

1.3. Erori în calculele numerice

Prof.dr.ing. Gabriela Ciuprina

Universitatea "Politehnica" București, Facultatea de Inginerie Electrică,
Departamentul de Electrotehnică

Suport didactic pentru disciplina *Metode numerice*, 2017-2018

Cuprins

- 1 Caracterizarea cantitativă a erorilor
 - În mod absolut
 - În mod relativ
- 2 Tipuri de erori
- 3 Analiza erorilor
 - Analiza erorilor de rotunjire
 - Analiza erorilor de trunchiere
 - Analiza erorilor inerente
- 4 Condiționare și stabilitate
 - Condiționare
 - Stabilitate

Eroarea absolută și marginea ei

Fie:

$\mathbf{x} \in \mathbb{R}^n$ - valoarea exactă a unei mărimi;

$\bar{\mathbf{x}}$ - valoarea aproximativă.

Eroarea absolută $\mathbf{e}_x \in \mathbb{R}^n$:

$$\mathbf{e}_x = \bar{\mathbf{x}} - \mathbf{x}. \quad (1)$$

Marginea erorii absolute $a_x \in \mathbb{R}$:

$$\|\mathbf{e}_x\| \leq a_x. \quad (2)$$

Dacă $n = 1$ rezultă

$$\bar{x} - a_x \leq x \leq \bar{x} + a_x. \quad (3)$$

Echivalentă cu: $x \in [\bar{x} - a_x, \bar{x} + a_x]$.

Scrisă pe scurt ca:

$$"x = \bar{x} \pm a_x". \quad (4)$$

Eroarea relativă și marginea ei

Eroarea relativă $\varepsilon_{\mathbf{x}} \in \mathbb{R}^n$:

$$\varepsilon_{\mathbf{x}} = \frac{\mathbf{e}_{\mathbf{x}}}{\|\mathbf{x}\|}. \quad (5)$$

Marginea erorii relative $r_{\mathbf{x}} \in \mathbb{R}$

$$\|\varepsilon_{\mathbf{x}}\| \leq r_{\mathbf{x}}. \quad (6)$$

Cel mai adesea, $r_{\mathbf{x}}$ se exprimă în procente.
Scriere pe scurt:

$$\mathbf{x} = \bar{\mathbf{x}} \pm r_{\mathbf{x}}\%. \quad (7)$$

Exemplu: π

$$x = 3.1415\dots$$

$$\bar{x} = 3.14$$

$$e_x = -0.0015\dots$$

$$a_x = 0.0016$$

$$\varepsilon_x = -0.0015\dots / 3.1415\dots$$

$$r_x = 0.0016/3 \leq 0.0006 = 0.06\%.$$

$$\pi = 3.14 \pm 0.0016 \quad \text{sau} \quad \pi = 3.14 \pm 0.06\%.$$

Concluzii

Relația " $\mathbf{x} = \bar{\mathbf{x}} \pm a_x$ "

unde $\mathbf{x}, \bar{\mathbf{x}} \in \mathbb{R}^n$ și $a_x \in \mathbb{R}$ se interpretează astfel:

$$(\exists) \mathbf{e}_x \in \mathbb{R}^n, \|\mathbf{e}_x\| \leq a_x, \text{ astfel încât } \bar{\mathbf{x}} = \mathbf{x} + \mathbf{e}_x, \quad (8)$$

Relația " $\mathbf{x} = \bar{\mathbf{x}} \pm r_x \%$ "

unde $\mathbf{x}, \bar{\mathbf{x}} \in \mathbb{R}^n$, $r_x \% = 100r_x$ și $r_x \in \mathbb{R}$ se interpretează astfel:

$$(\exists) \varepsilon_x \in \mathbb{R}^n, \|\varepsilon_x\| \leq r_x, \text{ astfel încât } \bar{\mathbf{x}} = \mathbf{x} + \|\mathbf{x}\|\varepsilon_x. \quad (9)$$

În cazul unei mărimi scalare ($n = 1$), relația (9) se scrie

$$\bar{x} = x(1 \pm \varepsilon_x), \quad (10)$$

semnul plus corespunzând unei valori x pozitive, iar semnul minus uneia negative.

Tipuri de erori

În funcție de tipul cauzelor care le generează:

- 1 **Erori de rotunjire** - datorate reprezentării finite a numerelor reale;
- 2 **Erori de trunchiere** - datorate reprezentării finite a algoritmului;
- 3 **Erori inerente** - datorate reprezentării imprecise a datelor de intrare.

