

**Gabriela
CIUPRINA**

**Algoritmi numerici
pentru calcule științifice
în ingineria electrică**

Algoritmi numerici
pentru calcule științifice în ingineria electrică

Gabriela Ciuprina

Partea I

MATRIX ROM
București 2013

MATRIX ROM
C.P. 16 – 162
062510 – BUCUREȘTI
tel. 021.4113617, 031.4012438, 0372743840
fax 021.4114280
email: office@matrixrom.ro
www.matrixrom.ro

Editura MATRIX ROM este acreditată de
CONSILIUL NAȚIONAL AL CERCETĂRII ȘTIINȚIFICE DIN ÎNVĂȚĂMÂNTUL SUPERIOR

Referenți științifici: Prof. univ. dr. ing. F.M.G. Tomescu

Prof. univ. dr. ing. Daniel Ioan

Universitatea « Politehnica » București

Descrierea CIP a Bibliotecii Naționale a României

CIUPRINA, GABRIELA

Algoritmi numerici pentru calcule științifice în ingineria electrică /
Gabriela Ciuprina. - București : Matrix Rom, 2013

Bibliogr.

ISBN 978-606-25-0008-5

519.6:621.3(075.8)

ISBN 978 – 606 – 25 – 0008 – 5

Algoritmi numerici
pentru calcule științifice în ingineria
electrică

Gabriela Ciuprina

Partea I

Referenți științifici:

Prof.dr.ing.F.M.G. Tomescu

Prof.dr.ing. Daniel Ioan

2013

Introducere

Cartea de față reprezintă o extindere a primei părți a suportului pentru cursul *Metode numerice în ingineria electrică*, predat studenților din anul II de licență ai Facultății de Inginerie Electrică din Universitatea Politehnica București. Lucrarea se adresează tuturor celor care nu au experiență în metode numerice și care doresc să se familiarizeze cu conceptele de bază și cu modul de gândire algoritmic, structurat.

Această carte este structurată în 5 capitole care tratează următoarele aspecte: descrierea și evaluarea algoritmilor (cap.1), rezolvarea numerică a sistemelor de ecuații algebrice liniare (cap.2) cu aplicație la conceperea unor algoritmi numerici pentru analiza circuitelor rezistive liniare (cap.3), metode de interpolare a funcțiilor (cap.4) și metode de derivare numerică (cap.5) aplicate pentru rezolvarea unor ecuații diferențiale.

Partea a II-a, ce va face obiectul unei alte cărți, va trata problema integrării numerice, a rezolvării numerice a sistemelor de ecuații algebrice neliniare, cu aplicație la algoritmi numerici pentru analiza circuitelor rezistive neliniare și rezolvarea numerică a sistemelor de ecuații cu derivate ordinare, cu aplicații la algoritmi numerici pentru analiza circuitelor în regim tranzitoriu.

Autoarea mulțumește referenților științifici Prof.dr.ing.F.M.G. Tomescu și Prof.dr.ing. Daniel Ioan pentru comentariile constructive și sugestiile de îmbunătățire, care au fost incluse în această carte.

Orice comentarii și sugestii în vederea îmbunătățirii acestui suport de curs sunt binevenite la adresa gabriela.ciuprina@upb.ro.

Studentilor mei, cu speranța că studiul acestei discipline le va face plăcere.

Copiilor mei, Andreea și Ștefan, cu toată dragostea.

Cuprins

1	Descrierea și evaluarea algoritmilor	1
1.1	Pseudocodul	2
1.1.1	Descrierea pseudocodului	2
1.1.2	Sugestii pentru utilizatorii de C	14
1.1.3	Sugestii pentru utilizatorii de Matlab	17
1.2	Complexitatea algoritmilor	19
1.2.1	Timp de calcul	19
	Algoritmi polinomiali	20
	Notații asimptotice	25
	Analiza complexității algoritmilor recursivi	28
1.2.2	Necesar de memorie	30
1.3	Erori în calculele numerice	30
1.3.1	Tipuri de erori	31
1.3.2	Analiza erorilor de rotunjire	33
1.3.3	Standardul IEEE pentru reprezentarea în virgulă mobilă	35
1.3.4	Analiza erorilor de trunchiere	37
1.3.5	Analiza erorilor inerente	38
1.3.6	Condiționare și stabilitate	43
2	Sisteme de ecuații algebrice liniare	51
2.1	Formularea problemei	51
2.2	Condiționarea problemei	52
2.3	Clasificarea metodelor	56
2.4	Metoda Gauss	58
2.4.1	Un exemplu simplu	59
2.4.2	Algoritm	59

2.4.3	Evaluarea algoritmului	63
2.4.4	Strategii de pivotare	64
2.4.5	Cazul sistemelor multiple	67
2.5	Metoda factorizării LU	70
2.5.1	Un exemplu simplu	71
2.5.2	Variante de factorizare	72
2.5.3	Algoritmul variantei Doolittle	73
2.5.4	Calculul soluției după factorizare	76
2.5.5	Evaluarea algoritmului. Strategii de pivotare.	77
2.5.6	Cazul sistemelor multiple	79
2.5.7	Varianta Cholesky	79
2.6	Cazul matricelor rare	82
2.6.1	Formate de memorare a matricelor rare	84
2.6.2	Metode directe pentru matrice rare	86
2.7	Metode iterative staționare - generalități	88
2.7.1	Ideea de bază a metodelor iterative staționare	89
2.7.2	Criteriul de oprire	91
2.7.3	Intermezzo despre vectori și valori proprii	92
2.7.4	Convergență	94
2.7.5	Algoritm general	95
2.8	Metoda Jacobi	96
2.8.1	Un exemplu simplu	97
2.8.2	Algoritmul metodei	97
2.8.3	Aspecte legate de convergență	99
2.9	Metoda Gauss-Seidel	100
2.9.1	Un exemplu simplu	100
2.9.2	Algoritmul metodei	100
2.9.3	Aspecte legate de convergență	102
2.9.4	Evaluarea algoritmului	103
2.10	Metoda suprarelaxării succesive (SOR)	104
2.11	Algoritmul general al metodelor staționare	105
2.12	Metoda gradientilor conjugați	106
2.12.1	Forma pătratică asociată	107
2.12.2	Metoda celei mai rapide coborâri (a gradientului)	108

2.12.3	Metoda direcțiilor conjugate	111
2.12.4	Metoda gradientilor conjugați	114
2.13	Precondiționare	118
3	Algoritmi numerici pentru analiza circuitelor rezistive	121
3.1	Tehnica nodală	122
3.1.1	Structuri de date	125
3.1.2	Etapa de preprocesare	126
3.1.3	Etapa de rezolvare	130
3.1.4	Etapa de postprocesare	130
3.2	Tratarea surselor ideale de curent	131
3.3	Tratarea surselor ideale de tensiune	135
3.4	Tratarea surselor comandate liniar	140
4	Interpolarea funcțiilor	143
4.1	Distanța dintre două funcții	145
4.2	Formularea problemei interpolării	147
4.3	Metode de interpolare globală	148
4.3.1	Metoda clasică	148
4.3.2	Metoda Lagrange	149
4.3.3	Eroarea de trunchiere	153
4.3.4	Metoda Newton	154
4.3.5	Un exemplu simplu	160
4.3.6	Interpolarea Chebyshev	161
4.4	Metode de interpolare polinomială pe porțiuni	162
4.4.1	Interpolarea liniară pe porțiuni	163
4.4.2	Interpolarea Hermite	164
5	Derivarea numerică	169
5.1	Formularea problemei derivării numerice	169
5.2	Formule de derivare numerică	170
5.2.1	Formule de derivare cu eroare de ordinul 1	170
5.2.2	Formule de derivare cu eroare de ordinul 2	172
5.3	Algoritmi numerici pentru derivarea funcțiilor	173
5.3.1	Cazul funcțiilor date tabelar	174

5.3.2	Cazul funcțiilor date prin cod	175
5.4	Derivate de ordin superior	178
5.5	Derivate parțiale	180
5.6	Metoda diferențelor finite	182
5.6.1	Rezolvarea unei ecuații diferențiale ordinare	182
5.6.2	Rezolvarea unei ecuații cu derivate parțiale de tip eliptic	188
5.6.3	Rezolvarea unei ecuații cu derivate parțiale de tip hiperbolic	194
A	Fișiere utile pentru implementarea în C	197
	Referințe	201

Capitolul 1

Descrierea și evaluarea algoritmilor

Un algoritm este o metodă de rezolvare a unei probleme corect formulate, constând într-un număr finit de etape simple, elementare, susceptibile de a fi implementate pe un calculator.

În consecință, pentru a putea descrie un algoritm trebuie cunoscute etapele elementare ce pot fi implementate pe un calculator. Implementarea se realizează cu ajutorul unui limbaj de programare (C, FORTRAN sau Matlab, Scilab, etc). Deși limbajele de programare au o sintaxă proprie ce trebuie respectată cu rigurozitate, acțiunile descrise sunt însă aceleași, în oricare dintre acestea. De aceea, este mult mai util și mai eficient ca atunci când se concepe un algoritm, acesta să fie descris într-un mod cât mai natural, neîncorsetat de rigorile vreunui limbaj de programare, eventual cu notații matematice compacte, astfel încât să poată fi ușor de citit și înțeles. Ca urmare, în ceea ce urmează, vom descrie algoritmi într-un așa numit *pseudolimbaj* sau *pseudocod*. Nu există niciun standard pentru sintaxa pseudocodului. Un pseudocod este un fals limbaj, nu poate fi compilat. El însă reprezintă punctul de plecare pentru scrierea programului. Există și așa numitele scheme logice care pot fi văzute ca reprezentări grafice ale pseudocodului. În această carte se va promova însă folosirea pseudocodului și nu a schemelor logice, pentru că trecerea de la pseudocod la limbajul de programare este mult mai naturală. În prima parte a acestui capitol este descris pseudocodul folosit pe parcursul acestei cărți.

Deoarece în general o problemă admite mai multe metode de rezolvare, înseamnă că pentru rezolvarea ei cu ajutorul calculatorului pot fi concepuți algoritmi diferiți. De aceea, în partea a doua a acestui capitol sunt prezentate criteriile după care pot fi evaluați algoritmi, criterii ce sunt independente de calculator, sistemul de operare și limbajul de programare folosit.

1.1 Pseudocodul

1.1.1 Descrierea pseudocodului

Pseudocodul este o metodă de descriere convențională, simplă, a algoritmilor. El nu are o sintaxă strictă, și poate chiar folosi cuvinte cheie din limba maternă a celui care îl concepe. Este important însă ca el să fie clar și neambiguu. *Cuvintele cheie* sunt acele cuvinte care au corespondențe în fiecare limbaj de programare. Pe parcursul acestei cărți cuvintele cheie sunt scrise cu un alt font și sunt subliniate.

Pseudocodul este alcătuit din linii care fie descriu acțiuni, fie descriu date. Liniile ce descriu *acțiuni* se vor traduce în limbajul de programare cu ajutorul unor *instrucțiuni*. Liniile ce descriu *date* se vor traduce în limbajul de programare cu ajutorul unor *declarații*.

Există limbaje de programare (de tip Matlab, Scilab, Octave) care nu necesită declararea datelor și limbaje de programare (C, FORTRAN) la care declarațiile sunt obligatorii. Folosirea unui limbaj din prima categorie facilitează implementarea prototipurilor de programe. Scrierea declarațiilor în pseudocod are avantajul că, pe de o parte, permite analiza tipurilor unor expresii ce pot apărea în pseudocod și, pe de altă parte, permite estimarea mai ușoară a necesarului de memorie, mărime utilă în evaluarea algoritmilor.

Un concept deosebit de important în descrierea algoritmilor îl reprezintă *variabila*.

Variabila este o zonă din memorie cu trei caracteristici: nume, valoare și tip.

Numele unei variabile este un șir de caractere alfa-numerice care permit identificarea zonei (adresei) de memorie în care se află variabila. Pseudocodul utilizat în această carte permite ca numele să conțină atât litere cât și cifre, dar cu următoarele două restricții, menite a evita confuziile: să nu existe spații goale în interiorul numelui și numele să înceapă obligatoriu cu o literă.

Valoarea unei variabile reprezintă conținutul zonei de memorie. Acest conținut este un număr, fără să aibă asociată o unitate de măsură. În pseudolimbaj nu vom face deosebire între nume și valoare.

Tipul este un atribut foarte important al variabilei care permite interpretarea conținutului zonei de memorie.

Un algoritm are nevoie de mai multe variabile, care vor ocupa în calculator o anumită *zonă de memorie*.

Declarațiile sunt linii de pseudocod ce descriu datele. Sintaxa declarațiilor diferă în funcție de tipul variabilelor. Pseudocodul utilizat admite două tipuri principale de variabile: fundamentale (simple) și agregate.

Tipurile fundamentale de date sunt: logic, întreg, real sau caracter. Cuvintele cheie folosite pentru declararea acestor tipuri în pseudocod sunt, respectiv, logic, întreg,

real, caracter.

Tipurile agregate pot fi tablouri sau înregistrări, pentru declararea cărora se vor folosi cuvintele cheie tablou și, respectiv, înregistrare.

Declararea unei variabile se face specificând cuvântul cheie asociat tipului, urmat de numele variabilei respective sau de lista numelor variabilelor de tipul respectiv.

Iată exemple de declarații ale unor variabile simple:

logic T, F ; adevărat (*true* - în engl.), fals
întreg N ; numărul de noduri
real Pc, Pg ; putere consumată, putere generată
caracter c, C

În pseudolimbaj se face deosebirea între literele mari și cele mici. De asemenea, este permis să se declare mai multe variabile de același tip în aceeași declarație. De asemenea, tot ce este după semnul ”;” va fi considerat comentariu. Este bine ca variabilele alese să aibă nume cât mai sugestive.

Deși în anumite limbaje de programare nu este permis, în pseudocod nu este grav dacă sunt omise declarațiile unor variabile. Este totuși recomandat să se declare toate variabilele folosite pe de o parte pentru analiza corectă a operațiilor de evaluare a unor expresii, și, pe de altă parte pentru analiza corectă a complexității algoritmului (a se vedea și paragraful 1.2).

Tabloul este o mulțime de elemente de același tip, cu nume comun. Numărul de elemente al tabloului este cunoscut anterior, chiar din declarație.

Declararea unui tablou se face astfel:

tablou tip nume[dimensiune]

unde ”tip” este un cuvânt cheie ce desemnează unul din tipurile fundamentale, ”nume” este numele variabilei de tip tablou, iar ”dimensiune” reprezintă numărul de elemente al tabloului.

În declarația următoare, vă puteți imagina că, în calculator se rezervă spațiu de memorie pentru un număr de 10 variabile de tip real, ce poartă numele V :

tablou real $V[10]$

În algoritmul propriu-zis pot fi folosite cel mult zece variabile, componente ale tabloului V .

O variantă mai elegantă, care conduce și la economie de memorie, este de a declara exact atâtea variabile cât este nevoie. De aceea, se va admite în pseudocod o declarație de tipul

întreg N

...

tablou real $V[N]$

Detaliile despre cum se implementează acest tip de declarație în limbajul C sunt descrise în paragraful 1.1.2.

Adresarea elementelor unui tablou se poate face cu paranteze sau cu indici. De exemplu, adresarea elementelor din tabloul declarat în exemplu anterior se poate face astfel $V(1)$, $V(5)$ sau V_1 , V_5 , etc.

Unele metode de rezolvare sunt descrise în mod natural cu ajutorul unor tablouri care au indici ce pornesc din 0. De aceea, se va accepta ca element posibil și $B(0)$ sau B_0 . La fiecare algoritm se va preciza dacă indicii sunt de la 1 la N sau de la 0 la $N - 1$.

Înregistrarea este o mulțime de elemente de tipuri diferite. Fiecare element al unei înregistrări poartă numele de *câmp* și este identificat printr-un nume.

Declarația unei înregistrări se face astfel:

înregistrare nume

tip1 nume1

tip2 nume2

unde "nume" este numele înregistrării, "tip1", "tip2", etc. sunt cuvinte cheie asociate tipurilor fundamentale, iar "nume1", "nume2", etc. sunt nume de câmpuri.

Iată un exemplu de înregistrare

înregistrare punct

logic polar

real coord1

real coord2

unde, dacă "polar" este fals atunci "coord1" reprezintă abscisa și "coord2" reprezintă ordonata unui punct în plan, iar dacă "polar" are valoarea adevărat atunci "coord1" reprezintă raza, și "coord2" reprezintă unghiul în sistemul de coordonate polar.

Adresarea elementelor unei înregistrări se face invocând numele înregistrării, urmat de separatorul ".", urmat de numele câmpului. Adresarea elementelor înregistrării declarate în exemplul anterior se face cu

punct.polar

punct.coord1

punct.coord2

Instrucțiunile sunt linii de pseudocod ce reprezintă acțiunile ce se execută asupra datelor.

Există două categorii de instrucțiuni: simple și structurate. *Instrucțiunile simple* pot fi la rândul lor: de intrare, de ieșire și de atribuire. *Instrucțiunile structurate* sunt: secvența, decizia fără sau cu alternativă, ciclul cu test inițial/final sau cu contor, rutinele (proceduri sau funcții). Descrierea fiecărui tip de instrucțiune este prezentată în cele ce urmează.

Cele mai simple instrucțiuni sunt cele de transfer de la dispozitivul de intrare la unitatea centrală de prelucrare (CPU) și cele de la CPU la dispozitivul de ieșire.

Instrucțiunea de intrare are sintaxa

citește nume_var

iar *instrucțiunea de ieșire* are sintaxa

scrie nume_var.

Iată un exemplu de utilizare a instrucțiunilor de intrare și ieșire.

citește N

citește q, p ; în pseudocod se admite citirea unei liste de variabile

scrie N

scrie q, p ; în pseudocod se admite scrierea unei liste de variabile

Instrucțiunea de atribuire are sintaxa

nume_variabilă = expresie

în care "nume_variabilă" este numele unei variabile, iar "expresie" este un șir de operatori și operanzi care respectă anumite reguli sintactice.

Este important de reamintit faptul că nu orice tip de operand poate fi folosit cu orice tip de operator și rezultatul evaluării expresiei poate fi de tip logic sau de tip aritmetic (real sau întreg) în funcție de tipul de operanzi și operatori folosiți.

O *expresie logică* (al cărei rezultat se atribuie unei variabile de tip logic) poate avea fie *operatori logici* și, sau, nu, caz în care operanzii trebuie să fie logici, fie *operatori de relație* $<$, $>$, \leq , \geq , $=$, \neq , caz în care operanzii trebuie să fie aritmetici (reali sau întregi), sau de tip caracter (caz în care se compară poziția lor în tabelul ASCII¹).

Să considerăm următorul exemplu:

¹ASCII - *American Standard Code for Information Interchange* - codificare de tip text a 128 de caractere, din care 26 reprezintă litere ale alfabetului latin.

<u>logic</u> a	; rezultatul evaluării expresiei logice
<u>logic</u> b, c	; operanzi logici
<u>real</u> x, y	; operanzi aritmetici
$a = (b \text{ sau } (\text{nu } c))$; expresie logică cu operatori logici
$a = (x \leq y)$; expresie logică cu operatori de relație
$a = (x = y)$; expresie logică cu operatori de relație

Ultima linie a acestui exemplu conține de două ori semnul ”=”. Primul reprezintă simbolul de atribuire, iar al doilea este un operator de relație. În pseudocod nu se va face deosebirea între modul de scriere al lor, dar în limbajele de programare nu se întâmplă la fel. De exemplu, în C sau în Matlab, această linie se scrie $a = (x == y)$, iar în Pascal $a := (x = y)$. În unele cărți se evită această confuzie folosind simbolul \leftarrow pentru a reprezenta operația de atribuire, caz în care linia de mai sus s-ar fi scris $a \leftarrow (x = y)$.

O *expresie aritmetică* (al cărei rezultat se atribuie unei variabile de tip real sau întreg) are operanzi aritmetici, iar operatorii pot fi: $+$, $-$, \cdot , $/$, dar și $\sqrt{\quad}$, funcții standard (de exemplu funcțiile trigonometrice).

Iată un exemplu de pseudocod în care apar expresii aritmetice.

```

întreg  $i$ 
real  $d, x, y$ 
 $i = i + 1$ 
 $d = xy + \sin(y)$ 
 $d = \sqrt{d}$ 
 $d = \frac{d^2 + d}{\frac{1}{3}}$ 

```

În pseudocod se admite scrierea simbolului radical, scrierea unui exponent și chiar a fracțiilor etajate. Se admite chiar să nu se pună simbolul explicit pentru înmulțire dacă acest lucru nu produce confuzii. Astfel, sintaxa nu este rigidă, ca cea a unui limbaj de programare. Totul este ca expresia să fie clară, astfel încât implementarea să fie imediată.

Secvența este un șir de instrucțiuni simple, scrise una sub alta, dar aliniate la stânga și indentate. Indentarea este deosebit de importantă pentru a obține un cod clar și ușor de înțeles². În exemplele următoare se folosește scrierea indentată.

Decizia fără alternativă (simplă) are sintaxa:

```

dacă condiție [atunci]
    secvență

```

²Mediile de programare oferă facilități pentru indentare. De exemplu, în Linux există comanda `indent`, iar în editorul din Matlab există comanda `Smart indent`.

Cuvântul cheie atunci poate lipsi. Pentru claritate, ”secvența” este scrisă indentat (deplasat la dreapta) față de cuvântul cheie dacă. Ea se va executa doar dacă rezultatul evaluării expresiei ”condiție” este adevărat.

Iată un exemplu de folosire a scrierii indentate.

```
real x
x = 2
dacă x < 0
    x = x + 3
    x = x/2
    x = x2
x = 2x
```

În acest caz, la sfârșitul execuției acestui cod, x va avea valoarea 4. O modificare ”minoră” a indentării, ca de exemplu

```
real x
x = 2
dacă x < 0
    x = x + 3
    x = x/2
x = x2
x = 2x
```

schimbă rezultatul final, care va fi în acest caz 8. Pentru mărirea clarității pseudocodului, uneori se va marca suplimentar sfârșitul instrucțiunii de decizie cu \bullet , ca de exemplu

```
real x
x = 2
dacă x < 0
    x = x + 3
    x = x/2
    x = x2
     $\bullet$ 
x = 2x
```

Este recomandat ca fiecare instrucțiune să fie scrisă separat pe o linie, utilizând scrierea indentată. În fond, codul de mai sus se putea scrie pe o singură linie astfel

```
real x  x = 2  dacă x < 0 atunci x = x + 3  x = x/2  x = x2   $\bullet$  x = 2x
```

Cu siguranța că s-ar fi făcut economie de spațiu, dar claritatea s-ar fi pierdut. Iar acesta nu a fost decât un exemplu banal. Algoritmii problemelor reale sunt cu mult mai complicați.

Decizia cu alternativă are sintaxa

```
dacă condiție [atunci]  
    secvența1  
altfel  
    secvența2
```

În acest caz, "condiție" este o expresie care este evaluată, iar dacă rezultatul este adevărat, atunci se execută "secvența1", în caz contrar executându-se "secvența2".

Iată un exemplu de pseudocod în care se citește un număr real și se afișează modulul lui. În acest pseudocod se folosește o decizie cu alternativă.

```
real  $x$            ; un număr real - dată de intrare  
real  $modul$        ; dată de ieșire - modulul numărului real  
citește  $x$   
dacă  $x \geq 0$   
     $modul = x$   
altfel  
     $modul = -x$   
•  
scrie  $modul$ 
```

Ciclul cu test inițial are sintaxa

```
cât timp condiție [repetă]  
    secvență
```

În acest caz se repetă corpul ciclului, adică "secvență" cât timp rezultatul evaluării expresiei logice "condiție" este adevărată. S-ar putea ca "secvență" să nu fie executată niciodată. Acest lucru se întâmplă dacă în momentul execuției ciclului rezultatul evaluării expresiei "condiție" este fals.

Iată un exemplu de folosire a unui ciclu cu test inițial.

```
întreg  $N$   
întreg  $i$   
tablou real  $x[N]$   
 $N = 3$   
 $x_1 = 1$ 
```

```

 $x_2 = -1$ 
 $x_3 = 2$ 
 $i = 1$ 
cât timp ( $i \geq 1$ ) și ( $i \leq N$ )
    dacă  $x_i \geq 0$ 
        scrie "Elementul ", i, " este pozitiv"
    •
     $i = i + 1$ 
•

```

Ciclul cu test final are sintaxa

```

repetă
    secvență
cât timp condiție

```

În acest caz se repetă corpul ciclului, adică "secvență" cât timp "condiție" este adevărată. Spre deosebire de ciclul cu test inițial, ciclul cu test final este executat cel puțin o dată.

Este posibil ca sintaxa ciclului cu test final să o scriem și astfel

```

repetă
    secvență
până când condiție

```

În acest caz corpul ciclului se repetă cât timp "condiție" este falsă. Când "condiție" devine adevărată, programul iese din ciclu. Logica este deci inversă primei variante.

Oricare din aceste două variante pot fi folosite în pseudocod. Ele doar trebuie traduse cu grijă în limbajul de programare ales³.

Iată un exemplu de folosire al unui ciclu cu test final:

```

întreg  $N$ 
tablou real  $x[N]$ 
real  $s, \varepsilon$ 
întreg  $k$ 
...
 $k = 0$ 
 $s = 0$ 

```

³În C ciclul cu test final folosește cuvântul cheie `while` (care în limba engleză înseamnă cât timp), iar în Pascal ciclul cu test final folosește cuvântul cheie `until` (până când).

repetă

$k = k + 1$

$s = s + x_k$

cât timp $|x_k| \geq \varepsilon$

Un astfel de pseudocod calculează suma unei serii, adunând termenii până când ultimul termen adunat este mai mic în modul decât o anumită valoare, suficient de mică.

Ciclul cu contor are sintaxa

pentru contor = val_in, val_fin[, pas] [repetă]
secvență

Într-un astfel de ciclu, variabila "contor" ia, pe rând, valorile întregi "val_in", "val_in" + "pas", ..., "val_fin". Pentru fiecare din aceste valori se execută "secvența". Valoarea "pas" poate lipsi, caz în care este considerată, în mod implicit, 1. Această valoare poate fi pozitivă sau negativă.

Ciclul cu contor se folosește atunci când se cunoaște de la început numărul de repetiții ale secvenței, pe când ciclurile cu test se folosesc atunci când acest număr nu este cunoscut.

Un exemplu simplu de folosire a unui ciclu cu contor îl reprezintă citirea elementelor unei matrice.

întreg N

citește N

tablou real $a[N, N]$

întreg i, j

pentru $i = 1, N$

pentru $j = 1, N$

citește $a_{i,j}$

Rutinele sunt folosite atunci când aceeași secvență apare de mai multe ori în același algoritm. În acest fel, secvența respectivă va fi memorată într-un singur loc în memoria calculatorului și va fi apelată ori de câte ori este nevoie. În general se preferă utilizarea rutinelor chiar dacă ele sunt apelate o singură dată. În acest fel, algoritmul se modularizează și, chiar dacă el face lucruri complicate, nu este lung ci face apel la rutine. Modularizarea are avantajul că permite și împărțirea problemei de rezolvat în sub-probleme, astfel încât conceperea unui algoritm pentru rezolvarea unei probleme complicate poate fi făcut în echipă. Se recomandă ca programul principal să fie scris cu cât mai puține linii de cod, el făcând apel la rutine. O dată ce o rutină este apelată, programul îi transferă ei execuția. După ce rutina a fost executată, întoarcerea se face imediat după punctul de apel.

După numărul de parametri de ieșire, rutinele se împart în proceduri și funcții.

Procedura este o rutină care poate avea oricâți parametri de ieșire.

Definiția unei proceduri se face astfel

```
procedură nume_proc (lista argumentelor formale de intrare și ieșire)
; comentarii ce descriu ce face procedura și parametrii acesteia
...
; declarații pentru argumente
...
; instrucțiuni
...
retur                ; cuvânt cheie ce comandă întoarcerea în punctul de apel
```

Apelul unei proceduri se face simplu, prin invocarea numelui ei, urmat de lista argumentelor actuale:

```
nume_proc (lista argumentelor actuale de intrare și ieșire)
```

Funcția este o rutină care are exact un singur parametru de ieșire.

Definiția unei funcții se face astfel

```
funcție nume_fct (lista argumentelor formale de intrare)
; comentarii ce descriu ce face funcția și parametrii acesteia
...
; declarații pentru argumente
...
; instrucțiuni
...
întoarce valoare     ; cuvântul cheie comandă întoarcerea în punctul de apel
```

Apelul unei funcții se face în instrucțiuni de atribuire, funcția putând fi înțeleasă ca un operator aplicat unor operanzi ce apar ca o listă de argumente actuale de intrare ale funcției:

```
val = nume_fct (lista argumentelor actuale de intrare)
```

Funcția poate fi privită ca un caz particular de procedură. Este foarte important ca în cazul rutinelor, să nu se confunde argumentele formale, folosite la definirea rutinei,

cu argumentele actuale, cele pentru care se face calculul propriu-zis. Există o restricție valabilă în toate limbajele de programare și anume ca argumentele formale să concorde cu cele actuale ca număr, ordine și tip. Nerespectarea acestui lucru generează neclarități în pseudocod și greșeli în implementarea propriu-zisă.

Ca exemplu, vom modulariza, cu ajutorul rutinelor, următorul algoritm ce calculează produsul scalar a doi vectori.

```

; program principal
întreg N                                ; dimensiunea vectorilor
citește N
tablou real a[N], b[N]                ; date de intrare: doi vectori de numere reale, indici de la 1
real p                                ; rezultatul: produsul scalar al vectorilor a și b
întreg i                                ; alte variabile utile
pentru i = 1, N
    citește ai
    •
    pentru i = 1, N
        citește bi
    •
    p = 0
    pentru i = 1, N
        p = p + aibi
    •
scrie p

```

Vom rescrie acest pseudocod folosind o procedură pentru citirea vectorilor și o funcție pentru calculul produsului scalar. Codul rezultat este următorul:

```

; program principal
întreg N                                ; dimensiunea vectorilor
citește N
tablou real a[N], b[N]                ; date de intrare: doi vectori de numere reale, indici de la 1
real p                                ; rezultatul: produsul scalar al vectorilor a și b
citește_vector(N, a)
citește_vector(N, b)
p = produs_scalar(N, a, b)
scrie p

procedură citește_vector(N, x)

```

```

întreg  $N$  ; dimensiunea vectorului - argument formal de intrare
tablou real  $x[N]$  ; vectorul propriu-zis - argument formal de ieşire
întreg  $i$ 
pentru  $i = 1, N$ 
    citeşte  $x_i$ 
•
retur

funcţie produs_scalar( $N, v, w$ )
întreg  $N$ 
tablou real  $v[N], w[N]$ 
întreg  $i$ 
real  $r$ 
 $r = 0$ 
pentru  $i = 1, N$ 
     $r = r + v_i w_i$ 
•
întoarce  $r$ 

```

În cazul în care algoritmiile sunt complicaţi, este permis ca primele versiuni ale pseudocodului să fie scrise într-un pseudocod simplificat, ce permite o scriere de tip matriceal. Într-un astfel de pseudocod un tablou unidimensional va fi notat cu literă mică scrisă cu font aldin, iar un tablou bidimensional va fi notat cu literă mare scrisă cu font aldin. Tabloul unidimensional poate fi primit ca un vector coloană, iar cel bidimensional ca o matrice. Astfel, funcţia ce calculează produsul scalar de mai sus ar putea fi scrisă simplificat astfel:

```

funcţie produs_scalar( $\mathbf{v}, \mathbf{w}$ )
întoarce  $\mathbf{v}^T \cdot \mathbf{w}$ 

```

În acest caz, nici nu mai este nevoie să se complice pseudocodul programului principal cu instrucţiunea

```
p = produs_scalar(a, b)
```

ci se poate scrie simplu

$$p = \mathbf{a}^T \cdot \mathbf{b}$$

Pe parcursul acestei cărți se vor folosi ambele variante de pseudocod, de fiecare dată fiind folosită varianta ce permite explicarea mai clară a noțiunilor.

În momentul în care un algoritm este definitivat și se dorește implementarea lui într-un limbaj de programare, atunci stilul de lucru va depinde și de limbajul ales. Următoarele două paragrafe oferă sugestii de lucru utilizatorilor de C și, respectiv, Matlab.

1.1.2 Sugestii pentru utilizatorii de C

Metodele care vor fi studiate în această carte conduc în exclusivitate la algoritmi numerici, care efectuează de cele mai multe ori calcule cu vectori și matrice. În acest paragraf se vor face câteva considerații asupra declarării și alocării vectorilor și matricelor în C.

În C există o strânsă corespondență între adrese (pointeri) și tablouri. De exemplu, în cazul unui tablou unidimensional declarat ca

```
float a[4]
```

valoarea reprezentată de `a[j]` este același lucru cu `*(a+j)` adică ”conținutul adresei obținute incrementând pointer-ul `a` cu `j`”. O consecință a acestei definiții este aceea că dacă `a` este adresa unei locații valide, atunci `a[0]` este întotdeauna definit. Tablourile unidimensionale sunt numerotate în C, în mod natural, începând de la 0. Șirul definit mai sus are referințele valide `a[0]`, `a[1]`, `a[2]` și `a[3]`, dar nu și `a[4]`.

Mulți algoritmi sunt însă descriși în mod natural cu indici care încep de la 1. Cu siguranță că acești algoritmi pot fi modificați, dar aceasta presupune o aritmetică suplimentară în operarea cu indici, lucru care nu este prea plăcut. Se poate folosi însă puterea limbajului C pentru ca această problemă să dispară. Ideea este simplă:

```
float a[4], *aa;
aa = a - 1;
```

Pointer-ul `aa` indică acum o locație înaintea lui `a`. În consecință, elementele `aa[1]`, `aa[2]`, `aa[3]` și `aa[4]` există și vectorul `aa` are indici ce pornesc de la 1.

Uneori este convenabil ca vectorii să aibă indici care pornesc de la 0, iar alteori este convenabil ca indicii să înceapă cu 1. De exemplu, coeficienții unui polinom $a_0 + a_1x + \dots + a_nx^n$ sunt descriși în mod natural de un vector cu indici ce încep cu 0, pe când termenul liber al unui sistem de ecuații $b_i, i = 1, \dots, n$ este descris în mod natural cu un vector cu indici ce încep de la 1.

Pentru a evita rescrierea algoritmilor ce sunt deduși în mod natural cu indici ce pornesc de la 1, se poate folosi o funcție cu următoarea definiție.

```

typedef float *VECTOR;
VECTOR vector (int nl, int nh)
{
    VECTOR v;
    v = (float *) malloc ((unsigned) (nh - nl + 1) * sizeof (float));
    if (!v)
        perror ("Eroare de alocare in vector()");
    return v - nl;
}

```

Această funcție alocă un vector de variabile de tip `float`, care vor fi accesate cu indici cuprinși între `nl` și `nh`. O utilizare tipică a acestei funcții este

```

VECTOR a;
a = vector(1,4);

```

Această funcție precum și funcții similare ce alocă vectori de întregi se găsesc în fișierul `nrutil.c` pe care îl găsiți⁴ în anexa A. Acest fișier conține și rutinele corespunzătoare de dealocare. De exemplu, pentru dealocarea memoriei ocupate de vectorul definit mai sus, instrucțiunea este

```

free_vector(a,1,4);

```

Problema indicilor ce pornesc de la 0 sau de la 1 apare și în cazul matricelor. În sintaxa C, lucrul cu tabele bidimensionale este puțin mai complicat. Fie o valoare reală `a[i][j]` unde `i` și `j` sunt întregi. Un compilator de C va genera coduri mașină diferite pentru această referință, aceasta depinzând de declarația pentru variabila `a`. Dacă `a` a fost declarată de dimensiune fixă, de exemplu `float a[2][4]` atunci codul mașină ar putea fi descris astfel: "la adresa `a` adună de 4 ori `i`, apoi adună `j` și întoarce valoarea astfel adresată. Observați că valoarea constantă 4 trebuie cunoscută pentru a efectua în mod corect calculele.

Dacă `a` a fost declarat ca `float **a`, atunci codul mașină `a[i][j]` este "la adresa `a` adună `i`, valoarea astfel adresată consider-o o nouă adresă, la care adună `j` și întoarce valoarea astfel adresată.

Se observă că în al doilea caz nu este necesară cunoașterea dimensiunii matricei și nu este nevoie de nici o înmulțire. O indirectare suplimentară înlocuiește aceste informații.

În concluzie, tablourile de dimensiune fixă trebuie evitate. Ele nu sunt structuri de date potrivite pentru reprezentarea vectorilor și matricelor în calculele științifice.

⁴Acesta este stilul de lucru al cărții *Numerical Recipes in C* [9], <http://www.nr.com/>.

Ca exemplificare a celor spuse mai sus, iată traducerea în C a pseudocodului algoritmului ce calculează produsul scalar a doi vectori, descris la sfârșitul paragrafului anterior.

```
#include "nrutil.h"
void citeste_vector (int N, VECTOR x);
float produs_sclar (int N, VECTOR v, VECTOR w);

void main (void)
{
/* program principal */
  int N; /* dimensiunea vectorilor */
  VECTOR a, b; /* vectorii de intrare */
  float p; /* produsul scalar dintre a si b */
  int i;
  printf ("\n Introduceti dimensiunea vectorilor ");
  scanf ("%d", &N);
  a = vector (1, N); /* aloca spatiu de memorie pentru vectori */
  b = vector (1, N);

  printf ("\n Introduceti componentele vectorului a ");
  citeste_vector (N, a);
  printf ("\n Introduceti componentele vectorului b ");
  citeste_vector (N, b);
  p = produs_sclar (N, a, b);
  printf ("\n Rezultatul este %f \n", p);

  free_vector (a, 1, N); /* elibereaza spatiul de memorie */
  free_vector (b, 1, N);
}

void citeste_vector (int N, VECTOR x)
{
/* citeste un vector x de dimensiune N, indicii incep de la 1 */
  int i;
  for (i = 1; i <= N; i++)
  {
    printf ("\n componenta %d = ", i);
    scanf ("%f", &x[i]);
  }
}
```

```
    }  
}  
  
float produs_sclar (int N, VECTOR v, VECTOR w)  
{  
/* calculeaza produsul scalar a doi vectori v si w,  
   ai caror indici incep de la 1 */  
  int i;  
  float r;  
  r = 0;  
  for (i = 1; i <= N; i++)  
    r = r + v[i] * w[i];  
  return r;  
}
```

1.1.3 Sugestii pentru utilizatorii de Matlab

În limbaje de programare de tip Matlab, Scilab, Octave, implementările pot arăta chiar mai simplu decât algoritmul prezentat în sintaxa pseudocodului descrisă anterior. Acest lucru se datorează pe de o parte faptului că în aceste limbaje nu sunt necesare declarațiile variabilelor, iar pe de altă parte în aceste limbaje se poate opera direct cu vectori și matrice.

Vom ilustra aceste afirmații folosind exemplul produsului scalar. O traducere exactă în Matlab a algoritmului prezentat este următoarea:

```
% program principal, fisier main_ps.m  
clear all;  
N = input('Introduceti N = ');  
a = zeros(N,1); % alocare de memorie pentru vectorul a  
b = zeros(N,1); % alocare de memorie pentru vectorul b  
a = citeste_vector(N,a);  
b = citeste_vector(N,b);  
p = produs_sclar(N,a,b);  
disp(sprintf('p = %g',p));  
  
% fisier citeste_vector.m  
function a = citeste_vector(N,a)  
for i = 1:N
```

```

    a(i) = input(sprintf('Introduceti componenta %d = ',i));
end

```

```

% fisier produs_sclar.m
function p = produs_sclar(N,a,b)
p = 0;
for i = 1:N
    p = p + a(i)*b(i);
end

```

Funcția de citire se poate simplifica având în vedere că sintaxa comenzii `input` permite chiar introducerea vectorului, nu numai a unei componente:

```

% fisier citeste_vector.m - varianta a doua
function a = citeste_vector(N,a)
a = input(sprintf('Introduceti un vector coloana de dimensiune %d :',N));
if or(size(a,1)~=N,size(a,2)~=1)
    error('Eroare');
end

```

De asemenea, calculul produsului scalar se poate face mult mai simplu, știind că se pot folosi operații cu vectori și matrice.

```

% fisier citeste_vector.m - varianta a doua
function p = produs_sclar(N,a,b)
p = a'*b;

```

În Matlab de altfel, codurile ce folosesc operații cu matrice sunt mai eficiente din punct de vedere al timpului de calcul decât folosirea unui cod echivalent ce detaliază toate operațiile, așa cum s-ar fi făcut în C. Operarea cu matrice ascunde însă operațiile elementare și, în consecință, va face mai dificilă estimarea complexității codului.

Exemplul de mai sus ilustrează doar câteva aspecte ale implementării în Matlab. De fapt, funcția `produs_sclar` nici nu mai are nevoie de N , primul argument de intrare⁵. Mai mult, este lipsit de sens și chiar ineficient să se scrie funcții care au o singură linie de cod deoarece se pierde timp și la transferul execuției de la programul principal la funcție și înapoi.

⁵Pentru a curăța programul de astfel de lucruri inutile, este deosebit de folositor să se folosească comanda Matlab `mlint` care permite identificarea variabilelor inutile și chiar a liniilor de cod cu probleme.

1.2 Complexitatea algoritmilor

În general, o problemă admite mai multe metode de rezolvare, și în consecință, vor exista mai mulți algoritmi diferiți pentru rezolvarea ei cu ajutorul calculatorului. De aceea, este absolut necesar să se facă o evaluare a algoritmilor pentru a stabili care dintre ei este cel mai bun pentru o problemă dată.

Nu există un singur criteriu de evaluare a algoritmilor. Algoritmul ideal trebuie să fie simplu, să dea o soluție corectă într-un timp scurt și să ocupe o zonă mică de memorie. Nu există însă nici un algoritm care să rezolve perfect o problemă, fără nici un fel de eroare și atunci se va urmări ca erorile pentru o anumită aplicație să fie rezonabile. O analiză a erorilor posibile dintr-un algoritm este prezentată în paragraful 1.3. De asemenea, criteriul referitor la timpul de calcul intră de multe ori în contradicție cu criteriul referitor la memoria necesară algoritmului.

În consecință, la alegerea unui algoritm potrivit pentru o aplicație dată, trebuie făcut un compromis între trei criterii: timpul de calcul necesar obținerii soluției, necesarul de memorie și acuratețea soluției. Analiza primelor două criterii furnizează informații despre complexitatea algoritmului.

1.2.1 Timp de calcul

Timpul de calcul al unui algoritm depinde de complexitatea problemei de rezolvat, de performanțele intrinseci ale calculatorului și limbajului de programare folosit și evident, de algoritm. Este util să poată fi apreciată doar calitatea algoritmului și nu a problemei sau a mediului în care se lucrează. De aceea s-a inventat conceptul de complexitate a unui algoritm din punct de vedere al timpului de calcul.

Complexitatea unui algoritm din punct de vedere al timpului de calcul este relația dintre timpul de calcul exprimat în număr de operații elementare și dimensiunea problemei.

Dimensiunea problemei depinde de problema studiată. De exemplu, pentru calculul produsului scalar a doi vectori, dimensiunea problemei este dimensiunea vectorilor. Pentru rezolvarea unui sistem de ecuații algebrice liniare, dimensiunea problemei este dimensiunea sistemului. Uneori, este potrivit să exprimăm dimensiunea problemei în funcție de două numere în loc de unul. De exemplu, dacă datele de intrare reprezintă graful unui circuit, dimensiunea problemei este reprezentată de perechea alcătuită din numărul de noduri și numărul de laturi.

Timpul de calcul este de fapt suma timpilor necesari pentru executarea tuturor instrucțiunilor algoritmului. Nu toate instrucțiunile durează la fel de mult, de aceea, estimarea complexității nu se poate face foarte precis, dar nici nu este necesar acest lucru. Se alege o

operație elementară, considerată a fi cea care durează cel mai mult (de exemplu evaluarea unei anumite funcții) și se numără câte astfel de operații elementare sunt executate.

Algoritmi polinomiali

Să considerăm următorul fragment de pseudocod, corespunzător calculului unui produs scalar $p = \sum_{i=1}^n a_i b_i$.

$p = 0$;

pentru $i = 1, n$

$p = p + a_i b_i$

•

Considerând ca operație elementară orice operație algebrică (adunare, scădere, înmulțire, împărțire) și neglijând timpul de calcul petrecut în declarații și atribuiri, rezultă că în algoritmul de mai sus se fac $2n$ operații elementare, câte două (o adunare și o înmulțire) pentru fiecare valoare a contorului i . Timpul de calcul este proporțional cu dimensiunea problemei, în acest caz n . Un astfel de algoritm se spune că este un *algoritm liniar, sau de ordinul 1* și se notează $T = O(n)$. În cazul produsului scalar a doi vectori de dimensiune n , putem scrie $T = O(2n) \approx O(n)$. Constanta 2 nu este atât de relevantă, ceea ce este important este dependența de dimensiunea problemei.

Fie acum cazul înmulțirii unei matrice pătrate a de dimensiune n cu un vector coloană x . Rezultatul este un vector coloană ale cărui componente se calculează ca $b_i = \sum_{j=1}^n a_{ij} x_j$.

pentru $i = 1, n$

$b_i = 0$

pentru $j = 1, n$

$b_i = b_i + a_{ij} x_j$

•

•

În acest caz, pentru fiecare i se fac $2n$ operații (raționând exact ca la produsul scalar), deci pentru toate cele n valori ale lui i se vor face $2n^2$ operații elementare. Un algoritm pentru care timpul de calcul T este proporțional cu pătratul dimensiunii problemei n se spune că este un *algoritm pătratic, sau de ordinul 2* și se notează $T = O(n^2)$. Înmulțirea dintre o matrice și un vector are deci complexitatea $T = O(2n^2) \approx O(n^2)$.

Să considerăm acum cazul înmulțirii a două matrice pătrate de dimensiune n . Rezultatul este o matrice pătrată ale cărei componente sunt $c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$.

```

pentru  $i = 1, n$ 
  pentru  $j = 1, n$ 
     $c_{ij} = 0$ 
    pentru  $k = 1, n$ 
       $c_{ij} = c_{ij} + a_{ik}b_{kj}$ 
    •
  •
•

```

În acest caz se fac $2n^3$ operații elementare. Un algoritm pentru care timpul de calcul T este proporțional cu cubul dimensiunii problemei n se spune că este un *algoritm cubic*, sau *de ordinul 3* și se notează $T = O(n^3)$. Înmulțirea a două matrice pătrate are complexitatea $T = O(2n^3) \approx O(n^3)$.

Se spune că un algoritm este *polinomial de ordin k* din punct de vedere al timpului de calcul, și se notează $T = O(n^k)$ dacă și numai dacă există o constantă $C > 0$ și un număr n_0 astfel încât $T \leq Cn^k$ pentru orice $n \geq n_0$.

Un algoritm pentru care timpul de calcul nu depinde de dimensiunea problemei se numește *algoritm de ordin zero* și pentru el se scrie $T = O(1)$.

Ordonarea algoritmilor s-ar putea face pe o scară astfel: $O(1) < O(\log n) < O(n) < O(n \log n) < O(n^2) < \dots < O(e^n) < O(n!)$. Între algoritmi de ordin zero și algoritmi liniari se afla algoritmi *logaritmici*, al căror timp de calcul este proporțional cu logaritmul dimensiunii problemei. Între algoritmi liniari și cei pătratici se afla algoritmi *liniaritmici* având timpul de calcul majorat de la un moment dat de o funcție proporțională cu $n \log n$. Cei mai costisitori, și de altfel lipsiți de importanță practică, sunt algoritmi nonpolinomiali. Rezolvarea unui sistem de ecuații algebrice liniare cu regula lui Cramer conduce de exemplu la un algoritm cu complexitate factorială. La fel ar fi și algoritmul unui jucător de șah perfect, care la orice mutare analizează toate combinațiile posibile de pe tabla de șah. În practică sunt utili cel mult algoritmi având o complexitate cubică din punct de vedere a timpului de calcul.

Dacă analizăm cele trei exemple anterioare, observăm că gradul de indentare este în strânsă corelare cu ordinul algoritmului. Cu cât gradul de indentare este mai mare, cu atât este de așteptat ca ordinul algoritmului să fie mai mare. Acest lucru este adevărat dacă operația elementară apare în instrucțiunea indentată cel mai mult.

Vom exemplifica acum evaluarea complexității timpului de calcul în cazul a trei algoritmi diferiți pentru evaluarea unui polinom de grad n , cu coeficienți reali

$$P(x) = a_0 + a_1x + \dots + a_nx^n. \quad (1.1)$$

Relația (1.1), scrisă compact

$$P(x) = a_0 + \sum_{i=1}^n a_i x^i \quad (1.2)$$

se poate implementa "natural" astfel (declarațiile au fost omise):

; Varianta 1

$P = a_0$;

pentru $i = 1, n$

$t = a_i$

pentru $j = 1, i$

$t = t * x$

•

$P = P + t$;

•

Considerând ca operație elementară orice operație algebrică, rezultă că pentru fiecare valoare i se fac $i + 1$ operații elementare (i înmulțiri și o adunare). Cum i ia valori de la 1 la n rezultă că numărul total de operații elementare este $\sum_{i=1}^n (i + 1) \approx n^2/2$. Ordinul de complexitate al acestui algoritm va fi deci $T_1 = O(n^2/2) \approx O(n^2)$, așteptat pentru că există două nivele de indentare.

