

Laboratorul de Metode Numerice

Universitatea "POLITEHNICA" București, Facultatea de Electrotehnică, <http://lmn.pub.ro>, lmn@lmn.pub.ro

Laborator - metode numerice în ingineria electrică

L2 - Erori în rezolvarea problemelor numerice

Autori:

Conf. Gabriela Ciuprina

Prof. Daniel Ioan

Prep. Marius Piper

As. Marius Rădulescu

As. Mihai Popescu

12 martie 2004

Cuprins

2	Erori în rezolvarea problemelor numerice	1
2.1	Caracterizarea lucrării	1
2.2	Principiul lucrării	1
2.2.1	Erori de rotunjire	2
2.2.2	Erori inerente	3
2.2.3	Erori de trunchiere	6
2.3	Chestiuni de studiat	8
2.4	Modul de lucru	8
2.4.1	Determinarea erorii relative de rotunjire a sistemului de calcul . . .	8
2.4.2	Analiza propagării erorilor inerente	8
2.4.3	Analiza erorii de trunchiere	9
2.4.4	Implementarea unor algoritmi cu controlul erorii	11
2.5	Întrebări și probleme	11

Capitolul 2

Erori în rezolvarea problemelor numerice

2.1 Caracterizarea lucrării

În majoritatea cazurilor, algoritmi cu caracter numeric, după implementarea lor pe un sistem de calcul determină nu soluția exactă a problemei ci o aproximare numerică a acesteia.

Scopul acestei lucrări este de a evidenția modul în care pot fi caracterizate erorile numerice, motivele apariției acestora și modul în care acestea se propagă. Se studiază *erorile inerente* (datorate datelor de intrare), *erorile de rotunjire* (datorate reprezentării finite, aproximative a numerelor reale) și *erorile de trunchiere* (datorate aproximării finite a unor procese teoretic infinite).

2.2 Principiul lucrării

Pentru a caracteriza abaterea unei variabile numerice de la valoarea sa exactă se poate folosi *eroarea absolută*, definită prin:

$$e_x = x - x^*, \quad (2.1)$$

în care x este valoarea exactă a variabilei, iar x^* este valoarea sa aproximativă. Deoarece în majoritatea cazurilor valoarea exactă nu este cunoscută se preferă utilizarea unei *marginii superioare a erorii absolute*:

$$a_x \geq |e_x| = |x - x^*|. \quad (2.2)$$

În acest caz, în locul valorii exacte x se operează cu intervalul:

$$x^* - a_x \leq x \leq x^* + a_x. \quad (2.3)$$

O altă modalitate de caracterizare a abaterii unei variabile de la valoarea exactă este eroarea relativă:

$$\varepsilon_x = \frac{e_x}{x} = \frac{x - x^*}{x}, \quad (2.4)$$

sau marginea superioară a acesteia:

$$\Gamma_x = \frac{a_x}{|x|} \geq |\text{eps}_x| = \frac{|x - x^*|}{|x|}. \quad (2.5)$$

În majoritatea calculelor tehnice, mărimile rezultate în urma măsurării sunt cunoscute cu 3 maxim 6 cifre semnificative exacte, ceea ce corespunde unor erori relative cuprinse în intervalul $10^{-2} \div 10^{-5}$.

2.2.1 Erori de rotunjire

Una din cauzele frecvente de eroare în calculele cu numere reale se datorează reprezentării finite a acestor numere în sistemele de calcul.

Modul uzual de reprezentare a unui număr real într-un sistem de calcul este de forma:

$$x^* = \pm 0, kkk \dots k \cdot 10^{\pm kk} = m \cdot 10^e, \quad (2.6)$$

unde k este cifră ($0, 1, \dots, 9$), m se numește mantisa, iar e exponent. Cu excepția numărului nul, la care $m = 0$, în rest mantisa satisface relația:

$$0.1 \leq |m| \leq 1. \quad (2.7)$$

Presupunând că mantisa este reprezentată cu n cifre, rezultă că prin această reprezentare cifrele "l" din reprezentarea exactă:

$$x = \pm 0, \underbrace{kk \dots k}_{n} ll \dots \cdot 10^{\pm kk} \quad (2.8)$$

sunt pierdute prin rotunjire.

