

LABORATOR 1 - ELEMENTE DE TEORIA ERORILOR ȘI ARITMETICA ÎN VIRGULĂ FLOTANTĂ.

1. SURSE ȘI CLASIFICAREA ERORILOR

Estimarea preciziei rezultatelor obținute reprezintă un aspect foarte important al Analizei Numerice (AN). Se disting mai multe tipuri de erori, care pot limita această precizie:

- erori în datele de intrare;
- erori de rotunjire;
- erori de aproximare.

Erorile în datele de intrare sunt în afara controlului calculelor. Aceste erori nu constituie obiectul AN. Ele se pot datora inexactității inerente a măsurătorilor fizice, introducerii eronate a datelor în calculator etc.

Erorile de rotunjire apar atunci când se operează cu numere a căror reprezentare conține un număr finit de cifre, așa cum se întâmplă, de exemplu, la reprezentarea numerelor în calculator.

Al treilea tip de erori, cele de aproximare, sunt datorate metodei numerice utilizate. În general, este necesar un număr mare de operații pentru a ajunge la soluția exactă, chiar infinit uneori, iar calculele sunt oprite în funcție de un anumit criteriu, după un număr finit de pași, atunci când se ajunge la o precizie acceptabilă (deci conștientizată).

2. ARITMETICA ÎN VIRGULĂ MOBILĂ

2.1. Numere în virgulă mobilă. Calculatorul poate lucra cu două categorii de numere: întregi și fracționare. Reprezentarea internă a acestor numere este în *virgulă fixă* și în *virgulă mobilă (flotantă)*. Pentru cazul reprezentării în virgulă fixă nu se pune problema aproximării, deoarece rezultatul adunării, scăderii sau înmulțirii numerelor întregi este tot un număr întreg, care se reprezintă în calculator dacă valoarea sa nu depășește posibilitățile de reprezentare a numerelor întregi în calculatorul respectiv. Astfel, dacă notăm cu a baza de reprezentare, iar cu t precizia mașinii (numărul de cifre în baza a care se pot reprezenta pe un cuvânt sau multiplu de cuvânt de memorie), numerele întregi reprezentabile formează mulțimea:

$$I := \{x \in \mathbb{Z} \mid -a^{t-1} \leq x \leq a^{t-1} - 1\}.$$

Pentru $a = 2$ și $t = 16$ cifre binare, $I = \{-2^{15}, \dots, 2^{15} - 1\}$. Dacă $t = 32$, domeniul de reprezentare este $I = \{-2^{31}, \dots, 2^{31} - 1\}$.

Problema erorilor de rotunjire se pune atunci când calculatorul operează cu numere fracționare, reprezentate în virgulă mobilă. Mulțimea F a numerelor cu virgulă flotantă este caracterizată de patru parametri: baza b , precizia mașinii t și

intervalul exponenților $[L, U]$. Fiecare număr x cu virgulă flotantă, care aparține lui F , are reprezentarea:

$$x = \pm \left(\frac{c_1}{b} + \frac{c_2}{b^2} + \dots + \frac{c_t}{b^t} \right) \cdot b^e,$$

unde numerele naturale c_1, \dots, c_t , numite *cifre*, satisfac $0 \leq c_i \leq b-1, i = 1, 2, \dots, t$, iar numărul întreg e , numit *exponent*, satisface $L \leq e \leq U$. Numărul

$$f = \left(\frac{c_1}{b} + \frac{c_2}{b^2} + \dots + \frac{c_t}{b^t} \right)$$

se numește *partea fracționară* sau *mantisă*. Dacă pentru $x \in F, x \neq 0$, este adevărată relația $c_1 \neq 0$, sistemul de reprezentare în virgulă mobilă se numește *normalizat*. Se remarcă faptul că mulțimea F nu este un continuu, ea fiind chiar finită, având în cazul reprezentării normalizate $2(b-1)b^{t-1}(U-L+1)+1$ elemente, care nu sunt reprezentate uniform.