O altă variantă de implementare poate calcula fiecare termen al sumei mai simplu, în funcție de termenul adunat anterior. Ea se bazează pe scrierea, echivalentă din punct de vedere matematic

$$P(x) = a_0 + a_1 t_1 + a_2 t_2 + \dots + a_n t_n = a_0 + \sum_{i=1}^n a_i t_i, \quad (1.3)$$

unde

$$t_1 = x, \quad t_i = t_{i-1} x, \quad i = 2, \dots, n. \quad (1.4)$$

Pentru această scriere matematică, implementarea naturală este

; Varianta 2

$P = a_0$;

$t = 1$;

pentru $i = 1, n$

$t = t * x$

$P = P + a_i * t$

•

Deoarece există un singur nivel de indentare, este de așteptat ca algoritmul să fie liniar. Într-adevăr, pentru fiecare i se execută 3 operații elementare (două înmulțiri și o adunare), deci complexitatea algoritmului este $T_2 = O(3n) \approx O(n)$.

O altă idee de implementare este sugerată de efectuarea calculelor conform relației

$$P(x) = a_0 + x(a_1 + x(a_2 + \dots x(a_{n-1} + a_n x) \dots)). \quad (1.5)$$

Implementarea corespunzătoare acestei expresii este

; Varianta 3

$P = a_n$;

pentru $i = n - 1, 0, -1$

$P = a_i + P * x$

•

Această variantă este cea mai eficientă, ordinul de complexitate fiind $T_3 = O(2n) \approx O(n)$. Din punct de vedere matematic, cele trei relații (1.1), (1.3) și (1.5) sunt echivalente, dar algoritmi corespunzători sunt diferiți din punct de vedere al complexității lor.

Un algoritm se poate îmbunătăți din punct de vedere al timpului de calcul dacă operația elementară este scoasă din instrucțiunea indentată cel mai mult. Un alt aspect important este acela de a nu repeta operațiile efectuate o dată. Pentru a clarifica aceste afirmații, să considerăm următorul exemplu.

; Varianta A

întreg i, j, n

real a, b

...

tablou real $c[n][n]$

pentru $i = 1, n$

pentru $j = 1, n$

$c_{ij} = f(i * a) + f(j * b)$; f este o funcție definită în altă parte

•

•

Considerând ca operație elementară evaluarea funcției f , rezultă o complexitate a algoritmului $T_A = O(2n^2)$.

Acest cod se poate simplifica observând că termenul $f(i * a)$ este evaluat de mai multe ori, chiar pentru aceeași valoare a lui i . De aceea, se poate scoate în afara ciclului în

j evaluarea funcției f pentru argumentul $i * a$ rezultând următoarea variantă de implementare. Pentru aceasta este însă nevoie de o variabilă suplimentară, notată mai jos cu p .

; Varianta B

întreg i, j, n

real a, b, p

...

tablou real $c[n][n]$

pentru $i = 1, n$

$p = f(i * a)$

pentru $j = 1, n$

$c_{ij} = p + f(j * b)$

•

•

Pentru această implementare, numărul de apeluri ale funcției f este $T_B = O(n(n+1)) = O(n^2 + n) \approx O(n^2)$. Pentru valori mari ale lui n această variantă de implementare este de două ori mai rapidă decât prima variantă, dar ambii algoritmi sunt pătratici.

Implementarea poate fi simplificată și mai mult, dacă se observă că, pentru fiecare valoare a contorului i , funcția f este evaluată din nou în $b, 2b, \dots, nb$. Se repetă astfel calcule inutile, consumatoare de timp. Este evident că este nevoie ca funcția f să fie evaluată în doar $2n$ puncte: $a, 2a, \dots, na, b, 2b, \dots, nb$. Următoarea variantă de implementare execută un număr minim de evaluări, cu prețul memorării lor anterioare asamblării rezultatului.

; Varianta C

întreg i, j, n

real a, b

...

tablou real $c[n][n]$

tablou real $p[n], q[n]$

pentru $i = 1, n$

$p_i = f(i * a)$

$q_i = f(i * b)$

•

pentru $i = 1, n$

pentru $j = 1, n$

$$c_{ij} = p_i + q_j$$

-
-

Într-adevăr, în acest caz complexitatea algoritmului este $T_C = O(2n)$. Algoritmul este acum liniar. Reducerea timpului de calcul s-a făcut pe baza alocării suplimentare de memorie pentru două tablouri unidimensionale.

Notații asimptotice

În cele de mai sus, am discutat și exemplificat cazul algoritmilor polinomiali pentru care timpul de calcul este mărginit superior de o funcție de forma Cn^k , cel puțin de la o anumită valoare a lui n în sus. Această definiție se poate generaliza astfel[10].

Un algoritm are *ordinul de complexitate* $O(g(n))$ din punct de vedere al timpului de calcul T (exprimat în unități de timp, de exemplu secunde) dacă și numai dacă **există** o constantă pozitivă $C > 0$ și un număr n_0 astfel încât $T \leq Cg(n)$ pentru orice $n \geq n_0$. *Notația* O descrie deci o margine superioară pentru timpul de execuție al unui program, aspect ilustrat în figura 1.1.

Notația Ω definește o delimitare asimptotică inferioară. Un algoritm are *ordinul de complexitate* $\Omega(g(n))$ din punct de vedere al timpului de calcul dacă și numai dacă **există** o constantă pozitivă $C > 0$ și un număr n_0 astfel încât $Cg(n) \leq T$ pentru orice $n \geq n_0$. Această definiție este ilustrată în figura 1.2.

Notația Θ delimitează asimptotic atât inferior cât și superior, timpul de execuție al unui algoritm. Se mai spune că notația Θ furnizează o margine asimptotică *strânsă*. Un algoritm are *ordinul de complexitate* $\Theta(g(n))$ din punct de vedere al timpului de calcul

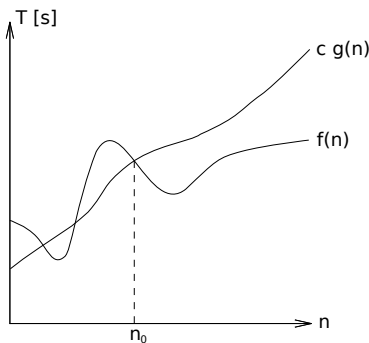


Figura 1.1: Notația O :
 $T = f(n) = O(g(n))$.

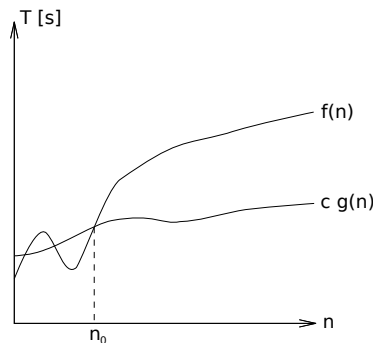


Figura 1.2: Notația Ω :
 $T = f(n) = \Omega(g(n))$.

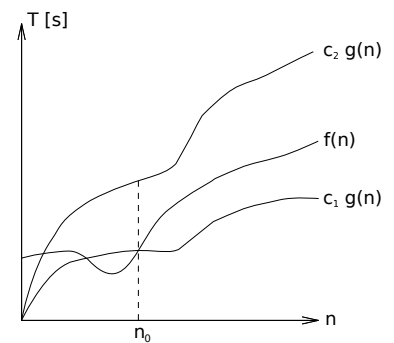


Figura 1.3: Notația Θ :
 $T = f(n) = \Theta(g(n))$.

dacă și numai dacă **există** două constante pozitive C_1 și C_2 și un număr n_0 astfel încât $C_1g(n) \leq T \leq C_2g(n)$ pentru orice $n \geq n_0$. Această definiție este ilustrată în figura 1.3.

Se poate demonstra următoarea echivalență

$$T = \Theta(g(n)) \Leftrightarrow T = O(g(n)) \text{ și } T = \Omega(g(n)).$$

De asemenea, se poate arăta că un algoritm pentru care timpul de calcul se exprimă printr-un polinom $T = \sum_{i=0}^k a_i n^i$ unde $a_k > 0$ este un algoritm de complexitate $T = \Theta(n^k)$, deci, conform echivalenței anterioare $T = O(n^k)$ și $T = \Omega(n^k)$.

O margine asimptotică superioară, dată de notația O poate să fie sau nu o margine asimptotic strânsă. Pentru a desemna o margine asimptotic superioară care nu este strânsă se folosește *notația o* .

Un algoritm are *ordinul de complexitate $o(g(n))$* din punct de vedere al timpului de calcul dacă și numai dacă pentru **orice** constantă pozitivă $C > 0$ există un număr n_0 astfel încât $T \leq Cg(n)$ pentru orice $n \geq n_0$. Definițiile pentru notația O și o sunt similare. Principala diferență este aceea că la notația O relația are loc pentru o anumită constantă C , în timp ce la notația o , relația are loc pentru toate constantele C . Intuitiv, în cazul notației o , timpul de calcul $T = T(n)$ devine neglijabil față de $g(n)$ atunci când n tinde la infinit: $\lim_{n \rightarrow \infty} T(n)/g(n) = 0$.

Similar, pentru marginea inferioară se folosește *notația ω* . Un algoritm are *ordinul de complexitate $\omega(g(n))$* din punct de vedere al timpului de calcul dacă și numai dacă pentru **orice** constantă pozitivă $C > 0$ există un număr n_0 astfel încât $Cg(n) \leq T$ pentru orice $n \geq n_0$. Aceasta implică faptul că T devine oricât de mare față de $g(n)$ atunci când n tinde la infinit: $\lim_{n \rightarrow \infty} T(n)/g(n) = \infty$.

Această analiză rafinată a ordinului de complexitate este importantă mai ales atunci când performanța algoritmului depinde și de modul în care sunt furnizate datele de intrare. De exemplu, pentru un algoritm de sortare a unui șir de numere reale, analiza cazului cel mai favorabil (în care vectorul de intrare este gata sortat) va furniza informații despre complexitatea notată cu Ω în timp ce analiza cazului cel mai defavorabil (în care vectorul este sortat în ordine inversă) va furniza informații despre complexitatea notată cu O .

De exemplu, fie un șir de n valori reale, ordonate crescător $x_1 < x_2 < \dots < x_n$ și un alt număr real x_{crt} . Dorim să știm în care subinterval $[x_k, x_{k+1}]$ se află valoarea x_{crt} .

Cu siguranță ideea de căutare naturală este de a parcurge sub-intervalele ca în pseudocodul următor care întoarce, în cazul în care x_{crt} este între x_1 și x_n , indicele k pentru care $x_{crt} \in [x_k, x_{k+1}]$.

funcție caută_v0(n, x, x_{crt})

; caută plasarea valorii x_{crt} în șirul ordonat x cu n valori

```

întreg n
tablou real x[n] ; şirul este presupus ordonat crescător
real xcrt
...
dacă (xcrt < x1)
    rez = 1
altfel dacă (xcrt > xn)
    rez = n - 1
altfel
    flag = 0
    k = 1
    cât timp ((k <= n - 1) și (flag = 0))
        dacă xcrt ≤ xk+1
            rez = k
            flag = 1
        altfel
            k = k + 1
    •
    •
    •
    întoarce rez

```

Considerând ca operație de referință compararea a două numere reale, este evident că ordinul de complexitate în cazul cel mai defavorabil este $T = O(n)$ și în cazul cel mai favorabil $T = \Omega(1)$. Căutarea se face mai eficient dacă se adoptă însă o strategie bazată pe înjumătățirea numărului de subintervale în care are loc căutarea. Astfel, la fiecare iterație se elimină jumătate din numărul de sub-intervale posibile, întorcându-se indexul capătului din stânga al subintervalului. O astfel de procedură este utilă de exemplu în cazul în care se dorește interpolarea pe porțiuni a unei funcții descrise cu ajutorul unui tabel de valori. În cazul în care valoarea *xcrt* nu se află în $[x_1, x_n]$ atunci rezultatul returnat este 1 dacă punctul se află în stânga domeniului și $n - 1$ dacă punctul se află în dreapta domeniului.

```

funcție caută(n, x, xcrt)
...
dacă (xcrt < x1)
    rez = 1
altfel dacă (xcrt > xn)

```



```

    rez = n - 1
altfel
    k1 = 1
    k2 = n
    cât timp k2 - k1 ≠ 1
        km = [(k1 + k2)/2]    ; [...] este partea întreagă
        dacă xcrt < xkm
            k2 = km
        altfel
            k1 = km
    •
•
rez = k1
întoarce rez

```

În acest caz, la fiecare iterație se elimină jumătate din intervale. Numărul de intervale eliminate după k iterații este $n/2 + n/4 + n/8 + \dots + n/2^k$. Căutarea durează cel mult până se elimină $n - 1$ subintervale. Din condiția $n/2 + n/4 + n/8 + \dots + n/2^k = n - 1$ rezultă $2^k = n$ de unde $k = \log_2(n)$. Algoritmul are în consecință o complexitate logaritmică.

Analiza complexității algoritmilor recursivi

Un algoritm recursiv este un algoritm care se apelează pe el însuși. Analiza complexității unui astfel de algoritm este mai dificilă decât în cazul algoritmilor nerecursivi. Ideea principală este aceea de a împărți problema în subprobleme de același tip. Pentru analiza complexității algoritmilor recursivi se scrie mai întâi o relație de recurență pentru timpul de calcul. Această relație leagă complexitatea problemei inițiale $T(n)$ de complexitatea unei subprobleme $T(n/b)$ unde b este numărul de subprobleme în care se descompune problema. Pentru aceasta trebuie evaluate numărul de apeluri recursive (fie acest număr a) precum și complexitatea instrucțiunilor de trecere la problemă la subproblemă, iar

$$T(n) = aT(n/b) + C(n).$$

Există mai multe tipuri de metode pentru rezolvarea evaluării acestei recurențe [10]. În cele ce urmează se va exemplifica *metoda iterației* pentru estimarea complexității algoritmului de căutare binară în varianta recursivă.

funcție caută_v2($n, x, xcrt$)

```

...
dacă ( $x_{crt} < x_1$ )
    rez = 1
altfel dacă ( $x_{crt} > x_n$ )
    rez =  $n - 1$ 
altfel
rez = binary_search( $x, x_{crt}, 1, n$ );
•
întoarce rez
;
funcție binary_search( $x, x_{crt}, k1, k2$ )
; presupuneri:
;  $x$  - ordonate crescător,  $k1 < k2$ 
;  $x_{crt}$  se afla în interiorul vectorului  $x$ 
...
 $km = \lfloor (k1 + k2)/2 \rfloor$ 
dacă  $k2 = k1 + 1$ 
    rez =  $km$ 
altfel dacă  $x_{km} > x_{crt}$ 
    rez = binary_search( $x, x_{crt}, k1, km$ )
altfel
    rez = binary_search( $x, x_{crt}, km, k2$ )
•
întoarce rez

```

Pentru acest algoritm recursiv, la fiecare iterație dimensiunea problemei se reduce la jumătate, iar timpul necesar pentru a decide care jumătate să fie eliminată este constant, nu depinde de dimensiunea problemei. De aceea putem scrie

$$T(n) = T(n/2) + O(1).$$

Pentru a simplifica raționamentul, vom presupune că n este o putere a lui 2: $n = 2^k$. Urmează că

$$T(n) \leq T(n/2) + C \leq T(n/4) + 2C \leq \dots \leq T(n/2^k) + kC = T(1) + C \log(n) = O(\log(n)).$$

Se poate arăta că acest algoritm are această complexitate pentru orice valoare a lui n , nu numai pentru valori de tipul $n = 2^k$.

1.2.2 Necesari de memorie

Complexitatea unui algoritm din punct de vedere al necesarului de memorie este dependența dintre necesarul de memorie exprimat în număr de locații elementare de memorie și dimensiunea problemei.

De obicei, o locație elementară de memorie este cea corespunzătoare unui număr real.

Se spune că un algoritm este *polinomial de ordin k* din punct de vedere al necesarului de memorie M , și se notează $M = O(n^k)$ dacă și numai dacă există o constantă C astfel încât $M \leq Cn^k$, unde n este dimensiunea problemei.

Necesarul de memorie se stabilește cu ușurință inspectând declarațiile făcute în pseudocod. Pentru cele trei variante (A, B, C) prezentate în paragraful 1.2.1 rezultă $M_A = O(n^2 + 2) \approx O(n^2)$ necesară pentru memorarea tabloului bidimensional c și a variabilelor reale a și b , $M_B = O(n^2 + 3) \approx O(n^2)$ care diferă de varianta A doar prin memorarea unei variabile suplimentare reale p și $M_C = O(n^2 + 2n + 2) \approx O(n^2)$ care diferă de varianta A prin memorarea a două tablouri unidimensionale de dimensiune n . Toate cele trei variante sunt, din punct de vedere al necesarului de memorie, de ordin 2. Câștigul variantei C este de aceea remarcabil, căci timpul de calcul a scăzut cu un ordin de mărime, pe seama unui necesar de memorie suplimentar, ne semnificativ la valori n mari ale dimensiunii problemei.

În concluzie, elaborarea unui algoritm înseamnă întotdeauna stabilirea unui compromis între timp de calcul și necesar de memorie.

1.3 Erori în calculele numerice

Unul din criteriile de evaluare a unui algoritm îl reprezintă eroarea cu care se obține rezultatul. Erorile nu pot fi înlăturate total dintr-un calcul deoarece, de exemplu, chiar numerele reale nu pot fi reprezentate în calculator cu o precizie infinită. Numărul real π are o infinitate de cifre semnificative, iar orice sistem de calcul lucrează cu un număr finit de cifre semnificative. Pe de altă parte, erorile se propagă în calcule, astfel încât, chiar dacă datele de intrare ale unui program sunt foarte precise, rezultatul poate avea doar câteva cifre semnificative corecte. De aceea este foarte important ca pentru orice algoritm să se estimeze eroarea rezultatului. În acest paragraf sunt prezentate tipurile de erori care pot afecta un algoritm și modul în care ele pot fi evaluate. În final sunt discutate două aspecte importante legate de erori: condiționarea unei probleme și stabilitatea unui algoritm.

1.3.1 Tipuri de erori

Erorile dintr-un algoritm se pot clasifica în funcție de tipul cauzelor care le generează în: erori inerente, erori de rotunjire și erori de trunchiere.

Erorile inerente sunt erorile datorate reprezentării imprecise a datelor de intrare. Datele de intrare ale unui algoritm pot proveni de exemplu din măsurători și de aceea ele sunt afectate de erori. Aceste erori nu pot fi eliminate și în consecință este important să se poată evalua modul în care se propagă ele în calculele numerice, astfel încât să poată fi făcută o estimare a erorii rezultatului în funcție de erorile datelor de intrare.

Erorile de rotunjire se datorează reprezentării finite a numerelor reale în calculator. Sistemele de calcul nu pot lucra decât cu aproximări raționale ale numerelor reale. Numerele reale sunt reprezentate cu un număr finit de cifre semnificative. Nici aceste erori nu pot fi eliminate și, ca și în cazul erorilor inerente, este importantă evaluarea modului în care ele afectează rezultatul final.

Erorile de trunchiere provin din finitudinea inerentă unui algoritm corect construit. Există metode matematice a căror soluție exactă ar necesita efectuarea unui număr infinit de calcule. Un astfel de exemplu este sumarea unei serii. Deoarece implementările nu pot fi făcute decât pentru algoritmi care fac un număr finit de calcule, nu are sens să permitem în pseudocod algoritmi infiniti. În acest caz este importantă evaluarea erorii care se face datorită trunchierii procesului infinit.

Pentru analiza cantitativă a acestor erori, este utilă definirea următoarelor mărimi.

Eroarea absolută \mathbf{e}_x a unei mărimi din \mathbb{R}^n este diferența dintre valoarea aproximativă \bar{x} și valoarea exactă \mathbf{x} a mărimii

$$\mathbf{e}_x = \bar{x} - \mathbf{x}. \quad (1.6)$$

Este evident că o astfel de mărime nu se poate calcula deoarece valoarea exactă nu este cunoscută. De aceea, este mai utilă aflarea unei *marginii a erorii absolute* a_x , adică a unei mărimi $a_x \in \mathbb{R}$ care majorează norma erorii absolute

$$\|\mathbf{e}_x\| \leq a_x. \quad (1.7)$$

Dacă presupunem că mărimea este scalară ($n = 1$), atunci din (1.6) și (1.7) rezultă că

$$\bar{x} - a_x \leq x \leq \bar{x} + a_x. \quad (1.8)$$

Altfel spus, cunoașterea marginii erorii absolute permite definirea intervalului în care este plasată soluția exactă. Relația (1.8), echivalentă cu $x \in [\bar{x} - a_x, \bar{x} + a_x]$, se scrie pe scurt sub forma

$$x = \bar{x} \pm a_x. \quad (1.9)$$

Dezavantajul folosirii erorii absolute este acela că valoarea numerică depinde de sistemul de unități de măsură folosit, făcând dificilă aprecierea gradului de acuratețe a soluției. De aceea, se preferă folosirea unei mărimi relative, invariantă la sistemul de unități de măsură.

Eroarea relativă ε_x a unei mărimi se definește ca fiind raportul dintre eroarea absolută și norma mărimii

$$\varepsilon_x = \frac{\mathbf{e}_x}{\|\mathbf{x}\|}. \quad (1.10)$$

Desigur, această mărime nu este definită pentru o valoare exactă nulă. Mai mult, ea nu se poate calcula deoarece nici eroarea absolută, nici mărimea exactă nu sunt cunoscute. De aceea se preferă folosirea unei *marginii a erorii relative* r_x , mărime care satisface

$$\|\varepsilon_x\| \leq r_x. \quad (1.11)$$

Cel mai adesea, marginea erorii relative se exprimă în procente, iar relația (1.11) se scrie sub forma

$$">\mathbf{x} = \bar{\mathbf{x}} \pm r_x \%$$
" . \quad (1.12)

Ca exemplu, valoarea exactă a numărului π este $x = 3.1415\dots$, iar valoarea aproximativă cea mai uzuală este $\bar{x} = 3.14$. Rezultă o eroare absolută $e_x = -0.0015\dots$ și deci o margine a ei $a_x = 0.0016$. Eroarea relativă $\varepsilon_x = -0.0015\dots/3.1415\dots$ este majorată de $r_x = 0.0016/3 \leq 0.0006 = 0.06\%$. Putem scrie deci

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

În concluzie, 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:

$$(\exists)\mathbf{e}_x \in \mathbb{R}^n, \|\mathbf{e}_x\| \leq a_x, \quad \text{astfel încât} \quad \bar{\mathbf{x}} = \mathbf{x} + \mathbf{e}_x, \quad (1.13)$$

iar 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:

$$(\exists)\varepsilon_x \in \mathbb{R}^n, \|\varepsilon_x\| \leq r_x, \quad \text{astfel încât} \quad \bar{\mathbf{x}} = \mathbf{x} + \|\mathbf{x}\|\varepsilon_x. \quad (1.14)$$

În cazul unei mărimi scalare, relația (1.14) se scrie

$$\bar{x} = x(1 \pm \varepsilon_x), \quad (1.15)$$

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

1.3.2 Analiza erorilor de rotunjire

Erorile de rotunjire se datorează reprezentării finite a numerelor reale în calculator. Numerele reale nu pot fi reprezentate în calculator decât cu un număr finit de cifre semnificative. În consecință, sistemele de calcul nu pot lucra cu numere reale exacte, ci doar cu aproximări raționale ale acestora.

În cele ce urmează este util să reprezentăm numerele reale în baza 10, cu ajutorul unei părți fracționare f și a unui exponent n :

$$\bar{x} = f \cdot 10^n. \quad (1.16)$$

Mai mult, pentru orice număr în afara lui 0, prin alegerea convenabilă a lui n , partea fracționară satisface $0.1 \leq |f| < 1$. De exemplu $3.14 = 0.314 \cdot 10^1$, $-0.007856 = -0.7856 \cdot 10^{-2}$.

În termenii convenției formulate, cifrele părții fracționare se numesc cifre semnificative. Astfel, numărul 3.14 are 3 cifre semnificative, iar numărul -0.007856 are patru cifre semnificative.

În calculator se pot memora un număr finit k de cifre semnificative. În consecință, un număr real oarecare x poate fi reprezentat doar ca

$$\bar{x} = \underbrace{0.***\dots*}_{k \text{ cifre}} \cdot 10^n, \quad (1.17)$$

unde fiecare simbol $*$ reprezintă o cifră, iar prima cifră nu este 0. Numărul exact ar fi avut o infinitate de cifre

$$x = 0. \underbrace{***\dots*}_{k \text{ cifre}} \#\#\#\dots \cdot 10^n, \quad (1.18)$$

cifrele marcate cu $\#$ fiind pierdute. Eroarea absolută datorată acestei reprezentări "rotunjite" este

$$e_x = \bar{x} - x = -0. \underbrace{000\dots0}_{k \text{ cifre}} \#\#\#\dots \cdot 10^n = -0. \#\#\#\dots \cdot 10^{n-k}, \quad (1.19)$$

ceea ce conduce la o eroare relativă

$$\varepsilon_x = \frac{e_x}{x} = \frac{-0. \#\#\#\dots \cdot 10^{n-k}}{0. \underbrace{***\dots*}_{k \text{ cifre}} \#\#\#\dots \cdot 10^n} = -\frac{0. \#\#\#\dots}{0. ***\dots} 10^{-k} \quad (1.20)$$

al cărei modul este majorat de

$$|\varepsilon_x| \leq \frac{1}{0.1} 10^{-k} = 10^{-k+1}. \quad (1.21)$$

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} .

Eroarea relativă de rotunjire nu depinde deci de cât de mare sau cât de mic este numărul (valoare care se reflectă în exponentul n), ci doar de numărul de cifre semnificative memorate.

Evident că dacă cifrele pierdute, notate cu # mai sus, sunt toate zero, atunci numărul poate fi reprezentat exact. Calculele în care intervin numere reale sunt afectate însă la rândul lor de procesul de rotunjire. De exemplu, adunarea a două numere reale se face adunând părțile fracționare după ce, în prealabil, numărul mai mic a fost rescris astfel încât să aibă exponentul numărului mai mare. La rescrierea numărului mai mic se pot pierde cifre semnificative. Pentru înțelegerea acestei afirmații să ne imaginăm că lucrăm pe un calculator ipotetic în care numărul de cifre semnificative este 3. Într-un astfel de calculator, să adunăm $x_1 = 3.73 = 0.373 \cdot 10^1$ cu $x_2 = 0.006 = 6 \cdot 10^{-3}$. Numărul x_2 , de modul mai mic, se rescrie astfel încât să aibă exponentul 1, ca al numărului mai mare. $x_2 = 6 \cdot 10^{-4} \cdot 10^1 = 0.0006 \cdot 10^1$. Deoarece calculatorul permite memorarea doar a 3 cifre după virgulă, iar cifra 6 ar fi a patra, înseamnă că de fapt cifra 6 este pierdută. La adunare x_2 este "văzut" a fi zero și rezultatul adunării este de fapt x_1 . Evident că un astfel de calculator nu există, în mod uzual se lucrează cu mai mult de 12 cifre semnificative, acest lucru depinzând de configurația hard, de limbajul de programare folosit și de declarația făcută pentru variabile, dar ideea este exact cea descrisă de acest exemplu simplu.

Se definește zeroul (acuratețea, precizia, "epsilon-ul") mașinii ca fiind cel mai mic număr real care adunat la unitate îi modifică valoarea. Am putea spune că zeroul mașinii `eps` este cel mai mic număr pentru care $1 + \text{eps} > 1$. Pentru orice număr a mai mic decât zeroul mașinii $1 + a = 1$, unde această din urmă relație o considerăm efectuată în calculator și nu în matematica exactă. În mod uzual⁶, zeroul mașinii are valoarea $2.22 \cdot 10^{-16}$. Există limbaje de programare în care această caracteristică a mediului de programare este predefinită: în Matlab există funcția `eps`, iar în Scilab variabila `%eps`. **Zeroul mașinii nu trebuie confundat cu cel mai mic număr reprezentabil în calculator** și care, în mod uzual are valoarea $2.23 \cdot 10^{-308}$.

Într-un mediu în care zeroul mașinii nu este cunoscut, el poate fi calculat cu ușurință cu ajutorul unui program ce implementează următorul pseudocod:

```
funcție zeroul_mașinii ()
real eps
```

⁶calculatoare ce folosesc standardul IEEE și variabile de tip `double`

```

eps = 1
cât timp (1 + eps > 1)
    eps = eps/2
•
eps = eps*2
întoarce eps

```

Ca o consecință a celor discutate mai sus, rezultă că, spre deosebire de matematica exactă, în calculator adunarea numerelor reale nu este asociativă. Dacă trebuie adunate mai multe numere reale, pentru a obține un rezultat afectat cât mai puțin de erori de rotunjire, trebuie ca numerele să fie adunate în ordinea crescătoare a valorii lor.

1.3.3 Standardul IEEE pentru reprezentarea în virgulă mobilă

În paragraful anterior, pentru analiza erorilor de rotunjire, s-a folosit notația (1.16) care a permis explicarea intuitivă a modului în care se fac adunările într-un calculator. Această notație, care folosește un exponent al lui 10 se numește *notație științifică*. Se mai întâlnește *notația științifică normalizată*, $x = a \cdot 10^b$ unde exponentul întreg b se alege astfel încât $1 \leq |a| < 10$, notație care pune în evidență ordinul de mărime al unui număr real. Mai există și *notația inginerescă* în care b se alege numai dintre multiplii lui 3, și, în consecință modului lui a poate avea valori între $1 \leq a < 1000$. Aceasta are avantajul că poate fi citită ușor cu ajutorul unor prefixe ca: mega, kilo, mili, micro etc.

Reprezentarea numerelor reale în calculator este un fel de echivalent hardware al notației științifice [7]. Calculatoarele folosesc un număr finit de biți pentru a reprezenta un număr real. Din acest motiv, în calculator poate fi reprezentată numai o mulțime finită de numere reale. Aceasta conduce la două probleme: pe de o parte numerele reale reprezentate nu pot fi oricât de mari sau oricât de mici, iar pe de altă parte, există spații între ele. Primul aspect pune rare probleme în practică (numite *overflow* în cazul în care s-ar dori reprezentarea unui număr prea mare sau *underflow* în cazul reprezentării unui număr prea mic). Problema a doua este mai delicată și a condus la elaborarea în 1985 de către *Institute of Electrical and Electronics Engineers (IEEE)* și *American National Standards Institute (ANSI)* a standardului 754 care reglementează reprezentarea în virgulă mobilă. Acest standard a reprezentat rezultatul muncii desfășurate pe parcursul unui deceniu de o echipă formată din aproape o sută de cercetători matematicieni, informaticieni și ingineri din universități și industrie (fabricanți de microprocesoare și calculatoare).

Toate calculatoarele construite după 1985 folosesc acest standard IEEE pentru reprezentarea în virgulă mobilă. Aceasta nu înseamnă că rezultatele obținute sunt identice pentru toate calculatoarele, deoarece există o oarecare flexibilitate în cadrul standardului.

Dar aceasta înseamnă că există un model independent de calculator pentru modul în care se fac calculele aritmetice cu numere reale.

În cadrul acestui standard, scrierea normalizată a unui număr real x este

$$\bar{x} = \pm(1 + f) \cdot 2^e, \quad (1.22)$$

unde f se numește *mantisă* și e *exponent*. Mai mult, mantisa f satisface $0 \leq f < 1$ iar exponentul e este un număr întreg cuprins în intervalul $-1022 \leq e \leq 1023$. Ambele se reprezintă în calculator ca numere binare (în baza 2). În calculator, orice număr real se aproximează cu un număr reprezentat ca în (1.22). De exemplu, numărul 0.1 nu poate fi reprezentat exact în acest standard deoarece reprezentarea sa în binar necesită un număr infinit de biți:

$$\frac{1}{10} = \frac{1}{2^4} + \frac{1}{2^5} + \frac{0}{2^6} + \frac{0}{2^7} + \frac{1}{2^8} + \frac{1}{2^9} + \frac{0}{2^{10}} + \frac{0}{2^{11}} + \frac{1}{2^{12}} + \dots$$

după primul termen secvența 1, 0, 0, 1 repetându-se la infinit. În consecință, reprezentarea numărului 0.1 este un pic mai mare decât valoarea exactă. **Iată deci că nu numai numerele iraționale sunt afectate de erori de rotunjire.**

Numerele în dublă precizie sunt reprezentate ca un cuvânt de 64 de biți din care 52 de biți pentru f , 11 biți pentru e și 1 bit pentru semnul numărului. În acest standard, intervalul $[1, 2]$ este reprezentat de mulțimea discretă $\{1, 1 + 2^{-52}, 1 + 2 \cdot 2^{-52}, 1 + 3 \cdot 2^{-52}, \dots, 2\}$. Spațiile dintre numerele acestui interval sunt egale cu $2^{-52} \approx 2.22 \cdot 10^{-16}$. Acesta este zeroul mașinii și corespunde la aproximativ 16 cifre semnificative în baza 10. Un interval $[2^j, 2^{j+1}]$ este reprezentat de mulțimea de mai sus înmulțită cu 2^j . Spațiile dintre numerele vecine sunt deci scalate în funcție de dimensiunea numerelor.

Zeroul mașinii reflectă rezoluția mulțimii reale discrete \mathcal{R} ce poate fi reprezentată în calculator. El are următoarea proprietate

$$(\forall)x \in \mathbb{R}, (\exists)\bar{x} \in \mathcal{R} \quad \text{astfel încât} \quad |x - \bar{x}| \leq \text{eps} |x|. \quad (1.23)$$

Pe scurt, eroarea relativă dintre numărul real x și reprezentarea lui discretă \bar{x} este mărginită superior de zeroul mașinii. Sau, conform relației (1.15) există ε , unde $|\varepsilon| \leq \text{eps}$ astfel încât

$$\bar{x} = x(1 + \varepsilon). \quad (1.24)$$

Restricțiunile impuse pentru mantisă și exponent conduc și la o restricție pentru cel mai mare/cel mai mic număr reprezentabil în calculator. Numărul cel mai mare care se poate reprezenta este $\text{realmax} = (2 - \text{eps}) \cdot 2^{1023}$ a cărui valoare zecimală este aproximativ $1.797 \cdot 10^{308}$. Dacă un calcul încearcă să producă o valoare decât aceasta, se spune că

apare o *depășire* în virgulă mobilă (*overflow*). Cel mai mic număr strict pozitiv care poate fi reprezentat în calculator este `realmin` = $2 \cdot 2^{-1022}$ a cărui valoare zecimală este aproximativ $2.2251 \cdot 10^{-308}$. Un calcul care încearcă să producă o valoare pozitivă mai mică decât aceasta reprezintă tot o depășire în virgulă mobilă (*underflow*).

Dacă un calcul încearcă să producă o valoare care este nedefinită și în matematica numerelor reale (cum ar fi $0/0$), atunci rezultatul este o valoare excepțională numită NaN (*not a number*).

1.3.4 Analiza erorilor de trunchiere

Erorile de trunchiere provin din caracterul finit al algoritmilor. Există metode matematice a căror soluție exactă ar necesita efectuarea unui număr infinit de calcule. Estimarea erorilor de trunchiere se poate face de multe ori apriori, folosind rezultatele teoretice ale matematicii.

Ca exemplu, să considerăm dezvoltarea în serie Taylor a unei funcții:

$$f(x) = f(x_0) + \frac{x - x_0}{1!} f'(x_0) + \frac{(x - x_0)^2}{2!} f''(x_0) + \dots \quad (1.25)$$

aplicată pentru funcția sinus și punctul de dezvoltare $x_0 = 0$:

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}. \quad (1.26)$$

Estimarea acestei sume cu ajutorul calculatorului se va face trunchiind-o și aproximând rezultatul cu

$$\bar{s} = \sum_{k=0}^n (-1)^k \frac{x^{2k+1}}{(2k+1)!}. \quad (1.27)$$

Conform teoriei seriilor alternate, se știe că eroarea absolută făcută prin această trunchiere este mai mică decât ultimul termen considerat:

$$|e_s| = |\bar{s} - s| \leq \frac{x^{2n+1}}{(2n+1)!}. \quad (1.28)$$

Cunoașterea acestui rezultat teoretic permite elaborarea unui algoritm de evaluare a funcției sinus care să stabilească singur când se va opri. Algoritmul va aduna termeni la suma parțială până când termenul curent sumat devine mai mic decât o anumită valoare. Această valoare nu trebuie să fie prea mare căci atunci soluția va fi nesatisfăcătoare. De asemenea, o valoare mai mică decât zeroul mașinii este lipsită de sens deoarece un termen curent cu o astfel de valoare, adunat la suma parțială acumulată până în acel moment, nu mai are nici o influență asupra valorii sumei parțiale. Într-un astfel de moment, continuarea sumării seriei este inutilă, ea nu mai îmbunătățește rezultatul, fenomen total diferit

de teoria matematică în care rezultatul este cu atât mai precis cu cât suma conține mai mulți termeni. Pseudocodul următor implementează calculul funcției sinus într-un punct x , cu o eroare impusă e .

```

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 ( $|t| > e$ )
     $k = k + 1$ 
     $t = (-1) * t * \frac{x^2}{(2k)(2k+1)}$ 
     $s = s + t$ 
•
intoarce  $s$ 

```

În figura 1.4 este reprezentat modulul termenului curent în cazul evaluării funcției sinus în punctul $x = 1$, iar în figura 1.5 modulul diferenței dintre sumele parțiale consecutive. Se observă că încă de la $k = 8$ adăugarea unui nou termen curent nu mai poate îmbunătăți rezultatul. Dacă, în matematică, adunarea a 30 de termeni ar fi însemnat obținerea unei erori de trunchiere mai mică decât 10^{80} , numeric, eroarea s-a saturat după adunarea a 8 termeni, la o valoare de ordinul zeroului mașinii.

Acest paragraf a exemplificat un algoritm în care se realizează controlul erorii de trunchiere. În general, algoritmi se concep astfel încât ei să ofere și informații despre eroare. Algoritmi care nu fac așa ceva sunt algoritmi lipsiți de valoare practică.

1.3.5 Analiza erorilor inerente

Erorile inerente sunt erorile datorate reprezentării imprecise a datelor de intrare. Datele de intrare ale unui algoritm pot proveni de exemplu din măsurători, și de aceea ele sunt afectate de erori. Aceste erori nu pot fi eliminate, fiind important să putem evalua modul în care ele se propagă în calculele numerice, astfel încât să poată fi estimată eroarea rezultatului în funcție de erorile datelor de intrare.

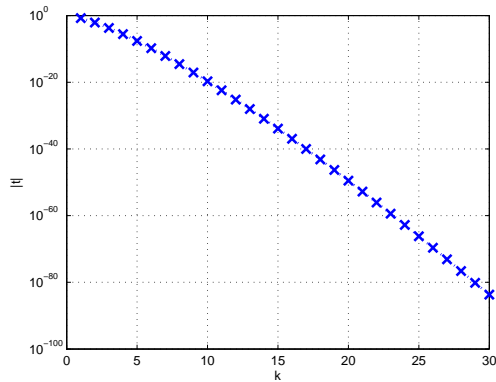


Figura 1.4: Modulus termenului curent al dezvoltării în serie Taylor a funcției sinus.

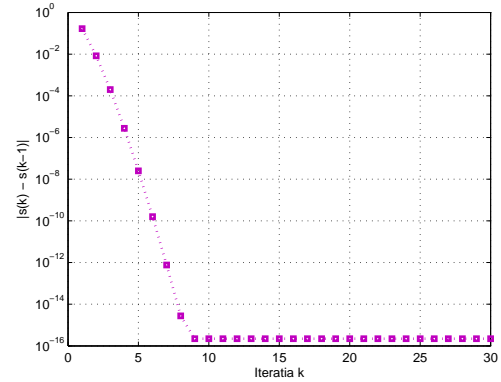


Figura 1.5: Modulus diferenței dintre sume parțiale consecutive la dezvoltarea în serie Taylor a funcției sinus.

Pentru a putea evalua care este efectul micilor perturbații ale datelor de intrare independente notate x_1, x_2, \dots, x_n , considerate a fi mărimi reale, asupra rezultatului y considerat și el un număr real, vom nota cu f dependența dintre date și rezultat

$$y = f(x_1, x_2, \dots, x_n). \quad (1.29)$$

Diferențiala rezultatului va fi în consecință

$$dy = \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \dots + \frac{\partial f}{\partial x_n} dx_n. \quad (1.30)$$

Relația (1.30) poate fi scrisă aproximativ ca

$$\Delta y \approx \frac{\partial f}{\partial x_1} \Delta x_1 + \frac{\partial f}{\partial x_2} \Delta x_2 + \dots + \frac{\partial f}{\partial x_n} \Delta x_n. \quad (1.31)$$

Dacă mărimile Δx_k reprezintă mici perturbații ale datelor de intrare:

$$\Delta x_k = \bar{x}_k - x_k = e_{x_k}, \quad (1.32)$$

unde cu e_{x_k} a fost notată eroarea absolută a mărimii x_k , atunci rezultă că *eroarea absolută a rezultatului* $e_y = \bar{y} - y = \Delta y$ se poate calcula ca o combinație liniară a erorilor absolute a datelor de intrare:

$$e_y = \sum_{k=1}^n \frac{\partial f}{\partial x_k} e_{x_k}. \quad (1.33)$$

Pentru a estima o margine a erorii absolute a rezultatului se majorează modulusul expresiei din membrul drept al relației (1.33):

$$\left| \sum_{k=1}^n \frac{\partial f}{\partial x_k} e_{x_k} \right| \leq \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} e_{x_k} \right| = \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| |e_{x_k}| \leq \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| a_{x_k}, \quad (1.34)$$

unde a_{x_k} este marginea erorii absolute a datei de intrare: $|e_{x_k}| \leq a_{x_k}$. Din (1.34) rezultă că *marginea erorii absolute a rezultatului* se poate calcula ca

$$a_y = \sum_{k=1}^n \left| \frac{\partial f}{\partial x_k} \right| a_{x_k}. \quad (1.35)$$

Folosind (1.33) și (1.29) se deduce că *eroarea relativă a rezultatului* $\varepsilon_y = e_y/|y|$ poate fi exprimată în funcție de erorile relative ale datelor de intrare $\varepsilon_{x_k} = e_{x_k}/|x_k|$ astfel:

$$\varepsilon_y = \frac{\sum_{k=1}^n \frac{\partial f}{\partial x_k} e_{x_k}}{|y|} = \sum_{k=1}^n \frac{\partial f}{\partial x_k} \frac{e_{x_k}}{|y|} = \sum_{k=1}^n \frac{\partial f}{\partial x_k} \frac{|x_k|}{|y|} \varepsilon_{x_k}. \quad (1.36)$$

Majorând expresia din membrul drept, rezultă că *marginea erorii relative a rezultatului* r_y , unde $|\varepsilon_y| \leq r_y$ poate fi calculată în funcție de marginile erorilor relative ale datelor de intrare r_{x_k} , unde $|\varepsilon_{x_k}| \leq r_{x_k}$:

$$r_y = \sum_{k=1}^n \left| \frac{\partial(\ln f)}{\partial x_k} \right| |x_k| r_{x_k}. \quad (1.37)$$

Este util să particularizăm aceste formule în cazul operațiilor algebrice elementare: adunare, scădere, înmulțire și împărțire. Rezultatele obținute în cazul adunării și scăderii a două numere reale sunt prezentate în tabelul 1.1.

Precizăm că operația $x_1 + x_2$ este considerată adunare dacă ambii operanzi au același semn. În caz contrar, operația este de fapt o scădere. Similar, $x_1 - x_2$ este o scădere dacă ambii operanzi au același semn, în caz contrar fiind de fapt efectuată o adunare.

Dacă analizăm marginea erorii relative de la adunare, se observă că erorile relative ale datelor de intrare sunt ponderate cu $|x_1/(x_1 + x_2)|$ și $|x_2/(x_1 + x_2)|$, ambele ponderi fiind subunitare. Aceasta înseamnă că eroarea rezultatului adunării a două numere reale nu este amplificată, ea rămâne de același ordin de mărime cu erorile datelor de intrare. Se spune că **adunarea este o operație bine condiționată**. Nu același lucru se poate spune despre

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 $\varepsilon_y =$	$\frac{x_1}{x_1+x_2}$	$\varepsilon_{x_1} +$	$\frac{x_2}{x_1+x_2}$	ε_{x_2}	$\frac{x_1}{x_1-x_2}$	$\varepsilon_{x_1} -$	$\frac{x_2}{x_1-x_2}$	ε_{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}

Tabelul 1.1: Erorile rezultatului adunării și scăderii a două numere reale în funcție de erorile datelor de intrare.

scădere. Ponderile în acest caz sunt $|x_1/(x_1 - x_2)|$ și $|x_2/(x_1 - x_2)|$, care pot fi oricât de mari pentru că diferența $x_1 - x_2$ poate fi oricât de mică. În consecință, rezultatul unei scăderi poate fi afectat de erori relative mult mai mari decât erorile relative ale datelor de intrare. Se spune că **scăderea este o operație prost condiționată**. Proasta condiționare la scădere apare mai ales atunci când numerele sunt foarte apropiate, efect cunoscut și sub numele de *efect de anulare prin scădere*.

De exemplu, dacă $x_1 = 1.23 \pm 1\%$ și $x_2 = 1.22 \pm 1\%$, rezultă atunci că marginea erorii relative la scădere este $r = |1.23/0.01 \cdot 1/100 + 1.22/0.01 \cdot 1/100| = 1.23 + 1.22 = 2.45 = 245\%$, deci $x_1 - x_2 = 0.01 \pm 245\%$. Eroarea relativă la scădere este de peste 200 de ori mai mare decât erorile relative ale datelor de intrare. La adunare în schimb eroarea relativă este $r = |1.23/2.45 \cdot 1/100 + 1.22/2.45 \cdot 1/100| \approx 0.5 \cdot 1/100 + 0.5 \cdot 1/100 = 1/100 = 1\%$. Deci $x_1 + x_2 = 2.45 \pm 1\%$.

Efectuarea unor calcule de acest tip, în care se urmărește nu numai valoarea rezultatului ci și modul în care acesta este afectat de erorile datelor de intrare se numește *calcul cu intervale*. Este acum mai evident faptul că adunarea numerelor reale în calculator nu este o operație asociativă, afirmație făcută cu ocazia discuției referitoare la erorile de rotunjire.

Tabelul 1.2 conține particularizarea formulelor de eroare în cazul operațiilor de înmulțire și împărțire. Se observă în acest caz faptul că **înmulțirea și împărțirea sunt operații bine condiționate**, eroarea rezultatului fiind de ordinul de mărime al datelor de intrare. Pentru exemplul numeric de mai sus, erorile relative vor fi majorate de 2%.

În concluzie, din cele patru operații algebrice elementare, numai scăderea este prost condiționată și ea trebuie evitată pe cât posibil în calculele numerice, mai ales atunci când numerele care se scad sunt apropiate. Un exemplu tipic în acest sens îl reprezintă algoritmul de calcul al soluțiilor unei ecuații de gradul doi $ax^2 + bx + c = 0$, $x_{1,2} = (-b \pm \sqrt{b^2 - 4ac})/(2a)$. Presupunând că $b > 0$ și că $b^2 \gg 4ac$, rezultă că rezultatul radicalului va fi foarte apropiat în modul de b , astfel încât, soluția cu plus reprezintă de fapt scăderea a două numere foarte apropiate și va fi afectată în consecință de erori

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 $r_y =$	$r_{x_1} + r_{x_2}$	$r_{x_1} + r_{x_2}$

Tabelul 1.2: Erorile rezultatului înmulțirii și împărțirii a două numere reale în funcție de erorile datelor de intrare.

foarte mari. În acest caz, este mult mai eficient să se calculeze doar soluția obținută prin adunare, $x_1 = (-b - \sqrt{b^2 - 4ac})/(2a)$, urmând ca soluția a doua să se obțină din relația lui Viète: $x_2 = c/(ax_1)$, relație care presupune doar folosirea înmulțirii și împărțirii, operații bine condiționate. Pe scurt, în pseudocod, aceasta se scrie astfel

dacă $b > 0$

$$x1 = (-b - \sqrt{b^2 - 4ac})/(2a) \quad ; \text{ soluția cu minus este o adunare dacă } b > 0$$

altfel

$$x1 = (-b + \sqrt{b^2 - 4ac})/(2a) \quad ; \text{ soluția cu plus este o adunare dacă } b < 0$$

•

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

Un alt exemplu interesant este cel al funcției $y = \sqrt{x}$ pentru care rezultă o eroare absolută

$$e_y = \frac{df}{dx} e_x = \frac{1}{2\sqrt{x}} e_x, \quad (1.38)$$

și o eroare relativă egală cu jumătate din eroarea relativă a datei x :

$$\varepsilon_y = \frac{e_y}{y} = \frac{1}{2\sqrt{x}\sqrt{x}} e_x = \frac{e_x}{2x} = \frac{\varepsilon_x}{2}. \quad (1.39)$$

Am putea deduce de aici, în mod eronat, că în acest caz, dacă aplicăm radicalul în mod repetat rezultatului, eroarea tinde către zero. Toate calculele prezentate în acest paragraf (inclusiv acesta din urmă) au presupus un calcul exact, adică un calcul fără erori de rotunjire. Rotunjirea nu poate fi însă ignorată. Nu putem separa proprietățile erorilor inerente de efectul rotunjirii. De aceea, vom accepta un fel de superpoziție a erorilor în sensul că eroarea relativă într-un calcul aproximativ este egală cu eroarea relativă produsă de calculul aproximativ cu numere exacte (adică eroarea de rotunjire) plus eroarea relativă produsă de calculul exact cu numere aproximative (afectate deci de erori inerente). Pentru a justifica această afirmație, să notăm cu \bar{y} valoarea rotunjită a rezultatului y_i (presupus pozitiv), în care s-au acumulat erori inerente. În consecință, folosind (1.15) rezultă

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

de unde $(\bar{y} - y)/y = \varepsilon_y + \mathbf{eps}$. Cu această precizare, rezultatului dat de relația (1.39) trebuie să i se adauge o eroare de rotunjire

$$\varepsilon_{\sqrt{x}} = \frac{\varepsilon_x}{2} + \mathbf{eps}. \quad (1.40)$$

În concluzie, eroarea relativă a oricărui rezultat numeric este cel puțin egală cu zeroul mașinii.