În consecință, eroarea relativă de rotunjire

$$\varepsilon_x = \frac{|x - x^*|}{|x|} = \frac{0, 0 \dots 0 ll \dots \cdot 10^{\pm kk}}{0, kk \dots ll \dots \cdot 10^{\pm kk}} = \frac{0, ll \dots}{0, kk \dots} 10^{-n} = 10^{-n+1} \quad (2.9)$$

depinde de numărul de cifre semnificative folosite în reprezentarea numărului real și nu de valoarea numărului. Această eroare relativă de rotunjire, specifică sistemului de calcul (calculator + mediu de programare) este cel mai mare număr real, care adăugat la unitate nu "modifică" valoarea acesteia. Ordinul de mărime al erorii relative de rotunjire în sistemele uzuale de calcul este $10^{-5} \div 10^{-20}$ și poate fi determinat pe fiecare sistem de calcul cu următorul algoritm:

```
; calculează eroarea relativă de rotunjire
real err
err=1
repetă
    err=err/2
până când (1 + err = 1)
scrie err
```

Eroarea relativă de rotunjire *err* este cunoscută sub numele de *zeroul mașinii* și nu trebuie confundată cu cel mai mic număr pozitiv, nenul, reprezentabil în calculator.

2.2.2 Erori inerente

Datele de intrare folosite în rezolvarea unei probleme tehnice provin în multe cazuri din determinări experimentale. Rezultatul oricărei măsurători este susceptibil de erori și chiar dacă datele de intrare sunt cunoscute cu maximă precizie ele pot fi afectate de erori de rotunjire.

Erorile datelor de intrare într-un algoritm se numesc erori inerente. Aceste erori se propagă în procesul de calcul și afectează în final soluția problemei. Chiar dacă algoritmul de calcul nu poate fi făcut responsabil de prezența erorilor în datele de intrare, el poate influența precizia soluției. Un algoritm la care eroarea relativă a soluției nu depășește erorile relative ale datelor de intrare este un *algoritm stabil din punct de vedere numeric*. În schimb, dacă erorile relative ale soluției sunt mult mai mari decât cele ale datelor de intrare se spune că algoritmul de rezolvare prezintă instabilități numerice. În acest caz este posibil ca abateri foarte mici ale datelor de intrare să determine abateri mari ale soluției numerice față de cea exactă, și să facă soluția numerică inutilizabilă. Dacă se notează cu y soluția problemei, iar cu x_1, x_2, \dots, x_n datele acesteia, procedeul de calcul va consta în fond în evaluarea funcției $y = f(x_1, x_2, \dots, x_n)$

Considerând erorile absolute suficient de mici, acestea se pot aproxima cu diferențiala:

$$\begin{aligned} dy &= \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \cdots + \frac{\partial f}{\partial x_n} dx_n, \\ e_y &= \frac{\partial f}{\partial x_1} e_{x_1} + \frac{\partial f}{\partial x_2} e_{x_2} + \cdots + \frac{\partial f}{\partial x_n} e_{x_n}. \end{aligned} \quad (2.10)$$

Erorile relative satisfac relația:

$$y \cdot \varepsilon_y = \frac{\partial f}{\partial x_1} x_1 \cdot \varepsilon_{x_1} + \frac{\partial f}{\partial x_2} x_2 \cdot \varepsilon_{x_2} + \cdots + \frac{\partial f}{\partial x_n} x_n \cdot \varepsilon_{x_n}. \quad (2.11)$$

Deoarece în sistemele numerice de calcul, cele mai complicate proceduri se reduc în final la operații aritmetice, va fi analizată propagarea erorilor în cazul celor patru operații aritmetice elementare:

Adunarea:

$$\begin{aligned} y &= f(x_1, x_2) = x_1 + x_2, \\ e_y &= e_{x_1} + e_{x_2}, \\ \varepsilon_y &= \frac{x_1}{x_1 + x_2} \varepsilon_{x_1} + \frac{x_2}{x_1 + x_2} \varepsilon_{x_2}. \end{aligned} \quad (2.12)$$

