10^{c-t}$ , atribuind rezultatului valoarea  $m_x \cdot 10^c$ .

În acest caz, modulul erorii de rotunjire maxime este:

$$\frac{e_x}{x} = \left| \frac{p_x \cdot 10^{c-t}}{m_x \cdot 10^c} \right| \leq \frac{1 \cdot 10^{c-t}}{0,1 \cdot 10^c} = 10^{-t+1} \quad (\text{II.7})$$

În consecință, maximul erorii relative de rotunjire prin tăiere nu depinde de mărimea numărului, ci de numărul de cifre semnificative reprezentate în mantisă ( $t$ ).

Un alt tip de rotunjire uzuală este cea simetrică, aproximarea numărului  $x$ , fiind:

$$\bar{x} = \begin{cases} |m_x| \cdot 10^c, & \text{daca } |p_x| < \frac{1}{2} \\ |m_x| \cdot 10^c + 10^{c-t}, & \text{daca } |p_x| \geq \frac{1}{2} \end{cases} \quad (\text{II.8})$$

unde  $\bar{x}$  și  $m_x$  sunt de același semn.

Eroarea absolută maximă este:

$$|e_x| = \begin{cases} |p_x| \cdot 10^{c-t} \leq \frac{1}{2} \cdot 10^{c-t}, & \text{daca } |p_x| < \frac{1}{2} \\ |1 - p_x| \cdot 10^{c-t} \leq \frac{1}{2} \cdot 10^{c-t}, & \text{daca } |p_x| \geq \frac{1}{2} \end{cases}$$

Valoarea absolută a erorii relative maxime de rotunjire simetrică este:

$$\left| \frac{e_x}{x} \right| \leq \left| \frac{\frac{1}{2} \cdot 10^{c-t}}{m_x \cdot 10^c} \right| \leq \frac{\frac{1}{2} \cdot 10^{c-t}}{0,1 \cdot 10^c} = 5 \cdot 10^{-t},$$

deci cel mult jumătate din cea obținută la rotunjirea prin tăiere.

## II.4. Propagarea erorilor

S-a constatat că, în cazuri extreme, precum scăderea a două numere foarte mici, eroarea relativă a diferenței poate fi destul de mare, propagându-se prin orice operație aritmetică ulterioară.

De aceea expresiile erorilor relative și absolute ale rezultatului unei operații în funcție de erorile operanzilor sunt foarte importante. Notându-se  $x = \bar{x} + e_x$  și  $y = \bar{y} + e_y$ , se determină erorile absolute pentru fiecare operație:

a) adunarea (scăderea)

$$e_{x \pm y} = (x \pm y) - (\bar{x} \pm \bar{y}) = e_x \pm e_y \quad (\text{II.9})$$

b) înmulțirea

$$e_{x \cdot y} = x \cdot y - \bar{x} \cdot \bar{y} = (\bar{x} + e_x) \cdot (\bar{y} + e_y) - \bar{x} \cdot \bar{y} = \bar{x} e_y + \bar{y} e_x + e_x \cdot e_y \cong \bar{x} e_y + \bar{y} e_x \quad (\text{II.10})$$

c) împărțirea

$$e_{x/y} = x/y - \bar{x}/\bar{y} = (\bar{x} + e_x)/(\bar{y} + e_y) - \bar{x}/\bar{y} = \frac{\bar{x} + e_x}{\bar{y}} \cdot \frac{1}{1 + \frac{e_y}{\bar{y}}} - \frac{\bar{x}}{\bar{y}} =$$

$$= \frac{\bar{x} + e_x}{\bar{y}} \cdot \left[ 1 - \frac{e_y}{\bar{y}} + \left( \frac{e_y}{\bar{y}} \right)^2 - \dots \right] - \frac{\bar{x}}{\bar{y}} \cong \frac{e_x}{\bar{y}} - \frac{\bar{x}}{\bar{y}^2} \cdot e_y \quad (\text{II.11})$$

Relațiile deduse specifică valoarea maximă posibilă a erorii rezultatului operației respective. Cum semnul erorilor operanzilor este rareori cunoscut, rezultă că nu întotdeauna adunarea mărește eroarea, iar scăderea o micșorează așa cum s-ar putea deduce la o analiză superficială a expresiilor.

Erorile relative derivă din cele absolute astfel:

a) adunare

$$\frac{e_{x+y}}{x+y} = \frac{\bar{x}}{x+y} \cdot \frac{e_x}{x} + \frac{\bar{y}}{x+y} \cdot \frac{e_y}{y} \quad (\text{II.12})$$

b) scădere

$$\frac{e_{x-y}}{x-y} = \frac{\bar{x}}{x-y} \cdot \frac{e_x}{x} - \frac{\bar{y}}{x-y} \cdot \frac{e_y}{y} \quad (\text{II.13})$$

c) înmulțire

$$\frac{e_{x \cdot y}}{x \cdot y} = \frac{e_x}{x} + \frac{e_y}{y} \quad (\text{II.14})$$

d) împărțire

$$\frac{e_{x/y}}{x/y} = \frac{e_x}{x} - \frac{e_y}{y} \quad (\text{II.15})$$

În formulele de propagare erorile pot fi de orice tip, termenii  $x$  și  $y$  putând fi obținuți în urma unor măsurători (cu erori inițiale) sau a unor calcule anterioare (cu erori de trunchiere, tăiere sau rotunjire). Chiar și în cazul operațiilor cu numere lipsite de erori rezultatul va fi afectat de rotunjire.

**Exemplu:**  $u = (x+y) \cdot z$ , unde  $x, y, z$  nu au erori.

În urma efectuării adunării va apare o eroare de rotunjire simetrică:

$$\left| \frac{e_{x+y}}{x+y} \right| \leq \frac{1}{2} \cdot 10^{-t+1}$$

La produs se va adăuga sumei erorilor relative ale factorilor și eroarea de rotunjire  $r_m$  specifică operației.

$$\left| \frac{e_u}{u} \right| = \left| \frac{e_{x+y}}{x+y} + \frac{e_z}{z} + r_m \right| \leq \left| \frac{e_{x+y}}{x+y} \right| + |r_m| \leq \frac{1}{2} \cdot 10^{-t+1} + \frac{1}{2} \cdot 10^{-t+1} = 10^{-t+1}$$

Unde  $e_z=0$ .

Marginea erorii absolute a rezultatului va fi:

$$|e_u| \leq |u| \cdot 10^{-t+1}$$

### Reguli pentru precizia calculelor

Sucesiunea operațiilor, adoptată prin algoritmul de calcul, poate avea o influență decisivă asupra mărimii erorii rezultatului final. Astfel:

- 1) Când se adună și/sau scad mai multe numere se va începe cu acelea mai mici
- 2) Scăderea a două numere aproximativ egale se va evita prin rescrierea expresiei (dacă este posibil)
- 3) Algoritmul trebuie conceput cu un număr minim de operații aritmetice.