Exemplul 1. Dacă F este normalizată, pentru $b = 2, t = 3, L = -1, U = 2, F$ are 33 de elemente:

$$0, \pm .100 \cdot 2^e, \pm .101 \cdot 2^e, \pm .110 \cdot 2^e, \pm .111 \cdot 2^e, \text{ cu } e \in \{-1, 0, 1, 2\}.$$



Figura 1 Elementele pozitive ale mulțimii F

În figura de mai sus s-au reprezentat numerele pozitive ale mulțimii F .

Mai remarcăm faptul că numerele nenule reprezentabile satisfac relația: $m \leq |x| \leq M$, unde $m = b^{L-1}$ se numește cel mai mic număr pozitiv reprezentabil, iar $M = b^U \left(1 - \frac{1}{b^t} \right)$ - cel mai mare număr pozitiv reprezentabil. În exemplul nostru $m = \frac{1}{4}, M = \frac{7}{2}$.

De asemenea, se observă faptul că numerele mulțimii F sunt mai "dense" în apropierea originii (puterile lui b scad) și mai "rare" spre extremitatea lui F (puterile lui b cresc).

Exemplul 2. Pentru $b = 2$ și 32 de poziții binare pentru un număr real, reprezentarea se realizează în **simplă precizie**, iar pentru 64 de poziții binare pentru un număr real, reprezentarea se realizează în **dublă precizie**.

După standardul IEEE (the "Institute of Electrical and Electronics Engineers" Inc. USA), pentru $b = 2$ rezultă $t = 24, L = -126$ și $U = 127$ pentru reprezentarea în simplă precizie și $t = 53, L = -1022$ și $U = 1023$ pentru reprezentarea în dublă precizie. Rezultă, de asemenea, pentru reprezentarea în simplă precizie, $m \approx 10^{-38}$ și $M \approx 10^{38}$, iar pentru reprezentarea în dublă precizie, $m \approx 10^{-308}$ și $M \approx 10^{308}$.

2.2. Reprezentarea aproximativă a numerelor. Scheme de rotunjire. Fie x un număr care nu depășește marginile mulțimii F ; în calculator acest număr este reprezentat de numărul cu virgulă mobilă notat $fl(x)$; în aceste condiții, spunem că x este reprezentat aproximativ (rotunjit).

Orice număr x ce satisface relația $m \leq |x| \leq M$ se poate scrie sub forma:

$$x = f \cdot b^e + g \cdot b^{e-t}.$$

În cazul fracțiilor normalizate, sunt îndeplinite relațiile:

$$\frac{1}{b} \leq |f| < 1, 0 \leq |g| < 1.$$

Există mai multe tipuri de rotunjire, și anume:

(a) *rotunjirea prin tăiere (trunchiere):*

$fl(x)$ este cel mai apropiat element din F față de x , cu proprietatea $|fl(x)| \leq |x|$.

În notația de mai sus, $fl(x) = f \cdot b^e$.

Exemplul 3. Dacă $b = 10, t = 4$ și $x = 12945.734$, putem scrie:

$$x = 0.12945734 \cdot 10^5 = 0.1294 \cdot 10^5 + 0.5734 \cdot 10.$$

În acest caz, $fl(x) = 0.1294 \cdot 10^5 = 12940 \neq x$.

(b) *rotunjirea uniformă (metoda cifrei pare):*

La rotunjirea uniformă, $fl(x)$ are următoarea expresie:

$$fl(x) = \begin{cases} f \cdot b^e, & \text{dacă } |g| < 0. \frac{b}{2} \\ f \cdot b^e \pm b^{e-t}, & \text{dacă } |g| > 0. \frac{b}{2} \\ f \cdot b^e \pm b^{e-t}, & \text{dacă } |g| = 0. \frac{b}{2}, \text{ ultima cifră } f - \text{impară} \\ f \cdot b^e, & \text{dacă } |g| = 0. \frac{b}{2}, \text{ ultima cifră } f - \text{pară} \end{cases}.$$

În această expresie, semnul "+" se consideră pentru $f > 0$ și semnul "-" se consideră pentru $f < 0$. Această modalitate de reprezentare e adoptată și de standardul IEEE.