1.3.6 Condiționare și stabilitate

Condiționarea și stabilitatea sunt două aspecte fundamentale în rezolvarea numerică a problemelor. Despre condiționare s-a discutat și în paragraful anterior, unde s-a demonstrat că, dintre operațiile algebrice, scăderea este prost condiționată.

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

Fie o problemă matematică, notată f , pentru care se dorește aflarea unei soluții \mathbf{x} pentru datele \mathbf{d} . Problema f , privită ca o aplicație definită pe mulțimea datelor \mathcal{D} cu valori în mulțimea soluțiilor \mathcal{X} , poate fi formulată astfel:

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

Să se găsească $\mathbf{x} \in \mathcal{X}$ astfel încât $f(\mathbf{d}) = \mathbf{x}$. (1.41)

O problemă este bine condiționată dacă perturbații mici ale datelor conduc la perturbații mici ale rezultatului. Reprezentări intuitive ale unei probleme bine condiționate sunt cele din figura 1.6. O problemă prost condiționată este cea pentru care perturbații mici ale datelor conduc la perturbații mari ale rezultatului (figura 1.7).

Dacă datele și rezultatele sunt numere reale, atunci o problemă f este prost condiționată dacă ea are alura curbei din figura 1.8 - stânga.

De multe ori, problema matematică nu este formulată explicit, ci implicit:

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

Să se găsească $\mathbf{x} \in \mathcal{X}$ astfel încât $g(\mathbf{x}) = \mathbf{d}$. (1.42)

În acest caz, problema g este prost condiționată dacă ea are alura curbei din figura 1.8, dreapta.

Pentru aprecierea cantitativă a bunei sau a proastei condiționări a unei probleme se folosește *numărul de condiționare*. El reprezintă o măsură a perturbațiilor rezultatului la perturbații mici ale datelor și poate fi definit atât în mod absolut cât și relativ.

Se definește *numărul de condiționare absolut* ca fiind

$$\hat{\kappa} = \lim_{\mu \rightarrow 0} \sup_{\|\delta \mathbf{d}\| < \mu} \frac{\|\delta f\|}{\|\delta \mathbf{d}\|}, \quad (1.43)$$

sau, într-o scriere simplificată

$$\hat{\kappa} = \sup_{\|\delta \mathbf{d}\|} \frac{\|\delta f\|}{\|\delta \mathbf{d}\|}, \quad (1.44)$$

unde $\delta f = f(\mathbf{d} + \delta \mathbf{d}) - f(\mathbf{d})$ și $\delta \mathbf{d}$ sunt mărimi infinitezimale. Dacă f este derivabilă, atunci perturbația δf se poate aproxima în condiția $\|\delta \mathbf{d}\| \rightarrow 0$ ca

$$\delta f = \mathbf{J}(\mathbf{d})\delta \mathbf{d}, \quad (1.45)$$

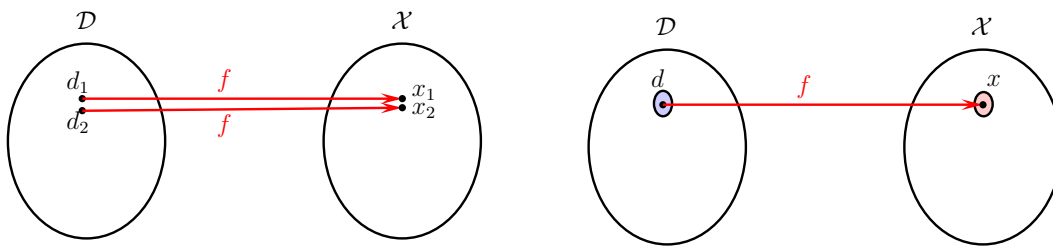


Figura 1.6: Reprezentări intuitive ale unei probleme bine condiționate.

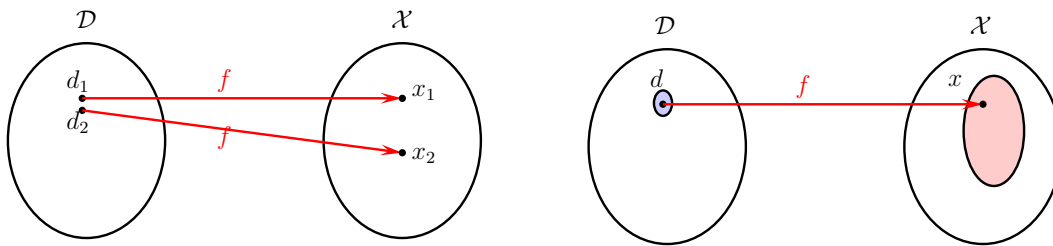


Figura 1.7: Reprezentări intuitive ale unei probleme prost condiționate.

unde $\mathbf{J}(\mathbf{d})$ este matricea Jacobian⁷ asociată funcției f . În consecință, numărul de condiționare absolut este norma matricea Jacobian:

$$\hat{\kappa} = \|\mathbf{J}(\mathbf{d})\|. \quad (1.46)$$

În cazul în care mulțimea de date și mulțimea de rezultate sunt incluse în \mathbb{R} , numărul de condiționare absolut este modulul derivatei funcției f (corespunzătoare problemei definite explicit). O problemă prost condiționată este o problema cu modulul derivatei mare, ceea ce corespunde unei pante mari a graficului funcției (figura 1.8 - stânga).

Atunci când interesează perturbațiile relative, se utilizează *numărul de condiționare relativ*

$$\kappa = \lim_{\mu \rightarrow 0} \sup_{\|\delta \mathbf{d}\| < \mu} \frac{\|\delta f\| / \|f(\mathbf{d})\|}{\|\delta \mathbf{d}\| / \|\mathbf{d}\|}, \quad (1.47)$$

sau, scris mai simplu în ipoteza unor variații infinitezimale

$$\kappa = \sup_{\|\delta \mathbf{d}\|} \frac{\|\delta f\| / \|f(\mathbf{d})\|}{\|\delta \mathbf{d}\| / \|\mathbf{d}\|}. \quad (1.48)$$

Dacă f este derivabilă, atunci

$$k = \frac{\|\mathbf{J}(\mathbf{d})\|}{\|f(\mathbf{d})\| / \|\mathbf{d}\|}. \quad (1.49)$$

⁷Dacă f este o funcție cu valori scalare, ce depinde de n date, atunci Jacobianul este un vector cu n componente, fiecare componentă fiind derivata parțială a funcției în raport cu o dată $\partial f / \partial d_j$, $j = 1, n$. Dacă f este o funcție cu valori vectoriale (fie de dimensiune m), fiecare componentă depinzând de n date, atunci matricea Jacobian are dimensiunea $m \times n$, un element al ei fiind $\partial f_i / \partial d_j$, $i = 1, m$, $j = 1, n$.

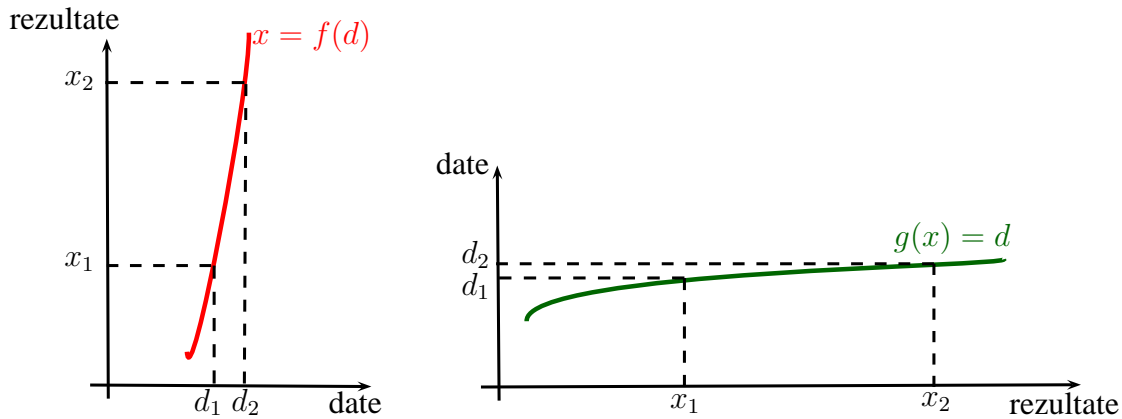


Figura 1.8: Problema prost condiționată: stânga - formulată explicit: $x = f(d)$; dreapta - formulată implicit: $g(x) = d$.

O problemă este bine condiționată dacă valoarea lui κ este mică și prost condiționată dacă valoarea lui κ este mare. Ce înseamnă mic sau mare, depinde de problemă.

Ca exemplu, să calculăm numărul de condiționare al scăderii a două numere reale $f(\mathbf{d}) = d_1 - d_2$, unde $\mathbf{d} = [d_1, d_2]^T$. Atunci, Jacobianul este $\mathbf{J} = [\partial f/\partial d_1, \partial f/\partial d_2]^T = [1, -1]^T$, iar numărul de condiționare calculat cu norme inf⁸ este

$$\kappa = \frac{\|\mathbf{J}(\mathbf{d})\|}{\|f(\mathbf{d})\|/\|\mathbf{d}\|} = \frac{1}{|d_1 - d_2|/\max\{|d_1|, |d_2|\}}. \quad (1.50)$$

Este evident că scăderea este prost condiționată dacă $d_1 \approx d_2$, rezultat justificat și în paragraful anterior.

Un alt exemplu este numărul de condiționare pentru înmulțirea dintre o matrice \mathbf{A} și un vector \mathbf{d} unde se presupun perturbate valorile lui \mathbf{d} . În acest caz $f(\mathbf{d}) = \mathbf{A}\mathbf{d}$. Se poate demonstra [16] că numărul de condiționare al acestei probleme este mărginit superior de

$$\kappa \leq \|\mathbf{A}\|\|\mathbf{A}^{-1}\|. \quad (1.51)$$

De asemenea, pentru problema înmulțirii inversei unei matrice \mathbf{A} cu un vector, numărul de condiționare satisface relația (1.51).

Produsul $\|\mathbf{A}\|\|\mathbf{A}^{-1}\|$ apare atât de des încât el a primit numele de *numărul de condiționare al unei matrice*:

$$\kappa(\mathbf{A}) = \|\mathbf{A}\|\|\mathbf{A}^{-1}\|. \quad (1.52)$$

În acest caz, el este atașat unei matrice și nu unei probleme. Dacă are valoare mică se spune că matricea e bine condiționată.

⁸Norma inf a unui vector este maximum dintre modulele componentelor sale: $\|\mathbf{v}\|_\infty = \max_{1 \leq i \leq n} |v_i|$.

Un alt exemplu este numărul de condiționare al rezolvării unui sistem de ecuații. El este discutat în detaliu în paragraful 2.1. Și acest număr este legat de numărul de condiționare al unei matrice.

Din (1.48) rezultă că, pentru o perturbație oarecare a datelor $\delta \mathbf{d}$, căreia îi corespunde o perturbație a rezultatului $\|\delta f\|$ rezultă

$$\frac{\|\delta f\|/\|f\|}{\|\delta \mathbf{d}\|/\|\mathbf{d}\|} \leq \kappa. \quad (1.53)$$

Perturbația rezultatului este o distanță în spațiul soluțiilor \mathcal{X} , deci reprezintă o eroare absolută $\mathbf{e}_x = \delta f$ sau relativă $\varepsilon_x = \delta f/\|f\|$. Perturbația datelor, numită și *reziduu* este o distanță în spațiul \mathcal{D} . Reziduuul poate fi absolut $\mathbf{e}_d = \delta \mathbf{d}$, unde $\delta \mathbf{d} = \bar{\mathbf{d}} - \mathbf{d}$, sau relativ $\varepsilon_d = \delta \mathbf{d}/\|\mathbf{d}\|$. Cu aceste notații, relația (1.53) devine

$$\|\mathbf{e}_x\| \leq \kappa \|\varepsilon_d\|. \quad (1.54)$$

Eroarea și reziduuul sunt legate prin numărul de condiționare. Din inegalitatea (1.54) rezultă că 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 se referă la eroarea soluției numerice. Interesează ca algoritmul să fie precis, adică eroarea rezultatului să fie mică (figura 1.9).

Un algoritm este precis dacă norma erorii relative

$$\|\varepsilon_x\| = \frac{\|\bar{f}(\mathbf{d}) - f(\mathbf{d})\|}{\|f(\mathbf{d})\|} \quad (1.55)$$

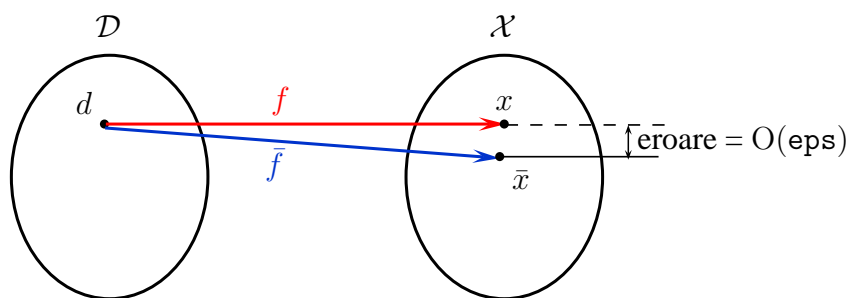


Figura 1.9: Reprezentarea intuitivă a unui algoritm a cărui precizie este ideală.

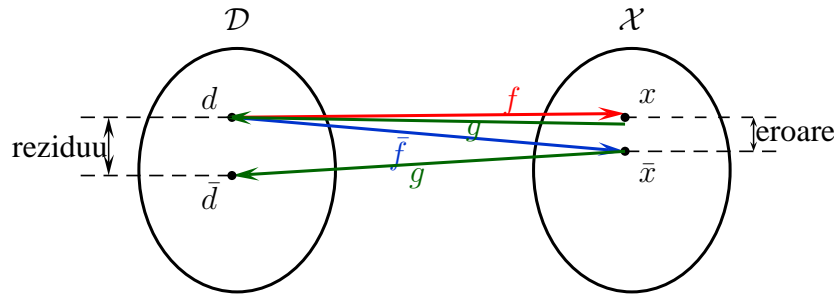


Figura 1.10: Reziduu este o distanță în spațiul datelor. Eroarea este o distanță în spațiul soluțiilor.

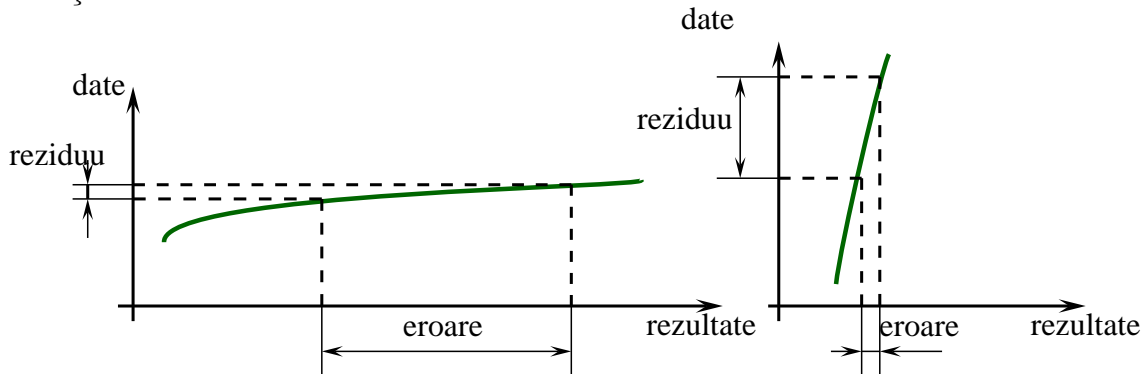


Figura 1.11: Reziduu nu oferă informații despre eroare: stânga - reziduu mic, eroare mare (problemă prost condiționată), dreapta - reziduu mare, eroare mică.

este mică. În mod ideal, un algoritm este precis dacă modulul erorii relative este de ordinul zeroului mașinii

$$\frac{\|\bar{f}(\mathbf{d}) - f(\mathbf{d})\|}{\|f(\mathbf{d})\|} = O(\text{eps}). \quad (1.56)$$

În cazul în care problema este formulată implicit ca (1.42) este ușor a se calcula cât de bine este satisfăcută relația ” $g(\mathbf{x}) = \mathbf{d}$ ” prin înlocuirea lui \mathbf{x} cu $\bar{\mathbf{x}}$. Mărimea

$$g(\bar{\mathbf{x}}) - \mathbf{d} = \bar{\mathbf{d}} - \mathbf{d} = \mathbf{e}_d \quad (1.57)$$

este o distanță în spațiul datelor, deci un reziduu (figura 1.10). Reziduu dă informații asupra problemei care s-a rezolvat de fapt. În general însă, el nu dă informații despre eroare. În cazul unei probleme prost condiționate reziduu poate fi mic deși eroarea este mare, iar în cazul unei probleme bine condiționate reziduu poate fi mare deși eroarea este mică (figura 1.11).

Pentru a ilustra această afirmație, să considerăm sistemul de două ecuații cu două necunoscute

$$\mathbf{Ax} = \mathbf{b}, \quad \text{unde} \quad \mathbf{A} = \begin{bmatrix} 0.78 & 0.563 \\ 0.913 & 0.659 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0.217 \\ 0.254 \end{bmatrix}.$$

Matricea coeficienților acestui sistem are determinantul egal cu $1e-6$, diferit de zero, ceea ce înseamnă că din punct de vedere matematic problema are o soluție unică. Se verifică ușor că soluția exactă a problemei este $\mathbf{x} = [1, -1]^T$. Dacă $\bar{\mathbf{x}}$ este o soluție aproximativă, atunci reziduul este

$$\varepsilon_{\mathbf{b}} = \bar{\mathbf{b}} - \mathbf{b} = \mathbf{A}\bar{\mathbf{x}} - \mathbf{b}.$$

Dacă analizăm posibilele soluții aproximative $\bar{\mathbf{x}}_1 = [0.999, -1.001]^T$ și $\bar{\mathbf{x}}_2 = [0.341, -0.087]^T$, se observă că reziduul primei mărimi $\bar{\mathbf{d}}_1 = [-0.0013, -0.0016]^T$ este mai mare în orice normă decât reziduul celei de a doua $\bar{\mathbf{d}}_2 = [-10^{-6}, 0]^T$, deși este evident că prima soluție are eroarea mai mică. Acest exemplu este o problemă slab condiționată. Intuitiv, cele două drepte în planul (x_1, x_2) corespunzătoare celor două ecuații, deși nu sunt identice, sunt aproape suprapuse. O discuție detaliată despre cum se poate estima cantitativ buna sau slaba condiționare a rezolvării unui sistem algebric de ecuații liniare este prezentată în paragraful 2.1. **Pentru probleme prost condiționate, reziduul poate furniza informații eronate despre eroare.**

În relația (1.55), mărimea $\bar{f}(\mathbf{d})$ semnifică "rezultatul algoritmului f aplicat datelor \mathbf{d} ". În realitate, rotunjirea datelor de intrare este inevitabilă, erorile se acumulează și perturbă rezultatul. În loc de a ținti o astfel de acuratețe, este 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

$$\frac{\|\bar{f}(\bar{\mathbf{d}}) - f(\mathbf{d})\|}{\|f(\mathbf{d})\|} = O(\text{eps}), \quad (1.58)$$

pentru $(\forall)\bar{\mathbf{d}}, \mathbf{d}$ care satisfac $\|\bar{\mathbf{d}} - \mathbf{d}\|/\|\mathbf{d}\| = O(\text{eps})$. Pe scurt, **un algoritm stabil dă răspunsul aproape corect pentru date reprezentate aproape precis.**

Pentru a ilustra noțiunea de stabilitate a unui algoritm, să considerăm sistemul de două ecuații cu două necunoscute

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

Sistemul, scris explicit ca

$$\begin{aligned} x_2 &= 1 \\ x_1 + x_2 &= 0 \end{aligned} \quad (1.59)$$

are soluția $x_1 = -1, x_2 = 1$. Cu notațiile anterioare $\mathbf{x} = f(\mathbf{d}) = [-1, 1]^T$.

Să considerăm acum că datele (numai matricea coeficienților) au fost perturbate:

$$\bar{\mathbf{A}} = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix},$$

cea ce corespunde scrierii explicite

$$\begin{aligned} 10^{-20}x_1 + x_2 &= 1 \\ x_1 + x_2 &= 0 \end{aligned} \tag{1.60}$$

Soluția exactă a acestei probleme perturbate este $x'_1 = -x'_2 = 1/(10^{-20} - 1) \approx -1$, extrem de puțin perturbată față de soluția exactă a problemei neperturbate. Într-adevăr, se poate demonstra că această problemă este bine condiționată. Intuitiv, soluția se găsește la intersecția celei de a doua bisectoare cu o dreaptă verticală.

Un posibil algoritm de rezolvare este algoritmul \bar{f}_1 , descris pe scurt astfel:

- **Pasul 1:** se înmulțește prima ecuație a sistemului (1.60) 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.

Un alt algoritm posibil de rezolvare este algoritmul \bar{f}_2 :

- **Pasul 1:** se înmulțește a doua ecuație a sistemului (1.60) 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}_1 este instabil, pe când algoritmul \bar{f}_2 este stabil. Stabilitatea lui este foarte puternică, el a dat răspunsul exact pentru date de intrare aproape precise.

În concluzie, pentru estimarea acurateții unei soluții numerice trebuie parcurși următorii pași:

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 rezidului relativ.

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

Capitolul 2

Sisteme de ecuații algebrice liniare

Rezolvarea cu ajutorul calculatorului a sistemelor de ecuații algebrice liniare este cea mai frecvent întâlnită problemă matematică ce apare pe parcursul rezolvării aplicațiilor de inginerie electrică. Un astfel de exemplu este analiza circuitelor rezistive liniare. Analiza circuitelor rezistive neliniare conduce la sisteme de ecuații algebrice neliniare care se rezolvă iterativ, fiecare iterație necesitând rezolvarea unui sistem algebric liniar. Circuitele liniare în regim tranzitoriu necesită și ele rezolvarea unui sistem de ecuații algebrice liniare la fiecare pas de timp. De asemenea, analiza oricărei probleme de câmp electromagnetic necesită rezolvarea unuia sau mai multor sisteme de ecuații algebrice liniare.

2.1 Formularea problemei

Problema matematică abordată în acest capitol este rezolvarea unui sistem de n ecuații algebrice liniare cu n necunoscute:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1, \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2, \\ \cdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n. \end{cases} \quad (2.1)$$

Altfel spus, dată fiind matricea coeficienților

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & & & \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (2.2)$$

și vectorul termenilor liberi

$$\mathbf{b} = \begin{bmatrix} b_1 & b_2 & \cdots & b_n \end{bmatrix}^T \in \mathbb{R}^n, \quad (2.3)$$

se cere să se rezolve sistemul

$$\mathbf{Ax} = \mathbf{b}, \quad (2.4)$$

unde \mathbf{x} este soluția

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}^T \in \mathbb{R}^n. \quad (2.5)$$

O astfel de problemă este bine formulată din punct de vedere matematic (soluția există și este unică) dacă și numai dacă matricea \mathbf{A} este nesingulară (are determinantul nenul). În aceste condiții, se poate scrie că $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$, unde \mathbf{A}^{-1} este inversa matricei \mathbf{A} . Vom vedea însă că este ineficient ca soluția \mathbf{x} să se determine numeric folosind inversa matricei calculată explicit. În contextul metodelor numerice, scrierea " $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ " este pur formală, ea trebuie citită ca " \mathbf{x} este soluția sistemului algebric liniar $\mathbf{Ax} = \mathbf{b}$ " și nu "se calculează inversa matricei \mathbf{A} care se înmulțește cu vectorul \mathbf{b} ".

2.2 Condiționarea problemei

În paragraful 1.3.6 s-a discutat problema condiționării în general, definindu-se numărul de condiționare ca (1.48) sau (1.49) în cazul în care funcția matematică ce descrie problema de rezolvat este derivabilă.

Nu orice problemă de rezolvare a unui sistem de ecuații algebrice liniare care este bine formulată matematic este și bine condiționată. Această afirmație poate fi înțeleasă intuitiv pentru cazul $n = 2$. Găsirea soluției unui sistem algebric liniar de două ecuații cu două necunoscute are ca semnificație geometrică găsirea punctului de intersecție dintre cele două drepte care reprezintă fiecare ecuație din sistem. În cazul unei bune formulări matematice, cele două drepte sunt concurente (fig. 2.1 a). Cazul unei matrice singulare are drept corespondent geometric fie două drepte paralele (caz în care sistemul nu are nici o soluție - fig. 2.1 b), fie două drepte confundate (caz în care sistemul are o infinitate de soluții - fig. 2.1 c).

În vecinătatea cazului de proastă formulare matematică există o mulțime de situații în care dreptele respective au pante foarte apropiate, fără să fie confundate (fig. 2.1 d). Știind că, de fapt, nici coeficienții matricei, nici termenii liberi nu sunt cunoscuți cu precizie infinită, aceasta înseamnă că, la o mică perturbație a datelor, soluția poate varia extrem de mult. Este ca și cum dreptele nu sunt desenate cu un creion bine ascuțit ci cu un creion cu un vârf mai gros, astfel încât zona în care se află soluția este mare. Dacă această zonă cuprinde mai multe ordine de mărime, atunci este dificil, dacă nu chiar imposibil de a se determina soluția. Mici perturbații ale datelor conduc la perturbații mari ale soluției și problema, deși bine formulată matematic, este slab condiționată.

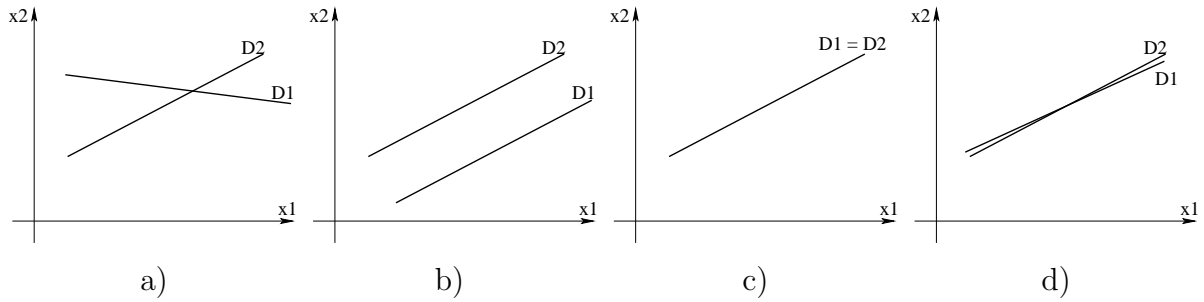


Figura 2.1: a) Problemă matematică bine formulată și bine condiționată. b) Problemă matematică prost formulată (nu există soluție). c) Problemă matematică prost formulată (are o infinitate de soluții). d) Problemă matematică bine formulată și slab condiționată.

În cele ce urmează se urmărește aflarea unui expresii cantitative pentru numărul de condiționare al acestei probleme.

Fie de rezolvat sistemul

$$\mathbf{Ax} = \mathbf{b}, \quad (2.6)$$

unde \mathbf{x} este soluția exactă și presupunem o perturbație $\mathbf{x} + \mathbf{e}_x$ a soluției, corespunzătoare unei perturbații $\mathbf{b} + \mathbf{e}_b$ a datelor:

$$\mathbf{A}(\mathbf{x} + \mathbf{e}_x) = \mathbf{b} + \mathbf{e}_b. \quad (2.7)$$

Din (2.6) și (2.7) rezultă că și perturbațiile satisfac un sistem de ecuații algebrice liniare, cu aceeași matrice a coeficienților:

$$\mathbf{A}\mathbf{e}_x = \mathbf{e}_b. \quad (2.8)$$

Vom nota eroarea relativă a soluției cu ε_x și eroarea relativă a datelor cu ε_b :

$$\varepsilon_x = \frac{\|\mathbf{e}_x\|}{\|\mathbf{x}\|}, \quad \varepsilon_b = \frac{\|\mathbf{e}_b\|}{\|\mathbf{b}\|}. \quad (2.9)$$

Din (2.8) rezultă

$$\mathbf{e}_x = \mathbf{A}^{-1}\mathbf{e}_b, \quad (2.10)$$

și, în consecință, norma perturbației soluției poate fi majorată de

$$\|\mathbf{e}_x\| = \|\mathbf{A}^{-1}\mathbf{e}_b\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{e}_b\|. \quad (2.11)$$

Din (2.6) rezultă un minorant pentru norma soluției

$$\|\mathbf{b}\| = \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \quad \Rightarrow \quad \|\mathbf{x}\| \geq \frac{\|\mathbf{b}\|}{\|\mathbf{A}\|}. \quad (2.12)$$

Din (2.11) și (2.12) rezultă un majorant pentru eroarea asupra soluției

$$\varepsilon_x = \frac{\|\mathbf{e}_x\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\| \|\mathbf{e}_b\|}{\frac{\|\mathbf{b}\|}{\|\mathbf{A}\|}} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \varepsilon_b. \quad (2.13)$$

Mărimea

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \quad (2.14)$$

se numește număr de condiționare la inversare al matricei \mathbf{A} .

Relația (2.13) devine

$$\varepsilon_x \leq \kappa(\mathbf{A}) \varepsilon_b, \quad (2.15)$$

ceea ce reprezintă de fapt relația (1.54) scrisă pentru problema rezolvării sistemelor algebrice liniare. Demonstrarea expresiei numărului de condiționare se poate face și pornind de la relația (1.49) astfel:

$$\begin{aligned} \kappa &= \frac{\|\mathbf{J}(\mathbf{b})\|}{\|\mathbf{f}(\mathbf{b})\|/\|\mathbf{b}\|} = \frac{\|\mathbf{A}^{-1}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|/\|\mathbf{b}\|} = \frac{\|\mathbf{A}^{-1}\| \|\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} = \frac{\|\mathbf{A}^{-1}\| \|\mathbf{A}\mathbf{A}^{-1}\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} \leq \\ &\leq \frac{\|\mathbf{A}^{-1}\| \|\mathbf{A}\| \|\mathbf{A}^{-1}\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} = \|\mathbf{A}^{-1}\| \|\mathbf{A}\| = \kappa(\mathbf{A}), \end{aligned} \quad (2.16)$$

unde s-au presupus doar datele \mathbf{b} perturbate, iar soluția problemei a fost scrisă formal ca $\mathbf{x} = \mathbf{f}(\mathbf{b}) = \mathbf{A}^{-1}\mathbf{b}$, având matricea Jacobian $\mathbf{J}(\mathbf{b}) = \mathbf{A}^{-1}$.

Un raționament similar permite găsirea unei margini inferioare pentru eroarea asupra soluției. Din (2.8) rezultă un majorant pentru eroarea în rezultat

$$\|\mathbf{e}_b\| = \|\mathbf{A}\mathbf{e}_x\| \leq \|\mathbf{A}\| \|\mathbf{e}_x\| \quad \Rightarrow \quad \|\mathbf{e}_x\| \geq \frac{\|\mathbf{e}_b\|}{\|\mathbf{A}\|}. \quad (2.17)$$

Din (2.6) rezultă un minorant pentru norma soluției

$$\|\mathbf{x}\| = \|\mathbf{A}^{-1}\mathbf{b}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{b}\|. \quad (2.18)$$

Din (2.17) și (2.18) rezultă un minorant pentru eroarea asupra soluției

$$\varepsilon_x = \frac{\|\mathbf{e}_x\|}{\|\mathbf{x}\|} \geq \frac{\|\mathbf{e}_b\|}{\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \|\mathbf{b}\|} = \frac{\varepsilon_b}{\kappa(\mathbf{A})}. \quad (2.19)$$

În concluzie, relația între eroarea relativă asupra soluției și eroarea relativă a termenilor liberi este dată de

$$\frac{\varepsilon_b}{\kappa(\mathbf{A})} \leq \varepsilon_x \leq \kappa(\mathbf{A}) \varepsilon_b. \quad (2.20)$$

Numărul de condiționare este întotdeauna supraunitar $\kappa(\mathbf{A}) \geq 1$, deoarece

$$1 = \|\mathbf{I}\| = \|\mathbf{A}\mathbf{A}^{-1}\| \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| = \kappa(\mathbf{A}). \quad (2.21)$$

Dacă numărul de condiționare $\kappa(\mathbf{A})$ are valori mici, inegalitatea (2.20) garantează că o mică perturbare a termenilor liberi nu este amplificată. În cazul cel mai favorabil,

$\kappa(\mathbf{A}) = 1$ și $\varepsilon_x = \varepsilon_b$. Acest caz corespunde unei matrice ortogonale (în cazul $n = 2$ cele două drepte ce reprezintă ecuațiile sistemului sunt perpendiculare).

Dacă numărul de condiționare este mare, atunci chiar și termeni liberi foarte puțin perturbați pot conduce la perturbații mari ale soluției numerice, problema fiind slab condiționată.

Numărul de condiționare este o proprietate a matricei și nu are legătură nici cu metoda de rezolvare propriu-zisă, nici cu erorile de rotunjire care apar în mediul de calcul. Cu toate acestea, în practică se consideră că o problemă este slab condiționată dacă numărul de condiționare este mai mare decât inversul zeroului mașinii.

Dacă $\kappa(\mathbf{A}) > 1/\text{eps}$ problema se consideră slab condiționată.

Demonstrația de mai sus a fost făcută considerând o perturbație în vectorul termenilor liberi, matricea coeficienților fiind considerată neperturbată. Același număr de condiționare este util și dacă se consideră o perturbație în matricea coeficienților, vectorul termenilor liberi fiind neafectat de erori. Soluția perturbată va satisface atunci relația

$$(\mathbf{A} + \mathbf{e}_A)(\mathbf{x} + \mathbf{e}_x) = \mathbf{b}. \quad (2.22)$$

Din (2.6) și (2.22) rezultă că perturbațiile satisfac

$$\mathbf{A}\mathbf{e}_x = -\mathbf{e}_A(\mathbf{x} + \mathbf{e}_x), \quad (2.23)$$

de unde rezultă că perturbația soluției este mărginită superior de

$$\|\mathbf{e}_x\| = \|\mathbf{A}^{-1}\mathbf{e}_A(\mathbf{x} + \mathbf{e}_x)\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{e}_A\| \|\mathbf{x} + \mathbf{e}_x\|. \quad (2.24)$$

Eroarea relativă, calculată aproximativ față de soluția numerică este majorată de

$$\varepsilon_x \approx \frac{\|\mathbf{e}_x\|}{\|\mathbf{x} + \mathbf{e}_x\|} \leq \|\mathbf{A}^{-1}\| \|\mathbf{e}_A\| = \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \frac{\|\mathbf{e}_A\|}{\|\mathbf{A}\|} = \kappa(\mathbf{A})\varepsilon_A. \quad (2.25)$$

Deoarece $\|\mathbf{x} + \mathbf{e}_x\| \leq \|\mathbf{x}\| + \|\mathbf{e}_x\|$, rezultă că $\|\mathbf{x}\| \geq \|\mathbf{x} + \mathbf{e}_x\| - \|\mathbf{e}_x\|$. Dacă presupunem că

$$\|\mathbf{x} + \mathbf{e}_x\| - \|\mathbf{e}_x\| > 0, \quad (2.26)$$

rezultă o majorare pentru eroarea relativă a soluției

$$\varepsilon_x = \frac{\|\mathbf{e}_x\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{e}_x\|}{\|\mathbf{x} + \mathbf{e}_x\| - \|\mathbf{e}_x\|} = \frac{\|\mathbf{e}_x\|/\|\mathbf{x} + \mathbf{e}_x\|}{1 - \|\mathbf{e}_x\|/\|\mathbf{x} + \mathbf{e}_x\|} \leq \frac{\kappa(\mathbf{A})\varepsilon_A}{1 - \kappa(\mathbf{A})\varepsilon_A}, \quad (2.27)$$

relație valabilă în ipoteza $\kappa(\mathbf{A})\varepsilon_A < 1$.

2.3 Clasificarea metodelor

Există mai multe metode numerice posibile pentru rezolvarea cu calculatorul a sistemelor de ecuații algebrice liniare. Ele se pot clasifica în metode *directe* și *iterative*. La rândul lor, metodele iterative pot fi *staționare* și *nestaționare* (numite și *semiiterative*).

Metodele directe sunt metodele care găsesc soluția teoretică a problemei într-un număr *finit* de pași. Astfel de metode sunt metoda Gauss și metoda factorizării LU. În realitate, calculele sunt afectate de erori de rotunjire și la terminarea algoritmului nu se obține soluția exactă, ci o aproximare a ei. Ambele metode se bazează pe transformarea sistemului de ecuații într-unul echivalent din punct de vedere al soluției, mult mai ușor de rezolvat. Metodele directe sunt:

- Gauss (paragraful 2.4) - În această metodă sistemul este adus la o formă triunghiular superioară, ușor de rezolvat prin substituție regresivă;
- Factorizare LU (paragraful 2.5) - În această metodă sistemul este adus la o formulare echivalentă, constând în două sisteme de ecuații ușor de rezolvat, unul triunghiular inferior, rezolvat prin substituție progresivă, și unul triunghiular superior, rezolvat prin substituție regresivă.

Dezavantajul acestor metode este acela că, în anumite situații, efortul de generare a problemei echivalente (factorizarea) este mare sau necesarul de memorie poate deveni extrem de mare.

Metodele iterative generează un șir de aproximații ale soluției care se dorește a fi convergent către soluția exactă. În consecință, din punct de vedere teoretic, soluția poate fi obținută cu o astfel de metodă numai într-un număr infinit de pași.

Metodele iterative se grupează în două categorii principale: staționare și nestaționare. Metodele staționare au apărut mai întâi, ele sunt simple de înțeles și de implementat dar nu sunt întotdeauna eficiente. Metodele nestaționare au apărut mai târziu, analiza lor este mai dificilă, dar ele pot fi foarte eficiente.

Metodele staționare se bazează pe partiționarea matricei coeficienților în două și construirea unei probleme echivalente de punct fix. Metodele staționare sunt:

- Jacobi (paragraful 2.8) - În această metodă fiecare variabilă se calculează local, în funcție de celelalte variabile. Metoda este ușor de înțeles și implementat dar convergența ei este lentă;
- Gauss-Seidel (paragraful 2.9) - Metoda seamănă cu metoda Jacobi, dar valorile calculate sunt imediat folosite în calculele următoare. În general, dacă pentru

o problemă metoda Jacobi converge, metoda Gauss-Seidel converge mai rapid. Convergența este însă tot lentă.

- Metoda suprarelaxării succesive (SOR - *Successive OverRelaxation*, paragraful 2.10) - este dedusă din metoda Gauss-Seidel prin introducerea unui parametru de extrapolare. Pentru o valoare optimă a acestuia, SOR poate converge într-un număr de iterații cu un ordin de mărime mai mic decât numărul de iterații al metodei Gauss-Seidel.
- Metoda SOR simetrică (SSOR) - este SOR pentru matrice simetrice, utilă ca metodă de preconditionare pentru metodele nestaționare.

Metodele nestaționare se bazează în general pe idea construirii unui șir de vectori ortogonali. Singura metodă care nu folosește această idee este metoda iterațiilor Chebyshev. Metodele iterative nestaționare se mai numesc și metode semiiterative deoarece, chiar dacă ele sunt metode iterative, într-o aritmetică exactă și pentru un anumit tip de probleme, garantează obținerea soluției exacte după un anumit număr de pași. Metodele nestaționare sunt

- Metoda gradientilor conjugați (CG, paragraful 2.12) - generează un șir de vectori conjugați (sau ortogonali). Acești vectori sunt reziduurile corespunzătoare iterațiilor. Ei sunt de asemenea vectorii gradient ai unei funcționale pătratice asociate problemei, a cărei minimizare este echivalentă cu rezolvarea sistemului de ecuații. CG este o metodă foarte eficientă dacă matricea coeficienților este simetrică și pozitiv definită, în acest caz fiind necesară doar memorarea unui număr limitat de vectori generați.
- Metoda reziduuului minimin (MINRES) și metoda LQ simetrică (SYMMLQ) - sunt alternative pentru CG atunci când matricea este simetrică dar posibil indefinită.
- CG pentru ecuația normală (CGNE sau CGNR). Dacă sistemul de rezolvat $\mathbf{Ax} = \mathbf{b}$ are matricea coeficienților \mathbf{A} nesimetrică și nesingulară, există două variante prin care putem folosi CG pentru rezolvare, ambele bazate pe faptul că $\mathbf{A}^T \mathbf{A}$ este simetrică și pozitiv definită. Varianta CGNE este CG aplicată sistemului $(\mathbf{A}\mathbf{A}^T)\mathbf{y} = \mathbf{b}$ urmând ca apoi $\mathbf{x} = \mathbf{A}^T \mathbf{y}$. Varianta CGNR este CG aplicată sistemului $(\mathbf{A}^T \mathbf{A})\mathbf{x} = \tilde{\mathbf{b}}$ unde $\tilde{\mathbf{b}} = \mathbf{A}^T \mathbf{b}$. Convergența acestor metode poate fi înceată, sistemul normal fiind mai prost condiționat decât sistemul inițial.
- Metoda reziduuului minim generalizat (GMRES) - poate fi aplicată matricelor generale, nesimetrice. Ea calculează o secvență de vectori ortogonali ca și MINRES, memorând întregul șir, necesitând astfel mai multă memorie. Din acest motiv se

practică restartarea metodei, la fiecare rulare numărul de iterații fiind limitat de un număr fixat de vectori care trebuie generați.

- Metoda gradientilor biconjugați (BiCG) și variante ale ei (CGS și Bi-CGSTAB) - pot fi folosite pentru matrice neregulate nesimetrice. Metoda generează două șiruri de vectori ca CG, unul pentru matricea originală a coeficienților \mathbf{A} și unul pentru matricea transpusă \mathbf{A}^T . În loc să ortogonalizeze fiecare șir, vectorii sunt construiți mutual ("bi-") ortogonali. O îmbunătățire a metodei, utilă atunci când BiCG nu converge, este metoda rezidului cvasiminimal (QMR).
- Metoda iterațiilor Chebyshev - determină în mod recursiv polinoame cu coeficienți aleși astfel încât norma rezidului să fie minimizată în sens minimax. Se poate aplica pentru sisteme cu matricea coeficienților pozitiv definită și necesită cunoașterea valorilor proprii extreme.

Nu există o "cea mai bună metodă" dintre cele enumerate mai sus. Pentru fiecare clasă de probleme, utilizatorul trebuie să decidă ce metodă de rezolvare va alege. Ideile principale ale metodelor iterative și criteriile după care pot fi alese aceste metode sunt prezentate în [17]. Aspecte algoritmice detaliate se găsesc în [1] iar aspecte legate de preconditionare se găsesc de exemplu în [11].

În cele ce urmează vor fi descrise câteva metode de rezolvare din fiecare categorie, precum și algoritmi corespunzători.

2.4 Metoda Gauss

Ideea metodei Gauss este de a rezolva sistemul de ecuații în două etape. Prima etapă, numită *eliminarea gaussiană (triangularizare)* urmărește transformarea sistemului de ecuații astfel încât soluția să rămână neschimbată, dar matricea coeficienților să aibă în final o structură specială, și anume triunghiular superioară. A doua etapă este *retrosubstituția (substituția regresivă)* în care se calculează pe rând componentele vectorului soluție, începând cu cel de indice cel mai mare. Pe scurt:

$$\mathbf{Ax} = \mathbf{b} \quad \underbrace{\Leftrightarrow}_{\text{eliminare}} \quad \mathbf{Ux} = \mathbf{b} \quad \underbrace{\Rightarrow}_{\text{subst.regresivă}} \quad \mathbf{x} = \mathbf{U}^{-1}\mathbf{b}. \quad (2.28)$$

2.4.1 Un exemplu simplu

Ideea metodei poate fi ușor înțeleasă pe un exemplu simplu. Fie sistemul de ecuații

$$\begin{cases} x_1 + 2x_2 - x_3 = -1 \\ -2x_1 + 3x_2 + x_3 = 0 \\ 4x_1 - x_2 - 3x_3 = -2. \end{cases} \quad (2.29)$$

Matricea coeficienților se va triangulariza în două sub-etape de eliminare. Mai întâi se vor anula coeficienții de pe prima coloană, aflați sub prima linie, după cum urmează. Pentru a anula primul coeficient din a doua ecuație se va înmulți prima ecuație cu 2 și se va aduna la a doua. Apoi se va anula primul coeficient din a treia ecuație înmulțind prima ecuație cu -4 și se va aduna la a treia. Valorile 2 și -4 se numesc *multiplicatori*. Sistemul rezultat, echivalent din punct de vedere al soluției cu cel inițial este

$$\begin{cases} x_1 + 2x_2 - x_3 = -1 \\ 7x_2 - x_3 = -2 \\ -9x_2 + x_3 = 2. \end{cases} \quad (2.30)$$

În a doua etapă de eliminare nu trebuie să se anuleze decât coeficientul -9 din ecuația a treia, și pentru aceasta se înmulțește ecuația a doua cu multiplicatorul $9/7$ și se adună la ecuația a treia. Sistemul rezultat după etapa de eliminare este unul superior triunghiular:

$$\begin{cases} x_1 + 2x_2 - x_3 = -1 \\ 7x_2 - x_3 = -2 \\ -2/7x_3 = -4/7. \end{cases} \quad (2.31)$$

Soluția se calculează ușor, pornind de la ultima ecuație, către prima (retrosubstituție), astfel:

$$\begin{aligned} x_3 &= (-4/7)/(-2/7) = 2, \\ x_2 &= (-2 + x_3)/7 = 0, \\ x_1 &= -1 - 2x_2 + x_3 = 1. \end{aligned} \quad (2.32)$$

2.4.2 Algoritm

Etapa de eliminare urmărește aducerea matricei coeficienților la forma triunghiular superioară. Mai precis, coeficienții aflați sub diagonală sunt anulați în sub-etape de eliminare, sugerate în figura 2.2.

În prima sub-etapă de eliminare se anulează toți coeficienții a_{i1} unde $i = 2, \dots, n$, astfel: se înmulțește prima ecuație cu elementul de multiplicare $-a_{21}/a_{11}$ și se adună la a doua ecuație, se înmulțește prima ecuație cu elementul de multiplicare $-a_{31}/a_{11}$ și se

adună la a treia ecuație, ș.a.m.d. La sfârșitul acestei prime sub-etape noul sistem de rezolvat, echivalent cu cel inițial din punct de vedere al soluției, este:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a'_{22}x_2 + \cdots + a'_{2n}x_n = b'_2 \\ \dots \\ a'_{n2}x_2 + \cdots + a'_{nn}x_n = b'_n, \end{cases} \quad (2.33)$$

unde, de exemplu, noii coeficienți ai celei de a doua ecuații sunt

$$a'_{2j} = a_{2j} + pa_{1j}, \quad (2.34)$$

unde $p = -a_{21}/a_{11}$ este primul multiplicator folosit pentru anularea coeficientului a_{21} , iar noul termen liber este

$$b'_2 = b_2 + pb_1. \quad (2.35)$$

În a doua sub-etapă de eliminare se urmărește anularea coeficienților matricei aflați pe coloana a doua, sub diagonală. Pentru aceasta se înmulțește ecuația a doua cu elementul de multiplicare $p = -a'_{32}/a'_{22}$ și se adună la ecuația a treia, apoi elementul de multiplicare este $p = -a'_{42}/a'_{22}$, ș.a.m.d.