Dacă cei doi termeni x_1, x_2 au același semn, adunarea este o operație stabilă numeric, deoarece

$$\left| \frac{x_1}{x_1 + x_2} \right| \leq 1, \quad \left| \frac{x_2}{x_1 + x_2} \right| \leq 1, \quad |\varepsilon_y| \leq |\varepsilon_{x_1}| + |\varepsilon_{x_2}|.$$

Scăderea:

$$\begin{aligned} y &= f(x_1, x_2) = x_1 - x_2, \\ e_y &= e_{x_1} - e_{x_2}, \\ \varepsilon_y &= \frac{x_1}{x_1 - x_2} \varepsilon_{x_1} - \frac{x_2}{x_1 - x_2} \varepsilon_{x_2}. \end{aligned} \quad (2.13)$$

Dacă x_1 și x_2 au același semn, operația de scădere este instabilă numeric. Eroarea relativă a diferenței poate depăși cu multe ordine de mărime eroarea termenilor, prin fenomenul numit de anulare prin scădere, ca în exemplul:

$$x_1 = 0,12345 \pm 1\% ; x_2 = 0,12344 \pm 1\% ; y = x_1 - x_2 = 0,00001 \pm 3 \cdot 10^4\%.$$

Înmulțirea:

$$\begin{aligned} y &= f(x_1, x_2) = x_1 x_2 \\ e_y &= x_2 e_{x_1} + x_1 e_{x_2} \\ \varepsilon_y &= \varepsilon_{x_1} + \varepsilon_{x_2} \end{aligned} \quad (2.14)$$

Înmulțirea este o operație stabilă din punct de vedere numeric deoarece

$$|\varepsilon_y| \leq |\varepsilon_{x_1}| + |\varepsilon_{x_2}|.$$

Împărțirea:

$$\begin{aligned} y &= f(x_1, x_2) = \frac{x_1}{x_2}, \\ e_y &= \frac{1}{x_2}e_{x_1} - \frac{x_1}{(x_2)^2}e_{x_2}, \\ \varepsilon_y &= \varepsilon_{x_1} - \varepsilon_{x_2}. \end{aligned} \tag{2.15}$$

Și împărțirea este o operație stabilă din punct de vedere numeric deoarece

$$|\varepsilon_y| \leq |\varepsilon_{x_1}| + |\varepsilon_{x_2}|.$$

În aceste evaluări ale erorilor s-a considerat că toate calculele se efectuează exact. În realitate, rezultatul este rotunjit (la numărul de cifre semnificative specifice sistemului de calcul) și eroarea relativă a rezultatului este majorată cu eroarea relativă de rotunjire.

În cazul în care prezintă importanță deosebită evaluarea erorii unui calcul numeric se poate utiliza următorul tip abstract de date care permite controlul propagării erorii. În acest tip de date numerele reale x sunt reprezentate prin valoarea x și marginea erorii relative Γ_x , deci prin intervalul:

$$x^*(1 - \Gamma_x) \leq x \leq x^*(1 + \Gamma_x). \tag{2.16}$$

real tip inregistrare interval

real val

;valoare

real er

;marginea erorii relative

procedura suma (x,y,s)

;calculează $s = x + y$

interval x, y, s

$s.val = x.val + y.val$

$s.er = err + (|x.val \cdot x.er| + |y.val \cdot y.er|)/|s.val|$

retur

procedura dif(x,y,d)

;calculează $d = x - y$

interval x, y, d

$d.val = x.val - y.val$

$d.er = err + (|x.val \cdot x.er| + |y.val \cdot y.er|)/|d.val|$

retur

procedura prod(x,y,p)

;calculează $p = x * y$


```

interval  $x, y, p$ 
 $p.val = x.val \cdot y.val$ 
 $p.er = err + |x.er| + |y.er|$ 
retur
procedura rap( $x, y, c$ ) ; calculează  $c = x / y$ 
interval  $x, y, c$ 
 $c.val = x.val / y.val$ 
 $c.er = err + |x.er| + |y.er|$ 
retur

```

2.2.3 Erori de trunchiere

În esență, metodele numerice constau în reducerea rezolvării unei probleme complicate la un număr finit de etape simple, care în fond sunt operații aritmetice elementare. Datorită caracterului finit al oricărui algoritm, încercarea de a rezolva probleme de analiză matematică, ce presupun teoretic o infinitate de pași (cum sunt de exemplu calculele limitelor) determină apariția unor erori de metodă numite erori de trunchiere.