Cifre semnificative

Reprezentarea unui număr real în baza 10:

$$\bar{x} = f \cdot 10^n. \quad (11)$$

unde $0.1 \leq |f| < 1$.

Cifrele părții fracționare se numesc **cifre semnificative**.

Exemple:

$$3.14 = 0.314 \cdot 10^1$$

$$-0.007856 = -0.7856 \cdot 10^{-2}.$$

Rotunjirea afectează reprezentarea numerelor reale

$$\bar{x} = 0.\overbrace{***\dots*}^f \cdot 10^n, \quad (12)$$

k cifre

$$x = 0.\underbrace{***\dots*}_{k \text{ cifre}} \#\#\#\dots \cdot 10^n, \quad (13)$$

$$e_x = \bar{x} - x = -0.\underbrace{000\dots0}_{k \text{ cifre}} \#\#\#\dots \cdot 10^n = -0.\#\#\#\dots \cdot 10^{n-k}, \quad (14)$$

Rotunjirea afectează reprezentarea numerelor reale

$$\varepsilon_x = \frac{e_x}{x} = \frac{-0.\#\#\#\dots 10^{n-k}}{0.\underbrace{*\#\#\#\dots*}_{k \text{ cifre}} \#\#\#\dots 10^n} = -\frac{0.\#\#\#\dots}{0.*\#\#\#\dots} 10^{-k} \quad (15)$$

$$|\varepsilon_x| \leq \frac{1}{0.1} 10^{-k} = 10^{-k+1}. \quad (16)$$

Marginea erorii relative de rotunjire a unui sistem de calcul depinde doar de numărul de cifre semnificative ce pot fi memorate. Pentru un sistem de calcul ce lucrează cu k cifre semnificative, marginea erorii relative de rotunjire este 10^{-k+1} .

Rotunjirea afectează calculele

Adunarea a două numere reale

Intuitiv: pp. $k = 3$, $x_1 + x_2 = ?$

$$x_1 = 3.73 = 0.373 \cdot 10^1$$

$$x_2 = 0.006 = 6 \cdot 10^{-3}$$

$$x_2 = 6 \cdot 10^{-4} \cdot 10^1 = 0.0006 \cdot 10^1 = 0.000 \cdot 10^1$$

Rezultat: $x_1 + x_2 = x_1$.

Zeroul mașinii

Zeroul (acuratețea, precizia, "epsilon-ul") mașinii = cel mai mic eps pentru care $1 + \text{eps} > 1$.

- $(\forall)a < \text{eps}, 1 + a = 1$ (în calculator)
- în mod uzual $\text{eps} = 2.22 \cdot 10^{-16}$.
- Matlab: `eps`
- Scilab `%eps`.
- Zeroul mașinii nu trebuie confundat cu cel mai mic număr reprezentabil în calculator și care, în mod uzual are valoarea $2.23 \cdot 10^{-308}$.

Consecință: adunarea numerelor reale în calculator nu este asociativă.

Determinarea eps într-un mediu de programare

funcție zeroul_mașinii ()

real eps

eps = 1

cât timp (1 + eps > 1)

eps = eps/2

•

eps = eps*2

întoarce eps

Exemplu

$$f(x) = f(x_0) + \frac{x - x_0}{1!} f'(x_0) + \frac{(x - x_0)^2}{2!} f''(x_0) + \dots \quad (17)$$

sinus, $x_0 = 0$:

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}. \quad (18)$$

$$\bar{s} = \sum_{k=0}^n (-1)^k \frac{x^{2k+1}}{(2k+1)!}. \quad (19)$$

$$|e_s| = |\bar{s} - s| \leq \frac{x^{2n+1}}{(2n+1)!}. \quad (20)$$

Algoritm cu controlul erorii de trunchiere

funcție sinus(x, e)

; întoarce valoarea funcției sinus în punctul x

; prin trunchierea seriei Taylor dezvoltată în 0

real x ; punctul în care se va evalua funcția sin

real e ; eroarea de trunchiere impusă

real t, s

întreg k

$t = x$

$s = t$

$k = 0$

cât timp ($|t| > e$)

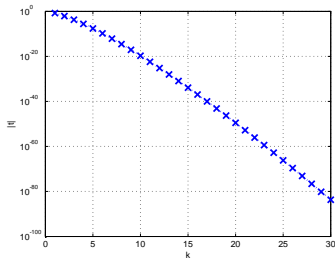
$k = k + 1$

$t = (-1) * t * \frac{x^2}{(2k)(2k+1)}$

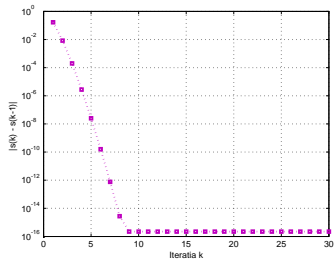
$s = s + t$

•
întoarce s

Rezultate numerice



Modulul termenului curent al dezvoltării în serie Taylor a funcției sinus.