În total sunt $n - 1$ sub-etape de eliminare la sfârșitul cărora sistemul echivalent obținut este de forma

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \dots \\ a_{nn}x_n = b_n. \end{cases} \quad (2.36)$$

În mod intenționat în sistemul (2.36) au fost omise notațiile de tip ”prim” deoarece operațiile se fac ”în loc”, peste valorile vechi sunt scrise valorile noi, adică relațiile de tipul (2.34) se vor implementa ca

$$a_{2j} = a_{2j} + pa_{1j}, \quad (2.37)$$

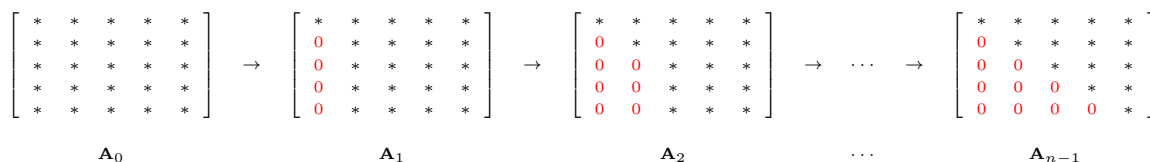


Figura 2.2: Eliminare în metoda Gauss: pentru un sistem de dimensiune n există $n - 1$ sub-etape de eliminare. În final matricea este superior triunghiulară. Matricea inițială este notată \mathbf{A}_0 iar matricea superior triunghiulară obținută este notată \mathbf{A}_{n-1} . În realitate, transformările sunt memorate ”în loc”, în același tablou bidimensional.

unde simbolul = reprezintă acum o atribuire.

Vom construi acum pseudocodul etapei de eliminare. Modificarea ecuației a doua în prima sub-etapă de eliminare poate fi descrisă astfel:

; anularea elementului a_{21}

$p = -a_{21}/a_{11}$; element de multiplicare

pentru $j = 1, n$; parcurge coloanele

$$a_{2j} = a_{2j} + pa_{1j}$$

•

$$b_2 = b_2 + pb_1$$

Dacă în secvența de mai sus înlocuim 2 cu i atunci se obține secvența ce permite anularea elementului a_{i1} . Inserată într-un ciclu cu contor în care i parcurge liniile de la 2 la n , se obține secvența ce descrie prima sub-etapă de eliminare:

; prima sub-etapă de eliminare

pentru $i = 2, n$; parcurge liniile

$p = -a_{i1}/a_{11}$; element de multiplicare

pentru $j = 2, n$; parcurge coloanele

$$a_{ij} = a_{ij} + pa_{1j}$$

•

$$b_i = b_i + pb_1$$

•

În ciclul în j contorul începe cu valoarea 2 pentru a evita calcule inutile pentru valorile a_{i1} , al căror rezultat este nul. Aceasta înseamnă că, de fapt, zerourile nu sunt memorate în matricea \mathbf{A} . În triunghiul inferior vor rămâne valorile avute exact înainte de anulare.

Dacă în pseudocodul de mai sus se înlocuiește 1 cu k și 2 cu $k + 1$ se obține secvența de cod corespunzătoare sub-etapei de eliminare k . Aceasta, repetată pentru cele $n - 1$ valori posibile conduce la următorul pseudocod pentru etapa de eliminare:

; etapa de eliminare din metoda Gauss

pentru $k = 1, n - 1$; parcurge sub-etape ale eliminării

pentru $i = k + 1, n$; parcurge liniile

$p = -a_{ik}/a_{kk}$; element de multiplicare

pentru $j = k + 1, n$; parcurge coloanele

$$a_{ij} = a_{ij} + pa_{kj}$$

•

$$b_i = b_i + pb_k$$

-
-

În *etapa de substituție regresivă* sistemul triunghiular obținut (2.36) se rezolvă ușor, în ordinea descrescătoare a indicelui componentelor vectorului soluție. Din ultima ecuație rezultă

$$x_n = b_n/a_{nn}, \quad (2.38)$$

iar apoi se calculează pe rând x_{n-1} din penultima ecuație, x_{n-2} din antepenultima ș.a.m.d.

În general, din ecuația i

$$a_{ii}x_i + a_{i,i+1}x_{i+1} + \dots + a_{in}x_n = b_i, \quad (2.39)$$

se calculează componenta x_i :

$$x_i = \frac{b_i - \sum_{j=i+1}^n a_{ij}x_j}{a_{ii}}. \quad (2.40)$$

Algoritmul etapei de substituție regresivă este descris de următorul pseudocod:

; etapa de retrosubstituție

$$x_n = b_n/a_{nn}$$

pentru $i = n - 1, 1, -1$

$$s = 0$$

pentru $j = i + 1, n$

$$s = s + a_{ij}x_j$$

-
-

$$x_i = (b_i - s)/a_{ii}$$

Algoritmul complet al procedurii Gauss este următorul.

procedură Gauss(n, a, b, x)

; rezolvă sistemul algebric liniar $ax = b$ prin metoda Gauss

întreg n ; dimensiunea sistemului

tablou real $a[n][n]$; matricea coeficienților - indici de la 1

tablou real $b[n]$; vectorul termenilor liberi

tablou real $x[n]$; vectorul soluție

întreg i, j, k

real p, s

; etapa de eliminare

pentru $k = 1, n - 1$; parcurge sub-etape ale eliminării

; aici se poate introduce pivotarea

pentru $i = k + 1, n$; parcurge liniile
 $p = -a_{ik}/a_{kk}$; element de multiplicare
pentru $j = k + 1, n$; parcurge coloanele
 $a_{ij} = a_{ij} + pa_{kj}$
 •
 $b_i = b_i + pb_k$

•
 ; etapa de retrosubstituție

$$x_n = b_n/a_{nn}$$

pentru $i = n - 1, 1, -1$

$$s = 0$$

pentru $j = i + 1, n$

$$s = s + a_{ij}x_j$$

•
 $x_i = (b_i - s)/a_{ii}$

•
retur

Algoritmul poate fi îmbunătățit prin folosirea unei strategii de pivotare la fiecare etapă de eliminare. Detalii despre pivotare sunt prezentate în paragraful 2.4.4.

2.4.3 Evaluarea algoritmului

Algoritmul prezentat anterior va fi analizat acum conform criteriilor prezentate în capitolul 1.

Din punct de vedere al timpului de calcul, complexitatea este cubică deoarece în etapa de eliminare există trei cicluri imbricate. Într-adevăr, considerând ca operații elementare orice operații algebrice (adunări, scăderi, înmulțiri, împărțiri), rezultă că în etapa de eliminare se execută un număr de operații egal cu

$$T_e = \sum_{k=1}^{n-1} [2(n-k) + 3](n-k) \approx \sum_{k=1}^{n-1} 2(n-k)^2 = 2 \frac{(n-1)n(2n-1)}{6} \approx \frac{2n^3}{3}. \quad (2.41)$$

În etapa de substituție se execută un număr de operații egal cu

$$T_s = \sum_{i=1}^{n-1} [2(n-i) + 2] \approx \sum_{i=1}^{n-1} [2(n-i)] = 2 \frac{n(n-1)}{2} \approx n^2. \quad (2.42)$$

Etapa de eliminare a algoritmului Gauss este cea mai costisitoare din punct de vedere al timpului de calcul, complexitatea ei fiind $T_e = O(2n^3/3)$, cu un ordin de mărime mai mare decât complexitatea etapei de substituție regresivă care este $T_s = O(n^2)$. În consecință, complexitatea algoritmului Gauss este cubică, dată de etapa de eliminare $T_{\text{Gauss}} = O(2n^3/3)$.

Din punct de vedere al necesarului de memorie, considerând locația ocupată de un număr real ca locație elementară de memorie, rezultă (din inspectarea declarațiilor din pseudocod) că este nevoie de $n^2 + 2n + 2$ locații de memorie. În consecință $M = O(n^2)$, algoritmul este pătratic din punct de vedere al necesarului de memorie.

Algoritmul Gauss prezentat în acest paragraf este un algoritm costisitor, el poate fi folosit până la valori de cel mult 1000 ale dimensiunii sistemului, foarte puțin pentru necesitățile problemelor reale. El poate deveni un algoritm eficient dacă este adaptat unor matrice cu structuri speciale (rare), caz analizat în paragraful 2.6.

Din punct de vedere al erorilor, în algoritmul Gauss apar, pe lângă erorile inerente, numai erori de rotunjire. Cu cât sistemul este de dimensiune mai mare, cu atât cresc numărul de operații efectuate și erorile acumulate datorită rotunjirii. O diminuare a erorilor de rotunjire se poate obține dacă se includ în algoritm strategii de pivotare.

Din punct de vedere al stabilității, algoritmul Gauss poate să nu fie stabil chiar dacă problema matematică este bine formulată și bine condiționată (numărul de condiționare al matricei \mathbf{A} este mic). Acest lucru se întâmplă atunci când numărul de condiționare al matricei \mathbf{U} este mare. Remediul îl constituie în acest caz pivotarea.

2.4.4 Strategii de pivotare

Elementele diagonale a_{kk} obținute în urma etapei de eliminare se numesc pivoți.

Deoarece transformările elementare făcute nu modifică determinantul sistemului și deoarece determinanul unei matrice tridiagonale este egal cu produsul elementelor de pe diagonală, înseamnă că determinatul matricei este egal cu produsul pivoților. Dacă problema este bine formulată matematic, atunci determinantul matricei coeficienților este diferit de zero, ceea ce înseamnă că toți pivoții sunt nenuli.

Pe parcursul algoritmului Gauss prezentat anterior, elementele de multiplicare se calculează ca $p = -a_{ik}/a_{kk}$. Valoarea a_{kk} folosită pentru calculul elementului de multiplicare este un pivot, ea rămânând în diagonala matricei triunghiulare. Se poate întâmpla însă ca pivotul să fie nul chiar dacă problema este bine formulată matematic, caz în care algoritmul Gauss prezentat anterior eșuează. În acest caz se poate face o operație de permutare (de linii sau/și de coloane), astfel încât pe poziția kk să fie adus un element nenul.

Pivotarea este o operație de permutare care urmărește obținerea unor valori nenule pentru pivotți. Ea este o operație care trebuie făcută înainte de calculul multiplicatorului.

Există mai multe strategii de pivotare, cele mai des întâlnite fiind: pivotarea pe linii, pe coloane, totală, diagonală.

Pivotarea pe linii (parțială) constă în schimbarea a două linii între ele (linia k și o linie aflata sub ea), ceea ce este echivalent cu scrierea ecuațiilor în altă ordine. Este cea mai simplă strategie de pivotare. Pivotul este căutat numai pe coloana de sub elementul diagonal.

Pivotarea pe coloane constă în schimbarea a două coloane între ele (coloana k și o coloană aflată la dreapta ei), ceea ce este echivalent cu renumerotarea variabilelor. Pivotul este căutat numai pe linie, la dreapta elementului diagonal. Permutările făcute trebuie memorate deoarece la finalul algoritmului soluția trebuie reconstituită în ordinea inițială.

Pivotarea totală (completă sau maximală) caută pivotul și pe linii și pe coloane, sub și la dreapta elementului diagonal. După aflarea lui vor fi permutate și două linii și două coloane. Algoritmul acestei permutări este cel mai costisitor și nu neaparat cel mai eficient.

Pivotarea diagonală caută pivotul pe diagonală. Vor fi permutate apoi două linii și două coloane. O astfel de pivotare poate fi utilă în cazul în care matricea este simetrică.

Pseudocodul strategiei de pivotare pe linii este următorul:

$p = 0$

pentru $i = k, n$; parcurge coloana k , de la elementul diagonal în jos

dacă $|a_{ik}| > p$ atunci

$l = i$; memorează poziția potențialului pivot

$p = |a_{ik}|$

•
•

dacă $p = 0$ atunci

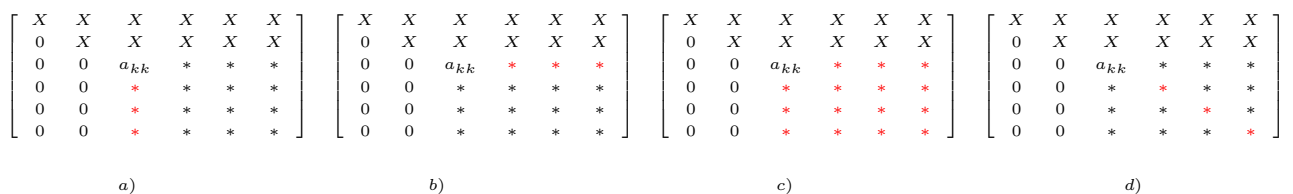


Figura 2.3: Zona de căutare a pivotului în pivotarea pe linii (a), pe coloane (b), totală (c), diagonală (d). Cu X sunt marcate elementele nenule ale căror valori nu se vor mai modifica.

scrie ”problema este prost formulată matematic”

altfel

pentru $j = k, n$; permută linia l cu linia k

$$p = a_{kj}$$

$$a_{kj} = a_{lj}$$

$$a_{lj} = p$$

•

$p = b_k$; permută termenii liberi

$$b_k = b_l$$

$$b_l = p$$

•

Pivotarea este absolut necesară dacă pe parcursul algoritmului Gauss se întâlnește un pivot nul. Ea are însă și un efect benefic asupra stabilității și acurateții soluției. Pentru a ilustra această afirmație, să reluăm exemplul discutat la paragraful 1.3.6.

$$\begin{cases} 10^{-20}x + y = 1 \\ x + y = 0, \end{cases} \quad (2.43)$$

având soluția corectă $(x, y) \approx (-1, 1)$. După eliminarea Gauss, sistemul triunghiular devine

$$\begin{cases} 10^{-20}x + y = 1 \\ (1 - 10^{20})y = -10^{20}. \end{cases} \quad (2.44)$$

Presupunând că zeroul mașinii este de ordinul 10^{-16} rezultă că sistemul ce se va rezolva la retrosubstituție va fi

$$\begin{cases} 10^{-20}x + y = 1 \\ 10^{20}y = -10^{20}. \end{cases} \quad (2.45)$$

Rezultatul final va fi $(x, y) = (0, 1)$, foarte departe de soluția adevărată. Deși sistemul inițial are matricea coeficienților bine condiționată $\kappa(\mathbf{A}) \approx 2.6$, matricea \mathbf{U} este extrem de prost condiționată $\kappa(\mathbf{U}) = 10^{40}$, ceea ce duce la eșecul rezolvării, în ciuda efectuării cu acuratețe a etapei de eliminare. Acest sistem nu este stabil la retrosubstituție și în consecință, întreaga rezolvare este compromisă. Deși problema este bine condiționată, algoritmul fără pivotare nu este stabil.

Aplicarea pivotării este o etapă esențială în rezolvarea directă a sistemelor de ecuații algebrice liniare. Pivotarea parțială asigură stabilitatea procedurii, cu un efort de implementare nesemnificativ. Pivotarea totală este rareori aplicată în practică deoarece căutarea pivoților și pe linii și pe coloane duce la o creștere semnificativă a timpului de calcul, nerealizând decât o îmbunătățire nesemnificativă a acurateții soluției.

2.4.5 Cazul sistemelor multiple

De multe ori este necesar să se rezolve mai multe sisteme de ecuații algebrice liniare care au aceeași matrice a coeficienților, dar termeni liberi diferiți. Un astfel de caz este cel al rezolvării circuitelor rezistive neliniare, în care pentru rezolvare se folosește metoda tangentelor paralele.

Fie m sisteme de ecuații algebrice liniare

$$\mathbf{Ax}^{(1)} = \mathbf{b}^{(1)}, \quad \mathbf{Ax}^{(2)} = \mathbf{b}^{(2)}, \quad \dots, \quad \mathbf{Ax}^{(m)} = \mathbf{b}^{(m)}, \quad (2.46)$$

care au aceeași matrice a coeficienților $\mathbf{A} \in \mathbb{R}^{n \times n}$ dar termeni liberi diferiți $\mathbf{b}^{(k)} \in \mathbb{R}^{n \times 1}$, unde $k = 1, m$. Se dorește rezolvarea celor m sisteme

$$\mathbf{Ax}^{(k)} = \mathbf{b}^{(k)}, \quad k = 1, \dots, m, \quad (2.47)$$

adică găsirea soluțiilor $\mathbf{x}^{(k)} \in \mathbb{R}^{n \times 1}$, în condițiile în care se cunoaște matricea \mathbf{A} și termenii liberi $\mathbf{b}^{(k)}$.

Dacă notăm cu

$$\mathbf{B} = [\mathbf{b}^{(1)} \quad \mathbf{b}^{(2)} \quad \dots \quad \mathbf{b}^{(m)}] \in \mathbb{R}^{n \times m} \quad (2.48)$$

matricea dreptunghiulară a termenilor liberi și cu

$$\mathbf{X} = [\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots \quad \mathbf{x}^{(m)}] \in \mathbb{R}^{n \times m} \quad (2.49)$$

matricea dreptunghiulară a necunoscutelor, atunci relațiile (2.47) se scriu compact

$$\mathbf{AX} = \mathbf{B}. \quad (2.50)$$

Varianta I - aplicarea succesivă a algoritmului Gauss

Cea mai proastă idee de a rezolva o astfel de problemă este de a apela de m ori succesiv algoritmul Gauss. Efortul total de calcul ar fi $m(2n^3/3 + n^2) \approx 2mn^3/3$. Deoarece matricea coeficienților este aceeași, etapa de eliminare este repetată inutil, de m ori.

Varianta II - adaptarea algoritmului Gauss pentru *rezolvarea simultană a sistemelor*

Presupunând că în momentul rezolvării toți termenii liberi sunt cunoscuți, atunci putem modifica algoritmul Gauss, astfel încât la eliminare să se prelucreze toți termenii liberi, iar la retrosubstituție, să se calculeze "simultan" componentele tuturor soluțiilor. Algoritmul complet al unei astfel de proceduri este următorul.

procedură Gauss_multiplu(n, m, a, B, X)

; rezolvă simultan sistemele algebrice liniare $aX = B$ prin metoda Gauss


```

întreg  $n$  ; dimensiunea sistemului
întreg  $m$  ; numărul de sisteme
tablou real  $a[n][n]$  ; matricea coeficienților - indici de la 1
tablou real  $B[n][m]$  ; matricea termenilor liberi
tablou real  $X[n][m]$  ; matricea soluție
întreg  $i, j, k$ 
real  $p, s$ 
; etapa de eliminare
pentru  $k = 1, n - 1$  ; parcurge sub-etape ale eliminării
    ; aici se poate introduce pivotarea
    pentru  $i = k + 1, n$  ; parcurge liniile
         $p = -a_{ik}/a_{kk}$  ; element de multiplicare
        pentru  $j = k + 1, n$  ; parcurge coloanele
             $a_{ij} = a_{ij} + pa_{kj}$ 
            •
        pentru  $j = 1, m$  ; parcurge coloanele termenilor liberi
             $b_{ij} = b_{ij} + pb_{kj}$ 
            •
        •
    •
; etapa de retrosubstituție
pentru  $k = 1, m$ 
     $x_{nk} = b_{nk}/a_{nn}$ 
    pentru  $i = n - 1, 1, -1$ 
         $s = 0$ 
        pentru  $j = i + 1, n$ 
             $s = s + a_{ij}x_{jk}$ 
            •
         $x_{ik} = (b_{ik} - s)/a_{ii}$ 
        •
    •
retur

```

În etapa de eliminare se face acum un număr de operații egal cu

$$\begin{aligned}
 T_e &= \sum_{k=1}^{n-1} [2(n-k) + 2m + 1](n-k) \approx \sum_{k=1}^{n-1} [2(n-k)^2 + 2m(n-k)] = \\
 &= 2 \frac{(n-1)n(2n-1)}{6} + 2m \frac{n(n-1)}{2} \approx \frac{2n^3}{3} + mn^2.
 \end{aligned} \tag{2.51}$$

În etapa de substituție se face un număr de operații egal cu

$$T_s = m \sum_{i=1}^{n-1} [2(n-i) + 2] \approx m \sum_{i=1}^{n-1} [2(n-i)] = 2m \frac{n(n-1)}{2} \approx mn^2. \quad (2.52)$$

În consecință, efortul total de calcul în aceasta variantă este de ordinul $O(2n^3/3 + 2mn^2)$, mai mic decât în cazul primei variante.

Varianta III - *rezolvarea succesivă* a sistemelor folosind calculul inversei

Dacă termenii liberi nu sunt cunoscuți simultan, atunci metoda anterioară nu se poate aplica. O nouă variantă ar fi calculul inversei \mathbf{A}^{-1} și aflarea fiecărei soluții $\mathbf{x}^{(k)} = \mathbf{A}^{-1}\mathbf{b}^{(k)}$ imediat ce este cunoscut termenul liber.

Calculul inversei se poate face folosind procedura Gauss pentru sisteme multiple rezolvate simultan, în care matricea termenilor liberi este chiar matricea unitate. Dacă $\mathbf{B} = \mathbf{I}$, atunci soluția sistemului $\mathbf{AX} = \mathbf{B}$ va fi chiar $\mathbf{X} = \mathbf{A}^{-1}$. Aceasta idee este descrisă de următorul pseudocod.

funcție invA(n, a)

; calculează inversa matricei a

întreg n

; dimensiunea matricei

tablou real $a[n][n]$

; matricea, indici de la 1

; alte declarații

....

pentru $i = 1, n$

pentru $j = 1, n$

$B_{ij} = 0$

•

$B_{ii} = 1$

•

Gauss_multiplu(n, n, a, B, X)

întoarce X

; X este inversa matricei

Complexitatea acestui calcul este obținută din complexitatea procedurii Gauss multiplu, pentru cazul particular $n = m$, deci $2n^3/3 + 2mn^2 = 8n^3/3$.

Calculul inversei este o operație foarte costisitoare din punct de vedere al timpului de calcul.

În final, pentru calculul fiecărei soluții, matricea inversă trebuie înmulțită cu vectorul termenilor liberi, algoritmul acestei înmulțiri având o complexitate de ordinul $2n^2$, așa cum poate fi estimat cu ușurință din pseudocodul următor:

```

funcție produs_Mv (n, M, v)
; calculează produsul dintre o matrice pătrată M și un vector coloană v
întreg n ; dimensiunea problemei
tablou real M[n][n] ; matricea, indici de la 1
tablou real v[n] ; vectorul
tablou real p[n] ; rezultatul p = Mv
; alte declarații
....
pentru i = 1, n
    p_i = 0
    pentru j = 1, n
        p_i = p_i + M_ij v_j
    •
•
întoarce p

```

Efortul total de calcul la rezolvarea succesivă a m sisteme algebrice liniare folosind calculul inversei, va fi $O(8n^3/3 + 2mn^2)$. Efortul de calcul este tot de ordin cubic, dar mai mare decât la rezolvarea simultană a sistemelor. Există însă și o variantă mai eficientă de rezolvare succesivă a sistemelor algebrice. Ea este bazată însă nu pe calculul inversei ci pe factorizarea matricei coeficienților.

2.5 Metoda factorizării LU

Metoda factorizării LU (*Lower-Upper*) este o metodă directă de rezolvare a unui sistem algebric liniar de dimensiune n

$$\mathbf{Ax} = \mathbf{b}, \quad (2.53)$$

bazată pe descompunerea matricei coeficienților într-un produs de două matrice

$$\mathbf{A} = \mathbf{LU}, \quad (2.54)$$

unde \mathbf{L} este o matrice inferior triunghiulară iar \mathbf{U} este o matrice superior triunghiulară (fig. 2.4).

Sistemul de rezolvat (2.53) este echivalent cu

$$\mathbf{LUx} = \mathbf{b}. \quad (2.55)$$

Dacă se notează

$$\mathbf{y} = \mathbf{Ux}, \quad (2.56)$$

$$\mathbf{A} = \begin{bmatrix} * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix} \quad \mathbf{L} = \begin{bmatrix} * & 0 & 0 & 0 & 0 & 0 \\ * & * & 0 & 0 & 0 & 0 \\ * & * & * & 0 & 0 & 0 \\ * & * & * & * & 0 & 0 \\ * & * & * & * & * & 0 \\ * & * & * & * & * & * \end{bmatrix} \quad \mathbf{U} = \begin{bmatrix} * & * & * & * & * & * \\ 0 & * & * & * & * & * \\ 0 & 0 & * & * & * & * \\ 0 & 0 & 0 & * & * & * \\ 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & * \end{bmatrix}$$

Figura 2.4: Structura termenilor în descompunerea $\mathbf{A} = \mathbf{LU}$. Cu * sunt marcate elementele nenule.

atunci rezolvarea sistemului (2.55) este echivalentă cu rezolvarea succesivă a sistemelor

$$\begin{aligned} \mathbf{L}\mathbf{y} &= \mathbf{b}, \\ \mathbf{U}\mathbf{x} &= \mathbf{y}. \end{aligned} \tag{2.57}$$

Rezolvarea sistemelor (2.57) va fi simplă deoarece ele au o structură particulară iar procedura de rezolvare va fi adaptată acestei structuri.

Pe scurt, rezolvarea sistemului se face prin *factorizarea* matricei coeficienților, urmată de rezolvarea a două sisteme de ecuații cu structuri particulare, prin *substituție progresivă și regresivă*.

2.5.1 Un exemplu simplu

Cea mai simplă metodă de factorizare provine din algoritmul Gauss. Matricea \mathbf{U} este chiar matricea coeficienților obținută la sfârșitul etapei de eliminare al algoritmului Gauss. Pentru exemplul prezentat în paragraful 2.4.1, matricea \mathbf{U} este

$$\mathbf{U} = \begin{bmatrix} 1 & 2 & -1 \\ 0 & 7 & -1 \\ 0 & 0 & -2/7 \end{bmatrix}. \tag{2.58}$$

Matricea \mathbf{L} are 1 pe diagonală, iar elementele nenule se obțin din ultimele valori nenule ale elementelor anulate în Gauss (-2 și 4 în relația (2.29) și -9 în relația (2.30)), împărțite la pivoți:

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ -2/1 & 1 & 0 \\ 4/1 & -9/7 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 4 & -9/7 & 1 \end{bmatrix}. \tag{2.59}$$

Se verifică ușor că $\mathbf{LU} = \mathbf{A}$, unde

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -1 \\ -2 & 3 & 1 \\ 4 & -1 & -3 \end{bmatrix}. \tag{2.60}$$

Rezolvarea sistemului $\mathbf{L}\mathbf{y} = \mathbf{b}$ devine

$$\begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 4 & -9/7 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ -2 \end{bmatrix}, \quad (2.61)$$

adică

$$\begin{cases} y_1 & = -1 \\ -2y_1 + y_2 & = 0 \\ 4y_1 - 9/7y_2 + y_3 & = -2 \end{cases} \quad (2.62)$$

de unde

$$\begin{aligned} y_1 &= -1 \\ y_2 &= 2y_1 = -2 \\ y_3 &= -2 - 4y_1 + 9/7y_2 = -4/7. \end{aligned} \quad (2.63)$$

Rezolvarea sistemului $\mathbf{U}\mathbf{x} = \mathbf{y}$ devine

$$\begin{bmatrix} 1 & 2 & -1 \\ 0 & 7 & -1 \\ 0 & 0 & -2/7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -2 \\ -4/7 \end{bmatrix}, \quad (2.64)$$

adică

$$\begin{cases} x_1 + 2x_2 - x_3 & = -1 \\ 7x_2 - x_3 & = -2 \\ -2/7x_3 & = -4/7. \end{cases} \quad (2.65)$$

de unde

$$\begin{aligned} x_3 &= (-4/7)/(-2/7) = 2, \\ x_2 &= (-2 + x_3)/7 = 0, \\ x_1 &= -1 - 2x_2 + x_3 = 1. \end{aligned} \quad (2.66)$$

2.5.2 Variante de factorizare

Nu există o descompunere unică a matricei \mathbf{A} în factorii \mathbf{L} și \mathbf{U} . Variantele considerate standard sunt: *Doolittle*, în care termenii diagonali ai matricei \mathbf{L} sunt unitari ($l_{ii} = 1$), *Crout*, în care termenii diagonali ai matricei \mathbf{U} sunt unitari ($u_{ii} = 1$) și *Cholesky* în care factorii sunt transpuși unul altuia ($\mathbf{L} = \mathbf{U}^T$).

Variantele Doolittle și Crout se aplică la orice matrice nesingulară, iar varianta Cholesky se aplică doar matricelor simetrice și pozitiv definite, la care se salvează simetria și după factorizare.

Pentru a exemplifica aceste variante, să considerăm cazul descompunerii în factori a unei matrice de dimensiune 2:

$$\begin{bmatrix} 3 & 2 \\ 6 & 1 \end{bmatrix} = \begin{bmatrix} l_{11} & 0 \\ l_{21} & l_{22} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} \\ 0 & u_{22} \end{bmatrix}. \quad (2.67)$$

Relațiile între componentele celor doi factori sunt

$$\begin{cases} l_{11}u_{11} = 3 \\ l_{11}u_{12} = 2 \\ l_{21}u_{11} = 6 \\ l_{21}u_{12} + l_{22}u_{22} = 1. \end{cases} \quad (2.68)$$

Un astfel de sistem este nedeterminat. Dar, dacă fixăm oricare două valori, restul de valori vor rezulta în mod unic. De aceea, pentru a găsi o factorizare LU unică, este necesar să se stabilească niște restricții suplimentare pentru matricele-factori. De exemplu, dacă alegem $l_{11} = l_{22} = 1$ obținem o soluție posibilă (numită varianta Doolittle), iar dacă alegem $u_{11} = u_{22} = 1$ obținem o altă soluție posibilă (numită varianta Crout). În cazul exemplului de mai sus

$$\begin{bmatrix} 3 & 2 \\ 6 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 3 & 2 \\ 0 & -3 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & -3 \end{bmatrix} \begin{bmatrix} 1 & 2/3 \\ 0 & 1 \end{bmatrix}. \quad (2.69)$$

Pentru o matrice simetrică și pozitiv definită se poate impune condiția de păstrare a simetriei după factorizare, ca de exemplu

$$\begin{bmatrix} 9 & 2 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 2/3 & \sqrt{5}/3 \end{bmatrix} \begin{bmatrix} 3 & 2/3 \\ 0 & \sqrt{5}/3 \end{bmatrix}. \quad (2.70)$$

2.5.3 Algoritmul variantei Doolittle

În cele ce urmează vom demonstra că algoritmul Gauss, extrem de puțin modificat, generează și factorizarea LU a matricei coeficienților, în varianta Doolittle.

Diferitele etape prin care trece matricea \mathbf{A} în algoritmul Gauss sunt etape de transformare liniară, respectiv de înmulțire cu o matrice. Mai precis, dacă notăm matricea inițială cu

$$\mathbf{A}_0 = \mathbf{A}, \quad (2.71)$$

atunci, după prima sub-etapă de eliminare se obține matricea \mathbf{A}_1 , care are zerouri pe prima coloană, sub diagonală (fig. 2.2). Această etapă de eliminare este echivalentă, din punct de vedere algebric, cu înmulțirea la stânga a matricei inițiale \mathbf{A} cu o matrice elementară notată \mathbf{E}_1 . Similar, efectuarea celei de a doua etape de eliminare este echivalentă cu înmulțirea cu o matrice elementară \mathbf{E}_2 ș.a.m.d:

$$\begin{aligned} \mathbf{A}_1 &= \mathbf{E}_1 \mathbf{A}_0, \\ \mathbf{A}_2 &= \mathbf{E}_2 \mathbf{A}_1 = \mathbf{E}_2 \mathbf{E}_1 \mathbf{A}_0, \\ &\dots \\ \mathbf{A}_{n-1} &= \mathbf{E}_{n-1} \mathbf{A}_{n-2} = \mathbf{E}_{n-1} \mathbf{E}_{n-2} \dots \mathbf{E}_2 \mathbf{E}_1 \mathbf{A}_0. \end{aligned} \quad (2.72)$$

Matricea obținută după ultima etapă de eliminare este superior triunghiulară, deci putem nota

$$\mathbf{U} = \mathbf{A}_{n-1}. \quad (2.73)$$

Dacă notăm produsul matricelor elementare de eliminare cu

$$\mathbf{E} = \mathbf{E}_{n-1}\mathbf{E}_{n-2}\cdots\mathbf{E}_2\mathbf{E}_1, \quad (2.74)$$

atunci relația (2.72) se scrie

$$\mathbf{U} = \mathbf{E}\mathbf{A}. \quad (2.75)$$

Pentru finalizarea demonstrației, vom arăta că matricea \mathbf{E} este nesingulară și, mai mult, inversa ei este triunghiular inferioară și, în consecință,

$$\mathbf{L} = \mathbf{E}^{-1}. \quad (2.76)$$

Pentru acesta trebuie să analizăm structura matricelor de eliminare. Vom scrie relațiile pentru cazul $n = 3$, dar generalizarea este imediată.

Matricea elementară \mathbf{E}_1 este cea care realizează prima sub-etapă de eliminare, deci

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a'_{22} & a'_{23} \\ 0 & a'_{32} & a'_{33} \end{bmatrix} = \mathbf{E}_1 \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}. \quad (2.77)$$

Se deduce ușor că \mathbf{E}_1 are 1 pe diagonală, iar pe prima coloană, sub diagonală, sunt exact elementele de multiplicare folosite pentru anularea coeficienților:

$$\mathbf{E}_1 = \begin{bmatrix} 1 & 0 & 0 \\ -a_{21}/a_{11} & 1 & 0 \\ -a_{31}/a_{11} & 0 & 1 \end{bmatrix}. \quad (2.78)$$

În general, o matrice de eliminare \mathbf{E}_k este inferior triunghiulară, are 1 pe diagonală, iar pe coloana k , sub diagonală, sunt exact elementele de multiplicare folosite pentru anularea coeficienților în sub-etapa k de eliminare. O astfel de matrice este inversabilă, iar inversa ei are 1 pe diagonală, iar pe coloana k , sub diagonală, sunt exact opusele elementelor de multiplicare folosite pentru anularea coeficienților în etapa k de eliminare. De exemplu

$$\mathbf{E}_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ a_{21}/a_{11} & 1 & 0 \\ a_{31}/a_{11} & 0 & 1 \end{bmatrix}, \quad \mathbf{E}_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & a'_{32}/a'_{22} & 1 \end{bmatrix}. \quad (2.79)$$

În consecință, matricea \mathbf{E} este inversabilă, iar inversa ei este

$$\mathbf{E}^{-1} = \mathbf{E}_1^{-1}\mathbf{E}_2^{-1}\cdots\mathbf{E}_{n-2}^{-1}\mathbf{E}_{n-1}^{-1}. \quad (2.80)$$

Prin înmulțire, se obține o matrice care are 1 pe diagonală, iar sub fiecare diagonală k sunt opusele elementelor de multiplicare folosite în etapa k de eliminare:

$$\mathbf{E}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ a_{21}/a_{11} & 1 & 0 \\ a_{31}/a_{11} & a'_{32}/a'_{22} & 1 \end{bmatrix}. \quad (2.81)$$

Matricea \mathbf{E}^{-1} este matricea inferior triunghiulară căutată (matricea \mathbf{L} în varianta Doolittle). Elementele nenule și diferite de 1 ale matricei \mathbf{L} pot fi memorate exact în triunghiul inferior al tabloului real bidimensional inițial, așa cum se face în secvența următoare.

; etapa de eliminare din metoda Gauss cu memorarea opuselor elementelor

; de multiplicare în triunghiul inferior al matricei

pentru $k = 1, n - 1$; parcurge sub-etape ale eliminării

pentru $i = k + 1, n$; parcurge liniile

$p = -a_{ik}/a_{kk}$; element de multiplicare

pentru $j = k + 1, n$; parcurge coloanele

$a_{ij} = a_{ij} + pa_{kj}$

•

$a_{ik} = -p$

•

•

De fapt, putem chiar să renunțăm la variabila p , obținând următorul pseudocod pentru factorizarea "în loc" a unei matrice.

procedură factorizare_LU(n, a)

; factorizează "in loc" matricea a

; varianta Doolittle

; declarații

...

pentru $k = 1, n - 1$; parcurge sub-etape ale eliminării

pentru $i = k + 1, n$; parcurge liniile

$a_{ik} = a_{ik}/a_{kk}$; element de multiplicare

pentru $j = k + 1, n$; parcurge coloanele

$a_{ij} = a_{ij} - a_{ik}a_{kj}$

•

•

•

retur

Prin factorizare "în loc" a matricei se înțelege faptul că elementele semnificative ale factorilor (adică elementele diferite de zero ale lui \mathbf{U} și elementele diferite de zero aflate sub diagonală ale matricei \mathbf{L}) sunt memorate în locul valorilor inițiale ale matricei \mathbf{A} . Se mai scrie " $\mathbf{A} = \mathbf{L} + \mathbf{U} - \mathbf{I}$ ", unde egalul este cel în sens aritmetic, dar matricea \mathbf{A} conține valorile după factorizare.

2.5.4 Calculul soluției după factorizare

Sistemul de rezolvat (2.53) este echivalent cu

$$\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{b}. \quad (2.82)$$

Dacă se notează

$$\mathbf{y} = \mathbf{U}\mathbf{x}, \quad (2.83)$$

atunci rezolvarea sistemului (2.82) este echivalentă cu rezolvarea succesivă a sistemelor

$$\mathbf{L}\mathbf{y} = \mathbf{b}, \quad (2.84)$$

$$\mathbf{U}\mathbf{x} = \mathbf{y}. \quad (2.85)$$

Sistemul (2.84), a cărui soluție o putem scrie formal $\mathbf{y} = \mathbf{L}^{-1}\mathbf{b}$, este un sistem cu structură inferior triunghiulară, a cărui rezolvare se face prin *substituție progresivă*:

$$\begin{cases} l_{11}y_1 = b_1, \\ l_{21}y_1 + l_{22}y_2 = b_2, \\ \dots \\ l_{n1}y_1 + l_{n2}y_2 + \dots + l_{nn}y_n = b_n, \end{cases} \Rightarrow \begin{cases} y_1 = b_1/l_{11}, \\ y_2 = (b_2 - l_{21}y_1)/l_{22}, \\ \dots \\ y_n = (b_n - l_{n1}y_1 - \dots - l_{n,n-1}y_{n-1})/l_{nn}. \end{cases} \quad (2.86)$$

Scris compact, soluția se calculează la substituție regresivă ca

$$y_1 = b_1/l_{11}, \quad (2.87)$$

$$y_i = \left(b_i - \sum_{j=1}^{i-1} l_{ij}y_j \right) / l_{ii}, \quad i = 2, \dots, n. \quad (2.88)$$

Sistemul (2.85), a cărui soluție o putem scrie formal $\mathbf{x} = \mathbf{U}^{-1}\mathbf{y}$, este un sistem cu structură superior triunghiulară, a cărui rezolvare se face prin *substituție regresivă*, ca la metoda Gauss:

$$x_n = y_n/u_{nn}, \quad (2.89)$$

$$x_i = \left(y_i - \sum_{j=i+1}^n u_{ij}x_j \right) / u_{ii}, \quad i = n-1, \dots, 1. \quad (2.90)$$

Pe baza formulelor (2.87) ÷ (2.90) putem scrie acum algoritmul rezolvării propriu-zise a sistemului, presupunând că matricea coeficienților este factorizată în loc în varianta Doolittle.

```

procedură rezolvă.LU( $n, a, b, x$ )
; rezolvă sistemul de ecuații  $ax = b$  prin factorizare LU
; matricea este presupusă a fi deja factorizată "în loc"
; varianta Doolittle
; declarații
...
; substituție progresivă
 $y_1 = b_1$  ; formula (2.87), unde  $l_{11} = 1$ 
pentru  $i = 2, n$ 
     $s = 0$ 
    pentru  $j = 1, i - 1$ 
         $s = s + a_{ij}y_j$  ; formula (2.88), unde  $\mathbf{L}$  este memorat în  $a$ 
    •
     $y_i = b_i - s$  ; deoarece  $l_{ii} = 1$ 
    •
; substituție regresivă
 $x_n = y_n/a_{nn}$  ; formula (2.89), unde  $\mathbf{U}$  este memorat în  $a$ 
pentru  $i = n - 1, 1, -1$ 
     $s = 0$ 
    pentru  $j = i + 1, n$ 
         $s = s + a_{ij}x_j$ 
    •
     $x_i = (y_i - s)/a_{ii}$ 
    •
retur

```

2.5.5 Evaluarea algoritmului. Strategii de pivotare.

Deoarece algoritmul de factorizare este derivat din Gauss, putem evalua complexitatea lui prin similitudine cu acesta. Factorizarea propriu-zisă are complexitatea etapei de eliminare Gauss, deci $T_f = O(2n^3/3)$, iar rezolvările prin substituție progresivă și regresivă au împreună $T_s = O(2n^2)$.

Din punct de vedere al necesarului de memorie, efortul este dat de memorarea tabloului bidimensional, deci $M = O(n^2)$.

Ca la metoda Gauss, în factorizarea LU nu există erori de trunchiere, iar erorile de rotunjire pot fi micșorate dacă se aplică strategii de pivotare de tipul celor descrise la paragraful 2.4.4. Acestea au ca efect și îmbunătățirea stabilității algoritmului.

Pivotarea poate fi descrisă algebric prin folosirea unor matrice de permutare. *O matrice de permutare este o matrice care are exact un element egal cu 1 pe fiecare linie și pe fiecare coloană, și 0 în rest.* Ea reprezintă o permutare a liniilor și/sau coloanelor matricei unitate. Atunci când o matrice de permutare înmulțește la stânga (respectiv la dreapta) o matrice oarecare, rezultatul este o matrice în care sunt permutate liniile (respectiv coloanele) matricei inițiale.

Pivotarea pe linie poate fi descrisă prin înmulțirea la stânga cu o matrice de permutare notată în cele ce urmează cu \mathbf{P} , iar pivotarea pe coloane poate fi descrisă prin înmulțirea la dreapta cu o matrice de permutare ce va fi notată cu \mathbf{Q} . Se poate demonstra că inversa unei matrice de permutare este o matrice de permutare și că produsul a două matrice de permutare este de asemenea o matrice de permutare [3].

Astfel, dacă înaintea primului pas al etapei de eliminare se efectuează o permutare parțială, relația (2.71) se scrie

$$\mathbf{A}_0 = \mathbf{P}_1 \mathbf{A}. \quad (2.91)$$

Presupunând că la fiecare etapă de eliminare se efectuează o permutare parțială, relațiile (2.72) se scriu

$$\begin{aligned} \mathbf{A}_1 &= \mathbf{E}_1 \mathbf{A}_0 = \mathbf{E}_1 \mathbf{P}_1 \mathbf{A}, \\ \mathbf{A}_2 &= \mathbf{E}_2 \mathbf{P}_2 \mathbf{A}_1 = \mathbf{E}_2 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 \mathbf{A}, \\ &\dots \\ \mathbf{A}_{n-1} &= \mathbf{E}_{n-1} \mathbf{P}_{n-1} \dots \mathbf{E}_2 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 \mathbf{A}. \end{aligned} \quad (2.92)$$

În argumentarea factorizării Doolittle, ultima relație se scrie

$$\mathbf{U} = \mathbf{E}_{n-1} \mathbf{P}_{n-1} \dots \mathbf{E}_2 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 \mathbf{A}. \quad (2.93)$$

Aceasta pune în evidență factorii L și U pentru o permutare pe linii a matricei \mathbf{A} . De exemplu (scriind pentru comoditate pentru $n = 4$):

$$\begin{aligned} \mathbf{U} &= \mathbf{E}_3 \mathbf{P}_3 \mathbf{E}_2 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_1 \mathbf{A} = \\ &= \underbrace{\mathbf{E}_3}_{\mathbf{E}'_3} \underbrace{\mathbf{P}_3 \mathbf{E}_2 \mathbf{P}_3^{-1}}_{\mathbf{E}'_2} \underbrace{\mathbf{P}_3 \mathbf{P}_2 \mathbf{E}_1 \mathbf{P}_2^{-1} \mathbf{P}_3^{-1}}_{\mathbf{E}'_1} \mathbf{P}_3 \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} = \\ &= \underbrace{\mathbf{E}'_3 \mathbf{E}'_2 \mathbf{E}'_1}_{\mathbf{L}^{-1}} \underbrace{\mathbf{P}_3 \mathbf{P}_2 \mathbf{P}_1}_{\mathbf{P}} \mathbf{A} \end{aligned} \quad (2.94)$$

Factorizarea cu pivotare pe linii se scrie deci compact ca

$$\mathbf{PA} = \mathbf{LU}, \quad (2.95)$$

unde toate elementele subdiagonale ale matricei \mathbf{L} sunt în modul mai mici sau egale cu 1, drept consecință a modului în care a fost ales pivotul. Interpretarea formulei (2.95) ar putea fi făcută astfel. Mai întâi se permută liniile matricei \mathbf{A} în conformitate cu matricea de permutare \mathbf{P} apoi se aplică algoritmul de rezolvare pentru matricea \mathbf{PA} . Am văzut că de fapt pivotarea parțială nu se realizează așa, pentru că matricea \mathbf{P} nu este cunoscută apriori. În implementarea practică \mathbf{P} nu este reprezentată explicit ca o matrice, ci liniile sunt permutate la fiecare etapă.

Permutarea coloanelor se reprezintă algebric prin înmulțirea la dreapta cu o matrice de permutare \mathbf{Q} . Factorizarea LU cu pivotare totală se poate scrie formal ca

$$\mathbf{PAQ} = \mathbf{LU}, \quad (2.96)$$

Pivotarea totală este rareori aplicată în practică deoarece căutarea pivotilor și pe linii și pe coloane duce la o creștere semnificativă a timpului de calcul, nerealizând decât o îmbunătățire nesemnificativă a stabilității.

2.5.6 Cazul sistemelor multiple

Cazul sistemelor multiple a fost prezentat mai întâi în paragraful 2.4.5. Metoda factorizării LU este utilă în acest caz deoarece factorizarea se face o singură dată, iar substituțiile se fac pentru cele m sisteme, pe măsură ce sunt cunoscuți termenii liberi. Efortul total de rezolvare prin factorizare a m sisteme algebrice de dimensiune n va fi deci $T = O(2n^3/3 + 2mn^2)$, mai mic decât cel necesar calculului inversei (tabelul 2.1).

2.5.7 Varianta Cholesky

Dacă matricea \mathbf{A} este simetrică, atunci este de dorit ca și factorizarea ei LU să păstreze această simetrie, adică

$$\mathbf{U} = \mathbf{L}^T. \quad (2.97)$$

Tabelul 2.1: Efort de calcul pentru rezolvarea sistemelor multiple.

Nr. sisteme	Metoda	Complexitate T
1	Gauss	$2n^3/3 + n^2$
	LU	$2n^3/3 + 2n^2$
m - simultan	Gauss	$2n^3/3 + 2mn^2$
m - succesiv	folosind inversa	$8n^3/3 + 2mn^2$
	LU	$2n^3/3 + 2mn^2$

Decompunerea LU care satisface condiția (2.97) se numește *factorizare Cholesky* și ea se poate realiza doar dacă \mathbf{A} este pozitiv definită.

Pentru a justifica această afirmație, să considerăm \mathbf{A} o matrice nesingulară și simetrică și presupunem că există \mathbf{L} o matrice triunghiular inferioară, nesingulară, astfel încât

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T. \quad (2.98)$$

Fie un vector coloană arbitrar \mathbf{x} nenul. Atunci $\mathbf{y} = \mathbf{L}^T\mathbf{x}$ va fi un vector coloană nenul (pentru că altfel $\mathbf{x} = (\mathbf{L}^T)^{-1}\mathbf{y}$ ar rezulta nul). Atunci

$$\mathbf{x}^T\mathbf{A}\mathbf{x} = \mathbf{x}^T\mathbf{L}\mathbf{L}^T\mathbf{x} = \mathbf{y}^T\mathbf{y} = \sum_{i=1}^n y_i^2 > 0.$$

În consecință, numai matricile care satisfac $\mathbf{x}^T\mathbf{A}\mathbf{x} > 0$ pentru vectori coloana \mathbf{x} nenuli pot admite o descompunere Cholesky. Aceste matrici sunt prin definiție matrici pozitiv definite.

Demonstrația de mai sus arată că pozitiv definirea matricei este o condiție necesară pentru factorizarea Cholesky. Se poate arăta însă că ea este și o condiție suficientă:

Teoremă: Dacă \mathbf{A} este o matrice simetrică și pozitiv definită, atunci factorizarea ei Cholesky există în mod unic, adică există în mod unic o matrice triunghiular inferioară \mathbf{L} cu elementele diagonale pozitive, astfel încât $\mathbf{A} = \mathbf{L}\mathbf{L}^T$.