Exemplul 4. Se consideră $b = 10, t = 4$ și o rotunjire uniformă.

(i) $x = 12944.9942 = 0.1294 \cdot 10^5 + 0.4994 \cdot 10^{5-4}$.

Se observă că $|g| < 0.5$, deci $fl(x) = 0.1294 \cdot 10^5 = 12940 \neq x$.

(ii) $x = 129551 = 0.1295 \cdot 10^6 + 0.51 \cdot 10^{6-4}$.

Se observă că $|g| > 0.5$, deci $fl(x) = 0.1295 \cdot 10^6 + 10^{6-4} = 129600 \neq x$.

(iii) $x = 1297.5 = 0.12975 \cdot 10^4 = 0.1297 \cdot 10^4 + 0.5 \cdot 10^{4-4}$.

Se observă că $|g| = 0.5$ și ultima cifră a lui f este impară, deci $fl(x) = 0.1297 \cdot 10^4 + 10^{4-4} = 1298 \neq x$.

(iv) $x = 1296.5 = 0.12965 \cdot 10^4 = 0.1296 \cdot 10^4 + 0.5 \cdot 10^{4-4}$.

Se observă că $|g| = 0.5$ și ultima cifră a lui f este pară, deci $fl(x) = 0.1296 \cdot 10^4 = 1296 \neq x$.

Având în vedere reprezentarea aproximativă a numerelor, se pot defini următoarele două tipuri de erori:

(a) *eroarea absolută*, notată e_x :

$$e_x = |x - fl(x)|;$$

(b) *eroarea relativă*, notată ε_x :

$$\varepsilon_x = \frac{e_x}{|x|} \cong \frac{e_x}{|fl(x)|} = \frac{|x - fl(x)|}{|fl(x)|}.$$

În expresia anterioară, $|x|$ s-a aproximat prin $|fl(x)|$ deoarece, în general, valoarea x nu se cunoaște în sensul că nu se reprezintă exact în calculator. Se demonstrează faptul că eroarea relativă are valoarea cea mai mare atunci când f are în modul valoarea cea mai mică, iar g valoarea cea mai mare în modul, adică $|f| = \frac{1}{b}, |g| = 1$. Avem:

$$\varepsilon_x \cong \frac{|x - fl(x)|}{|fl(x)|} = \frac{|f \cdot b^e + g \cdot b^{e-t} - fl(x)|}{|fl(x)|} \leq k \cdot \frac{1 \cdot b^{e-t}}{\left(\frac{1}{b}\right) \cdot b^e} = k \cdot b^{1-t},$$

unde b^{1-t} este o mărime specifică mașinii de calcul, ea caracterizând precizia relativă de reprezentare. Se observă imediat că la rotunjirea prin tăiere, $k = 1$.

2.3. Operații elementare în virgulă mobilă. 1. Adunarea

Oricare ar fi două numere x și y , pentru care există $fl(x)$ și $fl(y)$, numărului $x + y$ i se asociază numărul $fl(x + y)$, obținut în felul următor:

- (i) se reprezintă intern numerele x și y prin $fl(x)$ și, respectiv, $fl(y)$;
- (ii) dacă numerele au exponent diferit, atunci numărul cu exponent mai mic se aduce la o formă în care exponentul să fie egal cu cel al celuilalt termen, operație numită *denormalizare*. Acest lucru se realizează prin deplasarea mantisei spre dreapta, inserând zerouri după virgulă;
- (iii) se adună mantisele și din rezultat se păstrează t cifre;
- (iv) dacă este necesar, se normalizează rezultatul.

Observații:

1. Consecința principală a acestui mod de definire este aceea că, spre deosebire de aritmetica reală, adunarea nu este asociativă.
2. Scăderea se realizează la fel ca adunarea, cu deosebirea că mantisele se scad. De fapt, scăderea reprezintă o adunare în care scăzătorul are semn schimbat.