Modulul diferenței dintre sume parțiale consecutive la dezvoltarea în serie Taylor a funcției sinus.

Efectul perturbațiilor datelor de intrare

$$y = f(x_1, x_2, \dots, x_n). \quad (21)$$

$$dy = \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \dots + \frac{\partial f}{\partial x_n} dx_n. \quad (22)$$

$$\Delta y \approx \frac{\partial f}{\partial x_1} \Delta x_1 + \frac{\partial f}{\partial x_2} \Delta x_2 + \dots + \frac{\partial f}{\partial x_n} \Delta x_n. \quad (23)$$

$$\Delta x_k = \bar{x}_k - x_k = e_{x_k}, \quad (24)$$

Eroarea absolută a rezultatului și marginea ei

$$e_y = \bar{y} - y = \Delta y:$$

$$e_y = \sum_{k=1}^n \frac{\partial f}{\partial x_k} e_{x_k}. \quad (25)$$

$$\left| \sum_{k=1}^n \frac{\partial f}{\partial x_k} e_{x_k} \right| \leq \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} e_{x_k} \right| = \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| |e_{x_k}| \leq \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| a_{x_k}, \quad (26)$$

unde $|e_{x_k}| \leq a_{x_k}$.

Marginea erorii absolute a rezultatului

$$a_y = \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| a_{x_k}. \quad (27)$$

Eroarea relativă a rezultatului și marginea ei

$$\varepsilon_y = \mathbf{e}_y / |y|$$

$$\varepsilon_y = \frac{\sum_{k=1}^n \frac{\partial f}{\partial x_k} \mathbf{e}_{x_k}}{|y|} = \sum_{k=1}^n \frac{\partial f}{\partial x_k} \frac{\mathbf{e}_{x_k}}{|y|} = \sum_{k=1}^n \frac{\partial f}{\partial x_k} \frac{|x_k|}{|y|} \varepsilon_{x_k}. \quad (28)$$

Marginea erorii relative a rezultatului

$$r_y = \sum_{k=1}^n \left| \frac{\partial(\ln f)}{\partial x_k} \right| |x_k| r_{x_k}. \quad (29)$$

Cazuri particulare: +, -

Erori	Adunare $y = x_1 + x_2$	Scădere $y = x_1 - x_2$
Eroare absolută: $e_y =$	$e_{x_1} + e_{x_2}$	$e_{x_1} - e_{x_2}$
majorată de: $a_y =$	$a_{x_1} + a_{x_2}$	$a_{x_1} + a_{x_2}$
Eroare relativă: $\varepsilon_y =$	$\frac{x_1}{x_1+x_2} \varepsilon_{x_1} + \frac{x_2}{x_1+x_2} \varepsilon_{x_2}$	$\frac{x_1}{x_1-x_2} \varepsilon_{x_1} - \frac{x_2}{x_1-x_2} \varepsilon_{x_2}$
majorată de $r_y =$	$\frac{x_1}{x_1+x_2} r_{x_1} + \frac{x_2}{x_1+x_2} r_{x_2}$	$\frac{x_1}{x_1-x_2} r_{x_1} + \frac{x_2}{x_1-x_2} r_{x_2}$

Erorile rezultatului adunării și scăderii a două numere reale în funcție de erorile datelor de intrare.

NB! La adunare și scădere marginile erorilor absolute se adună.

- Adunarea este o operație bine condiționată.
- Scăderea este o operație prost condiționată.