Calculul numeric al limitelor de șiruri, serii sau funcții (inclusiv calculul numeric al derivatelor și integralelor, care se reduc în fond tot la calculul limitelor) presupune trunchierea unui proces numeric infinit. Pentru a evidenția eroarea de trunchiere se consideră un șir numeric $x_k = 1, 2, \dots, n$ convergent către limita $x = \lim_{k \rightarrow \infty} x_k$. Conform definiției convergenței, pentru orice ε există un n astfel încât $|x - x_k| \leq \varepsilon$, pentru orice $k > n$. În consecință, pentru o eroare ε impusă există un rang finit n , astfel încât x_n reprezintă o aproximare satisfăcătoare pentru limita x . Rezultă că, după evaluarea unui număr n suficient de mare de termeni, ultimul poate fi adoptat ca soluție numerică $x^* = x_n$, cu o eroare de trunchiere dependentă de acest n . În general, eroarea de trunchiere este cu atât mai mică cu cât numărul de termeni calculați este mai mare. Modul în care eroarea de trunchiere depinde de n caracterizează viteza de convergență a algoritmului. Se spune că un algoritm are *viteza de convergență* de ordinul p , dacă există o constantă $c > 0$, astfel încât eroarea de trunchiere satisface inegalitatea:

$$|\varepsilon_n| = |x - x_n| \leq \frac{c}{n^p}, \quad (2.17)$$

pentru orice $n > 1$. Se constată că ordinul vitezei de convergență a unui algoritm este dat de panta dreptei ce mărginește superior graficul funcției $|\varepsilon| = f(n)$ în scară dublu logaritmică.

O categorie importantă de procese numerice infinite o reprezintă seriile Taylor. Acestea, fiind serii de puteri, pot fi utilizate la evaluarea funcțiilor elementare (sinus, cosinus, exponențială, logaritm, etc.) și a celor speciale (Bessel, integrale eliptice) prin reducerea la operații aritmetice elementare (adunări și înmulțiri).

Pentru a exemplifica modul în care poate fi sumată numeric o serie cu controlul erorii de trunchiere se consideră dezvoltarea în serie Taylor a funcției $y = \sin x$:

$$y = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \quad (2.18)$$

Termenul general al seriei poate fi calculat recursiv:

$$T_k = -\frac{(-1)^k}{(2k-1)!} x^{2k-1} = -T_{k-1} \frac{x^2}{(2k-1)(2k-2)}. \quad (2.19)$$

Prin trunchierea seriei la rangul n se obține

$$y^* = \sum_{k=1}^n \frac{(-1)^{k+1}}{(2k-1)!} \cdot x^{2k-1}. \quad (2.20)$$

Deoarece seria este alternantă, eroarea de trunchiere este majorată de modulul primului termen neglijat:

$$|\varepsilon_y| = |y - y^*| \leq \frac{|x|^{2n+1}}{(2n+1)!}. \quad (2.21)$$

Rezultă următorul pseudocod pentru evaluarea funcției sinus cu o eroare de trunchiere mai mică decât eroarea de rotunjire.

```

funcția sinus( $x$ )
  real  $x, t, s$ 
  întreg  $k$ 
   $t = x$ 
   $s = t$ 
   $k = 1$ 
  repetă
     $k = k + 2$ 
     $t = -tx^2/k/(k-1)$ 
     $s = s + t$ 
  până când  $|t| < \text{err}$ 
întoarce  $s$ 

```

Efortul de calcul pentru evaluarea funcției crește odată cu scăderea erorii impuse, dar și cu creșterea modulului variabilei x . Deoarece funcția $\sin(x)$ este periodică, algoritmul poate fi îmbunătățit prin reducerea variabilei la primul cerc (sau chiar la primul cadran).