Exemplul 5. *Se consideră o aritmetică a virgulei mobile cu $b = 10, t = 3$, reprezentare normalizată și rotunjire prin tăiere. Fie calculul: $0.001 + 1 - 1$. Asociind primii doi termeni, se ajunge la rezultatul $fl(fl(10^{-3} + 1) - 1) = 0$, care este eronat. În cazul asocierii ultimilor doi termeni se obține $fl(10^{-3} + fl(1 - 1)) = 10^{-3}$, care este rezultatul corect.*

Exemplul 6. *Se consideră aritmetică a virgulei mobile cu $b = 10, t = 3$, reprezentare normalizată și rotunjire prin tăiere. Fie calculul: $1.001 - 1$. Se obține: $fl(1.001 - 1) = 0$, rezultat eronat.*

Exemplele 5 și 6 pun în evidență două fenomene nedorite și generatoare de erori, care pot apărea la efectuarea unei adunări în virgulă mobilă:

1. *omiterea catastrofală*: apare atunci când se adună doi termeni și valoarea absolută a unui termen este mai mică decât precizia de reprezentare a celuilalt termen; în acest caz, rezultatul este dat de termenul cu valoare absolută mai mare (situație ilustrată în Exemplul 5).
2. *neutralizarea termenilor*: apare când se adună numere de semne diferite și cu valori absolute apropiate; în acest caz, în mod eronat, rezultatul este nul (situație ilustrată în Exemplul 6).

Precizia calculului numerice este caracterizată de două mărimi constante a căror valoare este dependentă de tipul mașinii de calcul folosite. Cele două valori menționate sunt:

- *epsilonul mașină pentru adunare* (notat ε_m^+) - reprezintă cel mai mic număr real reprezentabil care schimbă, prin adunare, unitatea mașinii de calcul: $fl(1+\varepsilon_m^+) > 1$; el are valoarea b^{1-t} ;

- *epsilonul mașină pentru scădere* (notat ε_m^-) - reprezintă cel mai mic număr real reprezentabil care schimbă, prin scădere, unitatea mașinii de calcul: $fl(1-\varepsilon_m^-) < 1$; el are valoarea b^{-t} .

Exemplul 7. *Valorile celor două constante de mașină, în standardul IEEE, sunt următoarele:*

- pentru reprezentarea în simplă precizie:

$$\varepsilon_m^- = 5.96 \cdot 10^{-8}, \varepsilon_m^+ = 1.19 \cdot 10^{-7};$$

- pentru reprezentarea în dublă precizie:

$$\varepsilon_m^- = 1.11 \cdot 10^{-16}, \varepsilon_m^+ = 2.22 \cdot 10^{-16}.$$

2. Înmulțirea

Oricare ar fi două numere x și y , pentru care există $fl(x)$ și $fl(y)$, numărului $x \cdot y$ i se asociază numărul $fl(x \cdot y)$, obținut în felul următor:

- (i) se reprezintă intern numerele x și y prin $fl(x)$ și, respectiv, $fl(y)$;
- (ii) se înmulțesc fracțiile și se adună exponenții;
- (iii) din fracția rezultată se opresc t cifre;
- (iv) dacă este necesar, se normalizează rezultatul.

Observații:

1. Înmulțirea nu este asociativă.
2. Împărțirea se realizează în aceeași manieră ca și înmulțirea, cu deosebirea că la pasul (ii) mantisele se împart, iar exponenții se scad.

Exemplul 8. *Se consideră o aritmetică a virgulei mobile cu $b = 10, t = 3$, reprezentare normalizată și rotunjire prin tăiere. Fie $x = 22.547$ și $y = 0.43936$. Urmărind etapele descrise anterior, avem:*

- (i) $fl(x) = 0.225 \cdot 10^2, fl(y) = 0.439 \cdot 10^0$;
- (ii) $0.225 \cdot 0.439 = 0.098775, 10^2 \cdot 10^0 = 10^2$;
- (iii) $0.098775 \cdot 10^2 \rightarrow 0.098 \cdot 10^2$;
- (iv) $fl(x \cdot y) = 0.098 \cdot 10^2 = 9.8$.