Exemplu

$$x_1 = 1.23 \pm 1\% , x_2 = 1.22 \pm 1\%$$

- Scădere:

$$r = |1.23/0.01 \cdot 1/100 + 1.22/0.01 \cdot 1/100| = 1.23 + 1.22 = 2.45 = 245\%$$

$$x_1 - x_2 = 0.01 \pm 245\%.$$

- Adunare:

$$r = |1.23/2.45 \cdot 1/100 + 1.22/2.45 \cdot 1/100| \approx 0.5 \cdot 1/100 + 0.5 \cdot 1/100 = 1/100 = 1\%.$$

$$x_1 + x_2 = 2.45 \pm 1\%.$$

Cazuri particulare: *, /

Erori	Înmulțire $y = x_1 x_2$	Împărțire $y = \frac{x_1}{x_2}$
Eroare absolută: $e_y =$	$x_2 e_{x_1} + x_1 e_{x_2}$	$\frac{1}{x_2} e_{x_1} - \frac{x_1}{x_2^2} e_{x_2}$
majorată de: $a_y =$	$ x_2 a_{x_1} + x_1 a_{x_2}$	$\frac{1}{ x_2 } a_{x_1} + \frac{ x_1 }{x_2^2} a_{x_2}$
Eroare relativă: $\varepsilon_y =$	$\varepsilon_{x_1} + \varepsilon_{x_2}$	$\varepsilon_{x_1} - \varepsilon_{x_2}$
majorată de $r_y =$	$r_{x_1} + r_{x_2}$	$r_{x_1} + r_{x_2}$

Erorile rezultatului înmulțirii și împărțirii a două numere reale în funcție de erorile datelor de intrare.

NB! La înmulțire și împărțire marginile erorilor relative se adună.

- Înmulțirea și împărțirea sunt operații bine condiționate.

Scăderea trebuie evitată

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

$$x_{1,2} = (-b \pm \sqrt{b^2 - 4ac}) / (2a)$$

Ce se întâmplă dacă $b > 0$ și $b^2 \gg 4ac$?

Ce avantaj are următorul cod?

dacă $b > 0$

$$x1 = (-b - \sqrt{b^2 - 4ac}) / (2a)$$

altfel

$$x1 = (-b + \sqrt{b^2 - 4ac}) / (2a)$$

•

$$x2 = c / (a * x1)$$

Extragerea radicalului

$$y = \sqrt{x}$$

$$e_y = \frac{df}{dx} e_x = \frac{1}{2\sqrt{x}} e_x, \quad (30)$$

$$\varepsilon_y = \frac{e_y}{y} = \frac{1}{2\sqrt{x}\sqrt{x}} e_x = \frac{e_x}{2x} = \frac{\varepsilon_x}{2}. \quad (31)$$

Dar rotunjirea nu poate fi ignorata!

Superpoziția erorilor

eroarea relativă într-un calcul aproximativ

=

eroarea relativă produsă de calculul aproximativ cu numere exacte (eroarea de rotunjire)

+

eroarea relativă produsă de calculul exact cu numere aproximative (afectate deci de erori inerente).

$$\bar{y} = y_i(1 + \text{eps}) = y(1 + \varepsilon_y)(1 + \text{eps}) \approx y(1 + \varepsilon_y + \text{eps}),$$

de unde $(\bar{y} - y)/y = \varepsilon_y + \text{eps}$.

$$\varepsilon_{\sqrt{x}} = \frac{\varepsilon_x}{2} + \text{eps}. \quad (32)$$

Eroarea relativă a oricărui rezultat numeric este cel puțin egală cu zeroul mașinii.

Condiționare vs. stabilitate

Condiționarea

se referă la comportarea **problemei matematice** la perturbații ale datelor.

Stabilitatea

se referă la comportarea **algoritmului** la perturbații ale datelor.

Condiționare

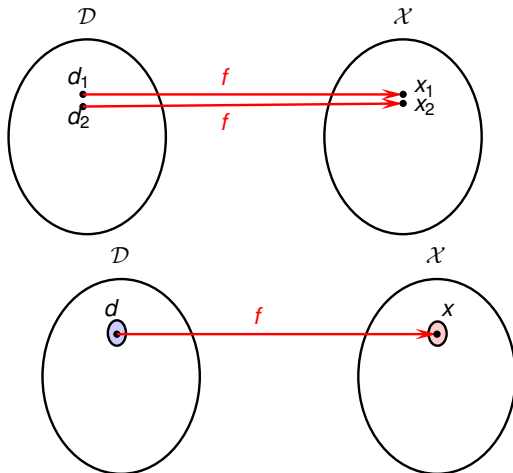
Problemă matematică f formulată explicit:

Fie $f : \mathcal{D} \rightarrow \mathcal{X}$ și $\mathbf{d} \in \mathcal{D}$.