Pentru a înțelege modul de generare al matricei \mathbf{L} , vom începe prin partiționarea matricelor astfel

$$\begin{bmatrix} \alpha & \mathbf{a}^T \\ \mathbf{a} & \hat{\mathbf{A}} \end{bmatrix} = \begin{bmatrix} \lambda & \mathbf{0} \\ \mathbf{l} & \hat{\mathbf{L}} \end{bmatrix} \begin{bmatrix} \lambda & \mathbf{l}^T \\ \mathbf{0} & \hat{\mathbf{L}}^T \end{bmatrix}. \quad (2.99)$$

Pentru intuirea mai clară a blocurilor, am notat scalarii cu litere grecești, vectorii cu litere mici aldine, iar matricile cu litere mari aldine. Din calculul elementului aflat pe poziția (1,1) rezultă $\alpha = \lambda^2$ de unde

$$\lambda = \sqrt{\alpha} \quad (2.100)$$

este primul element din \mathbf{L} . Din calculul blocului aflat pe poziția (2,1) rezultă $\mathbf{a} = \lambda\mathbf{l}$, de unde

$$\mathbf{l} = \mathbf{a}/\lambda \quad (2.101)$$

este prima coloană, aflată sub diagonală, a matricei \mathbf{L} . Din calculul blocului aflat pe poziția (2,2) rezultă $\hat{\mathbf{A}} = \mathbf{l}\mathbf{l}^T - \hat{\mathbf{L}}\hat{\mathbf{L}}^T$ de unde

$$\hat{\mathbf{L}}\hat{\mathbf{L}}^T = \hat{\mathbf{A}} - \mathbf{l}\mathbf{l}^T. \quad (2.102)$$

Matricea din membrul drept al acestei relații este complementul Schur al lui α :

$$\mathbf{S} = \hat{\mathbf{A}} - \mathbf{II}^T = \hat{\mathbf{A}} - \mathbf{aa}^T/\alpha. \quad (2.103)$$

Mai mult, se poate demonstra că \mathbf{S} este simetrică și pozitiv definită și, în consecință $\hat{\mathbf{L}}\hat{\mathbf{L}}^T$ este factorizarea ei Cholesky.

Astfel, echivalentul primei etape de eliminare Gauss, se poate scrie în cazul factorizării Cholesky ca

$$\mathbf{A} = \mathbf{A}_0 = \begin{bmatrix} \lambda & \mathbf{0} \\ \mathbf{1} & \hat{\mathbf{L}} \end{bmatrix} \begin{bmatrix} \lambda & \mathbf{I}^T \\ \mathbf{0} & \hat{\mathbf{L}}^T \end{bmatrix} = \begin{bmatrix} \lambda & \mathbf{0} \\ \mathbf{1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \lambda & \mathbf{I}^T \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \mathbf{L}_1 \mathbf{A}_1 \mathbf{L}_1^T. \quad (2.104)$$

Similar,

$$\begin{aligned} \mathbf{A}_1 &= \mathbf{L}_2 \mathbf{A}_2 \mathbf{L}_2^T, \\ &\vdots \\ \mathbf{A}_{n-1} &= \mathbf{L}_n \mathbf{A}_n \mathbf{L}_n^T, \end{aligned} \quad (2.105)$$

unde $\mathbf{A}_n = \mathbf{I}$.

Rezultă că, descompunerea finală a matricei \mathbf{A}

$$\mathbf{A} = \underbrace{\mathbf{L}_1 \mathbf{L}_2 \cdots \mathbf{L}_n}_{\mathbf{L}} \underbrace{\mathbf{L}_n^T \mathbf{L}_{n-1}^T \cdots \mathbf{L}_1^T}_{\mathbf{L}^T} \quad (2.106)$$

este unică deoarece la fiecare pas mărimile au fost determinate în mod unic.

În ce privește implementarea factorizării Cholesky, numai jumătate din matricea \mathbf{A} trebuie memorată. Nu numai necesarul de memorie este înjumătățit față de cazul general, dar și efortul de factorizare este mai mic. Pentru a justifica acest lucru, vom schița o variantă de algoritm care nu folosește o schemă particulară de memorarea a matricei.

Un mic exercițiu algebric arată că formulele de calcul pentru coloana k a matricei \mathbf{L} sunt

$$l_{kk} = \sqrt{a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2}, \quad k = 1, \dots, n \quad (2.107)$$

$$l_{ik} = (a_{ik} - \sum_{j=1}^{k-1} l_{ij} l_{kj}) / l_{kk}, \quad i = k + 1, \dots, n. \quad (2.108)$$

Acestea sunt implementate de pseudocodul următor:

procedură factorizare_LU_Cholesky(n, a, l)

; factorizează matricea a , presupusă simetrică și pozitiv definită

; întoarce matricea triunghiular inferioară l

; varianta Cholesky

; declarații

...

pentru $k = 1, n$

; parcurge sub-etape ale eliminării

pentru $i = k, n$

; calculează coloana, sub diagonală

$$s = a_{ik}$$

pentru $j = 1, k - 1$

$$s = s - l_{ij}l_{kj}$$

•

dacă $i = k$

$$l_{kk} = \sqrt{s}$$

altfel

$$l_{ik} = s/l_{kk}$$

•

•

•

retur

Efortul de calcul corespunzător este jumătate din efortul unei factorizări LU pentru matrice oarecare (egal cu efortul unei eliminări Gauss)

$$\begin{aligned} T_e &\approx \sum_{k=1}^n [2k(n-k)] = -2 \sum_{k=1}^n [(n-k-n)(n-k)] = -2 \left[\sum_{k=1}^n (n-k)^2 - n \sum_{k=1}^n (n-k) \right] = \\ &= -2 \left[\frac{(n-1)n(2n-1)}{6} - n \frac{n(n-1)}{2} \right] \approx -2 \left(\frac{2n^3}{6} - \frac{n^3}{2} \right) = \frac{n^3}{3}. \end{aligned} \quad (2.109)$$

În ce privește stabilitatea, algoritmul Cholesky este întotdeauna stabil și nu are nevoie de pivotare. Aceasta se datorează proprietăților speciale ale matricei \mathbf{A} , care fiind pozitiv definită este și diagonal dominantă.

2.6 Cazul matricelor rare

Algoritmii de rezolvare prezentați până acum (bazați în esență pe eliminarea Gauss) au complexitate destul de mare. Atât timpul de calcul este ridicat (ordin cubic), dar și spațiul de memorie este mare (ordin pătratic, datorită memorării matricei coeficienților).

În mod uzual, numerele reale se stochează în dublă precizie pe 8B, iar numerele întregi

pe $4B^1$. Un model mic al unei probleme, având de exemplu 1000 grade de libertate, va necesita o matrice a coeficienților cu 10^6 numere reale, deci 8 MB. Modele rafinate ale problemelor reale pot avea milioane de necunoscute. O problema cu un milion de necunoscute ar necesita un spațiu de memorare de 1 TB și nu ar încăpea în memoria calculatoarelor uzuale. Pe un PC obișnuit, nu s-ar putea lucra nici măcar cu o matrice de dimensiune 10000. Din fericire însă, formularea tuturor problemelor reale de inginerie poate fi făcută astfel încât ecuațiile care se scriu nu implică toate necunoscutele simultan. De fapt, pentru o ecuație, legăturile sunt chiar foarte puține comparate cu numărul total de necunoscute. De exemplu, în scrierea unei anumite ecuații nodale pentru un circuit, în ecuație apar doar termeni corespunzători laturilor care sunt incidente la nodul pentru care se scrie ecuația. Matricea conductanțelor nodale ar avea toate elementele nenule doar dacă graful circuitului ar fi un poligon complet, ceea ce nu se întâmplă în nici o aplicație reală.

Astfel de matrice nu se memorează în forma lor totală, ca tablouri bidimensionale, ci se vor memora doar elementele lor nenule. Mai mult, algoritmi de calcul se vor adapta acestei scheme de memorare, rezultând avantaje atât din punct de vedere al necesarului de memorie cât și al timpului de calcul.

O matrice care conține un număr foarte mare de elemente nule se numește matrice rară. Despre o matrice care nu este rară se spune că este *matrice densă sau plină*.

Se definește *densitatea unei matrice* ca fiind raportul dintre numărul de elemente nenule și numărul total de elemente al matricei. De obicei, algoritmi care exploatează raritatea matricelor devin avantajoși pentru valori mici ale densității. Nu se poate preciza o valoare exactă a densității care să demarceze matricile rare de matricile pline.

Dacă, pentru o anumită matrice care are și elemente nule, se poate elabora un algoritm care exploatează această structură și care este mai eficient decât algoritmul conceput pentru matricea plină, atunci aceasta este o matrice rară.

De exemplu, pentru un circuit cu $N = 1000$ noduri, presupunând că fiecare nod este conectat în medie cu alte 4, atunci matricea conductanțelor nodale va avea în medie 5 elemente nenule pe linie (4 termeni nediagonali și 1 termen diagonal), și deci va avea densitatea $d = 5N/N^2 = 5/N = 0.5\%$. Aceasta este cu siguranță o matrice pentru care merită aplicate tehnici de matrice rare.

¹1 Byte = 1 octet = 8 biți. Vezi și <http://en.wikipedia.org/wiki/Byte> pentru precizări asupra prefixului folosit pentru multipli de B.

2.6.1 Formate de memorare a matricelor rare

Există mai multe metode de a memora matricele rare. Cele mai generale nu fac nici o presupunere asupra structurii matricei, respectiv asupra pozițiilor în care se află elemente nenule.

Cea mai naturală metodă ar fi cea în care se memorează doar valorile nenule și poziția lor în matrice, dată prin indici de linie (r_{idx}) și de coloană (c_{idx}). Așa se face de exemplu în formatul de fișier ASCII al site-ului *Matrix Market*². Astfel, pentru un tablou bidimensional de dimensiune $m \times n$, în loc să se memoreze toate cele mn valori (inclusiv zerouri) într-un spațiu de $M_{plin} = 8mn$ B, sunt necesare doar n_{nz} locații de memorie pentru numere reale și $2n_{nz}$ locații de memorie pentru numere întregi, deci în total $M_{rar,coord} = 8 * n_{nz} + 4 * 2n_{nz} = 16n_{nz}$ B.

De exemplu, matricea \mathbf{M} de dimensiune 3×4 , având 6 elemente nenule, va fi memorată într-un vector val de dimensiune 6, și doi vectori de numere întregi, de dimensiune 6, astfel:

$$\mathbf{M} = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 0 & 5 & 1 \\ 2 & 3 & 0 & 7 \end{bmatrix} \Rightarrow \begin{cases} val & = [4 \ 5 \ 1 \ 2 \ 3 \ 7] \\ r_{idx} & = [1 \ 2 \ 2 \ 3 \ 3 \ 3] \\ c_{idx} & = [1 \ 3 \ 4 \ 1 \ 2 \ 4] \end{cases}$$

În acest exemplu, valorile nenule au fost citite pe linie, de la stânga la dreapta și de sus în jos, dar ele pot fi memorate în orice ordine, evident cu permutarea corespunzătoare a indicilor.

Memorarea poate fi și mai eficientă decât atât. Informația din vectorul r_{idx} poate fi comprimată, indicându-se doar locul unde se schimbă valoarea lui r_{idx} . Acesta este *formatul (vechi) Yale*, cunoscut și sub sigla *CRS - Compressed Row Storage*. O matrice \mathbf{M} , de dimensiuni $m \times n$ cu un număr de n_{nz} de elemente nenule, va fi stocată cu ajutorul a trei tablouri unidimensionale astfel:

- un tablou unidimensional val , de dimensiune n_{nz} , care conține toate valorile nenule, de la stânga la dreapta și de sus în jos;
- un talou unidimensional r_ptr , de dimensiune $m + 1$ (câte un element pentru fiecare linie, plus un element adițional care marchează sfârșitul tabloului), care conține indicii ce indică în tabloul val locurile în care încep liniile. Mai precis, linia i a matricei inițiale are valorile în val de la poziția $r_ptr(i)$ la poziția $r_ptr(i + 1) - 1$;
- un talou unidimensional r_idx , de dimensiune n_{nz} , care conține, pentru fiecare element din val , indicele coloanei pe care se află.

²Site-ul Matrix Market este accesibil la <http://math.nist.gov/MatrixMarket>.

Ordinul de complexitate din punct de vedere al necesarului de memorie este deci $M_{rar,CRS} = 8n_{nz} + 4(m + 1) + 4n_{nz} = 12n_{nz} + 4(m + 1)$ B.

De exemplu, aceeași matrice \mathbf{M} de dimensiune 3×4 , având 6 elemente nenule, va fi memorată într-un vector val de dimensiune 6, un vector r_ptr de dimensiune $3 + 1$ și un vector c_idx de dimensiune 6, astfel

$$\mathbf{M} = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 0 & 5 & 1 \\ 2 & 3 & 0 & 7 \end{bmatrix} \Rightarrow \begin{cases} val & = [4 \ 5 \ 1 \ 2 \ 3 \ 7] \\ r_ptr & = [1 \ 2 \ 4 \ 7] \\ c_idx & = [1 \ 3 \ 4 \ 1 \ 2 \ 3] \end{cases}$$

Un raționament similar se poate face pe coloane, formatul numindu-se *CCS - Compressed Column Storage*, fiind formatul folosit de colecția de matrice *Harwell - Boeing*, disponibilă de asemenea pe site-ul *Matrix Market*. Necesarul de memorie este $M_{rar,CCS} = 8n_{nz} + 4(n + 1) + 4n_{nz} = 12n_{nz} + 4(n + 1)$ B.

Același exemplu, memorat în format CCS este:

$$\mathbf{M} = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 0 & 5 & 1 \\ 2 & 3 & 0 & 7 \end{bmatrix} \Rightarrow \begin{cases} val & = [4 \ 2 \ 3 \ 5 \ 1 \ 7] \\ c_ptr & = [1 \ 3 \ 4 \ 5 \ 7] \\ r_idx & = [1 \ 3 \ 3 \ 2 \ 2 \ 3] \end{cases}$$

Dacă matricile au o structură particulară, atunci se poate alege o variantă și mai eficientă, care să țină cont de structură. Un exemplu ilustrativ este cel al *matricelor bandă* ale căror sub-diagonale se memorează în locații consecutive. De exemplu, o *matrice tridiagonală* are elemente nenule numai pe diagonala principală (notată q în fig. 2.5) și pe cele două subdiagonale vecine (p și r).

Așa cum sugerează notația din figura 2.5, o matrice tridiagonală poate fi memorată cu ajutorul a trei vectori, eventual grupați într-o matrice cu 3 linii și n coloane:

$$\mathbf{M}_{rar} = \begin{bmatrix} 0 & p_2 & p_3 & \cdots & p_{n-1} & p_n \\ q_1 & q_2 & q_3 & \cdots & q_{n-1} & q_n \\ r_1 & r_2 & r_3 & \cdots & r_{n-1} & 0 \end{bmatrix}$$

$$\mathbf{M} = \begin{bmatrix} q_1 & r_1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ p_2 & q_2 & r_2 & 0 & \cdots & 0 & 0 & 0 \\ 0 & p_3 & q_3 & r_3 & \cdots & 0 & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & 0 & \cdots & p_{n-1} & q_{n-1} & r_{n-1} \\ 0 & 0 & 0 & 0 & \cdots & 0 & p_n & q_n \end{bmatrix}$$

Figura 2.5: Structura unei matrice tridiagonale.

Acest mod de stocare este cunoscut și sub sigla *CDS - Compressed Diagonal Storage*. Se observă într-un astfel de mod de stocare vor fi incluse și elemente nule, fără semnificație, dar numărul lor este ne semnificativ.

2.6.2 Metode directe pentru matrice rare

Memorarea rară a matricelor nu are nici o valoare dacă algoritmi de calcul nu sunt adaptați structurii de date. De exemplu, algoritmul Gauss descris în paragraful 2.4.2 aplicat unei matrice tridiagonale ar face o mulțime de calcule inutile, mai precis operații algebrice cu zerouri, ale căror rezultate ar fi tot zerouri. Algoritmul poate fi regândit cu ușurință pentru matricea tridiagonală. Matricea coeficienților are la începutul etapei k de eliminare următoarea structură (exemplificată pentru $k = 4$):

$$\begin{bmatrix} * & * & 0 & 0 & 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & * & * & 0 & 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & * & * & 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & q_k & r_k & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & p_{k+1} & q_{k+1} & r_{k+1} & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & 0 & p_{k+2} & q_{k+2} & r_{k+2} & \cdots & 0 & 0 \\ \cdots & & & & & & & & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & p_n & q_n \end{bmatrix}$$

Se observă că singurul element de multiplicare ce trebuie calculat este $m = -p_{k+1}/q_k$, singura modificare suferind-o ecuația $k + 1$, mai precis doar un singur termen în membrul stâng al acesteia $q_{k+1} = q_{k+1} + m * r_k$, și termenul liber corespunzător.

În ceea ce privește etapa de retrosubstituție, matricea triangularizată va avea doar două diagonale, calculul se va face tot prin substituție regresivă, ultima componentă fiind calculată mai întâi

$$x_n = b_n/q_n, \quad (2.110)$$

iar o componentă oarecare x_i va rezulta din ecuația i după eliminare

$$q_i x_i + r_i x_{i+1} = b_i \quad \Rightarrow \quad x_i = (b_i - r_i x_{i+1})/q_i, \quad i = n - 1, \dots, 1. \quad (2.111)$$

Algoritmul următor implementează metoda Gauss adaptată unei matrice tridiagonale, memorată cu ajutorul a trei vectori, așa cum a fost descris mai sus.

procedură Gauss_tridiag(n, p, q, r, b, x)

; rezolvă sistemul algebric liniar $ax = b$ prin metoda Gauss

; matricea a este tridiagonală, memorată în p, q, r

întreg n

; dimensiunea sistemului

```

tablou real  $p[n], q[n], r[n]$  ; "matricea" coeficienților - indici de la 1
tablou real  $b[n]$  ; vectorul termenilor liberi
tablou real  $x[n]$  ; vectorul soluție
întreg  $i, k$ 
; etapa de eliminare din metoda Gauss
pentru  $k = 1, n - 1$  ; parcurge sub-etape ale eliminării
     $m = -p_{k+1}/q_k$  ; element de multiplicare
     $q_{k+1} = q_{k+1} + mr_k$  ; modifică element în linia  $k + 1$ 
     $b_{k+1} = b_{k+1} + mb_k$  ; modifică termenul liber al ecuației  $k + 1$ 
•
; etapa de retrosubstituție
 $x_n = b_n/q_n$ 
pentru  $i = n - 1, 1, -1$ 
     $x_i = (b_i - r_i x_{i+1})/q_i$ 
•
retur

```

Etapa de eliminare are acum o complexitate de ordinul $5n$, iar etapa de retrosubstituție de ordinul $3n$, complexitatea totală fiind liniară atât din punct de vedere al timpului de calcul $T = O(8n)$, cât și din punct de vedere al necesarului de memorie $T = O(5n)$.

În cazul unor matrice rare care nu sunt tridiagonale și nici nu au o structură particulară, algoritmul Gauss trebuie adaptat unor tehnici de memorare de tip CRS sau CCS. În timpul etapei de eliminare matricea se poate umple, astfel încât pivotarea urmărește nu numai stabilitatea numerică (pivoți nenuli de modul cât mai mare pentru acumularea unor erori de rotunjire cât mai mici), ci și minimizarea umplerilor, adică a elementelor nenule nou apărute. Găsirea pivotării optime care ar alege dintre toate pivotările posibile la un moment dat ar conduce la un algoritm non-polinomial, lipsit de importanță practică. Strategiile de pivotare care se implementează folosesc tehnici euristice, care dau o soluție pseudo-optimală dar au un grad de complexitate polinomial.

Motivul cel mai important pentru care calculul explicit al inversei unei matrice trebuie evitat este acela că, în marea majoritate a aplicațiilor practice, inversa unei matrice este plină, chiar dacă matricea este rară. Fenomenul prin care elemente inițial nule devin nenule prin calcul, se numește *fenomen de umplere*. La matrice rare de dimensiuni mari inversarea este practic imposibilă datorită acestui fenomen.

Din fericire, factorizarea unei matrice rare poate salva raritatea dacă matricea are o anumită structură. De exemplu o matrice cu structura \mathbf{A}_1 din figura 2.6 prin factorizare rămâne rară, în schimb matricea \mathbf{A}_2 se umple prin factorizare. Cele două matrice pot reprezenta același sistem, ultima ecuație din \mathbf{A}_1 fiind prima ecuație în \mathbf{A}_2 și ultima ne-

$$\mathbf{A}_1 = \begin{bmatrix} * & 0 & \cdots & 0 & 0 & * \\ 0 & * & \cdots & 0 & 0 & * \\ \cdots & & & & & \\ 0 & 0 & \cdots & * & 0 & * \\ 0 & 0 & \cdots & 0 & * & * \\ * & * & \cdots & * & * & * \end{bmatrix} \quad \mathbf{A}_2 = \begin{bmatrix} * & * & * & \cdots & * & * \\ * & * & 0 & \cdots & 0 & 0 \\ * & 0 & * & \cdots & 0 & 0 \\ \cdots & & & & & \\ * & 0 & \cdots & 0 & * & 0 \\ * & 0 & \cdots & 0 & 0 & * \end{bmatrix}$$

Figura 2.6: Matricea \mathbf{A}_1 are factorii LU rari, în timp ce matricea \mathbf{A}_2 are factorii LU plini.

cunoscută în \mathbf{A}_1 fiind prima necunoscută în \mathbf{A}_2 . Una din tehnicile euristice de pivotare este cea sugerată de aceste două cazuri și anume, pivotarea urmărește plasarea la sfârșit a ecuațiilor ”pline”.

Structura matricei joacă deci un rol important în conceperea algoritmului de rezolvare.

Deoarece factorizarea generează umpleri, un algoritm eficient va trebui să știe în avans unde apar umplerile astfel încât să aloce memorie pentru structurile de date ce memorează factorii L și U. Această etapă se numește *factorizare simbolică*. Abia după ce se efectuează factorizarea simbolică, are loc factorizarea numerică. Factorizarea simbolică este o etapă mai puțin costisitoare dar mai complicat de elaborat. Ea se bazează pe construirea unor grafuri asociate matricei și deducerea unor arbori de eliminare. Un document ce ilustrează în detaliu aceste concepte este [13], iar o carte excelentă dedicată exclusiv acestui subiect este [2].

În concluzie, pentru a obține algoritmi performanți pentru o anumită problemă, trebuie folosite structuri de date eficiente, care să valorifice eventualele proprietăți particulare ale datelor (de exemplu raritatea matricelor), iar algoritmi de calcul trebuie adaptați acestor structuri de date. Metodele directe de rezolvare a sistemelor rare de ecuații algebrice au patru pași principali: permutarea sistemului în vederea asigurării unei stabilități maxime, factorizarea simbolică, factorizarea numerică și rezolvarea celor două sisteme triunghiulare.

2.7 Metode iterative staționare - generalități

Metodele iterative pentru rezolvarea sistemelor algebrice liniare se bazează pe generarea unui șir de aproximații ale soluției, care se dorește a fi convergent către soluția exactă. În consecință, din punct de vedere teoretic, soluția poate fi obținută cu o astfel de metodă numai într-un număr infinit de pași. Cum algoritmul poate face numai un număr finit de pași, înseamnă că aceste metode sunt afectate și de erori de trunchiere,

nu numai de erori de rotunjire. Aparent, acesta ar fi un dezavantaj față de metodele directe care găsesc soluția, teoretic într-un număr finit de pași. În realitate, metodele iterative au proprietatea remarcabilă de a face autocorecția erorilor astfel încât se poate întâmpla ca soluția obținută printr-o metodă iterativă să fie mai precisă decât soluția obținută printr-o metodă directă. Metodele iterative sunt atractive și pentru faptul că oferă utilizatorului posibilitatea de a decide el însuși asupra compromisului dintre timpul de calcul și acuratețea soluției, calculând mai multe iterații dacă se dorește o soluție mai precisă, sau oprind algoritmul mai repede atunci când se dorește un timp mai scurt de rezolvare.

Un avantaj major al metodelor iterative este acela că, aplicate unor matrice rare, ele nu generează umpleri ale acestora, așa cum se întâmplă în cazul metodelor directe. Dacă memoria disponibilă pentru rezolvarea unei probleme printr-o metodă directă nu este suficientă, atunci o metodă iterativă este singura opțiune.

Dezavantajul metodelor iterative este acela că se poate întâmpla ca ele să nu fie convergente pentru o problemă dată, caz în care problema trebuie reformulată înainte de a fi rezolvată cu o metodă iterativă. Reformularea se poate face cu ajutorul unei matrice de transformare, numită *precondiționator*. Un preconditionator bun îmbunătățește convergența metodei iterative suficient de mult pentru a amortiza costul construirii și aplicării preconditionatorului. Fără preconditionare, o metodă iterativă s-ar putea chiar să eșueze.

2.7.1 Ideea de bază a metodelor iterative staționare

Ideea metodelor iterative staționare pentru rezolvarea unui sistem algebric liniar de dimensiune n

$$\mathbf{Ax} = \mathbf{b} \quad (2.112)$$

este de a construi un șir de aproximații $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}, \dots$ care să convergă către soluția exactă a ecuației (2.112), notată \mathbf{x}^* :

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*, \quad \text{unde } \mathbf{Ax}^* = \mathbf{b}. \quad (2.113)$$

Fiecare aproximație corespunzătoare unei iterații k este un vector cu n componente:

$$\mathbf{x}^{(k)} = \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \\ \vdots \\ x_n^{(k)} \end{bmatrix} \in \mathbb{R}^{n \times 1}.$$

Un astfel de algoritm are nevoie de o inițializare $\mathbf{x}^{(0)}$, un mod de generare a șirului de iterații și un criteriu de oprire.

Inițializarea șirului de iterații este, în principiu, arbitrară. Dacă însă sunt cunoscute apriori informații despre soluție, de exemplu o încadrare a ei, atunci este recomandat ca inițializarea să fie făcută cât mai aproape de soluție, algoritmul iterativ putând fi mai rapid convergent.

Șirul de iterații se generează *recursiv*, adică soluția la o nouă iterație se calculează exclusiv în funcție de soluția la iterația anterioară:

$$\mathbf{x}^{(k)} = F(\mathbf{x}^{(k-1)}), \quad (2.114)$$

unde F este o funcție continuă a cărei construcție o vom preciza în cele ce urmează. Trecând la limită pentru $k \rightarrow \infty$ relația (2.114) rezultă că soluția exactă \mathbf{x}^* va satisface relația

$$\mathbf{x}^* = F(\mathbf{x}^*), \quad (2.115)$$

ceea ce înseamnă că \mathbf{x}^* este punct fix pentru aplicația F .

În concluzie, soluția exactă a sistemului de ecuații este și punct fix pentru F . Rezolvarea sistemului de ecuații algebrice liniare se face prin căutarea unui punct fix pentru F .

Este evident atunci că trebuie să existe o legătură între sistemul algebric inițial, definit de matricea coeficienților \mathbf{A} împreună cu vectorul termenilor liberi \mathbf{b} și funcția F . Construcția aplicației F se bazează pe descompunerea matricei \mathbf{A} într-o diferență de două matrice

$$\mathbf{A} = \mathbf{B} - \mathbf{C}. \quad (2.116)$$

Sistemul de rezolvat (2.112) este echivalent cu

$$\mathbf{B}\mathbf{x} = \mathbf{C}\mathbf{x} + \mathbf{b} \quad (2.117)$$

sau

$$\mathbf{x} = \mathbf{M}\mathbf{x} + \mathbf{u}, \quad (2.118)$$

unde

$$\begin{aligned} \mathbf{M} &= \mathbf{B}^{-1}\mathbf{C}, \\ \mathbf{u} &= \mathbf{B}^{-1}\mathbf{b}. \end{aligned} \quad (2.119)$$

Matricea $\mathbf{M} \in \mathbb{R}^{n \times n}$ definită de (2.119) se numește *matrice de iterație*. Relația (2.118) sugerează modul de definire a aplicației F :

$$F(\mathbf{x}) = \mathbf{M}\mathbf{x} + \mathbf{u}, \quad (2.120)$$

iar calculul recursiv (2.114) poate fi scris explicit ca

$$\mathbf{x}^{(k)} = \mathbf{B}^{-1}\mathbf{C}\mathbf{x}^{(k-1)} + \mathbf{B}^{-1}\mathbf{b}. \quad (2.121)$$

În paragraful 2.4.5 s-a arătat că trebuie evitat calculul explicit al inversei unei matrice. De aceea, orice algoritm iterativ nu implementează relația (2.121) ci

$$\mathbf{B}\mathbf{x}^{(k)} = \mathbf{C}\mathbf{x}^{(k-1)} + \mathbf{b}, \quad (2.122)$$

care reprezintă rezolvarea unui nou sistem de ecuații algebrice liniare, având \mathbf{B} ca matrice a coeficienților.

Încă nu am precizat exact cum se descompune matricea coeficienților în (2.116). Există o infinitate de moduri de a face acest lucru. Ideea este de a avea o matrice \mathbf{B} cu o structură particulară, astfel încât sistemul (2.122) să poată fi rezolvat cu un efort de calcul mic, folosind o rezolvare adaptată structurii particulare a matricei \mathbf{B} . De exemplu, în metoda Jacobi \mathbf{B} este diagonală și rezolvarea sistemului (2.122) este imediată, iar în metoda Gauss-Seidel \mathbf{B} este inferior triunghiulară, rezolvarea sistemului (2.122) făcându-se prin substituție progresivă.

2.7.2 Criteriul de oprire

În ceea ce privește condiția de oprire, se pune problema dacă șirul de iterații este convergent sau nu. Convergența șirului de iterații către soluția exactă \mathbf{x}^* este echivalentă cu convergența șirului $\mathbf{x}^{(k)} - \mathbf{x}^*$ către zero. Ideal ar fi ca iterațiile să fie oprite atunci când eroarea $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \varepsilon$, unde ε este o valoare impusă de utilizator. Cum eroarea nu poate fi calculată deoarece nu se cunoaște soluția exactă, în practică se implementează o condiție de oprire bazată de criteriul de convergență Cauchy și anume construcția șirului iterațiilor este oprită atunci când

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \varepsilon, \quad (2.123)$$

iar soluția numerică aproximativă este cea obținută la ultima iterație calculată.

Se poate întâmpla însă ca șirul iterațiilor să nu fie convergent. Din acest motiv, criteriul de oprire va include și o condiție legată de numărul maxim de iterații admis. Dacă eroarea dorită ε nu se atinge într-un număr de iterații mai mic decât o limită maximă impusă, atunci procedeul iterativ se consideră neconvergent și rezultatul obținut nu este de încredere. S-ar putea întâmpla ca procedeul iterativ să fie convergent mai lent, astfel încât să fie nevoie de mai multe iterații pentru obținerea unei precizii dorite. Este responsabilitatea utilizatorului de a urmări convergența procesului iterativ și, în cazul unui proces iterativ convergent, de a stabili un compromis între timpul de calcul și acuratețe.

În consecință, procedurile iterative de rezolvare a sistemelor algebrice liniare vor avea ca parametri de intrare, pe lângă mărimile ce definesc sistemul (coeficienți și termeni liberi), o eroare ce reprezintă criteriul dorit de oprire a

iterațiilor și un număr maxim de iterații, util pentru a asigura oprirea procedurii în caz de neconvergență.

Este evident că valoarea actuală a erorii folosită în momentul apelului unei astfel de proceduri iterative nu are sens să fie mai mică decât zeroul mașinii.

2.7.3 Intermezzo despre vectori și valori proprii

Pentru înțelegerea mai ușoară a paragrafului următor, vom reaminti câteva noțiuni legate de vectori și valori proprii. Aceste noțiuni sunt folosite ca instrumente de analiză în metodele iterative.

Prin definiție, un vector propriu \mathbf{v} al unei matrice pătrate reale \mathbf{M} , de dimensiune n este acel vector nenul, pentru care există un scalar λ astfel încât

$$\mathbf{M}\mathbf{v} = \lambda\mathbf{v}. \quad (2.124)$$

Altfel spus, dacă ne imaginăm reprezentarea geometrică a unui vector în spațiu, prin aplicarea matricei \mathbf{M} asupra lui, vectorul nu se rotește (pentru că dacă λ este real, atunci $\lambda\mathbf{v}$ este colinar cu \mathbf{v}). El își poate schimba cel mult sensul, dar nu direcția. Se observă imediat că vectorii proprii ai unei matrice nu sunt unici. Dacă \mathbf{v} este un vector propriu, atunci și vectorul scalat $\alpha\mathbf{v}$ este de asemenea vector propriu. Mărimea λ se numește valoare proprie a matricei \mathbf{M} asociată vectorului propriu \mathbf{v} .

Metodele iterative depind adesea de înmulțirea dintre o matrice și un vector, operație care se face în mod repetat. Dacă \mathbf{v} este un vector propriu a lui \mathbf{M} , atunci înmulțirea repetată cu \mathbf{M} conduce la

$$\mathbf{M}^k\mathbf{v} = \lambda^k\mathbf{v}. \quad (2.125)$$

Dacă $|\lambda| < 1$, atunci norma vectorului rezultat $\|\mathbf{M}^k\mathbf{v}\|$ va tinde către 0 când k tinde la infinit (un exemplu ilustrativ este prezentat în fig. 2.7), în timp ce dacă $|\lambda| > 1$, atunci norma vectorului rezultat va tinde către infinit.

Dacă matricea \mathbf{M} este simetrică, atunci se poate demonstra că ea are n vectori proprii liniar independenți, notați $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(n)}$. Evident că această mulțime nu este unică, fiecare vector propriu putând fi scalat. Fiecare vector propriu are asociată o valoare proprie. Aceste valori proprii sunt unice pentru o matrice dată.

Dacă \mathbf{M} este aplicată unui vector \mathbf{u} care nu este vector propriu, atunci acesta poate fi scris ca o combinație liniară de alți vectori al căror comportament este înțeles. Dacă cei n vectori proprii ai lui \mathbf{M} formează o bază, atunci $\mathbf{u} = \sum_{i=1}^n \alpha_i \mathbf{v}^{(i)}$. Cum produsul matrice-vector este distributiv, atunci se poate urmări efectul aplicării lui \mathbf{M} asupra fiecărui vector separat:

$$\mathbf{M}^k\mathbf{u} = \alpha_1\lambda_1^k\mathbf{v}^{(1)} + \dots + \alpha_n\lambda_n^k\mathbf{v}^{(n)}. \quad (2.126)$$

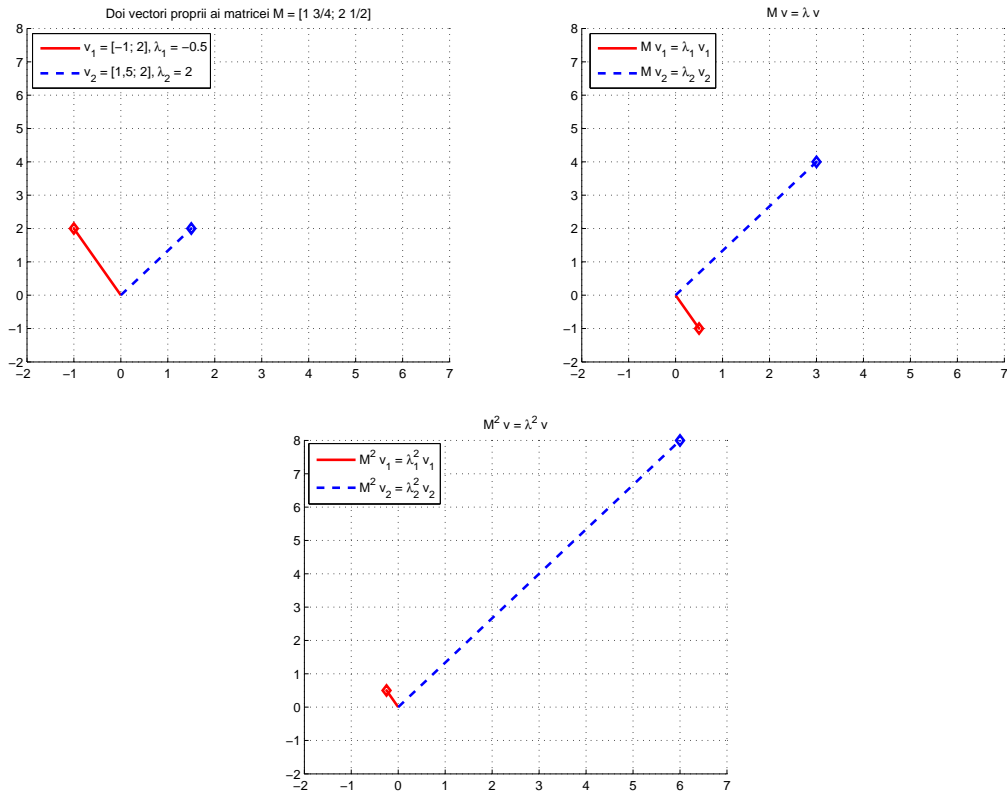


Figura 2.7: Înmulțirea repetată dintre matricea $\mathbf{M} = [1, 3/4; 2, 1/2]$ și vectorii ei proprii $\mathbf{v}_1 = [-1, 2]$ (corespunzător valorii proprii $\lambda_1 = -0.5$) și $\mathbf{v}_2 = [1.5, 2]$ ($\lambda_2 = 2$): sus stânga \mathbf{v}_1 și \mathbf{v}_2 ; sus dreapta $M\mathbf{v}_1 = \lambda_1\mathbf{v}_1$ și $M\mathbf{v}_2 = \lambda_2\mathbf{v}_2$; jos $M^2\mathbf{v}_1 = \lambda_1^2\mathbf{v}_1$ și $M^2\mathbf{v}_2 = \lambda_2^2\mathbf{v}_2$.

Dacă toate valorile proprii sunt subunitare în modul, atunci norma vectorului rezultat va tinde către zero. E suficient ca o singură valoare proprie să fie în modul mai mare ca 1, pentru ca norma vectorului rezultat să tindă către infinit.

De aceea, se acordă o deosebită importanță *razei spectrale a unei matrice*, definită ca modulul celei mai mari valori proprii

$$\rho(\mathbf{M}) = \max_i |\lambda_i|. \quad (2.127)$$

Valorile proprii sunt rădăcinile polinomului caracteristic. Pentru a justifica acest lucru, din (2.125) rezultă că valorile componentelor vectorilor proprii satisfac sistemul de ecuații algebrice liniare

$$(\lambda\mathbf{I} - \mathbf{M})\mathbf{v} = \mathbf{0} \quad (2.128)$$

Cum vectorii proprii sunt nenuli, rezultă că în mod necesar matricea coeficienților acestui sistem trebuie să fie singulară, altfel vectorii proprii ar fi nuli. Condiția de singularitate

este echivalentă cu anularea determinantului matricei, care este exact *polinomul caracteristic al matricei*:

$$\det(\lambda \mathbf{I} - \mathbf{M}) = 0. \quad (2.129)$$

Membrul stâng al relației (2.129) este un polinom de gradul n în λ , cu coeficienți reali. Conform teoremei fundamentale a algebrei, ecuația are exact n soluții (reale sau în perechi complex conjugate), care sunt valorile proprii ale matricei.

2.7.4 Convergență

Estimarea convergenței procesului iterativ poate fi făcută apriori, folosind următoarele teoreme.

Teorema 1: Condiția necesară și suficientă ca procesul iterativ să fie convergent este ca raza spectrală a matricei de iterație să fie strict subunitară

Intuitiv, această teoremă poate fi înțeleasă pe baza relației dintre eroarea soluției și matricea de iterație. Din (2.114), (2.115) și (2.120) rezultă că eroarea la o iterație k este:

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^* = F(\mathbf{x}^{(k-1)}) - F(\mathbf{x}^*) = \mathbf{M}\mathbf{x}^{(k-1)} + \mathbf{u} - \mathbf{M}\mathbf{x}^* - \mathbf{u} = \mathbf{M}\mathbf{e}^{(k-1)}, \quad (2.130)$$

și în consecință eroarea depinde de matricea de iterație ridicată la puterea k și de eroarea inițială $\mathbf{e}^{(0)}$:

$$\mathbf{e}^{(k)} = \mathbf{M}\mathbf{e}^{(k-1)} = \mathbf{M}^2\mathbf{e}^{(k-2)} = \dots = \mathbf{M}^k\mathbf{e}^{(0)}. \quad (2.131)$$

După cum a fost ilustrat în paragraful 2.7.3, dacă toate valorile proprii ale matricei \mathbf{M} au modulul strict subunitar, atunci norma erorii $\|\mathbf{e}^{(k)}\|$ va tinde către zero atunci când k tinde la infinit, procedeul iterativ fiind convergent.

Teorema 2: O condiție suficientă ca procesul iterativ să fie convergent este ca norma matricei de iterație să fie strict subunitară.

Această teoremă se poate demonstra cu ușurință pe baza primei teoreme și a relației

$$\rho(\mathbf{M}) \leq \|\mathbf{M}\|. \quad (2.132)$$

O altă demonstrație poate fi făcută pornind de la (2.131) și aplicând proprietățile normei:

$$\|\mathbf{e}^{(k)}\| \leq \|\mathbf{M}\|^k \|\mathbf{e}^{(0)}\|. \quad (2.133)$$

Când k tinde la infinit și norma matricei de iterație este subunitară, majorantul normei erorii tinde la zero. În consecință, norma erorii tinde la zero, procesul fiind convergent.

Mai mult, și diferența dintre două aproximații consecutive scade dacă norma matricei de iterație este strict subunitară:

$$\begin{aligned} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| &= \|\mathbf{M}\mathbf{x}^{(k-1)} + \mathbf{u} - \mathbf{M}\mathbf{x}^{(k-2)} - \mathbf{u}\| = \|\mathbf{M}(\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)})\| \leq \\ &\leq \|\mathbf{M}\| \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}\|, \end{aligned} \quad (2.134)$$

astfel încât utilizarea unui criteriu de oprire Cauchy este pe deplin justificată.

Fie o margine a erorii absolute notată cu a_k , unde $\|\mathbf{e}^{(k)}\| \leq a_k$. Din relația (2.133) rezultă că

$$a_k = \|\mathbf{M}\|^k a_0, \quad (2.135)$$

de unde

$$\log(a_k) = k \log \|\mathbf{M}\| + \log a_0. \quad (2.136)$$

Mărimea $R(\mathbf{M}) = -\log \|\mathbf{M}\|$ se numește rata de convergență. Rezultă că

$$\log(a_k) = -kR(\mathbf{M}) + \log a_0. \quad (2.137)$$

Într-un grafic ce ar ilustra dependența erorii în funcție de numărul de iterații, în care eroarea este reprezentată în scară logaritmică, rata de convergență este panta dreptei ce mărginește superior graficul.

Din (2.137) scrisă pentru două iterații consecutive, rezultă că

$$R(\mathbf{M}) = \log(a_{k-1}) - \log(a_k), \quad (2.138)$$

ceea ce înseamnă că rata de convergență reprezintă numărul de cifre semnificative corecte ce se câștigă la fiecare iterație. Inversa ratei de convergență reprezintă numărul de iterații necesar pentru câștigarea unei cifre semnificative corecte.

De exemplu, pentru o matrice având $\|\mathbf{M}\| = 10^{-3}$, rata de convergență este 3, deci la fiecare iterație numărul de cifre semnificative corecte crește cu 3. Aceasta este o convergență extrem de rapidă. Pentru o matrice cu $\|\mathbf{M}\| = 10^{-1}$, la fiecare iterație se câștigă o cifră semnificativă.

Este de remarcat că alegerea valorii inițiale nu are nici o influență asupra convergenței sau neconvergenței procesului iterativ. În cazul unui proces iterativ convergent, valoarea inițială afectează doar numărul de iterații necesar pentru atingerea unei erori impuse.

2.7.5 Algoritm general

Algoritmul general al unei metode iterative este dat de următorul pseudocod simplificat

procedură metodă_iterativă($n, B, C, b, x_0, er, maxit, x$)

...

$\mathbf{xv} = \mathbf{x0}$; inițializează șirul iterațiilor

$k = 0$; inițializare contor iterații

repetă

$\mathbf{t} = \mathbf{C} * \mathbf{xv} + \mathbf{b}$

```

metodă_directă (n, B, t, x)
d = ||xv - x||
xv = x ; actualizează soluția veche
k = k + 1
cât timp d > er și k ≤ maxit
retur

```

Acest pseudocod nu este scris conform regulilor din capitolul 1. El sugerează doar pașii importanți ai unei metode iterative. Instrucțiunile scrise cu litere aldine, de exemplu $\mathbf{xv} = \mathbf{x0}$ ascund operații de atribuire pe componente. De asemenea, instrucțiunea $\mathbf{t} = \mathbf{C} * \mathbf{xv} + \mathbf{b}$ reprezintă mai întâi înmulțirea dintre o matrice și un vector, urmată de adunarea cu un vector.

Se observă de asemenea că nu sunt memorate toate valorile aproximative din șirul de iterații. Deoarece o iterație se calculează exclusiv în funcție de iterația anterioară, este suficientă memorarea doar a unei iterații vechi (notate \mathbf{xv} în algoritm).

Estimarea gradului de complexitate a unui astfel de algoritm poate fi făcută doar pentru o singură iterație. Efortul de calcul depinde însă de structura matricelor în care a fost descompusă matricea coeficienților, fiind consumat mai ales în secvența în care se calculează noul termen liber \mathbf{t} și care implică produsul dintre o matrice și un vector (având o complexitate pătratică în cazul cel mai defavorabil, în care matricea \mathbf{C} este plină) și procedura de rezolvare directă, care în general are o complexitate liniară deoarece \mathbf{B} are o structură rară, particulară. Intuitiv, ne așteptăm ca procedeul iterativ să fie cu atât mai rapid convergent cu cât \mathbf{B} conține mai multă informație din \mathbf{A} . Dacă \mathbf{B} este mai rară (în cazul extrem \mathbf{B} reține doar elementele diagonale din \mathbf{A}), atunci rezolvarea cu metoda directă este foarte rapidă, deci efortul per iterație este mic, dar convergența este lentă.

Următoarele două paragrafe detaliază această procedură în funcție de alegerea precisă a descompunerii matricei coeficienților.

2.8 Metoda Jacobi

În metoda Jacobi descompunerea matricei \mathbf{A} se face astfel încât matricea \mathbf{B} este diagonală.

2.8.1 Un exemplu simplu

Să considerăm același exemplu de la paragraful 2.4.1. Ideea metodei Jacobi este de a păstra în membru stâng termenii diagonali.

$$\begin{cases} x + 2y - z = -1 \\ -2x + 3y + z = 0 \\ 4x - y - 3z = -2 \end{cases} \Leftrightarrow \begin{cases} x = -2y + z - 1 \\ 3y = 2x - z \\ -3z = -4x + y - 2 \end{cases} \quad (2.139)$$

Aceasta sugerează următorul calcul recursiv al șirului de iterații

$$\begin{aligned} x^{(n)} &= -2y^{(v)} + z^{(v)} - 1 \\ y^{(n)} &= 2x^{(v)} - z^{(v)} \\ z^{(n)} &= -4x^{(v)} + y^{(v)} - 2 \end{aligned} \quad (2.140)$$

Dacă inițializarea este nulă, atunci șirul iterațiilor Jacobi va fi $[0, 0, 0]^T$, $[-1, 0, -2]^T$, $[-3, 0, 2]^T$, etc.

2.8.2 Algoritmul metodei

Dacă notăm $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ (fig. 2.8), atunci descompunerea $\mathbf{A} = \mathbf{B} - \mathbf{C}$ în metoda Jacobi este

$$\mathbf{B} = \mathbf{D}, \quad \mathbf{C} = -(\mathbf{L} + \mathbf{U}), \quad (2.141)$$

iar relația (2.122) ce reprezintă calculul recursiv al noii iterații în funcție de cea veche prin rezolvarea unui sistem algebric liniar devine

$$\mathbf{D}\mathbf{x}^{(k)} = -(\mathbf{L} + \mathbf{U})\mathbf{x}^{(k-1)} + \mathbf{b}. \quad (2.142)$$

Ecuatia i a sistemului (2.142) este

$$a_{ii}\mathbf{x}_i^{(k)} = - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}x_j^{(k-1)} + b_i. \quad (2.143)$$

Se observă că fiecare ecuație este tratată izolat de celelalte.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \underset{\mathbf{A}}{=} \begin{bmatrix} 0 & 0 & 0 & 0 \\ a_{21} & 0 & 0 & 0 \\ a_{31} & a_{32} & 0 & 0 \\ a_{41} & a_{42} & a_{43} & 0 \end{bmatrix} \underset{\mathbf{L}}{=} + \begin{bmatrix} a_{11} & 0 & 0 & 0 \\ 0 & a_{22} & 0 & 0 \\ 0 & 0 & a_{33} & 0 \\ 0 & 0 & 0 & a_{44} \end{bmatrix} \underset{\mathbf{D}}{=} + \begin{bmatrix} 0 & a_{12} & a_{13} & a_{14} \\ 0 & 0 & a_{23} & a_{24} \\ 0 & 0 & 0 & a_{34} \\ 0 & 0 & 0 & 0 \end{bmatrix} \underset{\mathbf{U}}{=}$$

Figura 2.8: Partiționarea matricei în metodele iterative.

Rezolvarea sistemului (2.142) este imediată, soluția la o nouă iterație $\mathbf{x}^{(n)}$ calculându-se numai în funcție de soluția la iterația anterioară $\mathbf{x}^{(v)}$, nefiind necesară memorarea tuturor iterațiilor:

$$\mathbf{x}_i^{(n)} = (b_i - \sum_{\substack{j=1 \\ j \neq i}}^n x_j^{(v)}) / a_{ii} \quad i = 1, \dots, n. \quad (2.144)$$

Fiecare componentă nouă poate fi calculată independent de celelalte componente noi, motiv pentru care metoda Jacobi se mai numește și *metoda deplasărilor simultane*.