2.3 Chestiuni de studiat

1. Determinarea erorii relative de rotunjire a sistemului de calcul;
2. Analiza propagării erorilor inerente;
3. Analiza erorii de trunchiere;
4. Implementarea unor algoritmi numerici care controlează eroarea.

2.4 Modul de lucru

Pentru desfășurarea lucrării lansați lucrarea *Erori în rezolvarea problemelor numerice* din meniul general. Aceasta are ca urmare lansarea următorului meniu:

- Eroarea de rotunjire
- Calcule cu controlul erorii
- Erori de trunchiere

2.4.1 Determinarea erorii relative de rotunjire a sistemului de calcul

Se selectează opțiunea *Eroarea de rotunjire* care are ca efect calculul și afișarea erorii relative de rotunjire specifice sistemului.

2.4.2 Analiza propagării erorilor inerente

Se selectează opțiunea *Calcule cu controlul erorii* din meniul principal, care are ca efect lansarea unui program ce emulează un calculator de buzunar care folosește notația poloneză inversă (operand, operand, operator). Acest calculator nu operează cu numere reale ci cu intervale (valoare numerică și margine superioară a erorii relative).

Programul permite introducerea fie a operanzilor (valoare și eroare relativă) fie a operatorilor unari (I - invers, O - opus, V - radical) sau binari (+ - * /).

Operanzii sunt introduși într-o stivă vizualizată pe ecran de 5 coloane (număr de ordine, valoare numerică, eroare relativă, limita minimă și limita maximă a valorii exacte).

Operatorii unari acționează asupra ultimului element introdus în stivă, înlocuindu-l cu valoarea obținută în urma aplicării operatorului (stiva nu își modifică dimensiunea).

Operatori binari folosesc ca operanzi ultimele două elemente din stivă, iar rezultatul înlocuiește aceste elemente (stiva se micșorează cu o unitate).

Pentru a facilita manipularea stivei sunt disponibili încă doi operatori unari, cu caracter nenumeric (R - repetă vârful stivei și E - extrage element din stivă), care modifică (incrementează respectiv decrementează) înălțimea stivei.

Se propune rezolvarea ecuației de gradul doi:

$$ax^2 + bx + c = 0,$$

în care $a = 1$, $b = -100.01$, $c = 1$ (cu rădăcinile $x_1 = 100$; $x_2 = 0,01$). Valorile a , b , c se vor introduce cu erorile relative de $0.1\% = 0.001$. Se vor nota valorile intermediare:

$$\begin{aligned}\Delta &= b^2 - 4ac, \\ x_1 &= \frac{-b + \sqrt{\Delta}}{2a}, \\ x_2 &= \frac{-b - \sqrt{\Delta}}{2a},\end{aligned}$$

și erorile relative asociate.

Pentru a evita erorile datorate anulării prin scădere, se va calcula a doua rădăcină pe baza relației $x_1 x_2 = c/a$

$$x_2 = \frac{c}{ax_1}.$$

Se vor compara erorile relative și absolute ale rezultatelor obținute pe cele două căi.

2.4.3 Analiza erorii de trunchiere

Se selectează opțiunea *Erori de trunchiere* din meniul principal, care are ca efect lansarea unui program care realizează sumarea următoarelor serii:

- Taylor pentru funcția exponențială
- Taylor pentru funcția sinus
- Taylor pentru funcția logaritm natural
- Fourier pentru funcția crenel

- armonica alternantă
- armonica alternantă cu convergență îmbunătățită prin metoda Euler

După selectarea seriei și alegerea valorii variabilei independente x , programul calculează și afișează suma parțială a seriei și eroarea absolută de trunchiere pentru diferite ordine, apoi reprezintă grafic variația erorii în funcție de numărul de termeni sumați. Pentru a putea facilita comparațiile se pot reprezenta și mai multe grafice de erori simultan.