3. PROPAGAREA ERORILOR ÎN CALCULELE NUMERICE

Având în vedere cele prezentate mai sus, putem scrie relațiile:

$$x = fl(x) \pm e_x, y = fl(y) \pm e_y, \varepsilon_x = \frac{e_x}{|fl(x)|}, \varepsilon_y = \frac{e_y}{|fl(y)|}.$$

Se numește *calcul aproximativ* un calcul efectuat într-o aritmetică a virgulei mobile. Vom accepta ca **postulat** următoarea afirmație: *eroarea relativă într-un calcul cu numere aproximative este egală cu suma dintre eroarea relativă produsă de calculul exact respectiv cu numere aproximative (T_1) și eroarea relativă produsă de calculul aproximativ cu numerele exacte corespunzătoare (T_2).*

Fie $*$ una din operațiile descrise anterior. Postulatul de mai sus, descris de relația:

$$\varepsilon_{x*y} = T_1 + T_2,$$

permite determinarea modalității de propagare a erorilor relative pentru operațiile în virgulă mobilă, după cum urmează.

Pentru *calculul exact cu numere aproximative* (T_1) avem:

$$\begin{aligned} e_x &\neq 0, e_y \neq 0, fl(x * y) = fl(x) * fl(y), \\ T_1 &= k_1 \cdot \varepsilon_x + k_2 \cdot \varepsilon_y. \end{aligned}$$

Constantele k_1 și k_2 pot avea valori diferite, după cum urmează:

- pentru adunare:

$$k_1 = \frac{|fl(x)|}{fl(x) + fl(y)}, k_2 = \frac{|fl(y)|}{fl(x) + fl(y)};$$

- pentru scădere:

$$k_1 = \frac{|fl(x)|}{fl(x) - fl(y)}, k_2 = -\frac{|fl(y)|}{fl(x) - fl(y)};$$

- pentru înmulțire:

$$k_1 = k_2 = 1;$$

- pentru împărțire:

$$k_1 = 1, k_2 = -1.$$

Pentru *calculul aproximativ cu numere exacte* (T_2) avem:

$$e_x = 0, e_y = 0, fl(x * y) \neq fl(x) * fl(y)$$

și atunci:

$$T_2 \leq k \cdot b^{1-t},$$

ca pentru orice număr real care nu se reprezintă exact.

Exemplul 9. *Se consideră următorul polinom:*

$$y(x) = x^7 - 7x^6 + 21x^5 - 35x^4 + 35x^3 - 21x^2 + 7x - 1,$$

pentru $x \in [0.998, 1.012]$, cu pasul de evaluare 0.0001. Figura de mai jos ilustrează cu linie punctată rezultatul calculelor în dublă precizie (standardul IEEE) utilizând formula de mai sus, iar cu linie continuă sunt reprezentate cele corecte. Acestea din urmă sunt obținute evaluând polinomul $y(x) = (x - 1)^7$.