Pseudocodul procedurii Jacobi, în care pentru calculul normei se folosește norma euclidiană, este următorul.

procedură Jacobi($n, a, b, x0, er, maxit, x$)

; rezolvă sistemul algebric liniar $ax = b$, de dimensiune n prin metoda Jacobi

întreg n ; dimensiunea sistemului

tablou real $a[n][n]$; matricea coeficienților, indici de la 1

tablou real $b[n]$; vectorul termenilor liberi

; mărimi specifice procedurilor iterative

tablou real $x0[n]$; inițializarea soluției

real er ; eroarea folosită de criteriul de oprire

întreg $maxit$; număr maxim de iterații admis

tablou real $xv[n]$; aproximația anterioară ("veche")

pentru $i = 1, n$

$xv_i = x0_i$

•

$k = 0$; inițializare contor iterații

repetă

$d = 0$

pentru $i = 1, n$

$s = 0$

pentru $j = 1, n$

dacă $j \neq i$

$s = s + a_{ij} * xv_j$

•

•

$x_i = (b_i - s) / a_{ii}$

$d = d + (x_i - xv_i)^2$

•

$$d = \sqrt{d}$$

pentru $i = 1, n$

$xv_i = x_i$; actualizează soluția veche

•

$$k = k + 1$$

cât timp $d > er$ și $k \leq maxit$

retur

2.8.3 Aspecte legate de convergență

Convergența procedurii Jacobi depinde de matricea de iterație (2.119) care este în acest caz

$$\mathbf{M} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}). \quad (2.145)$$

Pentru o problemă bine formulată matematic și bine condiționată, algoritmul Jacobi poate fi sau nu convergent. Acest lucru poate fi ilustrat grafic pentru un sistem de două ecuații cu două necunoscute. Din punct de vedere matematic nu contează ordinea în care sunt scrise ecuațiile. Această ordine poate avea însă efect asupra convergenței sau neconvergenței procedurii iterativ. Dacă notăm cu Δ_1 dreapta corespunzătoare primei ecuații și cu Δ_2 dreapta corespunzătoare celei de a doua ecuații, atunci deplasările simultane din metoda Jacobi pot decurge ca în fig. 2.9, fiind convergente în cazul din stânga și neconvergente în cazul din dreapta.

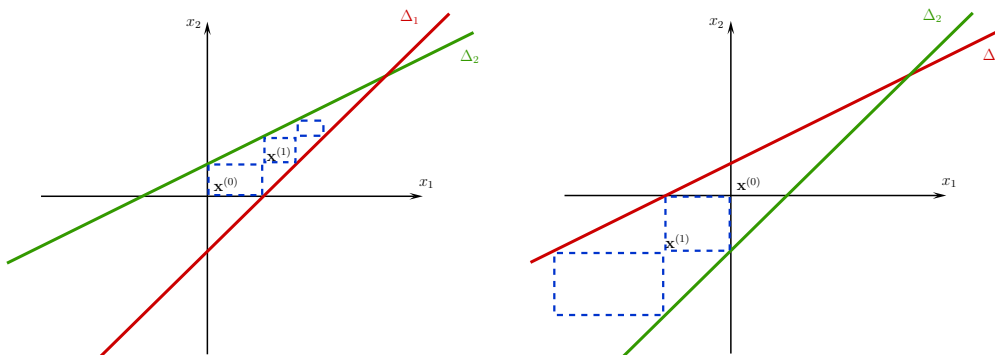


Figura 2.9: Schimbarea ordinii ecuațiilor din sistem înseamnă schimbarea matricei de iterație, deci a proprietăților de convergență ale metodei Jacobi.

În general, condițiile de convergență sunt date de teoremele prezentate în paragraful 2.7.4, unde matricea de iterație este (2.145). Un alt rezultat util este următorul. **Dacă matricea coeficienților este diagonal dominantă, atunci condiția suficientă de**

convergență este satisfăcută și algoritmul Jacobi este convergent. Pentru justificarea acestei afirmații este suficient să se explicitizeze componentele matricei (2.145). Rezultă imediat că $\|\mathbf{M}\| = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}/a_{ii}|$, valoare mai mică decât 1 dacă matricea este diagonal dominantă.

O analiză riguroasă a convergenței metodelor Jacobi și Gauss-Seidel poate fi făcută pentru sisteme care provin din rezolvarea numerică a ecuațiilor cu derivate parțiale (de exemplu ecuația Laplace) [8].

2.9 Metoda Gauss-Seidel

În metoda Gauss-Seidel descompunerea matricei \mathbf{A} se face astfel încât matricea \mathbf{B} este triunghiular inferioară.

2.9.1 Un exemplu simplu

Să considerăm același exemplu de la paragraful 2.4.1. Ideea metodei Gauss-Seidel este de a păstra în membru stâng termenii diagonali și sub-diagonali.

$$\begin{cases} x + 2y - z = -1 \\ -2x + 3y + z = 0 \\ 4x - y - 3z = -2 \end{cases} \Leftrightarrow \begin{cases} x & & & = -2y + z - 1 \\ -2x + 3y & & & = -z \\ 4x - y - 3z & & & = -2 \end{cases} \quad (2.146)$$

ceea ce sugerează următorul calcul recursiv al șirului de iterații

$$\begin{aligned} x^{(n)} &= -2y^{(v)} + z^{(v)} - 1 \\ -2x^{(n)} + 3y^{(n)} &= -z^{(v)} \\ 4x^{(n)} - y^{(n)} - 3z^{(n)} &= -2 \end{aligned} \quad (2.147)$$

Dacă inițializarea este nulă, atunci șirul iterațiilor Gauss-Seidel va fi $[0, 0, 0]^T$, $[-1, -2, 4]^T$, $[7, 10, -40]^T$, etc.

2.9.2 Algoritmul metodei

Dacă notăm $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ (fig.2.8), atunci descompunerea $\mathbf{A} = \mathbf{B} - \mathbf{C}$ în metoda Gauss-Seidel este

$$\mathbf{B} = \mathbf{L} + \mathbf{D}, \quad \mathbf{C} = -\mathbf{U}, \quad (2.148)$$

iar relația (2.122) ce reprezintă calculul recursiv al noii iterații în funcție de cea veche prin rezolvarea unui sistem algebric liniar devine

$$(\mathbf{L} + \mathbf{D})\mathbf{x}^{(k)} = -\mathbf{U}\mathbf{x}^{(k-1)} + \mathbf{b}. \quad (2.149)$$

Ecuția i a sistemului (2.149) este

$$\sum_{j=1}^{i-1} a_{ij}x_j^{(k)} + a_{ii}x_i^{(k)} = - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} + b_i. \quad (2.150)$$

Deoarece calculul la o iterație nouă depinde exclusiv de iterația anterioară, nu este necesară memorarea tuturor iterațiilor, iar relația (2.150) se poate scrie ca

$$\sum_{j=1}^{i-1} a_{ij}x_j^{(n)} + a_{ii}x_i^{(n)} = - \sum_{j=i+1}^n a_{ij}x_j^{(v)} + b_i, \quad (2.151)$$

în care (k) s-a înlocuit cu (n) -nou și $(k-1)$ cu (v) -vechi. Rezolvarea sistemului se va face prin substituție progresivă, conform formulei:

$$x_i^{(n)} = (b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(n)} - \sum_{j=i+1}^n a_{ij}x_j^{(v)})/a_{ii}. \quad (2.152)$$

Se observă că este respectat *principiul lui Seidel*, conform căruia o valoare nouă a unei necunoscute trebuie folosită imediat în calcule. O componentă nouă nu poate fi calculată independent de celelalte componente noi, motiv pentru care metoda Gauss-Seidel se mai numește și *metoda deplasărilor succesive*. Altfel spus, ecuațiile sunt examinate secvențial, iar rezultatele obținute sunt imediat folosite în calcule.

Pseudocodul procedurii Gauss-Seidel, în care pentru calculul normei se folosește norma Chebyshev, este următorul.

procedură Gauss-Seidel($n, a, b, x_0, er, maxit, x$)

; rezolvă sistemul algebric linear $ax = b$, de dimensiune n prin metoda Gauss-Seidel

; declarații ca la procedura Jacobi

...

pentru $i = 1, n$

$xv_i = x_0$

•

$k = 0$; inițializare contor iterații

repetă

$d = 0$

pentru $i = 1, n$

$s = b_i$

pentru $j = 1, n$

dacă $j \neq i$

$s = s + a_{ij} * xv_j$

```

      •
    •
    s = s/aii
    p = |xi - s|
    dacă p > d
      d = p
    •
    xi = s
  •
  k = k + 1
  cât timp d > er și k ≤ maxit
  retur

```

Spre deosebire de algoritmul Jacobi, în algoritmul Gauss-Seidel nu este necesară memorarea soluției anterioare. O dată calculată noua valoare, vechea valoare nu mai este folosită decât la calculul erorii.

2.9.3 Aspecte legate de convergență

Convergența procedurii Gauss-Seidel depinde de matricea de iterație (2.119) care este în acest caz

$$\mathbf{M} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}. \quad (2.153)$$

Pentru o problemă bine formulată matematic și bine condiționată, algoritmul Gauss-Seidel poate fi sau nu convergent. Acest lucru poate fi ilustrat grafic pentru un sistem de două ecuații cu două necunoscute. Din punct de vedere matematic nu contează ordinea în care sunt scrise ecuațiile. Această ordine poate avea însă efect asupra convergenței sau neconvergenței procedurii iterativ. Dacă notăm cu Δ_1 dreapta corespunzătoare primei ecuații și cu Δ_2 dreapta corespunzătoare celei de a doua ecuații, atunci deplasările succesive din metoda Gauss-Seidel pot decurge ca în fig. 2.10, fiind convergente în cazul din stânga și neconvergente în cazul din dreapta.

Condițiile generale de convergență sunt date de teoremele prezentate în paragraful 2.7.4, unde matricea de iterație este (2.153). Ca și în cazul metodei Jacobi, metoda Gauss-Seidel este convergentă dacă matricea coeficienților \mathbf{A} este diagonal dominantă. O altă condiție suficientă de convergență este ca matricea \mathbf{A} să fie simetrică și pozitiv definită.

În general, dacă metoda Jacobi este convergentă, metoda Gauss-Seidel este mai rapid convergentă. În cazul matricilor simetrice și pozitiv definite, Gauss-Seidel este de aproximativ două ori mai rapid convergentă decât Jacobi [4].

2.9.4 Evaluarea algoritmului

Efortul de calcul atât al algoritmului Gauss-Seidel cât și al algoritmului Jacobi depinde de numărul de iterații m , număr care depinde de rata de convergență a algoritmului, deci de matricea de iterație.

Analizând pseudocodurile celor doi algoritmi, constatăm că efortul de calcul per iterație este $O(2n^2)$. **Efortul total de calcul al algoritmilor Jacobi și Gauss-Seidel este $T = O(2mn^2)$.**

Comparând acest efort de calcul cu cel al metodei Gauss, rezultă că metoda Jacobi sau Gauss-Seidel este mai eficientă decât metoda Gauss dacă $2mn^2 < 2n^3/3$, deci dacă $m < n/3$. Practic, dacă algoritmi nu converg în $n/3$ iterații atunci este mai bine ca ei să fie opriți și rezolvarea să se facă cu o metodă directă.

Din punct de vedere al necesarului de memorie, algoritmul Gauss-Seidel necesită doar structurile de date necesare pentru memorarea sistemului (matricea coeficienților și termenii liberi) și a vectorului soluție (ultima iterație), în total $M_{GS} = O(n^2 + 2n) \approx O(n^2)$, fiind mai mic decât necesarul de memorie pentru algoritmul Jacobi $M_J = O(n^2 + 3n) \approx O(n^2)$, diferența nefiind semnificativă dacă matricile sunt pline.

Dacă însă matricea coeficienților este rară, atunci efortul de calcul pe iterație se poate diminua, calculul $s = s + a_{ij} * x_j$ făcându-se doar dacă a_{ij} este diferit de zero. Într-un astfel de caz, memorarea matricei coeficienților nu se va face însă într-un tablou bidimensional ci în structuri de date de tipul celor descrise în paragraful 2.6.1, deci nici necesarul de memorie nu va mai fi pătratic. Ceea ce este important de remarcat este faptul că nu putem vorbi de umpleri ale matricei coeficienților așa cum se întâmplă în cazul algoritmului Gauss aplicat pentru matrice rare.

Metodele iterative sunt mai potrivite decât metodele directe pentru re-

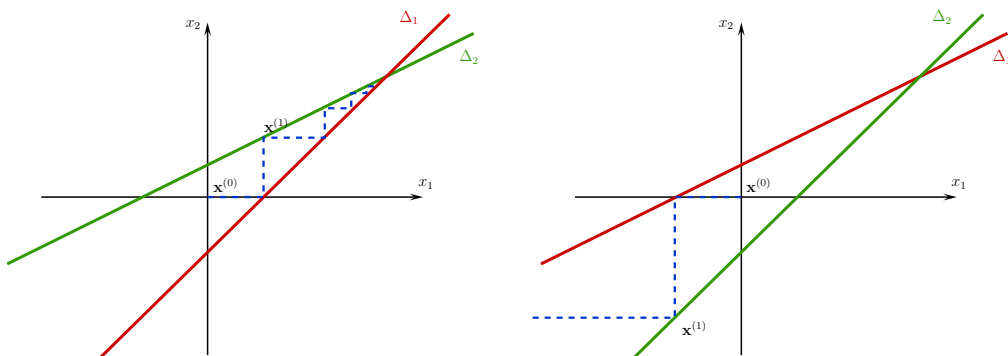


Figura 2.10: Schimbarea ordinii ecuațiilor din sistem înseamnă schimbarea matricei de iterație, deci a proprietăților de convergență ale metodei Gauss-Seidel.

zolvarea sistemelor cu matrice rare, într-un context hardware în care memoria disponibilă nu este suficientă rezolvării prin metode directe.

2.10 Metoda suprarelaxării succesive (SOR)

Interpretarea geometrică a iterațiilor Gauss-Seidel (fig.2.10) sugerează următorul procedeu de accelerare a convergenței. Pasul prescris de metoda Gauss-Seidel se amplifică cu un factor supraunitar ω care se numește *factor de suprarelaxare*, soluția nouă calculându-se ca

$$x_i^{(n)} = x_i^{(v)} + \omega(x_i^{(n-GS)} - x_i^{(v)}). \quad (2.154)$$

Cazul $\omega = 1$ corespunde metodei Gauss-Seidel. Pseudocodul algoritmului suprarelaxării succesive se obține înlocuind în pseudocodul algoritmului Gauss-Seidel instrucțiunea $x_i = s$ cu $x_i = x_i + \omega(s - x_i)$.

Relația (2.154) în care înlocuim formula de calcul a soluției Gauss-Seidel, dată de (2.151) devine

$$\begin{aligned} x_i^{(n)} &= \omega x_i^{(n-GS)} + (1 - \omega)x_i^{(v)} = \\ &= \omega(b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(n)} - \sum_{j=i+1}^n a_{ij}x_j^{(v)})/a_{ii} + (1 - \omega)x_i^{(v)}. \end{aligned} \quad (2.155)$$

Aceasta corespunde rezolvării rezolvării sistemului

$$(\mathbf{D} + \omega\mathbf{L})\mathbf{x}^{(n)} = [(1 - \omega)\mathbf{D} - \omega\mathbf{U}]\mathbf{x}^{(v)} + \omega b. \quad (2.156)$$

Matricea de iterație a metodei SOR este în consecință

$$\mathbf{M} = (\mathbf{D} + \omega\mathbf{L})^{-1} [(1 - \omega)\mathbf{D} - \omega\mathbf{U}], \quad (2.157)$$

iar proprietățile ei de convergență depind de valoarea factorului de relaxare ω . S-a demonstrat (teorema lui Kahan) că metoda nu converge dacă factorul de relaxare ω este ales în afara intervalului $(0, 2)$ [1]. Cazul $\omega \in (0, 1)$ corespunde unei subrelaxări și ea se folosește atunci când metoda Gauss-Seidel nu converge. Dacă matricea este simetrică și pozitiv definită, SOR este garantat convergentă pentru orice valoare $\omega \in (0, 2)$. Alegerea valorii lui ω poate afecta în mod semnificativ rata de convergență. În general, este dificil să se calculeze apriori valoarea optimală a factorului de suprarelaxare. Chiar atunci când un astfel de calcul este posibil, efortul de calcul este atât de mare încât este rareori folosit.

Pentru o clasă specială de matrici, numite ”consistent ordonate”, obținute prin discretizarea unor ecuații cu derivate parțiale prin parcurgerea sistematică a nodurilor rețelei

de discretizare se poate demonstra că valoarea optimă pentru ω este

$$\omega_{opt} = \frac{2}{1 + \sqrt{1 - \rho^2}}, \quad (2.158)$$

unde ρ este raza spectrală a matricei de iterație Jacobi. O astfel de estimare este rareori folosită deoarece estimarea lui ρ este costisitoare. În practică se folosesc tehnici euristice, ce iau în considerare dimensiunea rețelei de discretizare a problemei [1, 8].

În cazul matricelor simetrice, o variantă a metodei SOR, numită SSOR, este folosită ca procedură de preconditionare pentru metode nestaționare.

2.11 Algoritm general al metodelor staționare

Metodele iterative discutate până acum pentru rezolvarea sistemului $\mathbf{Ax} = \mathbf{b}$ au la bază descompunerea matricei coeficienților în $\mathbf{A} = \mathbf{B} - \mathbf{C}$, iterațiile construindu-se pe baza relației

$$\mathbf{B}\mathbf{x}^{(k+1)} = \mathbf{C}\mathbf{x}^{(k)} + \mathbf{b}. \quad (2.159)$$

Această relație a condus la algoritmul general prezentat într-un pseudocod simplificat în paragraful 2.7.5. Dacă se scade $\mathbf{B}\mathbf{x}^{(k)}$ din ambii membri ai relației (2.159), se obține

$$\mathbf{B}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}. \quad (2.160)$$

Dacă se definește reziduul la iterația k

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}, \quad (2.161)$$

atunci membrul drept al relației (2.160) reprezintă reziduul la iterația k . În consecință, calculul soluției la o nouă iterație poate fi pus sub forma

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{B}^{-1}\mathbf{r}^{(k)}, \quad (2.162)$$

adică o nouă iterație se calculează din iterația anterioară, adunând o corecție calculată în funcție de reziduul iterației anterioare. Acest reziduu este întotdeauna înmulțit cu matricea \mathbf{B}^{-1} , aceeași pe parcursul tuturor iterațiilor. La acest aspect se referă ”staționari-tatea” metodelor Jacobi unde $\mathbf{B} = \mathbf{D}$, Gauss-Seidel unde $\mathbf{B} = \mathbf{D} + \mathbf{L}$, SOR unde $\mathbf{B} = \omega^{-1}(\mathbf{D} + \omega\mathbf{L})$. Ca de obicei, nu se face inversarea propriu zisă, ci se rezolvă un sistem algebric liniar, cu matricea coeficienților \mathbf{B} .

Pseudocodul simplificat care implementează acest mod de efectuare a calculelor este următorul.

procedură metodă_iterativă_v2($n, B, A, b, x_0, er, maxit, x$)

...

$\mathbf{xv} = \mathbf{x0}$; inițializează șirul iterațiilor

$k = 0$; inițializare contor iterații

repetă

$\mathbf{r} = \mathbf{b} - \mathbf{A} \cdot \mathbf{xv}$; calculează reziduu

metodă_directă (n, B, r, z)

$d = \|\mathbf{z}\|$

$\mathbf{x} = \mathbf{xv} + \mathbf{z}$

$\mathbf{xv} = \mathbf{x}$; actualizează soluția veche

$k = k + 1$

cât timp $d > er$ și $k \leq maxit$

retur

Algoritmii detaliați prezentați pentru metodele Jacobi și Gauss-Seidel au presupus că matricea coeficienților este plină. Metodele iterative sunt eficiente însă pentru matrici rare. În cazul matricilor pline, timpul de rezolvare cu metode iterative poate fi comparabil cu timpul de factorizare. Într-o astfel de situație, factorizarea este mai utilă deoarece, o dată ce factorii L și U sunt calculați, rezolvarea sistemului poate fi făcută oricând pentru alt termen liber. De aceea un pseudocod simplificat, în care sunt scrise operații cu matrice, este mai general, putând fi adaptat unor matrice rare.

2.12 Metoda gradientilor conjugați

Metoda gradientilor conjugați este o metodă iterativă nestaționară pentru rezolvarea unor sisteme de ecuații algebrice de forma

$$\mathbf{Ax} = \mathbf{b}, \quad (2.163)$$

unde matricea \mathbf{A} este *simetrică și pozitiv definită*³. Ca orice metodă iterativă, algoritmul este eficient pentru matrice rare, și de aceea pseudocodurile descrise vor fi prezentate simplificat, evidențiind operații de algebră liniară, presupuse a fi implementate ținând cont de raritatea structurilor de date.

³O matrice reală $\mathbf{A} \in \mathbb{R}^{n \times n}$ este pozitiv definită dacă ea este simetrică și dacă $\mathbf{x}^T \mathbf{Ax} > 0$ pentru orice vector real, nenul $\mathbf{x} \in \mathbb{R}^{n \times 1}$.

2.12.1 Forma pătratică asociată

Metoda gradientilor conjugați este simplu de înțeles dacă ea este formulată ca o problemă de minimizare [12]. Pentru acesta este util să se considere o formă pătratică $f : \mathbb{R}^n \rightarrow \mathbb{R}$ definită de

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} + c, \quad (2.164)$$

unde $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{x}, \mathbf{b} \in \mathbb{R}^{n \times 1}$, $c \in \mathbb{R}$.

Se poate demonstra următoarea afirmație **Dacă \mathbf{A} este simetrică și pozitiv definită, atunci funcția (2.164) este minimizată de soluția ecuației (2.163).**

Pentru a justifica această proprietate, este util să se considere gradientul funcției f

$$\nabla f(\mathbf{x}) = \frac{1}{2} \mathbf{A}^T \mathbf{x} + \frac{1}{2} \mathbf{A} \mathbf{x} - \mathbf{b}. \quad (2.165)$$

În punctul de minim gradientul este zero

$$\nabla f(\mathbf{x}) = 0 \quad \Rightarrow \quad \frac{1}{2} (\mathbf{A}^T + \mathbf{A}) \mathbf{x} = \mathbf{b}. \quad (2.166)$$

Dacă matricea este simetrică, atunci $\mathbf{A}^T = \mathbf{A}$ și condiția (2.166) este echivalentă cu (2.163). În consecință, soluția ecuației (2.163) este punct critic pentru f . Dacă, în plus, matricea este pozitiv definită, atunci punctul critic este un punct minim.

Cea mai simplă înțelegere intuitivă a unei matrice pozitiv definite este obținută prin considerarea variației funcționalei în cazul particular al unei matrice de dimensiune 2. Fig.2.11 reprezintă forma pătratică asociată unei matrice pozitiv definite. Minimul ei coincide cu soluția sistemului de rezolvat și metoda gradientilor conjugați va funcționa. Fig.2.12 reprezintă forma pătratică asociată unei matrice negativ definite. Maximul ei coincide cu soluția sistemului de rezolvat și metoda gradientilor conjugați ar putea fi ușor modificată pentru a maximiza f . Curbele de nivel ale funcționalei (forme) pătratice asociate unei matrice simetrice pozitiv sau negativ definite de ordin 2 sunt elipse; similar, în general, hipersuprafețele de nivel în cazul unor sisteme de ordin superior sunt hiperelipsoizi. Fig.2.13 reprezintă forma pătratică asociată unei matrice singulare, pozitiv semidefinite. Minimul ei este atins într-o infinitate de valori, aflate pe o dreaptă. Sistemul de ecuații are o infinitate de soluții, problema nefiind bine formulată matematic. Fig.2.14 reprezintă forma pătratică asociată unei matrice indefinite. Soluția sistemului de ecuații este punct șarpe pentru această funcțională, fiind punct de minim pentru o direcție și punct de maxim pentru altă direcție. Pentru o astfel de problema, metoda gradientilor conjugați nu poate fi aplicată.

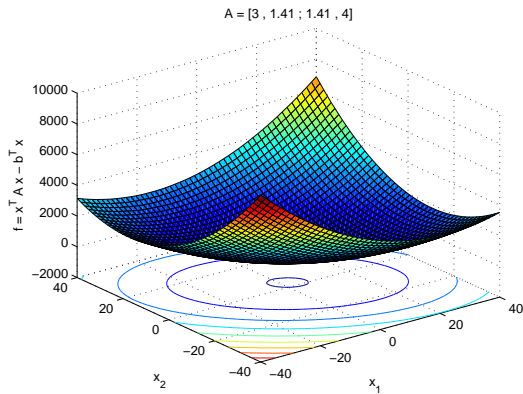


Figura 2.11: Funcția pătratică asociată unei matrice pozitiv definite.

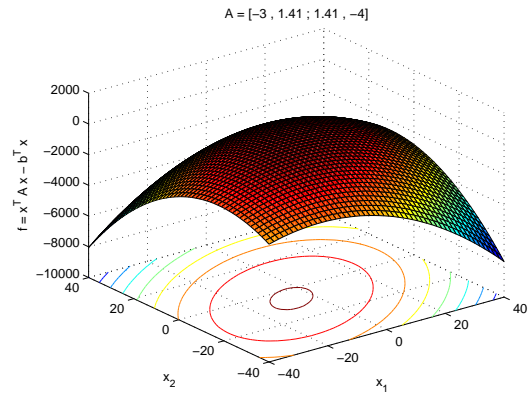


Figura 2.12: Funcția pătratică asociată unei matrice negativ definite.

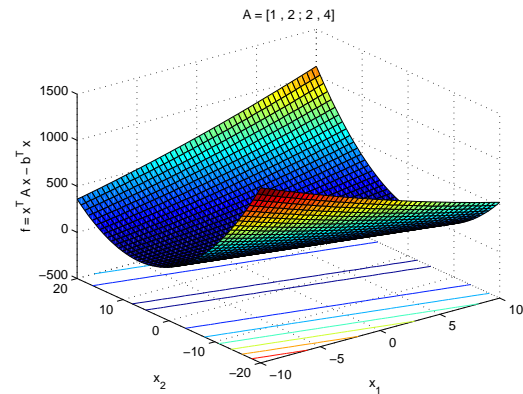


Figura 2.13: Funcția pătratică asociată unei matrice singulare, pozitiv semidefinite.

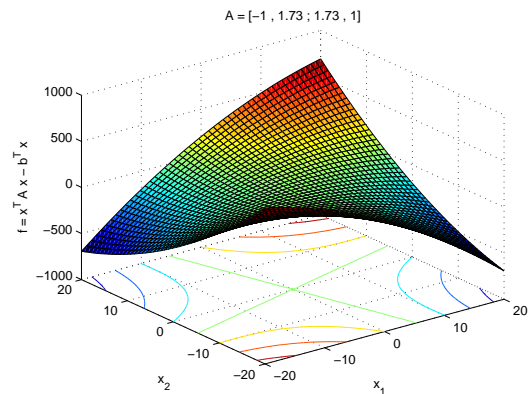


Figura 2.14: Funcția pătratică asociată unei matrice indefinite.

2.12.2 Metoda celei mai rapide coborâri (a gradientului)

Prima idee de a găsi minimul funcției f este de a merge în căutarea lui pe direcții care corespund ratei celei mai mari de schimbare a funcției. Dacă la un moment dat aproximația soluției este $\mathbf{x}^{(k)}$, atunci direcția celei mai rapide coborâri este direcția opusă gradientului în acest punct

$$-\nabla f(\mathbf{x}^{(k)}) = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}. \quad (2.167)$$

Vom defini la fiecare iterație eroarea $\mathbf{e}^{(k)}$ și reziduul $\mathbf{r}^{(k)}$ ca

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}, \quad (2.168)$$

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}. \quad (2.169)$$

Din (2.168) și (2.163) rezultă că reziduul este eroarea transformată de \mathbf{A} în același

spațiu ca al lui \mathbf{b} :

$$\mathbf{r}^{(k)} = -\mathbf{A}\mathbf{e}^{(k)}. \quad (2.170)$$

De asemenea, din (2.169) și (2.167) rezultă că reziduul este direcția celei mai rapide coborâri:

$$\mathbf{r}^{(k)} = -\nabla f(\mathbf{x}^{(k)}). \quad (2.171)$$

Corecția primei iterații se va face după direcția reziduului corespunzător inițializării:

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha \mathbf{r}^{(0)}. \quad (2.172)$$

Valoarea scalară α se va alege astfel încât pe direcția respectivă funcția să fie minimă. Mai precis $g(\alpha) = f(\mathbf{x}^{(1)})$ să fie minimă. Aceasta implică anularea derivatei de ordinul 1 a funcției g :

$$\frac{dg}{d\alpha}(\mathbf{x}^{(1)}) = 0 \quad \Rightarrow \quad (\nabla f(\mathbf{x}^{(1)}))^T \cdot \frac{d\mathbf{x}^{(1)}}{d\alpha} = 0. \quad (2.173)$$

Din (2.171) rezultă că $\nabla f(\mathbf{x}^{(1)}) = -\mathbf{r}^{(1)}$, iar din (2.172) rezultă că $d\mathbf{x}^{(1)}/d\alpha = \mathbf{r}^{(0)}$. În consecință, relația (2.173) este echivalentă cu faptul că reziduul de la prima iterație este ortogonal reziduului inițial. Urmează că:

$$\begin{aligned} (\mathbf{r}^{(1)})^T \cdot \mathbf{r}^{(0)} &= 0, \\ (\mathbf{b} - \mathbf{A}\mathbf{x}^{(1)})^T \cdot \mathbf{r}^{(0)} &= 0, \\ [\mathbf{b} - \mathbf{A}(\mathbf{x}^{(0)} + \alpha \mathbf{r}^{(0)})]^T \mathbf{r}^{(0)} &= 0, \\ (\mathbf{b} - \mathbf{A}\mathbf{x}^{(0)} - \alpha \mathbf{A}\mathbf{r}^{(0)})^T \cdot \mathbf{r}^{(0)} &= 0, \\ (\mathbf{b} - \mathbf{A}\mathbf{x}^{(0)})^T \cdot \mathbf{r}^{(0)} - \alpha (\mathbf{A}\mathbf{r}^{(0)})^T \mathbf{r}^{(0)} &= 0, \\ (\mathbf{r}^{(0)})^T \mathbf{r}^{(0)} &= \alpha (\mathbf{r}^{(0)})^T \mathbf{A}^T \mathbf{r}^{(0)}. \end{aligned}$$

În final, folosind faptul că \mathbf{A} este simetrică, rezultă

$$\alpha = \frac{(\mathbf{r}^{(0)})^T \mathbf{r}^{(0)}}{(\mathbf{r}^{(0)})^T \mathbf{A} \mathbf{r}^{(0)}}. \quad (2.174)$$

Acest raționament este valabil la fiecare iterație, corecția făcându-se după direcția ultimului reziduu:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{r}^{(k)}, \quad (2.175)$$

mărimea scalară α_k fiind dedusă din condiția de ortogonalitate dintre reziduul la iterația k și cel la iterația $k + 1$:

$$\alpha_k = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k)})^T \mathbf{A} \mathbf{r}^{(k)}}. \quad (2.176)$$

Acest mod de calcul este descris de următorul algoritm

procedură metoda_gradientului_v0($n, A, b, x0, er, maxit, x$)

...

$\mathbf{xv} = \mathbf{x0}$; inițializează șirul iterațiilor

$k = 0$; inițializare contor iterații

$\mathbf{r} = \mathbf{b} - \mathbf{A} \cdot \mathbf{xv}$; calculează reziduu

cât timp $\|\mathbf{r}\| > er$ și $k \leq maxit$

$$\alpha = \mathbf{r}^T \mathbf{r} / (\mathbf{r}^T \mathbf{A} \mathbf{r})$$

$$\mathbf{x} = \mathbf{xv} + \alpha \mathbf{r}$$

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \cdot \mathbf{xv}$$

$\mathbf{xv} = \mathbf{x}$; actualizează soluția veche

$$k = k + 1$$

•

retur

În cadrul unei iterații, efortul cel mai mare de calcul este dat de produsul dintre o matrice și un vector. În această variantă, sunt două astfel de produse la fiecare iterație.

Dacă înmulțim relația (2.175) la stânga cu $-\mathbf{A}$ și adunăm \mathbf{b} obținem

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{r}^{(k)}, \quad (2.177)$$

O variantă îmbunătățită a algoritmului este cea care folosește această relație și calculează un singur produs matrice vector în cadrul unei iterații:

procedură metoda_gradientului($n, A, b, x0, er, maxit, x$)

...

$\mathbf{xv} = \mathbf{x0}$; inițializează șirul iterațiilor

$k = 0$; inițializare contor iterații

$\mathbf{r} = \mathbf{b} - \mathbf{A} \cdot \mathbf{xv}$; calculează reziduu

$$e = \|\mathbf{r}\|$$

cât timp $e > er$ și $k \leq maxit$

$$\mathbf{m} = \mathbf{A} \cdot \mathbf{r}$$

$$\alpha = \mathbf{r}^T \mathbf{r} / (\mathbf{r}^T \mathbf{m})$$

$$e = \|\mathbf{r}\|$$

$$\mathbf{x} = \mathbf{xv} + \alpha \mathbf{r}$$

$\mathbf{xv} = \mathbf{x}$; actualizează soluția veche

$$k = k + 1$$

$$\mathbf{r} = \mathbf{r} - \alpha \mathbf{m}$$

•

retur

În această variantă de pseudocod, reziduul se calculează recursiv. Dezavantajul este acela că în valoarea reziduului se pot acumula erori de rotunjire. Pentru a corecta acest lucru, este bine ca, periodic, reziduul să se calculeze cu relația sa de definiție, dată de (2.169).

Analiza convergenței acestei metode se face pe baza *numărului de condiționare spectrală* care se definește ca raportul dintre cea mai mare și cea mai mică valoare proprie

$$\kappa = \frac{\max \lambda_i}{\min \lambda_i} \geq 1. \quad (2.178)$$

Metoda poate converge rapid, într-un singur pas, dacă punctul de start este ales pe axele elipsoidului sau dacă forma pătratică este sferică (ceea ce corespunde cazului în care toate valorile proprii sunt egale). În rest, convergența depinde de numărul de condiționare κ și de inițializare (fig.2.15). Se poate demonstra că eroarea la fiecare iterație i este mărginită de

$$\|\mathbf{e}^{(i)}\| \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^i \|\mathbf{e}^{(0)}\| \quad (2.179)$$

2.12.3 Metoda direcțiilor conjugate

Metoda gradientului converge atât de încet deoarece direcțiile de căutare sunt ortogonale și se întâmplă ca pe parcursul iterațiilor direcțiile de căutare să se repete. Ar fi de dorit să existe exact n direcții de căutare, $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n-1)}$ astfel încât soluția să fie găsită după exact n pași. La fiecare pas nou

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad (2.180)$$

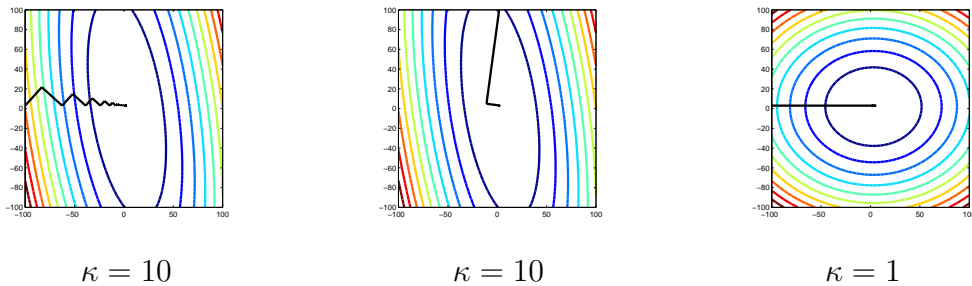


Figura 2.15: Convergența metodei gradientului depinde de inițializare și de numărul de condiționare spectrală. O problemă care are numărul de condiționare spectrală $\kappa = 1$ (toate valorile proprii sunt egale, iar forma pătratică este sferică) converge într-un singur pas.

iar α_k se va alege astfel încât eroarea $\mathbf{e}^{(k+1)}$ să fie ortogonală pe $\mathbf{d}^{(k)}$ și, în consecință, nici o altă corecție să nu se mai facă în direcția lui $\mathbf{d}^{(k)}$:

$$(\mathbf{d}^{(k)})^T \mathbf{e}^{(k+1)} = 0. \quad (2.181)$$

Dacă în (2.180) se scade soluția exactă din ambii termeni, rezultă

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad (2.182)$$

relație care înlocuită în (2.181) conduce la

$$\alpha_k = -\frac{(\mathbf{d}^{(k)})^T \mathbf{e}^{(k)}}{(\mathbf{d}^{(k)})^T \mathbf{d}^{(k)}}. \quad (2.183)$$

Relația (2.183) nu este utilă deoarece eroarea nu este cunoscută.

Soluția constă în a face *direcțiile de căutare A-ortogonale* și nu ortogonale, adică

$$(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(j)} = 0, \quad j < k. \quad (2.184)$$

Noua cerință este ca $\mathbf{e}^{(k+1)}$ să fie A-ortogonal pe $\mathbf{d}^{(k)}$, această condiție fiind echivalentă cu găsirea unui punct de minim de-a lungul direcției $\mathbf{d}^{(k)}$, așa cum se făcea la metoda gradientului. Într-adevăr, sunt valabile următoarele echivalențe

$$\begin{aligned} \frac{d}{d\alpha} f(\mathbf{x}^{(k+1)}) &= 0, \\ (\nabla f(\mathbf{x}^{(k+1)}))^T \frac{d}{d\alpha} (\mathbf{x}^{(k+1)}) &= 0, \\ -(\mathbf{r}^{(k+1)})^T \mathbf{d}^{(k)} &= 0, \\ (\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{e}^{(k+1)} &= 0. \end{aligned} \quad (2.185)$$

Rezultă următoarea expresie pentru α_k

$$\alpha_k = -\frac{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{e}^{(k)}}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}} = \frac{(\mathbf{d}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}}, \quad (2.186)$$

expresie ce se poate calcula. În cazul particular în care direcțiile $\mathbf{d}^{(k)}$ sunt reziduurile, se obține exact metoda celei mai rapide coborâri.

Această procedură calculează soluția în exact n pași. Pentru demonstrație este util să se exprime eroarea ca o combinație liniară a direcțiilor de căutare. Astfel, eroarea inițială se notează

$$\mathbf{e}^{(0)} = \sum_{j=0}^{n-1} \delta_j \mathbf{d}^{(j)}. \quad (2.187)$$

Valorile δ_j pot fi calculate prin următorul artificiu matematic. Înmulțim la stânga relația (2.187) cu $(\mathbf{d}^{(k)})^T \mathbf{A}$ și folosim faptul că direcțiile sunt A-ortogonale:

$$(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{e}^{(0)} = \sum_{j=0}^{n-1} \delta_j (\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(j)} = \delta_k (\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}. \quad (2.188)$$

În consecință

$$\delta_k = \frac{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{e}^{(0)}}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}} = \frac{(\mathbf{d}^{(k)})^T \mathbf{A} (\mathbf{e}^{(0)} + \sum_{j=0}^{k-1} \alpha_j \mathbf{d}^{(j)})}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}} = \frac{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{e}^{(k)}}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}}, \quad (2.189)$$

unde am folosit din nou A-ortogonalitatea vectorilor pentru a introduce în paranteză termeni suplimentari ce fac ca expresia finală a lui δ_k să fie exprimată în funcție de eroarea $\mathbf{e}^{(k)}$. Din (2.186) și (2.189) rezultă că $\alpha_k = -\delta_k$. În consecință eroarea la iterația k este

$$\mathbf{e}^{(k)} = \mathbf{e}^{(0)} + \sum_{j=0}^{k-1} \alpha_j \mathbf{d}^{(j)} = \sum_{j=0}^{n-1} \delta_j \mathbf{d}^{(j)} - \sum_{j=0}^{k-1} \delta_j \mathbf{d}^{(j)} = \sum_{j=k}^{n-1} \delta_j \mathbf{d}^{(j)}. \quad (2.190)$$

Este ca și cum la fiecare iterație eroarea are un termen mai puțin, astfel încât după n iterații $\mathbf{e}^{(n)} = \mathbf{0}$.

Ceea ce rămâne de făcut acum este găsirea unui set de direcții A-ortogonale $\mathbf{d}^{(k)}$. Acest lucru se realizează ușor cu ajutorul procedurii Gram-Schmidt conjugate. Această procedură pornește de la o mulțime de n vectori liniar independenți $\mathbf{u}^{(0)}, \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n-1)}$ și construiește direcțiile A-conjugate astfel. Prima direcție este chiar primul vector

$$\mathbf{d}^{(0)} = \mathbf{u}^{(0)}. \quad (2.191)$$

A doua direcție pornește de la al doilea vector la care se adaugă un vector orientat pe direcția anterioară, scalat astfel încât rezultatul să fie A-ortogonal pe direcția anterioară. Mai precis

$$\mathbf{d}^{(1)} = \mathbf{u}^{(1)} + \beta_{10} \mathbf{d}^{(0)}, \quad (2.192)$$

unde β_{10} se alege astfel încât

$$(\mathbf{d}^{(1)})^T \mathbf{A} \mathbf{d}^{(0)} = 0. \quad (2.193)$$

Această relație este echivalentă cu

$$(\mathbf{u}^{(1)} + \beta_{10} \mathbf{d}^{(0)})^T \mathbf{A} \mathbf{d}^{(0)} = 0, \quad (2.194)$$

de unde

$$\beta_{10} = -\frac{(\mathbf{u}^{(1)})^T \mathbf{A} \mathbf{d}^{(0)}}{(\mathbf{d}^{(0)})^T \mathbf{A} \mathbf{d}^{(0)}}. \quad (2.195)$$

Similar, următoarea direcție este

$$\mathbf{d}^{(2)} = \mathbf{u}^{(2)} + \beta_{20} \mathbf{d}^{(0)} + \beta_{21} \mathbf{d}^{(1)}, \quad (2.196)$$

unde β_{20} și β_{21} se aleg astfel încât

$$(\mathbf{d}^{(2)})^T \mathbf{A} \mathbf{d}^{(0)} = 0 \quad \text{și} \quad (\mathbf{d}^{(2)})^T \mathbf{A} \mathbf{d}^{(1)} = 0, \quad (2.197)$$

de unde rezultă

$$\beta_{20} = -\frac{(\mathbf{u}^{(2)})^T \mathbf{A} \mathbf{d}^{(0)}}{(\mathbf{d}^{(0)})^T \mathbf{A} \mathbf{d}^{(0)}}, \quad \beta_{21} = -\frac{(\mathbf{u}^{(2)})^T \mathbf{A} \mathbf{d}^{(1)}}{(\mathbf{d}^{(1)})^T \mathbf{A} \mathbf{d}^{(1)}}. \quad (2.198)$$

În general

$$\mathbf{d}^{(k)} = \mathbf{u}^{(k)} + \sum_{j=0}^{k-1} \beta_{kj} \mathbf{d}^{(j)}, \quad (2.199)$$

unde β_{kj} , definiți pentru $k > j$, sunt deduși din condițiile de A-ortogonalitate impuse direcțiilor, rezultând

$$\beta_{kj} = -\frac{(\mathbf{u}^{(k)})^T \mathbf{A} \mathbf{d}^{(j)}}{(\mathbf{d}^{(j)})^T \mathbf{A} \mathbf{d}^{(j)}}. \quad (2.200)$$

Dificultatea care apare este aceea că toate direcțiile de căutare trebuie memorate pentru a construi o direcție de căutare nouă.

2.12.4 Metoda gradientilor conjugați

Metoda gradientilor conjugați este metoda direcțiilor conjugate în care direcțiile de căutare sunt construite prin conjugarea reziduurilor:

$$\mathbf{u}^{(k)} = \mathbf{r}^{(k)}. \quad (2.201)$$

Această alegere este justificată din mai multe motive. Unul din motive este intuitiv, inspirat de metoda celei mai rapide coborâri, în care direcțiile de căutare erau direcțiile reziduurilor. Un al doilea motiv este dat de proprietatea reziduuului de a fi ortogonal pe direcțiile de căutare anterioare. Într-adevăr, dacă înmulțim la stânga relația (2.190) cu $-(\mathbf{d}^k)^T \mathbf{A}$ obținem

$$-(\mathbf{d}^k)^T \mathbf{A} \mathbf{e}^{(i)} = -\sum_{j=i}^{n-1} \delta_j (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^{(j)}, \quad (2.202)$$

de unde

$$(\mathbf{d}^k)^T \mathbf{r}^{(i)} = 0, \quad k < i, \quad (2.203)$$

folosind A-ortogonalitatea direcțiilor și relația (2.170). De asemenea, fiecare reziduu este ortogonal pe reziduurile anterioare

$$(\mathbf{r}^k)^T \mathbf{r}^{(i)} = 0, \quad k < i. \quad (2.204)$$

În plus, reziduuul \mathbf{r}^k este o combinație liniară dintre reziduuul anterior și produsul $\mathbf{A} \mathbf{d}^{(k-1)}$:

$$\mathbf{r}^{(k)} = -\mathbf{A} \mathbf{e}^{(k)} = -\mathbf{A}(\mathbf{e}^{(k-1)} + \alpha_{k-1} \mathbf{d}^{(k-1)}) = \mathbf{r}^{(k-1)} - \alpha_{k-1} \mathbf{A} \mathbf{d}^{(k-1)}. \quad (2.205)$$

Calculând coeficienții β în acest caz, rezultă

$$\beta_{kj} = -\frac{(\mathbf{r}^{(k)})^T \mathbf{A} \mathbf{d}^{(j)}}{(\mathbf{d}^{(j)})^T \mathbf{A} \mathbf{d}^{(j)}}. \quad (2.206)$$

Pentru a explicita termenul de la numărător, vom înmulți relația (2.205) la stânga cu $(\mathbf{r}^{(i)})^T$:

$$(\mathbf{r}^{(i)})^T \mathbf{r}^{(k+1)} = (\mathbf{r}^{(i)})^T \mathbf{r}^{(k)} - \alpha_k (\mathbf{r}^{(i)})^T \mathbf{A} \mathbf{d}^{(k)}, \quad (2.207)$$

de unde

$$\alpha_k (\mathbf{r}^{(i)})^T \mathbf{A} \mathbf{d}^{(k)} = (\mathbf{r}^{(i)})^T \mathbf{r}^{(k)} - (\mathbf{r}^{(i)})^T \mathbf{r}^{(k+1)}. \quad (2.208)$$

Pentru $i \neq k$, membrul drept al acestei relații este nenul doar pentru cazul $i = k + 1$, când valoarea lui este $-(\mathbf{r}^{(i)})^T \mathbf{r}^{(i)} / \alpha_{i-1}$. Rezultă că valorile β sunt nenule doar dacă $i = k + 1$:

$$\beta_{kj} = \begin{cases} \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{\alpha_{k-1} (\mathbf{d}^{(k-1)})^T \mathbf{A} \mathbf{d}^{(k-1)}} & k = i - 1, \\ 0 & k < i - 1 \end{cases} \quad (2.209)$$

Din acest motiv, nu mai este necesar ca valorile β să fie notate cu doi indici. Vom renota

$$\beta_k = \beta_{k,k-1} = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{\alpha_{k-1} (\mathbf{d}^{(k-1)})^T \mathbf{A} \mathbf{d}^{(k-1)}} = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{d}^{(k-1)})^T \mathbf{r}^{(k-1)}} = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}. \quad (2.210)$$

În concluzie, metoda gradientilor conjugați se bazează pe următoarele șase relații:

$$\begin{aligned} \mathbf{d}^{(0)} &= \mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(0)} \\ \alpha_k &= \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{d}^{(k)})^T \mathbf{A} \mathbf{d}^{(k)}} \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)} \\ \mathbf{r}^{(k+1)} &= \mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{d}^{(k)} \\ \beta_{k+1} &= \frac{(\mathbf{r}^{(k+1)})^T \mathbf{r}^{(k+1)}}{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}} \\ \mathbf{d}^{(k+1)} &= \mathbf{r}^{(k+1)} + \beta_{k+1} \mathbf{d}^{(k)} \end{aligned}$$

Algoritmul este complet după n iterații (fig.2.16). În practică însă metoda este folosită pentru probleme atât de mari încât nu ar fi fezabil să se execute toate cele n iterații. De aceea, iterațiile în algoritmul de mai jos sunt executate într-un ciclu cu test și nu într-un ciclu cu contor. Din acest motiv, se mai spune că *metoda este semi-iterativă*, șirul iterațiilor nu tinde către soluția exactă atunci când numărul iterațiilor tinde la infinit, ci, într-o aritmetică exactă, soluția exactă se obține după exact n iterații.

Pseudocodul de mai jos implementează algoritmul descris.