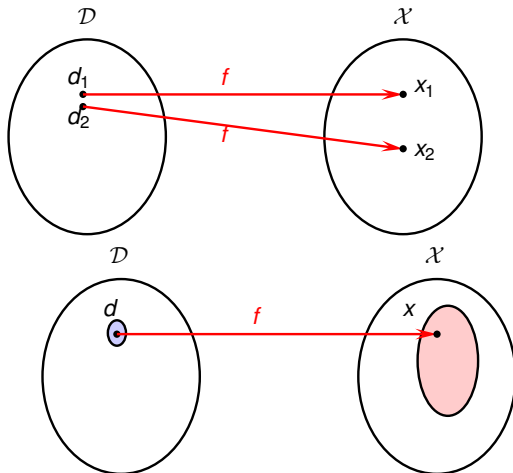
Să se găsească $\mathbf{x} \in \mathcal{X}$ astfel încât $f(\mathbf{d}) = \mathbf{x}$. (33)

O problemă este bine condiționată dacă perturbații mici ale datelor conduc la perturbații mici ale rezultatului.

Reprezentări intuitive - problemă bine condiționată



Reprezentări intuitive - problemă prost condiționată



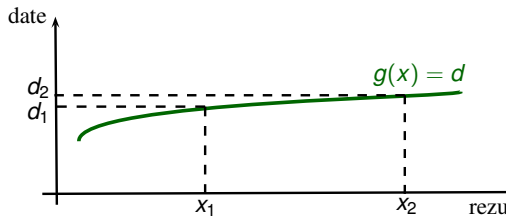
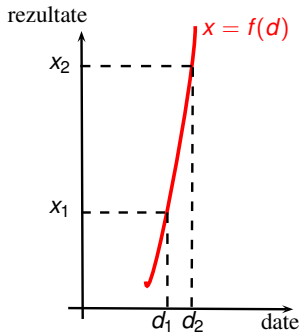
Condiționare

Problemă matematică poate fi formulată și implicit:

Fie $g : \mathcal{X} \rightarrow \mathcal{D}$ și $\mathbf{d} \in \mathcal{D}$.

Să se găsească $\mathbf{x} \in \mathcal{X}$ astfel încât $g(\mathbf{x}) = \mathbf{d}$. (34)

Reprezentări intuitive - problemă prost condiționată



Condiționare - rezultat important

Se demonstrează că între perturbația în date (reziduu) și perturbația în rezultat (eroare) există următoarea relație:

$$\|\mathbf{e}_x\| \leq \kappa \|\varepsilon_d\|, \quad (35)$$

unde κ este un scalar numit *număr de condiționare*, care depinde de problema numerică abordată. (Vom reveni asupra lui la cursul următor).

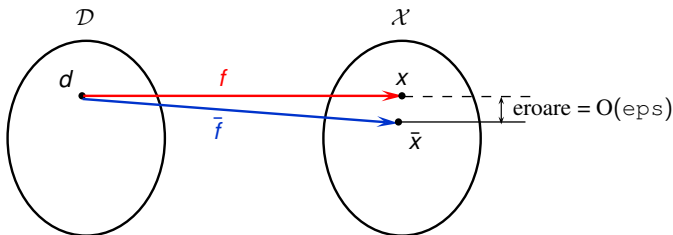
Condiționare - concluzii

- Reziduul nu dă informații despre eroare.
- Eroarea și reziduul sunt legate prin numărul de condiționare.
- Pentru o problemă cu număr de condiționare mic, o perturbație mică în date va duce la o perturbație mică a rezultatului.
- Problemele matematice care au κ mare sunt prost condiționate și ele nu pot fi rezolvate cu ajutorul calculatorului. Pentru astfel de probleme, trebuie găsită o formulare matematică echivalentă din punct de vedere al rezultatului, dar bine condiționată.

În cele ce urmează vom presupune că problema f este bine condiționată și pentru rezolvarea ei a fost conceput un algoritm \bar{f} .

Acuratețea unui algoritm

Acuratețea unui algoritm se referă la eroarea soluției numerice.



Reprezentarea intuitivă a unui algoritm a cărui precizie este ideală.

În mod ideal, un algoritm este precis dacă:

$$\frac{\|\bar{f}(\mathbf{d}) - f(\mathbf{d})\|}{\|f(\mathbf{d})\|} = O(\epsilon). \quad (36)$$

$\bar{f}(\mathbf{d})$ = "rezultatul algoritmului \bar{f} aplicat datelor \mathbf{d} ".