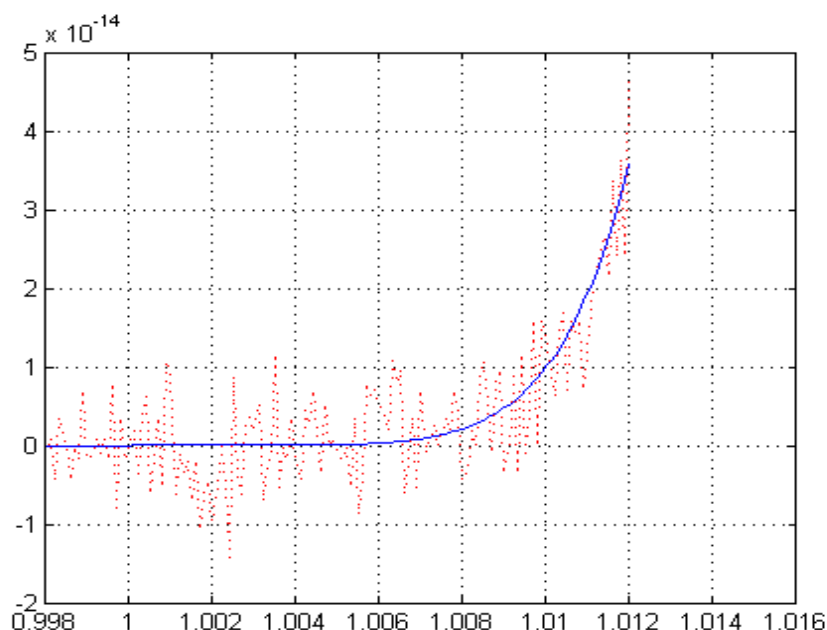


Figura 2 Ilustrarea propagării erorilor într-un calcul numeric

Analizând **Figura 2**, se remarcă faptul că valorile calculate sunt de ordinul 10^{-14} , ceea ce evidențiază erori relativ mici între rezultatele celor două maniere de calcul. Aceste diferențe se explică prin fenomenele de neutralizarea termenilor și omitere catastrofală.

Cele expuse până acum demonstrează faptul că într-un calcul numeric erorile se propagă de la o operație la alta. Pe măsură ce numărul operațiilor dintr-un calcul crește, pot apărea situații în care erorile se acumulează excesiv de mult, fapt care determină obținerea unei valori total incorecte a rezultatului final. Ca urmare, la întocmirea unui algoritm de calcul și apoi la implementarea acestuia, utilizatorul trebuie să se asigure că soluțiile nu vor fi afectate de erori care să depășească anumite limite admisibile.

Ca urmare, se pot formula următoarele reguli generale pentru mărirea preciziei calculelor:

1. când se adună sau se scad numere, este recomandabil să se înceapă cu cele mai mici în valoare absolută, separat pentru cele pozitive și separat pentru cele negative;
2. dacă este posibil, este recomandabil să se evite scăderea a două numere aproximativ egale; o expresie care conține o astfel de scădere poate fi rescrisă;
3. o expresie de forma $(a - b) \cdot c$ poate fi rescrisă sub forma $a \cdot c - a \cdot c$, iar o expresie de forma $\frac{a - b}{c}$ poate fi rescrisă sub forma $\frac{a}{c} - \frac{b}{c}$; dacă numerele a și b sunt aproximativ egale, este recomandabil să se efectueze mai întâi scăderea și apoi înmulțirea și împărțirea.
4. dacă regulile generale enunțate mai sus nu se pot aplica, se va urmări minimizarea numărului de operații aritmetice implicate.

Exemplul 10. *Se consideră următoarele relații de calcul:*

$$h = \frac{1}{2}, x = \frac{2}{3} - h, y = \frac{3}{5} - h, e = x + x + x - h,$$

$$f = y + y + y + y + y - h, q = \frac{f}{e}.$$

Efectuând calculele manual, se obțin rezultatele:

$$h = \frac{1}{2}, x = \frac{1}{6}, y = \frac{1}{10}, e = 0, f = 0, q = ?$$

Implementând aceste relații de calcul într-un program scris în MATLAB obținem rezultatele:

$$e \neq 0, f \neq 0, q = \text{valoare finită}.$$

Explicația acestor rezultate constă în următoarele:

- numărul $\frac{2}{3} = 0.(6)$ are un număr infinit de cifre în baza de numerație zece și deci nu va fi reprezentat exact;

- $\frac{3}{5} = 0.6$ are un număr infinit de cifre în baza de numerație doi și deci nici acesta nu va fi reprezentat exact: $0.6_{[10]} = 0.(1001)_{[2]}$.

Aceste erori se vor propaga în calculul valorilor e și f rezultând valori de ordinul epsilonului mașină.