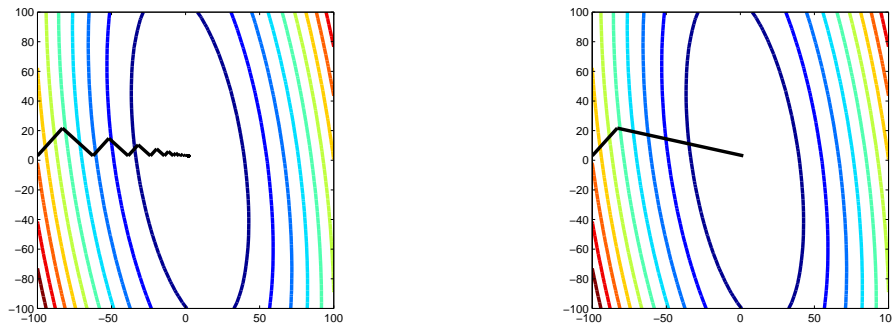


Figura 2.16: Convergența metodei gradientului depinde de inițializare și de numărul de condiționare spectrală (stânga). Metoda gradientilor conjugați converge în exact n pași, indiferent de inițializare și de numărul de condiționare spectrală (dreapta).

procedură metoda_gradientilor_conjugați($n, A, b, x_0, er, maxit, x$)

...

$\mathbf{xv} = \mathbf{x0}$; inițializează șirul iterațiilor

$k = 0$; inițializare contor iterații

$\mathbf{r} = \mathbf{b} - \mathbf{A} \cdot \mathbf{xv}$; calculează reziduu

cât timp $\|\mathbf{r}\| > er$ și $k \leq maxit$

dacă $k = 0$

$\mathbf{d} = \mathbf{r}$

altfel

$\beta = \mathbf{r}^T \mathbf{r} / (\mathbf{r} \mathbf{v}^T \cdot \mathbf{r} \mathbf{v})$

$\mathbf{d} = \mathbf{r} + \beta \mathbf{d}$

•

$\alpha = \mathbf{r}^T \mathbf{r} / (\mathbf{d}^T \mathbf{A} \mathbf{d})$

$\mathbf{x} = \mathbf{xv} + \alpha \mathbf{d}$

$\mathbf{rv} = \mathbf{r}$

$\mathbf{r} = \mathbf{rv} - \alpha \mathbf{A} \mathbf{d}$

$\mathbf{xv} = \mathbf{x}$

$\mathbf{dv} = \mathbf{d}$

$k = k + 1$

•

retur

Evident că algoritmul de mai sus se poate îmbunătăți, calculând la o iterație o singură dată produsul matrice - vector $\mathbf{A} \mathbf{d}$ și produsul scalar $\mathbf{r}^T \mathbf{r}$.

Faptul că metoda gradientilor conjugați necesită un efort mult mai mic decât metoda

direcțiilor conjugate poate fi justificat și făcând un raționament bazat pe spații vectoriale. Să notăm cu \mathcal{D}_k subspațiul vectorial generat de direcțiile de căutare până la iterația k :

$$\mathcal{D}_k = \text{span}\{\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k-1)}\}. \quad (2.211)$$

În metoda gradientilor conjugați, direcțiile de căutare sunt construite din reziduuri:

$$\begin{aligned} \mathbf{d}^{(0)} &= \mathbf{r}^{(0)}, \\ \mathbf{d}^{(1)} &= \mathbf{r}^{(1)} + \beta_1 \mathbf{d}^{(0)} = \mathbf{r}^{(1)} + \beta_1 \mathbf{r}^{(0)}, \\ \mathbf{d}^{(2)} &= \mathbf{r}^{(2)} + \beta_2 \mathbf{d}^{(1)} = \mathbf{r}^{(2)} + \beta_2 \mathbf{r}^{(1)} + \beta_2 \beta_1 \mathbf{r}^{(0)} \\ &\dots \end{aligned} \quad (2.212)$$

și în consecință subspațiul \mathcal{D}_k este în același timp și subspațiul generat de reziduuri

$$\mathcal{D}_k = \text{span}\{\mathbf{r}^{(0)}, \mathbf{r}^{(1)}, \dots, \mathbf{r}^{(k-1)}\}. \quad (2.213)$$

Pe de altă parte, din relația (2.205) se observă că reziduul la iterația k este o combinație liniară dintre reziduul la iterația $k-1$ și produsul dintre matricea \mathbf{A} și direcția de căutare la iterația $k-1$. Aceasta înseamnă că subspațiul \mathcal{D}_k este reuniunea dintre subspațiul \mathcal{D}_{k-1} și subspațiul $\mathbf{A}\mathcal{D}_{k-1}$.

$$\mathcal{D}_k = \mathcal{D}_{k-1} \cup \mathbf{A}\mathcal{D}_{k-1}. \quad (2.214)$$

De aici rezultă că

$$\mathcal{D}_k = \mathcal{D}_0 \cup \mathbf{A}\mathcal{D}_0 \cup \mathbf{A}^2\mathcal{D}_0 \cup \dots \cup \mathbf{A}^{k-1}\mathcal{D}_0, \quad (2.215)$$

deci \mathcal{D}_k este un *subspațiu Krylov*, creat prin aplicarea repetată a unei matrice unui vector:

$$\begin{aligned} \mathcal{D}_k &= \text{span}\{\mathbf{d}^{(0)}, \mathbf{A}\mathbf{d}^{(0)}, \mathbf{A}^2\mathbf{d}^{(0)} \dots, \mathbf{A}^{k-1}\mathbf{d}^{(0)}\} = \\ &= \text{span}\{\mathbf{r}^{(0)}, \mathbf{A}\mathbf{r}^{(0)}, \mathbf{A}^2\mathbf{r}^{(0)} \dots, \mathbf{A}^{k-1}\mathbf{r}^{(0)}\}. \end{aligned} \quad (2.216)$$

Aceasta înseamnă că subspațiul $\mathbf{A}\mathcal{D}_k$ este inclus în subspațiul \mathcal{D}_{k+1} și deoarece reziduul $\mathbf{r}^{(k+1)}$ este ortogonal pe \mathcal{D}_{k+1} (conform relației (2.203)) rezultă că $\mathbf{r}^{(k+1)}$ este deja A-ortogonal pe toate direcțiile de căutare cu excepția direcției $\mathbf{d}^{(k)}$. De aceea procedura Gram-Schmidt necesită la fiecare iterație calculul unui singur coeficient.

Algoritmul gradientilor conjugați poate fi dedus și din algoritmul Lanczos care este dedicat calculului valorilor proprii ale unei matrice simetrice. Mărimile deja calculate în algoritm (reziduurile, direcțiile de căutare, coeficienții β) pot fi folosiți pentru estimarea numărului de condiționare spectrală a matricei \mathbf{A} . Detalii legate de acest aspect pot fi găsite în [3].

2.13 Precondiționare

Convergența metodelor nestaționare depinde de numărul de condiționare spectrală al matricei coeficienților. Una dintre marile descoperiri în domeniul algoritmilor numerici a fost faptul că, în multe cazuri, problema inițială poate fi transformată într-una echivalentă care are proprietăți îmbunătățite considerabil. Problema echivalentă se obține de exemplu prin *precondiționarea la stânga* a problemei inițiale

$$\mathbf{Ax} = \mathbf{b} \quad (2.217)$$

adică prin rezolvarea problemei echivalente din punct de vedere al soluției

$$\mathbf{M}^{-1}\mathbf{Ax} = \mathbf{M}^{-1}\mathbf{b}, \quad (2.218)$$

unde \mathbf{M} este o matrice nesingulară, numită *matrice de precondiționare* sau *precondiționator*. Dacă se rezolvă iterativ sistemul (2.218), convergența depinde de proprietățile matricei $\mathbf{M}^{-1}\mathbf{A}$ și nu de proprietățile matricei \mathbf{A} a sistemului inițial. Dacă precondiționatorul \mathbf{M} este ales corespunzător, atunci sistemul (2.218) ar putea fi rezolvat mult mai repede decât sistemul (2.217).

Așa cum a mai fost explicat și cu alte ocazii, calculul matricei $\mathbf{M}^{-1}\mathbf{A}$ nu se face explicit deoarece nu ar fi eficient, ci prin rezolvarea unui sistem de ecuații de forma

$$\mathbf{My} = \mathbf{A} \quad (2.219)$$

Este evident că în cazurile extreme $\mathbf{M} = \mathbf{I}$ și $\mathbf{M} = \mathbf{A}$ nu există nici un câștig în rezolvare. Între aceste două situații pot exista precondiționatori utili, suficient de apropiați de \mathbf{A} într-un anumit sens astfel încât iterațiile sistemului (2.218) să convergeze mai repede decât iterațiile sistemului (2.217). O condiție puternică pentru un bun precondiționator este ca valorile proprii ale matricei $\mathbf{M}^{-1}\mathbf{A}$ să fie apropiate de 1 și ca norma $\|\mathbf{M}^{-1}\mathbf{A} - \mathbf{I}\|_2$ să fie mică [16].

O problemă echivalentă cu (2.217) se poate obține și prin *precondiționarea la dreapta*

$$\mathbf{AM}^{-1}\mathbf{y} = \mathbf{b}, \quad (2.220)$$

unde $\mathbf{x} = \mathbf{M}^{-1}\mathbf{y}$. Uneori se folosesc ambele tipuri de precondiționare simultan.

Dacă matricea \mathbf{A} este simetrică și pozitiv definită, atunci precondiționarea se face astfel încât să se păstreze această proprietate. De exemplu, alegem și matricea de precondiționare \mathbf{M} simetrică și pozitiv definită, scrisă cu ajutorul unei matrice \mathbf{C} astfel încât $\mathbf{M} = \mathbf{CC}^T$. Atunci, sistemul precondiționat de rezolvat este

$$[\mathbf{C}^{-1}\mathbf{AC}^{-T}] \mathbf{C}^T\mathbf{x} = \mathbf{C}^{-1}\mathbf{b}, \quad (2.221)$$

unde $\mathbf{C}^{-T} = (\mathbf{C}^{-1})^T = (\mathbf{C}^T)^{-1}$, iar matricea din paranteza dreapta este simetrică și pozitiv definită.

Există o mare varietate de idei de precondiționare pentru rezolvarea iterativă a sistemelor algebrice liniare. Printre cele mai cunoscute metode de precondiționare sunt:

- *Precondiționarea Jacobi (sau scalarea diagonală)* - în care $\mathbf{M} = \text{diag}(\mathbf{A})$ este o matrice diagonală care are pe diagonală valorile matricei inițiale. Matricea de precondiționare astfel construită trebuie să fie nesingulară. O generalizare a ei este alegerea $\mathbf{M} = \text{diag}(\mathbf{c})$, unde \mathbf{c} este un vector ales convenabil.
- *Factorizarea incompletă Cholesky sau LU* în care se aplică procedura de factorizare, dar factorii nu sunt calculați complet, valorile care ar umple matricea nefiind luate în considerare. Matricea de precondiționare se obține ca produsul acestor factori incompleți.

Aceste două exemple de precondiționatori nu fac nicio referire la problema inițială care a generat sistemul (2.217). Cel mai bun sfat general care poate fi dat pentru construcția precondiționatorilor este de a examina problema și de a încerca să se construiască precondiționatorul pe baza unei versiuni mai simple a problemei. Pe această idee se bazează *metodele multigrid*, care vor fi discutate ulterior. Metodele Jacobi și SSOR furnizează de multe ori precondiționatori foarte buni pentru metodele de tip multigrid.

Capitolul 3

Algoritmi numerici pentru analiza circuitelor electrice rezistive liniare

Cel mai simplu exemplu din ingineria electrică ce conduce la un sistem de ecuații algebrice liniare îl reprezintă analiza circuitelor electrice rezistive liniare. Un circuit rezistiv liniar este un circuit ce conține rezistoare, surse ideale de tensiune și curent precum și surse comandate liniar. Problema fundamentală a analizei acestor circuite are ca date topologia circuitului și valorile parametrilor (rezistențele, valorile surselor), și urmărește calculul curenților și tensiunilor din fiecare latură.

Există mai multe metode de rezolvare sistematică a unor astfel de circuite: metoda ecuațiilor Kirchhoff, metoda potențialelor nodurilor, metoda curenților ciclici. Atât metoda potențialelor nodurilor cât și metoda curenților ciclici generează un sistem de ecuații mai mic decât metoda Kirchhoff, având o matrice a coeficienților cu proprietăți mai bune (de exemplu simetrie, diagonal dominantă). Metoda curenților ciclici necesită și alegerea unui sistem convenabil de bucle independente, fiind mai dificil de transpus într-un algoritm decât metoda potențialelor nodurilor (numită mai scurt *metoda* sau *tehnica nodală*).

De aceea, în acest capitol se va explica modul în care se poate concepe un algoritm pentru rezolvarea problemei analizei circuitelor rezistive liniare folosind tehnica nodală. În prima parte se consideră cazul cel mai simplu, în care fiecare latură a circuitului conține un rezistor și eventual o sursă ideală de tensiune. Apoi, se discută modul în care algoritmul trebuie modificat pentru a putea fi aplicat unui circuit general, care conține și surse ideale de curent, surse comandate liniar și laturi alcătuite numai din surse ideale de tensiune.

3.1 Tehnica nodală

Vom considera un circuit electric alcătuit numai din laturi standard de tip sursă reală de tensiune (fig.3.1). În particular, aceste laturi pot fi rezistoare ideale, dar nu surse ideale de tensiune.

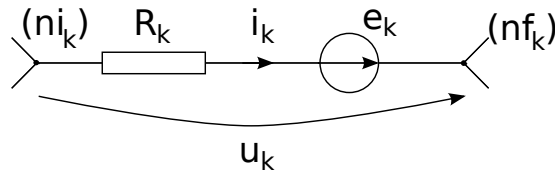


Figura 3.1: Latura standard.

Datele problemei sunt

- topologia: numărul de noduri N , numărul de laturi L și modul în care sunt conectate laturile, adică graful circuitului,
- toate rezistențele laturilor R_k , $k = 1, \dots, L$, presupuse nenule,
- toate tensiunile electromotoare e_k , $k = 1, \dots, L$

Se cere să se calculeze

- toate tensiunile u_k la bornele laturilor,
- toți curenții i_k ce străbat laturile,
- puterea consumată și puterea generată în circuit.

Deși alegerea sensurilor de referință poate fi arbitrară, în vederea transformării metodei nodale într-un algoritm, este mult mai ușor dacă se aleg sensuri de referință în mod sistematic, ca de exemplu: sensul de referință al curentului este dat de sensul de referință al sursei de tensiune (indicat de săgeata interioară) și sensul de referință al tensiunii se alege astfel încât să se respecte convenția de la receptoare (așa cum sunt reprezentate în fig. 3.1). Sensul de referință al curentului ce străbate o latură k stabilește și o relație de ordonare între cele două noduri ale laturii respective, el fiind sensul de la nodul inițial notat (ni_k) la nodul final notat (nf_k) .

Conform teoriei circuitelor electrice, fenomenele sunt complet descrise de un număr de $N - 1$ ecuații¹ Kirchhoff I

$$\sum_{k \in (n)}^A i_k = 0, \quad n = 1, \dots, N, \quad (3.1)$$

¹Notăția \sum^A reprezintă o sumă alegebrică, adică $\sum^A x_k = \sum \varepsilon_k x_k$, unde $\varepsilon_k = \pm 1$.

un număr de $L - N + 1$ ecuații Kirchoff II scrie pentru un sistem de bucle independente

$$\sum_{k \in [b]}^A u_k = 0, \quad b = 1, \dots, L - N + 1, \quad (3.2)$$

și L relații Joubert, scrise pentru toate laturile

$$u_k = R_k i_k - e_k, \quad k = 1, \dots, L, \quad (3.3)$$

în total un număr de $2L$ ecuații cu $2L$ necunoscute (toți curenții și toate tensiunile).

În tehnica nodală necunoscutele principale ale problemei (adică necunoscutele care sunt calculate mai întâi) sunt potențialele nodurilor v_k , $k = 1, \dots, N$, unde unul dintre noduri are, în mod convențional, potențialul nul. Vom presupune că numerotarea nodurilor este făcută astfel încât nodul de indice maxim este cel de referință, $v_N = 0$. Prin exprimarea tensiunilor ca diferențe de potențial, teorema Kirchoff II este identic satisfăcută. Altfel spus, relația (3.2) este satisfăcută dacă se scriu toate relațiile

$$u_k = v_{ni_k} - v_{nf_k}, \quad k = 1, \dots, L. \quad (3.4)$$

Pentru a face scrierea mai compactă este util să folosim următoarele notații:

$$\begin{aligned} \mathbf{u} &= [u_1 \ u_2 \ \dots \ u_L]^T \in \mathbb{R}^{L \times 1} && \text{vectorul tensiunilor laturilor,} \\ \mathbf{i} &= [i_1 \ i_2 \ \dots \ i_L]^T \in \mathbb{R}^{L \times 1} && \text{vectorul curenților prin laturi,} \\ \mathbf{v} &= [v_1 \ v_2 \ \dots \ v_{N-1}]^T \in \mathbb{R}^{(N-1) \times 1} && \text{vectorul potențialelor nodurilor,} \\ \mathbf{e} &= [e_1 \ e_2 \ \dots \ e_L]^T \in \mathbb{R}^{L \times 1} && \text{vectorul tensiunilor electromotoare,} \\ \mathbf{R} &= \text{diag}([R_1 \ R_2 \ \dots \ R_L]) \in \mathbb{R}^{L \times L} && \text{matricea diagonală a rezistențelor laturilor.} \end{aligned} \quad (3.5)$$

Cu aceste notații, relațiile Kirchoff I (3.1) se scriu în mod compact

$$\mathbf{A}\mathbf{i} = \mathbf{0}, \quad (3.6)$$

unde $\mathbf{A} = (a_{ij})_{i=1, N-1; j=1, L}$ este *matricea incidentelor laturi-noduri*, o matrice topologică de dimensiune $(N - 1) \times L$, un element al ei fiind definit astfel

$$a_{ij} = \begin{cases} 0 & \text{dacă nodul } i \text{ nu aparține laturii } j; \\ +1 & \text{dacă nodul } i \text{ este nod inițial pentru latura } j; \\ -1 & \text{dacă nodul } i \text{ este nod final pentru latura } j. \end{cases}$$

Relația Kirchoff II în forma (3.4) se scrie

$$\mathbf{u} = \mathbf{A}^T \mathbf{v}, \quad (3.7)$$

iar relațiile lui Joubert (3.3) se scriu compact

$$\mathbf{u} = \mathbf{R}\mathbf{i} - \mathbf{e}. \quad (3.8)$$

Dacă matricea \mathbf{R} este inversabilă (lucru adevărat dacă toate rezistențele sunt presupuse nenule) atunci din (3.8) rezultă că

$$\mathbf{i} = \mathbf{R}^{-1}(\mathbf{u} + \mathbf{e}). \quad (3.9)$$

Înlocuind (3.9) și (3.7) în (3.6), rezultă ecuația matriceală satisfăcută de vectorul potențialelor

$$\mathbf{A}\mathbf{R}^{-1}\mathbf{A}^T\mathbf{v} = -\mathbf{A}\mathbf{R}^{-1}\mathbf{e}. \quad (3.10)$$

Matricea coeficienților sistemului algebric liniar de rezolvat

$$\mathbf{G} = \mathbf{A}\mathbf{R}^{-1}\mathbf{A}^T \in \mathbb{R}^{(N-1) \times (N-1)} \quad (3.11)$$

este numită *matricea conductanțelor nodale*. Termenii ei au următoarea semnificație: orice termen diagonal G_{ii} reprezintă suma conductanțelor laturilor care concură la nodul i , iar orice termen nediagonal G_{ij} cu $i \neq j$ reprezintă suma conductanțelor laturilor care unesc direct nodul i cu nodul j , luată cu semnul minus. Altfel scris,

$$G_{ii} = \sum_{k \in (i)} \frac{1}{R_k}, \quad (3.12)$$

$$G_{ij} = - \sum_{k \in (i); k \in (j)} \frac{1}{R_k} \quad \text{pentru } i \neq j. \quad (3.13)$$

Termenul liber al sistemului de ecuații

$$\mathbf{t} = -\mathbf{A}\mathbf{R}^{-1}\mathbf{e} \in \mathbb{R}^{(N-1) \times 1} \quad (3.14)$$

este numit *vectorul injecțiilor de curent*. Un termen al acestui vector t_k reprezintă suma algebrică a unor termeni de tipul e_m/R_m pentru toate laturile m care concură la nodul k , cu plus dacă săgeata internă a sursei de tensiune electromotoare de pe latura m intră în nodul k , și cu semnul minus în caz contrar:

$$t_k = \sum_{m \in (k)} \frac{A e_m}{R_m}. \quad (3.15)$$

Este important să remarcăm că matricea \mathbf{G} este simetrică și pozitiv definită dacă rezistențele sunt pozitive. Deoarece \mathbf{R} este o matrice diagonală, inversa ei este tot o matrice diagonală, având pe diagonală valorile conductanțelor laturilor

$$\mathbf{R}^{-1} = \text{diag}([1/R_1 \quad 1/R_2 \quad \dots \quad 1/R_L]). \quad (3.16)$$

Transpusa matricei conductanțelor nodale va fi

$$\mathbf{G}^T = (\mathbf{A}\mathbf{R}^{-1}\mathbf{A}^T)^T = (\mathbf{A}^T)^T (\mathbf{R}^{-1})^T (\mathbf{A})^T = \mathbf{A}\mathbf{R}^{-1}\mathbf{A}^T = \mathbf{G}, \quad (3.17)$$

deci ea este o matrice simetrică. Relațiile (3.12) și (3.13) indică faptul că ea este și diagonal dominantă.

Pentru demonstrarea proprietății de pozitiv definire, să considerăm un vector coloană \mathbf{x} arbitrar, nenul. Atunci

$$\mathbf{x}^T \mathbf{G} \mathbf{x} = \mathbf{x}^T \mathbf{A} \mathbf{R}^{-1} \mathbf{A}^T \mathbf{x} = \mathbf{y}^T \mathbf{R}^{-1} \mathbf{y} = \sum_{k=1}^L \frac{y_k^2}{R_k} > 0, \quad (3.18)$$

unde am notat $\mathbf{y} = \mathbf{A}^T \mathbf{x}$ un vector coloană de componente y_k , $k = 1, \dots, L$. Inegalitatea demonstrată în (3.18) este strictă, ea ar putea fi zero doar dacă vectorul \mathbf{y} , și în consecință vectorul \mathbf{x} este nul, ceea ce contrazice ipoteza făcută. În concluzie, matricea conductanțelor nodale este pozitiv definită dacă rezistențele laturilor sunt strict pozitive.

Există mai multe variante posibile de concepere a algoritmului acestei metode. Toate au trei etape principale: *etapa de preprocesare* în care se descrie problema și se assemblează sistemul de ecuații de rezolvat, *etapa de rezolvare* în care se apelează o procedură propriu-zisă de rezolvare a sistemului de ecuații rezultat (procedură numită foarte adesea "solver") și *etapa de postprocesare* în care se calculează alte mărimi de interes.

3.1.1 Structuri de date

În conceperea unui algoritm, stabilirea structurilor de date ce vor fi folosite este de asemenea o etapă foarte importantă. Numele datelor ce le vom folosi în algoritm vor fi alese în concordanță cu teoria prezentată în paragraful anterior.

Declarații posibile pentru datele problemei sunt

; declaratii date - varianta A

<u>întreg</u> N	; număr de noduri
<u>întreg</u> L	; număr de laturi
<u>tablou</u> <u>întreg</u> $ni[L]$; noduri inițiale ale laturilor
<u>tablou</u> <u>întreg</u> $nf[L]$; noduri finale ale laturilor
<u>tablou</u> <u>real</u> $R[L]$; rezistențe
<u>tablou</u> <u>real</u> $e[L]$; tensiuni electromotoare

Pentru a evita transmiterea unui număr mare de parametri ca argumente de funcții, se recomandă agregarea datelor, ca de exemplu

; declarații date - varianta B

înregistrare circuit

întreg N ; număr de noduri
întreg L ; număr de laturi
tablou întreg $ni[L]$; noduri inițiale ale laturilor
tablou întreg $nf[L]$; noduri finale ale laturilor
tablou real $R[L]$; rezistențe
tablou real $e[L]$; tensiuni electromotoare

•

Un alt aspect important este acela că matricea conductanțelor nodale și vectorul injecțiilor de curent sunt rare. De exemplu, presupunând că în medie sunt 4 laturi care concură la un nod, rezultă o densitate a matricei egală cu $5n/n^2 = 5/n$, foarte mică chiar și pentru probleme mici (pentru $N \approx 1000$ rezultă $d = 0.5\%$). De aceea, pentru a concepe un algoritm eficient (atât din punct de vedere al timpului de calcul cât și al memoriei, vor trebui folosite tehnici de matrice rare pentru stocarea ei. Altfel, programul rezultat poate deveni ineficient chiar la valori mici (≈ 1000) ale numărului de noduri. Totuși, pentru a păstra simplă descrierea algoritmului, în cele ce urmează vom păstra același tip de declarații și pentru matricile rare și vectorii rari. În consecință, variabile mai importante care vor apărea în algoritm sunt:

; declarații variabile utile
tablou real $G[N, N]$; matricea conductanțelor nodale (stocată rar)
tablou real $t[N]$; vectorul injecțiilor de curent (stocat rar)
tablou real $v[N]$; vectorul potențialelor

3.1.2 Etapa de preprocesare

Etapa de preprocesare constă în citirea datelor (de exemplu de la tastatură sau din fișiere) și în asamblarea sistemului de ecuații de rezolvat.

În cazul variantei A de declarații, procedura de citire a datelor poate fi

procedură citire_date_A (N, L, ni, nf, R, e)
 ; declarații
 ...
citește N, L
pentru $k = 1, L$
 citește ni_k, nf_k, R_k, e_k

•

retur

	ni _k				nf _k					
	*	*	*	*	*	*	*	*	*	
ni _k	*	+1/R _k	*	*	-1/R _k	*	*	*	ni _k	-e _k /R _k
	*	*	*	*	*	*	*	*	*	
	*	*	*	*	*	*	*	*	*	
nf _k	*	-1/R _k	*	*	+1/R _k	*	*	*	nf _k	+e _k /R _k
	*	*	*	*	*	*	*	*	*	
	*	*	*	*	*	*	*	*	*	

Figura 3.2: Contribuția unei laturi k la matricea conductanțelor nodale (stânga) și la vectorul injecțiilor de curent (dreapta).

În cazul variantei B de declarații, rutina de citire a datelor poate fi o funcție, de tipul

```

funcție citire_date_B ()
; declarații
...
citește circuit.N, circuit.L
pentru k = 1,circuit.L
    citește circuit.nik, circuit.nfk, circuit.Rk, circuit.ek
•
întoarce circuit

```

Există mai multe variante de a asambla sistemul de ecuații. O variantă ar fi cea în care se parcurg nodurile și se scriu ecuațiile una câte una. Așa ar face și o persoană care ar aplica această metodă, cu creionul pe hârtie, pentru a rezolva o problemă de dimensiuni extrem de mici. În această abordare trebuie identificate care sunt laturile ce ating un nod. Numărul lor variază de la nod la nod, iar algoritmul corespunzător este destul de costisitor necesitând analiza grafului circuitului. O altă abordare se bazează pe parcurgerea laturilor și adunarea contribuțiilor acestora la sistem. Această variantă este mult mai simplă de conceput deoarece fiecare latură are exact două noduri. Din același motiv și descrierea circuitului a fost orientată pe laturi și nu pe noduri. Contribuția fiecărei laturi la matricea conductanțelor nodale și la vectorul injecțiilor de curent este ilustrată în fig. 3.2.

Algoritmul metodei nodale este următorul

```

procedură nodalRE_v1 (circuit, G, t)
; assemblează sistemul de ecuații pentru un circuit cu laturi de tip
; sursă reală de tensiune, folosind tehnica nodală

```

```

; parametri de intrare:
;          circuit - structură de date ce descrie circuitul
; parametri de ieșire:
;          G - matricea conductanțelor nodale și
;          t - vectorul injecțiilor de curent
; declarații
....
L = circuit.L ; pentru simplificarea scrierii algoritmului
N = circuit.N
ni = circuit.ni
nf = circuit.nf
R = circuit.R
e = circuit.e
; anulează componentele matricei G și ale vectorului termenilor liberi t
G = 0
t = 0
; assemblează sistem
pentru k = 1, L ; parcurge laturi
    i = nik ; nodul inițial al laturii k
    j = nfk ; nodul final al laturii k
    Gii = Gii + 1/Rk
    Gjj = Gjj + 1/Rk
    Gij = Gij - 1/Rk
    Gji = Gji - 1/Rk
    ti = ti - ek/Rk
    tj = tj + ek/Rk

```

•

retur

O primă observație este aceea că, pentru simplificarea pseudocodului, am preferat să nu scriem explicit modul în care se anulează componentele matricei coeficienților și ale vectorului termenilor liberi, înainte de asamblarea propriu-zisă. Pentru a sugera însă că se fac mai multe operații asupra componentelor am folosit litere aldine. Atragem atenția însă că o variantă de tipul

```

pentru i = 1,N
    pentru j = 1,N
        Gij = 0

```

•

•

scrisă pentru ”instrucțiunea” $\mathbf{G} = \mathbf{0}$ va face ca matricea coeficienților să devină o matrice plină. Această instrucțiune trebuie scrisă explicit în funcție de modul de memorare al matricei rare. De asemenea, pentru a evita repetarea unor calcule, se pot memora valorile $1/R_k$ și e_k/R_k .

O altă variantă de implementare poate folosi asamblarea matricei de incidență. Calculul matricei conductanțelor nodale și a vectorului injecțiilor de curent se va face apoi folosind produse de matrice sau produse matrice-vectori, scrise pentru matrice rare.

procedură nodalRE_v2 (circuit, G , t)

; assemblează sistemul de ecuații pentru un circuit cu laturi de tip

; sursă reală de tensiune, folosind tehnica nodală

; parametri de intrare:

; circuit - structură de date ce descrie circuitul

; parametri de ieșire:

; G - matricea conductanțelor nodale și

; t - vectorul injecțiilor de curent

; declarații

....

$L = \text{circuit}.L$; pentru simplificarea scrierii algoritmului

$N = \text{circuit}.N$

$ni = \text{circuit}.ni$

$nf = \text{circuit}.nf$

$R = \text{circuit}.R$

$e = \text{circuit}.e$

; anulează componentele:

$\mathbf{A} = \mathbf{0}$; matricei incidente laturi noduri

$\mathbf{Glat} = \mathbf{0}$; matricei diagonale \mathbf{R}^{-1}

; assemblează sistem

pentru $k = 1, L$; parcurge laturi

$i = ni_k$; nodul inițial al laturii k

$j = nf_k$; nodul final al laturii k

$A_{ik} = -1$

$A_{jk} = +1$

$\text{Glat}_{kk} = 1/R_k$

•

$\mathbf{G} = \mathbf{A} * \mathbf{Glat} * \mathbf{A}^T$; apel proceduri speciale pentru matrice rare

```
t = -A * Glat * e
```

```
retur
```

Varianta a doua a algoritmului se poate dovedi mult mai eficientă decât prima variantă dacă implementarea se face în medii de calcul precum Matlab.

3.1.3 Etapa de rezolvare

În etapa de preprocesare asamblarea sistemului nu a fost făcută la dimensiunea $(N - 1) \times (N - 1)$ ci la dimensiunea $N \times N$, nodul de referință nefiind tratat special. Acest lucru face ca scrierea pseudocodului pentru asamblarea matricei să fie foarte ușoară. Pentru rezolvare însă, sistemul trebuie să fie de dimensiune $N - 1$. Presupunând că nodul N este nodul de referință ($v_N = 0$), atunci matricea coeficienților se obține eliminând ultima linie și ultima coloană din matricea asamblată, iar vectorul termenilor liberi se obține eliminând ultima componentă.

Pentru rezolvare avem la dispoziție algoritmi pregătiți în capitolul 1. De exemplu, dacă vrem să folosim metoda Gauss, atunci apelul ei se face astfel

```
Gauss (N - 1, G, t, v)
```

```
vN = 0
```

Vectorul soluție întors reprezintă potențialele a $N - 1$ noduri. Potențialul ultimului nod a fost considerat zero și de aceea, pentru completitudine, este adăugată atribuirea $v_N = 0$ imediat după apelul procedurii de rezolvare. O altă observație este aceea că matricea conductanțelor nodale este diagonal dominantă, ceea ce înseamnă că, în cazul folosirii metodei Gauss, pivotarea nu este necesară nici pentru diminuarea erorilor de rotunjire.

Algoritmul Gauss prezentat în paragraful 2.4.2 nu este însă conceput pentru matrice rare. Într-o astfel de înlănțuire a pseudocodului, chiar dacă sistemul a fost generat cu tehnici de matrice rare, rezolvarea lui cu această procedură strică raritatea. De fapt, este mai mult decât atât, matricea devine plină cu zerouri, dar memorată cu structuri de matrice rare, ceea ce este mult mai costisitor decât memorarea matricei pline.

În conceperea unui algoritm pentru probleme reale trebuie acordată o atenție deosebită memorării cu eficiență a structurilor de date și efectuarea unor instrucțiuni adecvate cu acestea.

3.1.4 Etapa de postprocesare

După rezolvarea sistemului putem calcula orice alte mărimi de interes: tensiuni, curenți prin laturi, puteri. Următorul pseudocod ilustrează acest calcul

```

procedură postprocesare_circuitRE (circuit, v)
; declarații
...
L = circuit.L
ni = circuit.ni
nf = circuit.nf
R = circuit.R
e = circuit.e
Pc = 0 ; puterea consumată
Pg = 0 ; puterea generată
pentru k = 1, L ; parcurge laturi
    u = vnik - vnfk ; tensiunea laturii
    c = (u + ek)/Rk ; curentul prin latură
    scrie "Latura" k "are tensiunea" u "și curentul" c
    Pc = Pc + Rkc2 ; adaugă contribuția laturii la puterea consumată
    Pg = Pg + ekc ; adaugă contribuția laturii la puterea generată
•
scrie Pc, Pg
retur

```

Este posibil ca, datorită erorilor de rotunjire, valorile calculate ale puterilor consumată și generată să nu fie identice. Acest lucru se întâmplă mai ales dacă matricea sistemului este prost condiționată numeric, caz ce poate apărea dacă valorile rezistențelor sunt foarte diferite.

3.2 Tratarea surselor ideale de curent

Tratarea surselor ideale de curent se face cu ușurință dacă se consideră latura standard de tipul celei din fig. 3.3.

Și în acest caz vom presupune că toate rezistențele R_k sunt nenule. În particular, această latură poate fi o sursă reală de tensiune (dacă $J_k = 0$), un rezistor (dacă $J_k = 0$ și $E_k = 0$), o sursă ideală de curent (dacă $1/R_k = 0$).

Ecuatiile Kirchhoff I și II se scriu sub aceeași formă:

$$\mathbf{A}\mathbf{i} = \mathbf{0}, \quad (3.19)$$

$$\mathbf{u} = \mathbf{A}^T\mathbf{v}. \quad (3.20)$$

Relațiile între tensiunile și curenții unei laturi sunt acum

$$u_k = R_k(i_k - J_k) - e_k, \quad k = 1, \dots, L. \quad (3.21)$$

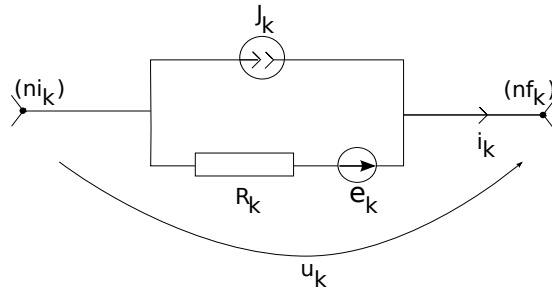


Figura 3.3: Latura standard.

Folosind aceleași notații (3.5), relațiile de mai sus se scriu compact

$$\mathbf{u} = \mathbf{R}(\mathbf{i} - \mathbf{J}) - \mathbf{e}, \quad k = 1, \dots, L, \quad (3.22)$$

unde $\mathbf{J} = (J_k)_{k=1,L}$ este vectorul curenților electromotori prin laturile circuitului. Rezultă curenții prin laturile standard

$$\mathbf{i} = \mathbf{R}^{-1}(\mathbf{u} + \mathbf{e}) + \mathbf{J}, \quad (3.23)$$

și ecuația matriceală în potențiale:

$$\mathbf{A}\mathbf{R}^{-1}\mathbf{A}^T\mathbf{v} = -\mathbf{A}\mathbf{R}^{-1}\mathbf{e} - \mathbf{A}\mathbf{J}. \quad (3.24)$$

Comparând ecuația rezultată cu ecuația (3.10) obținută în cazul în care latura standard conținea numai o sursă reală de tensiune, se observă că matricea coeficienților este neschimbată, iar în membrul drept apare un termen suplimentar, corespunzător sursei ideale de curent. Acum, un element al termenului liber

$$\mathbf{t} = -\mathbf{A}\mathbf{R}^{-1}\mathbf{e} - \mathbf{A}\mathbf{J} \in \mathbb{R}^{(N-1) \times 1} \quad (3.25)$$

se calculează ca suma algebrică a unor termeni de tipul e_m/R_m pentru toate laturile m care concură la nodul k , cu plus dacă săgeata internă a sursei de tensiune electromotoare de pe latura m intra în nodul k , și cu semnul minus în caz contrar și suma algebrică a unor termeni de tipul J_m cu plus dacă săgeata internă a sursei de curent de pe latura m intra în nodul k , și cu semnul minus în caz contrar.

Pentru a simplifica algoritmul, vom presupune că tensiunea electromotoare e_k și curențul electromotor J_k sunt au sensurile de referință la fel față de nod. Modificările algoritmului vor afecta etapa de preprocesare unde vor trebui citite și valorile curenților electromotori și vor trebui corecțati termenii liberi ai sistemului de rezolvat.

funcție citire_date_REJ ()

...

citește circuit. N , circuit. L

pentru $k = 1, \text{circuit}.L$

citește circuit. ni_k , circuit. nf_k , circuit. R_k , circuit. e_k , **circuit. J_k**

•

întoarce circuit

procedură nodalREJ_v1 (circuit, G , t)

; asamblează sistemul de ecuații pentru un circuit cu laturi de tip

; sursă reală de tensiune **în paralel cu sursă ideală de curent**,

; folosind tehnica nodală

; parametri de intrare:

; circuit - structură de date ce descrie circuitul

; parametri de ieșire:

; G - matricea conductanțelor nodale și

; t - vectorul injecțiilor de curent

; declarații

....

$L = \text{circuit}.L$; pentru simplificarea scrierii algoritmului

$N = \text{circuit}.N$

$ni = \text{circuit}.ni$

$nf = \text{circuit}.nf$

$R = \text{circuit}.R$

$e = \text{circuit}.e$

$J = \text{circuit}.J$

; anulează componentele matricei G și a vectorului termenilor liberi t

$G = 0$

$t = 0$

; asamblează sistem

pentru $k = 1, L$

; parcurge laturi

$i = ni_k$

; nodul inițial al laturii k

$j = nf_k$

; nodul final al laturii k

$G_{ii} = G_{ii} + 1/R_k$

$G_{jj} = G_{jj} + 1/R_k$

$G_{ij} = G_{ij} - 1/R_k$

$G_{ji} = G_{ji} - 1/R_k$

$t_i = t_i - e_k/R_k - J_k$

$t_j = t_j + e_k/R_k + J_k$

•

retur

procedură nodalREJ_v2 (circuit, G , t)

; asamblează sistemul de ecuații pentru un circuit cu laturi de tip

; sursă reală de tensiune în paralel cu sursă ideală de curent,

; folosind tehnica nodală

; parametri de intrare:

; circuit - structură de date ce descrie circuitul

; parametri de ieșire:

; G - matricea conductanțelor nodale și

; t - vectorul injecțiilor de curent

; declarații

....

$L = \text{circuit}.L$; pentru simplificarea scrierii algoritmului

$N = \text{circuit}.N$

$ni = \text{circuit}.ni$

$nf = \text{circuit}.nf$

$R = \text{circuit}.R$

$e = \text{circuit}.e$

$J = \text{circuit}.J$

; anulează componentele:

$\mathbf{A} = \mathbf{0}$; matricei incidentelor laturi noduri

$\mathbf{Glat} = \mathbf{0}$; matricei diagonale \mathbf{R}^{-1}

; asamblează sistem

pentru $k = 1, L$; parcurge laturi

$i = ni_k$; nodul inițial al laturii k

$j = nf_k$; nodul final al laturii k

$A_{ik} = -1$

$A_{jk} = +1$

$\text{Glat}_{kk} = 1/R_k$

•

$\mathbf{G} = \mathbf{A} * \mathbf{Glat} * \mathbf{A}^T$; apel proceduri speciale pentru matrice rare

$\mathbf{t} = -\mathbf{A} * \mathbf{Glat} * \mathbf{e} - \mathbf{A} * \mathbf{J}$

retur

Etapa de postprocesare trebuie adaptată noii configurații a laturii standard:

procedură postprocesare_circuitREJ (circuit, v)

; declarații

...

 $L = \text{circuit}.L$ $ni = \text{circuit}.ni$ $nf = \text{circuit}.nf$ $R = \text{circuit}.R$ $e = \text{circuit}.e$ $J = \text{circuit}.J$ $P_c = 0$; puterea consumată $P_g = 0$; puterea generatăpentru $k = 1, L$; parcurge laturi $u = v_{ni_k} - v_{nf_k}$; tensiunea laturii $c = (u + e_k)/R_k + J_k$; curentul prin laturăscrie "Latura" k "are tensiunea" u "și curentul" c $P_c = P_c + R_k(c - J_k)^2$; adaugă contribuția laturii la puterea consumată $P_g = P_g + e_k(c - J_k) - uJ_k$; adaugă contribuția laturii la puterea generată

•

scrie P_c, P_g retur

3.3 Tratarea surselor ideale de tensiune

Vom presupune acum că circuitul are și surse ideale de tensiune. Laturile circuitului sunt fie de tip sursă reală de tensiune în paralel cu sursă ideală de curent (vom numi această latură de tip "REJ"), fie de tip sursă ideală de tensiune (tip "E") (fig.3.4). Pentru a simplifica scrierea algoritmului, vom presupune că sunt numerotate mai întâi laturile de tip REJ și apoi laturile de tip E. Vom nota cu L_{REJ} și cu L_E laturile din fiecare categorie și cu $L = L_{REJ} + L_E$ numărul total de laturi.

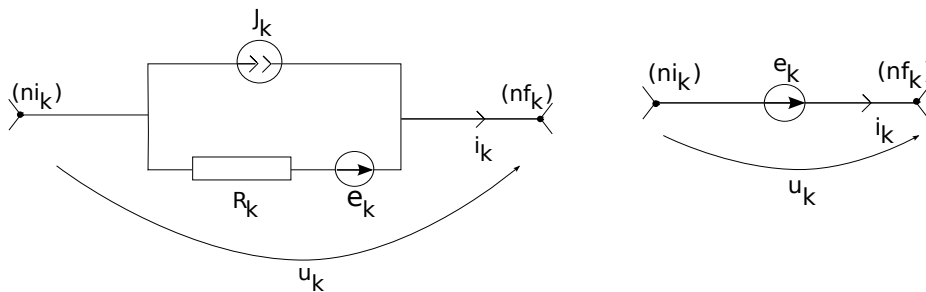


Figura 3.4: Laturi standard.

Ecuatiile Kirchhoff I și II se scriu sub aceeași formă:

$$\mathbf{A}\mathbf{i} = \mathbf{0}, \quad (3.26)$$

$$\mathbf{u} = \mathbf{A}^T \mathbf{v}. \quad (3.27)$$

Relațiile între tensiunile și curenții unei laturi sunt acum

$$u_k = R_k(i_k - J_k) - e_k, \quad k = 1, \dots, L_{REJ}, \quad (3.28)$$

$$u_k = -e_k, \quad k = L_{REJ} + 1, \dots, L_{REJ} + L_E. \quad (3.29)$$

Pentru scrierea compactă a acestor relații, vom partiționa vectorii în două, corespunzător celor două tipuri de laturi:

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_{REJ} \\ \mathbf{u}_E \end{bmatrix}, \quad \mathbf{i} = \begin{bmatrix} \mathbf{i}_{REJ} \\ \mathbf{i}_E \end{bmatrix}, \quad \mathbf{e} = \begin{bmatrix} \mathbf{e}_{REJ} \\ \mathbf{e}_E \end{bmatrix}, \quad (3.30)$$

unde vectorii cu indice REJ au dimensiunea L_{REJ} , iar vectorii cu indice E au dimensiunea L_E .

Formele matriceale ale relațiilor (3.28) și (3.29) sunt

$$\mathbf{u}_{REJ} = \mathbf{R}(\mathbf{i}_{REJ} - \mathbf{J}) - \mathbf{e}_{REJ}, \quad (3.31)$$

$$\mathbf{u}_E = -\mathbf{e}_E, \quad (3.32)$$

unde \mathbf{R} este acum o matrice diagonală, de dimensiune L_{REJ} , având pe diagonală rezistențele laturilor de tip REJ, iar \mathbf{J} este un vector de dimensiune L_{REJ} .

Partiționarea în două a vectorului curenților face utilă partiționarea în două blocuri a matricei incidențelor laturi-noduri, astfel încât relația (3.26) devine

$$\begin{bmatrix} \mathbf{A}_{REJ} & \mathbf{A}_E \end{bmatrix} \begin{bmatrix} \mathbf{i}_{REJ} \\ \mathbf{i}_E \end{bmatrix} = \mathbf{0}, \quad (3.33)$$

unde \mathbf{A}_{REJ} este o matrice de dimensiune $(N - 1) \times L_{REJ}$ ce conține incidența laturilor de tip REJ la noduri, iar matricea \mathbf{A}_E este o matrice de dimensiune $(N - 1) \times L_E$ ce conține incidența laturilor de tip E la noduri, de unde rezultă, după realizarea operațiilor cu blocuri de matrice că

$$\mathbf{A}_{REJ}\mathbf{i}_{REJ} + \mathbf{A}_E\mathbf{i}_E = \mathbf{0}. \quad (3.34)$$

Relația (3.27) se explicitază și ea ca

$$\begin{bmatrix} \mathbf{u}_{REJ} \\ \mathbf{u}_E \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{REJ}^T \\ \mathbf{A}_E^T \end{bmatrix} \mathbf{v}, \quad (3.35)$$

de unde rezultă că

$$\mathbf{u}_{REJ} = \mathbf{A}_{REJ}^T \mathbf{v}, \quad (3.36)$$

$$\mathbf{u}_E = \mathbf{A}_E^T \mathbf{v}. \quad (3.37)$$

Recapitulând, relațiile matriceale utile sunt:

$$\mathbf{A}_{REJ} \mathbf{i}_{REJ} + \mathbf{A}_E \mathbf{i}_E = \mathbf{0}, \quad (3.38)$$

$$\mathbf{u}_{REJ} = \mathbf{A}_{REJ}^T \mathbf{v}, \quad (3.39)$$

$$\mathbf{u}_E = \mathbf{A}_E^T \mathbf{v}, \quad (3.40)$$

$$\mathbf{u}_{REJ} = \mathbf{R}(\mathbf{i}_{REJ} - \mathbf{J}) - \mathbf{e}_{REJ}, \quad (3.41)$$

$$\mathbf{u}_E = -\mathbf{e}_E. \quad (3.42)$$

Din (3.41) rezultă curenții prin laturile de tip REJ:

$$\mathbf{i}_{REJ} = \mathbf{R}^{-1}(\mathbf{u}_{REJ} + \mathbf{e}_{REJ}) + \mathbf{J}. \quad (3.43)$$

Această expresie împreună cu (3.39) se substituie în (3.38). Sistemul de rezolvat devine

$$\mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{A}_{REJ}^T \mathbf{v} + \mathbf{A}_E \mathbf{i}_E = -\mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{e}_{REJ} - \mathbf{A}_{REJ} \mathbf{J}, \quad (3.44)$$

$$\mathbf{A}_E^T \mathbf{v} = -\mathbf{e}_E. \quad (3.45)$$

Acest sistem este un sistem algebric liniar, de dimensiune $N - 1 + L_E$. Necunoscutele acestui sistem sunt atât potențialele nodurilor cât și valorile curenților prin sursele ideale de tensiune:

$$\mathbf{x} = \begin{bmatrix} \mathbf{v} \\ \mathbf{i}_E \end{bmatrix}. \quad (3.46)$$

Sistemul de rezolvat

$$\mathbf{M} \mathbf{x} = \mathbf{p} \quad (3.47)$$

are matricea coeficienților

$$\mathbf{M} = \begin{bmatrix} \mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{A}_{REJ}^T & \mathbf{A}_E \\ \mathbf{A}_E^T & \mathbf{0} \end{bmatrix} \quad (3.48)$$

și vectorul termenilor liberi

$$\mathbf{p} = \begin{bmatrix} -\mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{e}_{REJ} - \mathbf{A}_{REJ} \mathbf{J} \\ -\mathbf{e}_E \end{bmatrix} \quad (3.49)$$

Matricea coeficienților rămâne simetrică dar își pierde pozitiv definirea. O prelucrare suplimentară a acestui sistem poate conduce însă la un sistem mai mic și cu proprietăți

mai bune. Pentru aceasta, vom presupune că sursele ideale de tensiune formează un subgraf conex. Această ipoteză nu este restrictivă deoarece se știe că, indiferent cum ar fi plasate sursele ideale de tensiune, se poate găsi un circuit, echivalent cu circuitul inițial din punct de vedere al grafului de curenți, care să satisfacă această proprietate. Mai mult, subgraful conex nu are bucle, deoarece se știe că un model de circuit în care apar bucle formate numai din surse ideale de tensiune nu este o problemă bine formulată (nu are soluție unică). Vom presupune că nodurile sunt numerotate astfel încât sunt lăsate la sfârșit nodurile care aparțin subgrafului conex al surselor ideale de tensiune. Ideea ce va fi explicitată în relații matematice în cele ce urmează este aceea că nodurile care aparțin acestui subgraf conex au potențiale cunoscute. Pentru fiecare astfel de nod, există o cale de la acest nod la nodul de referință (care este și el un nod al acestui subgraf conex), formată numai din surse ideale de tensiune. Aceste potențiale vor fi exprimate în funcție de tensiunile electromotoare ale surselor ideale de tensiune și apoi substituite în relațiile Kirchhoff I aduse la forma (3.44). Sistemul rezultat va fi de dimensiune $N - 1 - L_E$ iar matricea coeficienților va reprezenta conductanțe nodale, deci va fi simetrică, pozitiv definită și diagonal dominantă.