Se recomandă observarea convergenței

- seriilor Taylor pentru evaluarea expresiilor:
 - $\exp(0)$, $\exp(1)$, $\exp(5)$, $\exp(-10)$, $\exp(10)$,
 - $\sin(0)$, $\sin(1)$, $\sin(10)$,
 - $\ln(1)$, $\ln(0.1)$, $\ln(1.1)$.

Reamintim că dezvoltările în serie Taylor sunt:

$$\begin{aligned} e^x &= 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots \\ \sin(x) &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \\ \ln(x) &= (x-1) - \frac{(x-1)^2}{2} + \frac{(x-1)^3}{3} - \frac{(x-1)^4}{4} + \dots \end{aligned}$$

- seriei Fourier pentru evaluarea expresiilor:
 - $\text{crenel}(1)$, $\text{crenel}(0.1)$.

Funcția crenel este definită astfel

$$\text{crenel}(x) = \begin{cases} \frac{\pi}{4} & \text{dacă } \sin(x) > 0 \\ -\frac{\pi}{4} & \text{dacă } \sin(x) \leq 0 \end{cases}$$

iar dezvoltarea ei în serie Fourier conduce la

$$\text{crenel}(x) = \sin(x) + \frac{\sin(3x)}{3} + \frac{\sin(5x)}{5} + \dots$$

- armonicii alternante

$$A = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{1}{k}$$

- armonicii alternante cu convergență îmbunătățită

$$A = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots = 0.5 + \sum_{k=1}^{\infty} \frac{\left(\frac{1}{2k-1} - \frac{1}{2k+1}\right)}{4k}.$$

Se vor comenta rezultatele și se va aprecia viteza de convergență.

Indicație: Pentru estimarea vitezei de convergență se calculează panta dreptei ce mărginește graficul erorii, în scara dublu logaritmică.

2.4.4 Implementarea unor algoritmi cu controlul erorii

Se va genera pseudocodul funcției cosinus, determinată prin serie Taylor trunchiată până la atingerea unei erori de trunchiere impuse. Se va implementa acest algoritm într-un limbaj de programare sub forma unei funcții $\text{cosinus}(x, ert)$, în care x este variabilă independentă iar ert este eroarea de trunchiere impusă. Se va implementa algoritmul de determinare al erorii relative de rotunjire err , iar în final se va implementa algoritmul de evaluare a funcției cosinus cu $ert = err$ și cu controlul propagării erorii, sub forma unei proceduri care să calculeze $\cos(x)$ și eroarea relativă asociată.

2.5 Întrebări și probleme

1. Cum poate fi caracterizată eroarea unei variabile vectoriale x ?
2. Implementați algoritmul de determinare a erorii relative de rotunjire în diferite limbaje de programare (în simplă și dublă precizie) și apoi comparați rezultatele.
3. Ce modificări trebuie aduse algoritmului de determinare a erorii relative de rotunjire pentru ca acesta să determine nu numai ordinul de mărime al erorii ci și valoarea sa exactă?
4. Definiți și implementați un tip abstract de date, care să reprezinte fiecare număr real ca o pereche de numere reale ce îl încadrează (valoare maximă, valoare minimă).
5. Este adunarea numerelor reale rotunjite o operație asociativă? Pentru a aduna mai multe numere reale diferite cu eroare minimă ele trebuie sortate în ordine crescătoare sau descrescătoare?
6. Pentru a micșora erorile de rotunjire la sumarea unei serii cum trebuie adunați termenii, în ordine crescătoare sau descrescătoare?

7. Folosiți opțiunea *Calcule cu controlul erorii* pentru a determina rezistența de șunt ce trebuie folosită la transformarea unui galvanometru de 1mA (în clasa 2%) a cărui rezistență $R = 57.5\Omega$ a fost determinată cu precizie de 0,5(%) într-un ampermetru de 1A.
8. Generați un algoritm pentru rezolvarea unei ecuații de gradul doi, care evită fenomenul de anulare prin scădere.
9. Generați și implementați un algoritm, pentru evaluarea funcțiilor Bessel.