Stabilitatea unui algoritm

Dar, rotunjirea datelor este inevitabilă, erorile se acumulează și perturbă rezultatul. Este mai util să se țintească **stabilitatea algoritmului**.

Stabilitatea unui algoritm se referă la comportarea algoritmului atunci când datele de intrare sunt perturbate.

Un algoritm \bar{f} folosit pentru rezolvarea unei probleme f este stabil dacă

$$\frac{\|\bar{f}(\bar{\mathbf{d}}) - f(\mathbf{d})\|}{\|f(\mathbf{d})\|} = O(\epsilon_{\text{ps}}), \quad (37)$$

pentru $(\forall)\bar{\mathbf{d}}, \mathbf{d}$ care satisfac $\|\bar{\mathbf{d}} - \mathbf{d}\|/\|\mathbf{d}\| = O(\epsilon_{\text{ps}})$.

Pe scurt, **un algoritm stabil dă răspunsul aproape corect pentru date reprezentate aproape precis**.

Ilustrarea stabilității unui algoritm - problema

$$\mathbf{Ax} = \mathbf{b}, \quad \text{unde} \quad \mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

$$\begin{aligned} x_2 &= 1 \\ x_1 + x_2 &= 0 \end{aligned} \tag{38}$$

$$x_1 = -1, \quad x_2 = 1. \quad \mathbf{x} = f(\mathbf{d}) = [-1, 1]^T.$$

Să considerăm acum că datele au fost perturbate:

$$\begin{aligned} \bar{\mathbf{A}} &= \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix}, \\ 10^{-20}x_1 + x_2 &= 1 \\ x_1 + x_2 &= 0 \end{aligned} \tag{39}$$

$x'_1 = -x'_2 = 1/(10^{-20} - 1) \approx -1$. Se poate demonstra că această problemă este bine conditionată.

Ilustrarea stabilității unui algoritm - algoritmul \bar{f}_1

- **Pasul 1:** se înmulțește prima ecuație a sistemului cu (-10^{20}) și se adună cu a doua, rezultând x_2 ;
- **Pasul 2:** se calculează x_1 din prima ecuație.

La pasul 1 se ajunge la ecuația $(1 - 10^{20})x_2 = -10^{20}$ care, în calculator devine datorită rotunjirilor $-10^{20}x_2 = -10^{20}$, de unde va rezulta $x_2 = 1$, ceea ce este corect.

La pasul 2 ecuația de rezolvat devine $10^{-20}x_1 + 1 = 1$, de unde va rezulta $x_1 = 0$, ceea ce este greșit, foarte departe de valoarea adevărată.

Acest algoritm este instabil.

Ilustrarea stabilității unui algoritm - algoritmul \bar{f}_2

- **Pasul 1:** se înmulțește a doua ecuație a sistemului cu (-10^{-20}) și se adună cu prima, rezultând x_2 ;
- **Pasul 2:** se calculează x_1 din a doua ecuație.

La pasul 1 se ajunge la ecuația $(1 - 10^{-20})x_2 = 1$, care în calculator devine $x_2 = 1$.

La pasul 2 ecuația de rezolvat este $x_1 + 1 = 0$, de unde $x_1 = -1$, ceea ce este corect.

Algoritmul \bar{f}_2 este stabil. Stabilitatea lui este foarte puternică, el a dat răspunsul exact pentru date de intrare aproape precise.

Concluzii - estimarea acurateții unei soluții numerice

- 1 Se estimează numărul de condiționare al problemei. Se continuă numai dacă problema matematică este bine condiționată.
- 2 Se investighează stabilitatea algoritmului. Cel mai simplu este ca acest lucru să se realizeze experimental, rulându-se algoritmul pentru date perturbate. Dacă dispersia rezultatelor este mare atunci algoritmul este instabil și trebuie schimbat.
- 3 Dacă algoritmul este stabil, atunci acuratețea finală (modulul erorii relative) este majorată de produsul dintre numărul de condiționare și modulul reziduuului relativ.

Despre un algoritm stabil care generează erori mici pentru probleme bine condiționate se spune că este **robust**.

Lectura obligatorie pentru această săptămână

● Erori - Cap.2 din

[1] Gabriela Ciuprina, Mihai Rebican, Daniel Ioan - Metode numerice in ingineria electrica - Îndrumar de laborator pentru studenții facultății de Inginerie electrică, Editura Printech, 2013, disponibil la

http://mn.lmn.pub.ro/indrumar/IndrumarMN_Printech2013.pdf