Vom nota partiționarea vectorului potențialelor în

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}_* \\ \mathbf{v}_c \end{bmatrix}, \quad (3.50)$$

unde \mathbf{v}_c este vectorul potențialelor nodurilor subgrafului conex format din surse ideale de tensiune (în care nu este inclus nodul de referință), având dimensiunea L_E , iar \mathbf{v}_* este vectorul potențialelor restului de noduri, de dimensiune $N - 1 - L_E$. Această partiționare impune o partiționare suplimentară a matricei \mathbf{A}_E , matrice ce conține incidența laturilor de tip E la noduri. Deoarece laturile de tip E sunt acum conectate numai la noduri ale subgrafului conex, înseamnă că blocul superior, care reprezintă incidența laturilor de tip E la nodurile de potențial necunoscut este nul:

$$\mathbf{A}_E = \begin{bmatrix} \mathbf{0} \\ \mathbf{A}_{Ec} \end{bmatrix}, \quad (3.51)$$

unde \mathbf{A}_{Ec} este o matrice pătrată de dimensiunea numărului de noduri de potențial cunoscut (exceptând nodul de referință) ce reprezintă incidența laturilor de tip E la nodurile de potențial cunoscut. Relația (3.45) devine

$$\begin{bmatrix} \mathbf{0} & \mathbf{A}_{Ec}^T \end{bmatrix} \begin{bmatrix} \mathbf{v}_* \\ \mathbf{v}_c \end{bmatrix} = -\mathbf{e}_E, \quad (3.52)$$

de unde

$$\mathbf{A}_{Ec}^T \mathbf{v}_c = -\mathbf{e}_E. \quad (3.53)$$

Potențialele nodurilor subgrafului conex format numai din surse ideale de tensiune sunt

$$\mathbf{v}_c = -(\mathbf{A}_{Ec}^T)^{-1} \mathbf{e}_E. \quad (3.54)$$

Pentru o problemă de circuit bine formulată, matricea \mathbf{A}_{Ec}^T este inversabilă, fiecare linie a inversei ei reprezentând calea de la un nod la nodul de referință, cale aleasă doar din laturi ale subgrafului conex al surselor ideale de tensiune. Această cale este unică deoarece laturile acestui subgraf conex nu formează bucle.

Pentru a înlocui (3.54) în (3.44) este util să partiționăm matricile ce intervin astfel:

$$\mathbf{G}' = \mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{A}_{REJ}^T = \begin{bmatrix} \mathbf{G}_{**} & \mathbf{G}_{*c} \\ \mathbf{G}_{*c}^T & \mathbf{G}_{cc} \end{bmatrix} \in \mathbb{R}^{(N-1) \times (N-1)} \quad (3.55)$$

$$\mathbf{t}' = -\mathbf{A}_{REJ} \mathbf{R}^{-1} \mathbf{e}_{REJ} - \mathbf{A}_{REJ} \mathbf{J} = \begin{bmatrix} \mathbf{t}_* \\ \mathbf{t}_{*c} \end{bmatrix} \in \mathbb{R}^{(N-1) \times 1}. \quad (3.56)$$

Matricea \mathbf{G}_{**} este o matrice pătrată de dimensiune $N - 1 - L_E$, reprezentând conductanțe nodale pentru subgraful care conține laturi de tip REJ și nodurile de potențial necunoscut. Matricea \mathbf{G}_{*c} este o matrice de dimensiune $(N - 1 - L_E) \times L_E$ care conține conductanțe nodale nediagonale între nodurile cu potențial cunoscut și cele cu potențial necunoscut. \mathbf{G}_{cc} este o matrice pătrată de dimensiune L_E care conține conductanțe nodale pentru subgraful care are laturi de tip RE dar noduri de potențial cunoscut.

Relația (3.44) devine

$$\begin{bmatrix} \mathbf{G}_{**} & \mathbf{G}_{*c} \\ \mathbf{G}_{*c}^T & \mathbf{G}_{cc} \end{bmatrix} \begin{bmatrix} \mathbf{v}_* \\ \mathbf{v}_c \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{A}_{Ec} \end{bmatrix} \mathbf{i}_E = \begin{bmatrix} \mathbf{t}_* \\ \mathbf{t}_{*c} \end{bmatrix}, \quad (3.57)$$

relație echivalentă cu următoarele două relații matriciale:

$$\mathbf{G}_{**} \mathbf{v}_* + \mathbf{G}_{*c} \mathbf{v}_c = \mathbf{t}_*, \quad (3.58)$$

$$\mathbf{G}_{*c}^T \mathbf{v}_* + \mathbf{G}_{cc} \mathbf{v}_c + \mathbf{A}_{Ec} \mathbf{i}_E = \mathbf{t}_{*c}. \quad (3.59)$$

Înlocuind acum expresia (3.54), rezultă că potențialele necunoscute se determină prin rezolvarea sistemului

$$\mathbf{G}_{**} \mathbf{v}_* = \mathbf{t}_* + \mathbf{G}_{*c} (\mathbf{A}_{Ec}^T)^{-1} \mathbf{e}_E, \quad (3.60)$$

care este un sistem de dimensiune $N - 1 - L_E$, cu matricea coeficienților simetrică și pozitiv definită, urmând ca valorile curenților prin sursele ideale de tensiune să se determine din a două relație matriceală.

Algoritmul corespunzător acestui caz este ușor de conceput dacă el este orientat pe asamblarea matricelor și partiționarea lor.

3.4 Tratarea surselor comandate liniar

Un circuit rezistiv liniar poate conține și surse comandate liniar. Există patru feluri de surse comandate (fig. 3.5), notate pe scurt "SUCU" (sursă de tensiune comandată în tensiune), "SUCI" (sursă de tensiune comandată în curent, "SICU" (sursă de curent comandată în tensiune) și "SICI" (sursă de curent comandată în curent).

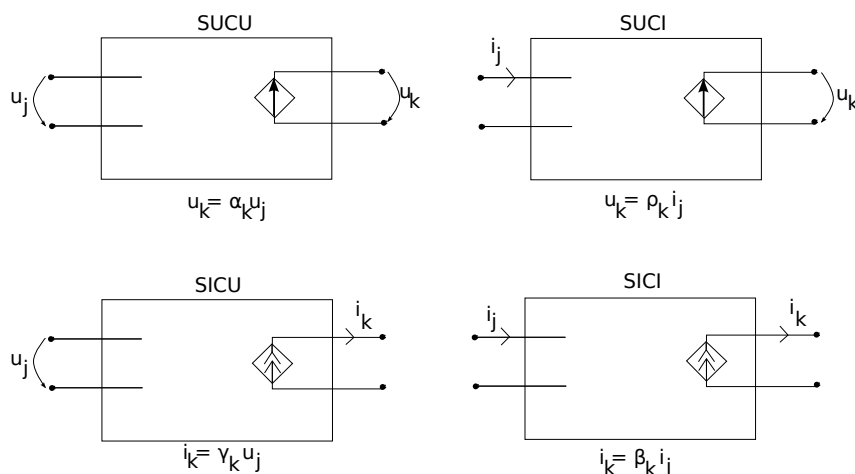


Figura 3.5: Tipuri de surse comandate liniar.

Vom sugera pașii care trebuie făcuți în asamblarea unui sistem de ecuații în cazul în care pe lângă laturi de tip REJ și E există numai laturi de tip SUCU, presupuse a fi numerotate la sfârșit. Tratarea celorlalte tipuri de surse comandate se face într-o manieră similară. Cu un raționament similar celui din paragraful anterior, ecuațiile matriceale sunt:

- Kirchhoff I:

$$\mathbf{A}_{REJ} \mathbf{i}_{REJ} + \mathbf{A}_E \mathbf{i}_E + \mathbf{A}_{SUCU} \mathbf{i}_{SUCU} = \mathbf{0}, \quad (3.61)$$

- Kirchhoff II:

$$\mathbf{u}_{REJ} = \mathbf{A}_{REJ}^T \mathbf{v}, \quad (3.62)$$

$$\mathbf{u}_E = \mathbf{A}_E^T \mathbf{v}, \quad (3.63)$$

$$\mathbf{u}_{SUCU} = \mathbf{A}_{SUCU}^T \mathbf{v}, \quad (3.64)$$

- relații constitutive:

$$\mathbf{u}_{REJ} = \mathbf{R}(\mathbf{i}_{REJ} - \mathbf{J}) - \mathbf{e}_{REJ}, \quad (3.65)$$

$$\mathbf{u}_E = -\mathbf{e}_E, \quad (3.66)$$

$$\mathbf{u}_{SUCU} = \alpha \mathbf{S}_{SUCU} \mathbf{v}, \quad (3.67)$$

unde α este o matrice diagonală, de dimensiune egală cu numărul de surse comandate, având pe diagonală parametrii surselor comandate, iar \mathbf{S}_{SUCU} este o matrice de selecție (topologică), de dimensiune $L_{SUCU} \times (N-1)$ care selectează pentru fiecare latură de tip SUCU, perechea de noduri care determină tensiunea de comandă.

În acest caz, sistemul de rezolvat va avea un număr de $N-1 + L_E + L_{SUCU}$ necunoscute, reprezentând potențialele nodurilor, curenții prin sursele ideale de tensiune și curenții prin sursele comandate:

$$\mathbf{x} = \begin{bmatrix} \mathbf{v} \\ \mathbf{i}_E \\ \mathbf{i}_{SUCU} \end{bmatrix}. \quad (3.68)$$

Sistemul de rezolvat

$$\mathbf{M}\mathbf{x} = \mathbf{p} \quad (3.69)$$

are matricea coeficienților

$$\mathbf{M} = \begin{bmatrix} \mathbf{A}_{REJ}\mathbf{R}^{-1}\mathbf{A}_{REJ}^T & \mathbf{A}_E & \mathbf{A}_{SUCU} \\ \mathbf{A}_E^T & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{SUCU}^T - \alpha\mathbf{S}_{SUCU} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (3.70)$$

iar vectorul termenilor liberi

$$\mathbf{p} = \begin{bmatrix} -\mathbf{A}_{REJ}\mathbf{R}^{-1}\mathbf{e}_{REJ} - \mathbf{A}_{REJ}\mathbf{J} \\ -\mathbf{e}_E \\ \mathbf{0} \end{bmatrix}. \quad (3.71)$$

În concluzie, sursele comandate complică problema. Nu numai că relațiile matematice devin mai complicate, conducând și la creșterea complexității algoritmului, dar sistemul de rezolvat are o matrice a coeficienților care nu mai este simetrică.

Capitolul 4

Interpolarea funcțiilor

Formularea celor mai multe probleme de inginerie poate fi scrisă formal

$$\mathbf{y} = f(\mathbf{x}), \quad (4.1)$$

unde \mathbf{x} sunt datele problemei, care sunt din punct de vedere matematic parametri independenți, iar \mathbf{y} sunt mărimile de interes ce se doresc a fi estimate. În aplicațiile reale f nu este definită explicit în funcție de datele problemei. În notația de mai sus, f poate reprezenta de exemplu chiar un proces de măsurare a mărimilor \mathbf{y} pentru o anumită stare sau configurație complet caracterizată de mărimile \mathbf{x} , sau f poate reprezenta programe software complicate, capabile să analizeze configurația caracterizată complet de datele \mathbf{x} și să calculeze printr-un algoritm de postprocesare mărimile \mathbf{y} .

Se ajunge astfel de multe ori ca, pentru un set de date (obținute de exemplu prin măsurători), să se dorească găsirea unei funcții care să aproximeze aceste date. Dacă se cunoaște un astfel de set de date se spune că *funcția este reprezentată prin date*. De exemplu, dacă setul de date este notat $(\mathbf{x}_k, \mathbf{y}_k)$, $k = 0, \dots, n$, iar $\mathbf{y}_k = f(\mathbf{x}_k)$, atunci se dorește găsirea unei expresii analitice pentru o funcție g care să aproximeze aceste date adică $g(\mathbf{x}_k) \approx \mathbf{y}_k$ sau chiar $g(\mathbf{x}_k) = \mathbf{y}_k$. În cele ce urmează, mulțimea (setul) de date \mathbf{x}_k îl vom numi și rețea (grid) de discretizare.

Pe scurt, *interpolarea unui set de date înseamnă găsirea unei funcții care să treacă prin punctele mulțimii de date* (fig.4.1). Dacă setul de date are foarte multe puncte, rezultatul interpolării poate reprezenta o funcție cu multe oscilații (fig. 4.1 - dreapta), care se datorează nu vreunui fenomen fizic, ci faptului că mulțimea de date inițiale este afectată de erori inerente, datorate de exemplu procedurii de măsurare. În acest caz, dacă se știe că funcția este netedă, se preferă *aproximarea (regresia) setului de date, adică găsirea unei funcții care să treacă printre punctele mulțimii de date* (fig.4.2), abordare ce are avantajul că diminuează erorile de măsurare din rezultatul final.

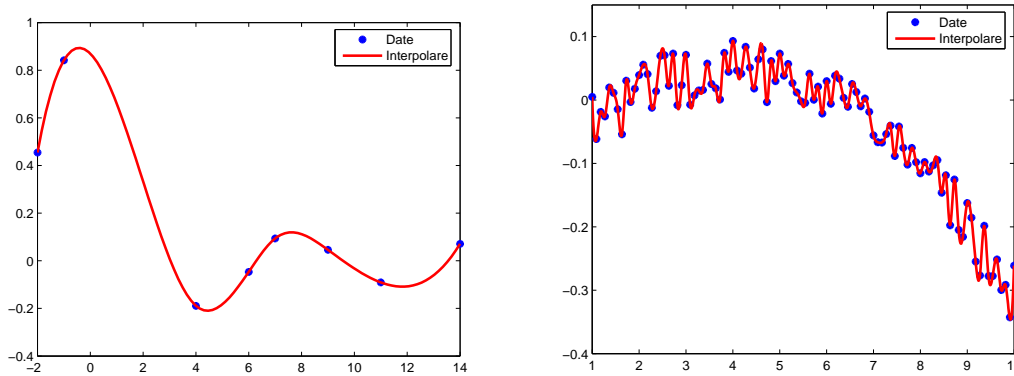


Figura 4.1: Interpolarea unui set de date. În cazul în care setul de date are foarte multe valori, interpolarea poate genera oscilații nedorite.

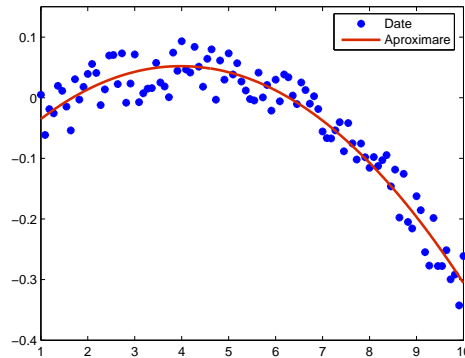


Figura 4.2: Aproximarea unui set de date.

Setul de date nu este neapărat obținut numai prin măsurători. Să presupunem că există un software capabil să calculeze $f(\mathbf{x})$ pentru orice \mathbf{x} dorit. În acest caz se spune că *funcția este reprezentată prin cod*. Dacă efortul de evaluare al funcției este mare, atunci se preferă evaluarea numerică doar în câteva puncte relevante, obținându-se astfel un set finit de date, iar pentru acesta se determină o aproximare sau interpolare g ce va fi folosită pentru estimări ulterioare ale funcției f .

Până acum nu am precizat nimic despre domeniul de definiție și domeniul de valori al funcției f . Cazul cel mai simplu este cel în care datele sunt mărimi scalare, reale, iar rezultatul evaluării funcției este tot un număr real, adică

$$f, g : [a, b] \rightarrow \mathbb{R}. \quad (4.2)$$

Acest caz este numit *cazul scalar unidimensional* (1D). Cazul vectorial unidimensional, în care valorile funcției sunt într-un spațiu multidimensional $f : [a, b] \rightarrow \mathbb{R}^m$, $m > 1$ se poate rezolva pe componente, reducându-se la m interpolări/aproximări 1D.

Problema se complică dacă domeniul de definiție are mai multe dimensiuni. Cazul

$$f, g : [a, b] \times [c, d] \rightarrow \mathbb{R} \quad (4.3)$$

îl vom numi *cazul scalar bidimensional* (2D). El este particularizarea cazului scalar n -dimensional (n D) $f, g : \mathcal{D} \subset \mathbb{R}^n \rightarrow \mathbb{R}$. Cazul cel mai general este $f, g : \mathcal{D} \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ care se reduce însă la m situații de tip n D.

Ideile metodelor de interpolare/aproximare pot fi cel mai ușor explicate pe cazul 1D, de aceea paragrafele vor presupune că are loc (4.2).

4.1 Distanța dintre două funcții

Evident se dorește ca funcția $g : [a, b] \rightarrow \mathbb{R}$ să aproximeze/interpoleze cât mai bine funcția $f : [a, b] \rightarrow \mathbb{R}$. Din punct de vedere matematic, aceasta înseamnă că distanța dintre cele două funcții, calculată ca norma diferenței lor

$$d(f, g) = \|f - g\| \quad (4.4)$$

să fie cât mai mică. Există mai multe procedee de definire a normei din relația (4.4).

O primă posibilitate este aceea de a estima aria cuprinsă între graficele celor două funcții, adică de a calcula distanța dintre funcții ca

$$d_1(f, g) = \frac{1}{b-a} \int_a^b |f(x) - g(x)| dx. \quad (4.5)$$

Dezavantajul acestei abordări este acela că, local, pot exista diferențe foarte mari între cele două funcții, fără ca valoarea rezultatului formulei (4.5) să reflecte acest lucru¹.

O a doua variantă, care nu elimină însă dezavantajul menționat mai sus, calculează abaterea pătratică medie dintre cele două funcții:

$$d_2(f, g) = \sqrt{\frac{1}{b-a} \int_a^b (f(x) - g(x))^2 dx}. \quad (4.6)$$

Norma (4.6) este norma L_2 din spațiul funcțiilor de pătrat integrabil și generalizează norma euclidiană.

O a treia posibilitate este de a evalua norma Chebyshev, care estimează abaterea maximă dintre cele două funcții:

$$d_3(f, g) = \max_{x \in [a, b]} |f(x) - g(x)|. \quad (4.7)$$

¹Această afirmație se bazează pe proprietatea că, dacă valoarea unei funcții se modifică într-o mulțime numărabilă de puncte, atunci valoarea integralei ei nu se modifică.

Din punctul de vedere al acurateții rezultatului aproximării, această normă este cea mai avantajoasă deoarece nu suferă de dezavantajul normelor ce folosesc integrale. Valoarea ei furnizează informații despre diferența locală maximă dintre cele două funcții.

Oricare din normele (4.5), (4.6), (4.7) presupun evaluarea funcției f în toate punctele domeniului de definiție, ceea ce nu este posibil într-un algoritm. Funcția f este cunoscută sau cel mult poate fi evaluată numai într-un număr discret de puncte x_k , $k = 0, \dots, n$. De aceea, în implementarea numerică a procedurilor de interpolare sau aproximare se vor folosi normele discrete:

$$d_{1d}(f, g) = \sum_{k=0}^n |g(x_k) - f(x_k)|, \quad (4.8)$$

$$d_{2d}(f, g) = \sqrt{\sum_{k=0}^n (g(x_k) - f(x_k))^2}, \quad (4.9)$$

$$d_{3d}(f, g) = \max_{k=0, n} |g(x_k) - f(x_k)|. \quad (4.10)$$

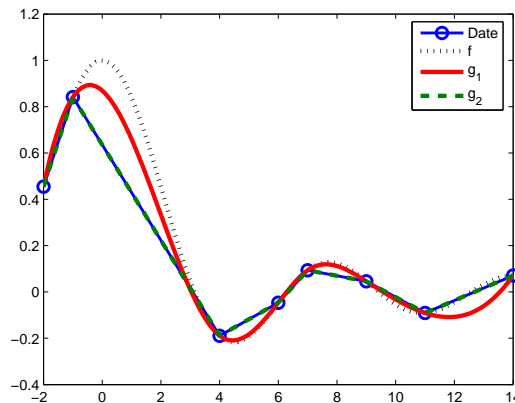


Figura 4.3: Oricare normă discretă este zero pentru diferența $f - g_1$ sau $f - g_2$, unde f este funcția adevărată, iar g_1, g_2 sunt posibile interpolări.

Normele discrete au avantajul că pot fi evaluate cu ușurință. Dezavantajul lor este acela că se pierde posibilitatea evaluării acurateții interpolării/aproximării între nodurile rețelei de interpolare. Mai mult, pentru orice funcție f există o infinitate de funcții g care nu sunt identice cu f și pentru care oricare din normele discrete de mai sus sunt zero dacă $g(x_k) = f(x_k)$ pentru toate valorile $k = 0, \dots, n$ (fig.4.3). Deci, problema interpolării/aproximării devine în acest moment prost formulată. Corectarea formulării se face impunând condiții suplimentare pentru funcția g , așa cum se prezintă în paragraful următor.

4.2 Formularea problemei interpolării

Problema interpolării este aceea de a găsi o funcție g pentru care distanța discretă față de funcția f , cunoscută într-un număr finit de puncte $f(x_j) = y_j$, este zero. Această condiție este echivalentă cu a impune egalitatea valorilor funcțiilor f și g în toate punctele (nodurile) rețelei de discretizare:

$$g(x_j) = f(x_j), \quad j = 0, \dots, n, \quad (4.11)$$

echivalent cu a scrie

$$g(x_j) = y_j, \quad j = 0, \dots, n. \quad (4.12)$$

Relațiile (4.12) se numesc *condiții de interpolare*. Pentru a face ca problema să fie bine formulată matematic (soluția să existe și să fie unică) funcția g se caută în spațiul polinoamelor generalizate, adică se caută de forma unei combinații liniare de m funcții φ_k , $k = 1, \dots, m$ numite *funcții de bază*:

$$g(x) = \sum_{k=0}^m c_k \varphi_k(x). \quad (4.13)$$

Funcțiile de bază pot fi de diferite tipuri și ele se aleg înainte de rezolvarea propriu-zisă a problemei interpolării. De exemplu, dacă se știe că funcția de interpolat are un caracter periodic cu perioada 2π atunci o alegere potrivită a funcțiilor de bază este $\varphi_0(x) = 1$, $\varphi_1(x) = \sin x$, $\varphi_2(x) = \cos x$, $\varphi_3(x) = \sin(2x)$, etc. Alegerea $\varphi_k(x) = x^k$ corespunde cazului polinoamelor scrise în mod clasic.

Astfel, problema interpolării se reduce la problema determinării celor m coeficienți c_k .

Condițiile de interpolare (4.12) scrise pentru polinoame de interpolare de tipul (4.13) devin

$$\sum_{k=0}^m c_k \varphi_k(x_j) = y_j, \quad j = 0, \dots, n, \quad (4.14)$$

ceea ce înseamnă că valorile coeficienților c_k sunt soluțiile unui sistem de ecuații algebrice liniare cu $n+1$ ecuații și $m+1$ necunoscute. Pentru ca sistemul (4.14) să fie bine formulat matematic vom impune ca numărul de funcții de bază să fie egal cu numărul de puncte al rețelei de discretizare

$$m = n \quad (4.15)$$

și determinantul matricei coeficienților să fie nenul:

$$\Delta = \begin{vmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \cdots & \varphi_n(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \cdots & \varphi_n(x_1) \\ \cdots & \cdots & \cdots & \cdots \\ \varphi_0(x_n) & \varphi_1(x_n) & \cdots & \varphi_n(x_n) \end{vmatrix} \neq 0. \quad (4.16)$$

Determinantul Δ este nenul dacă și numai dacă punctele x_j sunt distincte și funcțiile de bază φ_k sunt liniar independente.

Rezumând cele de mai sus, problema interpolării funcțiilor se formulează astfel.

Date:

- un tabel de valori (x_k, y_k) , $k = 0, \dots, n$, unde punctele rețelei de discretizare x_k sunt distincte două câte două;
- $n + 1$ funcții de bază liniar independente $\varphi_k(x)$, $k = 0, \dots, n$.

Se cer:

- coeficienții c_k , $k = 0, \dots, n$ pentru care sunt satisfăcute condițiile de interpolare $g(x_j) = y_j$, $j = 0, \dots, n$ unde $g(x) = \sum_{k=0}^n c_k \varphi_k(x)$ este funcția (în particular polinomul) de interpolare al datelor din tabel.

4.3 Metode de interpolare globală

Metodele de interpolare globală sunt metodele în care funcțiile de bază se definesc compact, printr-o singură expresie pe întreg domeniul de definiție al funcției de interpolat. În aceste metode, gradul polinomului de interpolare este în strânsă legătură cu numărul de puncte din rețeaua de discretizare, mai precis gradul polinomului de interpolare este cu unu mai mic decât numărul de puncte din tabelul de date (de exemplu prin două puncte trece în mod univoc o dreaptă, prin trei puncte o parabolă, etc.). În funcție de cum se aleg funcțiile de bază se obțin diferite metode de interpolare globală.

4.3.1 Metoda clasică

În metoda clasică funcțiile de bază sunt polinoamele algebrice $\varphi_k(x) = x^k$, $k = 0, \dots, n$. Atunci polinomul de interpolare (4.13), cu condiția (4.15) devine

$$g(x) = \sum_{k=0}^n c_k x^k = c_0 + c_1 x + c_2 x^2 + \dots + c_n x^n, \quad (4.17)$$

iar condițiile de interpolare (4.12) conduc la rezolvarea sistemului de ecuații algebrice liniare

$$\begin{cases} c_0 + c_1 x_0 + c_2 x_0^2 + \dots + c_n x_0^n = y_0 \\ c_0 + c_1 x_1 + c_2 x_1^2 + \dots + c_n x_1^n = y_1 \\ \dots \\ c_0 + c_1 x_n + c_2 x_n^2 + \dots + c_n x_n^n = y_n \end{cases} \quad (4.18)$$

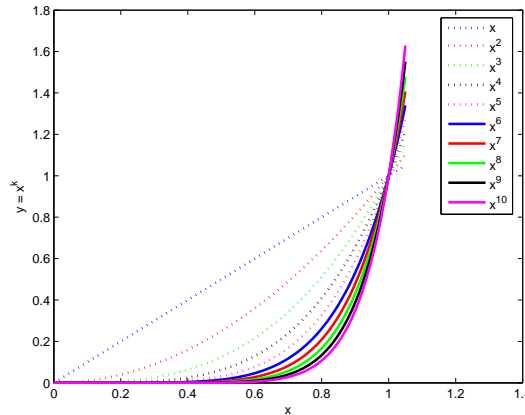


Figura 4.4: Funcțiile de bază x^k devin aproape liniar dependente pentru valori ale lui $k > 5$.

pentru determinarea celor $n + 1$ coeficienți.

Calculul coeficienților polinomului de interpolare se numește *etapa de pregătire* iar evaluarea propriu-zisă a polinomului de interpolare se numește *etapa de evaluare*.

În metoda clasică, etapa de pregătire constă în asamblarea și rezolvarea sistemului (4.18). Această etapă necesită un efort mare de calcul, corespunzător rezolvării unui sistem cu matrice plină, de ordinul $O(2n^3/3)$. Etapa de evaluare presupune evaluarea expresiei (4.17), deci un efort de calcul de ordinul $O(2n)$.

Dezavantajul major al acestei metode nu stă în primul rând în efortul mare de calcul necesar etapei de pregătire ci, mai ales în aceea că pentru valori mari ale lui n matricea coeficienților sistemului (4.18) este slab condiționată, pentru că funcțiile de bază devin aproape liniar independente (fig.4.4). Din acest motiv această metodă nu se recomandă să fie aplicată pentru polinoame de interpolare de grad mai mare decât 5.

O metodă mai eficientă de interpolare trebuie să folosească funcții de bază care să fie cât mai deosebite între ele. Mai mult, ar fi convenabil ca sistemul de ecuații de rezolvat în etapa de pregătire să fie cât mai simplu de rezolvat. Așa se întâmplă de exemplu în metoda Lagrange.

4.3.2 Metoda Lagrange

În metoda Lagrange funcțiile de bază sunt polinoamele Lagrange

$$\varphi_k(x) = l_k(x) = \frac{\prod_{i=0, i \neq k}^n (x - x_i)}{\prod_{i=0, i \neq k}^n (x_k - x_i)}, \quad (4.19)$$

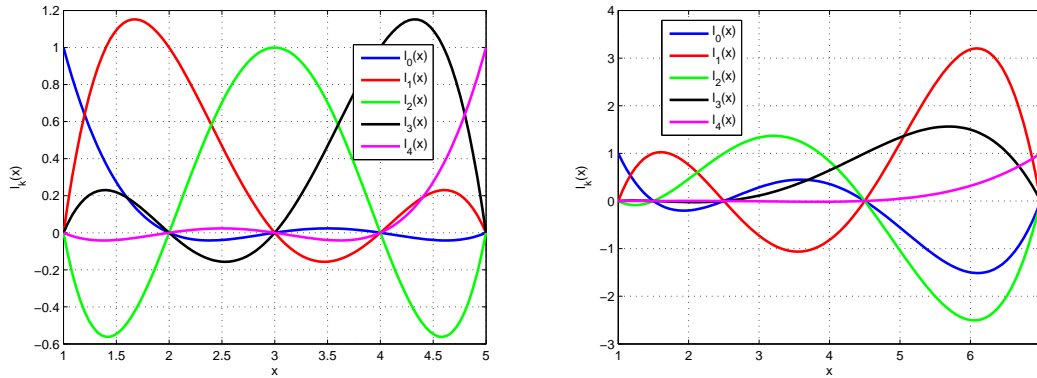


Figura 4.5: Funcțiile Lagrange pentru o rețea de discretizare uniformă (stânga), respectiv neuniformă (dreapta).

iar polinomul de interpolare este

$$g(x) = \sum_{k=0}^m c_k l_k(x). \quad (4.20)$$

Fiecare polinom Lagrange $l_k(x)$ poate fi considerat asociat unui nod k , și are proprietatea că ia valoarea 1 când este evaluat în acel nod și valoarea 0 când este evaluat în toate celelalte noduri:

$$l_k(x_j) = \begin{cases} 1 & \text{dacă } j = k, \\ 0 & \text{dacă } j \neq k. \end{cases} \quad (4.21)$$

Figura 4.5 ilustrează polinoamele Lagrange pentru o rețea de discretizare cu 5 puncte, uniformă și, respectiv, neuniformă. Între nodurile rețelei de discretizare polinoamele Lagrange pot avea valori oricât de mari.

Cu această alegere a funcțiilor de bază, condițiile de interpolare $g(x_j) = y_j$, $j = 0, \dots, n$ devin

$$\sum_{k=0}^m c_k l_k(x_j) = y_j, \quad (4.22)$$

de unde, aplicând (4.21)

$$c_j = y_j, \quad j = 0, \dots, n. \quad (4.23)$$

Practic, sistemul algebric de rezolvat pentru determinarea coeficienților este unul în care matricea coeficienților este matricea unitate, și etapa de pregătire așa cum a fost definită anterior dispare.

Expresia polinomului Lagrange poate fi scrisă în mod compact ca

$$g(x) = \sum_{k=0}^n y_k l_k(x) = \sum_{k=0}^n y_k \frac{\prod_{i=0, i \neq k}^n (x - x_i)}{\prod_{i=0, i \neq k}^n (x_k - x_i)}. \quad (4.24)$$

În particular, polinomul de interpolare linear pentru o rețea de discretizare cu două puncte este

$$g(x) = y_0 \frac{x - x_1}{x_0 - x_1} + y_1 \frac{x - x_0}{x_1 - x_0}, \quad (4.25)$$

iar parabola interpolatoare pentru o rețea de discretizare cu trei puncte este

$$g(x) = y_0 \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + y_1 \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} + y_2 \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}. \quad (4.26)$$

Implementarea metodei Lagrange se poate face în două variante. Prima din ele este o variantă fără pregătire, așa cum rezultă în mod natural din prezentarea de mai sus. Pentru orice valoare nouă din domeniul de definiție, evaluarea polinomului de interpolare se face cu expresia (4.24). Efortul de calcul poate fi estimat direct inspectând expresia (4.24) și este de ordinul $T_e = O(4n^2)$ pentru fiecare evaluare, rezultând un efort de calcul total pentru evaluarea în m puncte de $T_{L-fp} = O(4mn^2)$.

O variantă de implementare mai eficientă se bazează pe rescrierea relației (4.24) prin scoaterea forțată a factorului comun

$$p = \prod_{k=0}^n (x - x_k), \quad (4.27)$$

polinomul de interpolare fiind

$$g(x) = p \sum_{k=0}^n \frac{\alpha_k}{x - x_k}, \quad (4.28)$$

unde coeficienții notați

$$\alpha_k = \frac{y_k}{\prod_{j=0, j \neq k}^n (x_k - x_j)} \quad (4.29)$$

vor fi calculați în etapa de pregătire. Pseudocodul acestei variante este următorul.

procedură Lagrange_pregătire(n, x, y, α)

; pregătește coeficienții din metoda Lagrange

; declarații - parametri de intrare

întreg n ; dimensiunea problemei - numărul de intervale

tablou real $x[n], y[n]$; tabelul de valori, indici de la zero

; declarații - parametri de ieșire

tablou real $\alpha[n]$; coeficienții

pentru $k = 0, n$

$\alpha_k = y_k$

pentru $j = 0, n$

dacă $j \neq k$ atunci $\alpha_k = \alpha_k / (x_k - x_j)$

```

•
•
retur

funcție Lagrange_evaluare( $n, x, y, \alpha, xcrt$ )
; evaluează polinomul de interpolare Lagrange în punctul  $xcrt$ 
; declarații - parametri de intrare
întreg  $n$  ; dimensiunea problemei - numărul de intervale
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la zero
tablou real  $\alpha[n]$  ; coeficienții
real  $xcrt$  ; punctul de evaluat
; alte declarații
real  $p, s$ 
 $p = 1$ 
pentru  $k = 0, n$ 
    ... ; vezi comentariul de mai jos
     $p = p * (xcrt - x_k)$ 
•
 $s = 0$ 
pentru  $k = 0, n$ 
     $s = s + \alpha_k / (xcrt - x_k)$ 
•
întoarce  $s * p$ 

```

În acest caz etapa de pregătire are ordinul de complexitate $T_p = O(2n^2)$, iar etapa de evaluare ordinul $T_e = O(5n)$, rezultând un efort de calcul total $T_{L-cp} = O(2n^2 + 5mn)$.

Așa cum este scris, pseudocodul etapei de evaluare eșuează dacă punctul de evaluare coincide cu unul din punctele din tabel, din cauza împărțirii la zero. Un algoritm înțelept tratează și astfel de cazuri particulare. În acest caz se poate adăuga următoarea linie de cod în spațiul lăsat liber al pseudocodului de mai sus:

dacă $|xcrt - x_k| < \text{zeroul_mașinii}()$ atunci întoarce y_k

Reamintim că testarea egalității numerelor reale nu se poate face decât cel mult până la zeroul mașini (vezi discuția și algoritmul de la pagina 34).

4.3.3 Eroarea de trunchiere

În cazul în care funcția de interpolat f este continuă și derivabilă de un număr finit de ori, se poate obține o margine a erorii de interpolare

$$e(x) = f(x) - g(x), \quad (4.30)$$

care poate fi privită ca o eroare de trunchiere deoarece interpolarea caută o aproximare într-un spațiu finit dimensional de funcții.

Datorită condițiilor de interpolare (4.11), eroarea de interpolare este nulă în nodurile de interpolare

$$e(x_k) = 0, \quad k = 0, \dots, n \quad (4.31)$$

și, în consecință, putem considera o aproximare a erorii de interpolare de tipul

$$E(x) = a \prod_{k=0}^n (x - x_k), \quad (4.32)$$

unde a ține seama de abaterile dintre funcțiile f și g . Notăm diferența între eroarea de interpolare $e(x)$ și aproximarea ei prin polinomul $E(x)$ cu:

$$h(x) = e(x) - E(x) = f(x) - g(x) - a \prod_{k=0}^n (x - x_k). \quad (4.33)$$

Este evident că $h(x)$ se anulează în toate nodurile rețelei de interpolare

$$h(x_k) = 0, \quad k = 0, \dots, n. \quad (4.34)$$

Este interesant că, alegând convenabil constanta a , de exemplu

$$a = \frac{f(\alpha) - g(\alpha)}{\prod_{k=0}^n (\alpha - x_k)}, \quad (4.35)$$

unde α este ales arbitrar în intervalul $[a, b]$, $h(x)$ se mai anulează într-un punct și anume

$$h(\alpha) = 0. \quad (4.36)$$

În consecință, h are $n + 2$ rădăcini: α, x_0, \dots, x_n . Dacă presupunem f continuă și derivabilă, atunci și h este continuă și derivabilă, iar funcția

$$h'(x) = f'(x) - g'(x) - aP'_n(x), \quad (4.37)$$

unde $P_n(x) = \prod_{k=0}^n (\alpha - x_k)$ are, conform teoremei lui Rolle, $n + 1$ rădăcini. Dacă h' este continuă și derivabilă, atunci h'' are n rădăcini. Continuând raționamentul, va rezulta că $h^{(n+1)}$ va avea o singură rădăcină $\xi \in [a, b]$:

$$h^{(n+1)}(\xi) = 0. \quad (4.38)$$

Dar deoarece $g^{(n+1)}(x) = 0$ rezultă că

$$h^{(n+1)}(x) = f^{(n+1)}(x) - a(n+1)!, \quad (4.39)$$

de unde

$$a = \frac{f^{(n+1)}(\xi)}{(n+1)!}. \quad (4.40)$$

Deoarece $h(\alpha) = 0$, rezultă că

$$e(\alpha) = E(\alpha) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{k=0}^n (\alpha - x_k). \quad (4.41)$$

Dar α a fost ales arbitrar, astfel încât, în relația de mai sus se poate substitui α cu x .

În concluzie, pentru orice x există un $\xi \in [a, b]$ astfel încât eroarea de interpolare este

$$e(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{k=0}^n (x - x_k). \quad (4.42)$$

Să considerăm o margine M_{n+1} a derivatei de ordinul $n+1$ a funcției f :

$$|f^{(n+1)}(x)| \leq M_{n+1}, \quad (\forall)x \in [a, b]. \quad (4.43)$$

Rezultă atunci că eroarea de interpolare este mărginită de

$$|e(x)| \leq \frac{M_{n+1}}{(n+1)!} \prod_{k=0}^n |x - x_k|. \quad (4.44)$$

Eroarea de interpolare depinde deci de marginea derivatei de un ordin egal cu numărul de puncte de tabel. Această afirmație este ilustrată în fig.4.6.

4.3.4 Metoda Newton

În metoda Newton funcțiile de bază sunt alese astfel încât matricea sistemului de rezolvat în etapa de pregătire a coeficienților să fie triunghiular inferioară:

$$\begin{aligned} \varphi_0(x) &= 1, \\ \varphi_1(x) &= (x - x_0), \\ \varphi_2(x) &= (x - x_0)(x - x_1), \\ &\dots \\ \varphi_n(x) &= (x - x_0)(x - x_1) \cdots (x - x_{n-1}). \end{aligned} \quad (4.45)$$

Polinomul de interpolare (4.13) este atunci

$$g(x) = c_0 + c_1(x - x_0) + c_2(x - x_0)(x - x_1) + \cdots + c_n(x - x_0)(x - x_1) \cdots (x - x_{n-1}), \quad (4.46)$$

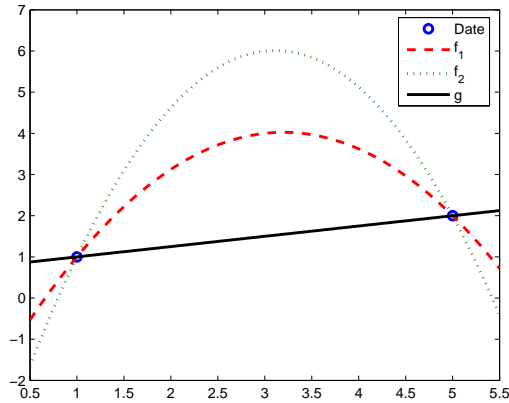


Figura 4.6: Ilustrarea formulei (4.44) în cazul unui tabel de valori cu două puncte. Presupunând că funcția adevărată este o parabolă, atunci eroarea de interpolare este mai mică pentru parabola f_1 , care are derivata de ordinul doi (curbura) mai mică.

iar condițiile de interpolare (4.14) devin

$$\begin{cases} c_0 = y_0, \\ c_0 + c_1(x_1 - x_0) = y_1, \\ c_0 + c_1(x_2 - x_0) + c_2(x_2 - x_0)(x_2 - x_1) = y_2, \\ \dots \\ c_0 + c_1(x_n - x_0) + c_2(x_n - x_0)(x_n - x_1) + \dots + c_n(x_n - x_0)(x_n - x_1) \dots (x_n - x_{n-1}) = y_n. \end{cases} \quad (4.47)$$

Sistemul de rezolvat în etapa de pregătire are structura inferior triunghiulară și se va rezolva prin substituție progresivă. Coeficienții reprezintă diferențe divizate ale funcției față de submulțimi ale nodurilor de discretizare:

$$\begin{cases} c_0 = y_0 = f[x_0] \\ c_1 = (y_1 - y_0)/(x_1 - x_0) = f[x_0, x_1] \\ c_2 = [y_2 - y_0 - (y_1 - y_0)/(x_1 - x_0)(x_2 - x_0)]/(x_2 - x_0)/(x_2 - x_1) = f[x_0, x_1, x_2], \\ \vdots \\ c_n = \dots = f[x_0, x_1, x_2, \dots, x_n]. \end{cases} \quad (4.48)$$

În general, *diferența divizată față de o submulțime a nodurilor rețelei de discretizare* se definește în mod recursiv astfel:

$$f[x_i, x_{i+1}, \dots, x_{i+m}] = \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+m}] - f[x_i, x_{i+1}, \dots, x_{i+m-1}]}{x_{i+m} - x_i}, \quad (4.49)$$

unde diferența divizată față de o submulțime cu un punct este valoarea funcției în acel punct $f[x_i] = f(x_i)$.

Următoarele două proprietăți ale diferențelor divizate sunt importante în contextul interpolării cu metoda Newton:

- Variabilele independente pot fi permutate într-o diferență divizată fără ca valoarea funcției să se modifice. Acest lucru este evident pentru diferența divizată față de o submulțime cu două noduri:

$$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0} = \frac{f[x_0] - f[x_1]}{x_0 - x_1} = f[x_1, x_0] \quad (4.50)$$

și, pe baza relației de recurență (4.49), această proprietate poate fi generalizată pentru o mulțime oarecare de noduri.

- Diferența divizată față de o submulțime cu două puncte identice este derivata funcției:

$$f[x_0, x_0] = \lim_{x_1 \rightarrow x_0} \frac{f(x_1) - f(x_0)}{x_1 - x_0} = f'(x_0). \quad (4.51)$$

Legătura dintre diferențele divizate ale unei funcții și derivatele sale poate fi generalizată, așa cum se va vedea în cele ce urmează.

Polinomul de interpolare (4.46) cu coeficienții dați de (4.48) este

$$\begin{aligned} g(x) &= f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \cdots \\ &+ \cdots + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1) \cdots (x - x_{n-1}). \end{aligned} \quad (4.52)$$

Dacă nodurile de interpolare tind toate către același punct x_0 , atunci polinomul de interpolare $g(x)$ tinde către seria Taylor a funcției f în x_0 :

$$f[x_0] + f[x_0, x_0](x - x_0) + f[x_0, x_0, x_0](x - x_0)^2 + \cdots + \underbrace{f[x_0, x_0, \dots, x_0]}_{n+1}(x - x_0)^n. \quad (4.53)$$

În consecință, termeni succesivi ai polinomului interpolării Newton aproximează termeni corespunzători din dezvoltarea în serie Taylor. Acest lucru conferă un avantaj deosebit metodei Newton față de metoda Lagrange, pentru că ea permite controlul erorii de trunchie-re, efortul de calcul putând fi adaptat preciziei impuse soluției. Mai mult, atunci când se adaugă un punct suplimentar în rețeaua de interpolare, se poate porni de la interpolarea cu un grad mai scăzut și trebuie doar adăugat un singur termen în sumă, nefiind necesară refacerea totală a calculelor ci doar evaluarea unui singur termen suplimentar.

Dacă $g(x)$ reprezintă interpolarea funcției $f(x)$ pe o rețea cu nodurile x_0, x_1, \dots, x_n , atunci la adăugarea unui nod nou x_{n+1} (care poate avea orice poziție față de punctele anterioare), noul polinom de interpolare va fi

$$h(x) = g(x) + f[x_0, x_1, \dots, x_n, x_{n+1}](x - x_0)(x - x_1) \cdots (x - x_n). \quad (4.54)$$

Condiția de interpolare pentru noul punct este $h(x_{n+1}) = f(x_{n+1})$, deci

$$f(x_{n+1}) = g(x_{n+1}) + f[x_0, x_1, \dots, x_n, x_{n+1}](x_{n+1} - x_0)(x_{n+1} - x_1) \cdots (x_{n+1} - x_n). \quad (4.55)$$

Deoarece x_{n+1} poate fi ales arbitrar, în (4.55) putem înlocui x_{n+1} cu x , deci

$$f(x) = g(x) + f[x_0, x_1, \dots, x_n, x](x - x_0)(x - x_1) \cdots (x - x_n), \quad (4.56)$$

de unde rezultă că eroarea de trunchiere făcută la interpolarea funcției f pe rețeaua de noduri x_0, x_1, \dots, x_n este

$$e(x) = f(x) - g(x) = f[x_0, x_1, \dots, x_n, x] \prod_{k=0}^n (x - x_k). \quad (4.57)$$

Am demonstrat însă și că eroarea de interpolare este dată de formula (4.42):

$$e(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{k=0}^n (x - x_k). \quad (4.58)$$

Cum pentru un tabel de valori dat polinomul de interpolare global este unic, rezultă că cele două erori de trunchiere (4.57) și (4.58) sunt același, de unde rezultă că, pentru orice diviziune, există ξ astfel încât

$$f[x_0, x_1, \dots, x_n, x] = \frac{f^{(n+1)}(\xi)}{(n+1)!}. \quad (4.59)$$

Iată deci că, diferențele divizate pe o rețea cu mai mult de două noduri reprezintă aproximări ale derivatelor de ordin superior, ordinul derivatei fiind cu unu mai mic decât numărul de noduri ale diviziunii. În particular, dacă în relația (4.59) toate nodurile tind către x_0 , rezultă că are loc:

$$f[\underbrace{x_0, x_0, \dots, x_0, x_0}_{n+2}] = \frac{f^{(n+1)}(x_0)}{(n+1)!}. \quad (4.60)$$

Algoritmul metodei Newton se bazează pe calculul tabelului diferențelor divizate în etapa de pregătire, calcul ilustrat de următoarea schemă:

x_0	x_1	x_2	\dots	x_{n-1}	x_n
$y_0 = f[x_0]$	$y_1 = f[x_1]$	$y_2 = f[x_2]$	\dots	$y_{n-1} = f[x_{n-1}]$	$y_n = f[x_n]$
	$f[x_0, x_1]$	$f[x_1, x_2]$	\dots	$f[x_{n-2}, x_{n-1}, x_n]$	$f[x_{n-1}, x_n]$
	$f[x_0, x_1, x_2]$		\dots		
		\vdots			
		\vdots			
			\vdots		
			$f[x_0, x_1, \dots, x_n]$		

În total sunt calculate $n(n+1)/2$ diferențe divizate, din care utile pentru polinomul de interpolare sunt doar n , și anume valorile marcate cu roșu în schema de mai sus. Pseudocodul metodei Newton este următorul.

```

procedură Newton_pregătire( $n, x, y, dd$ )
; pregătește diferențele divizate din metoda Newton
; declarații - parametri de intrare
întreg  $n$  ; dimensiunea problemei - numărul de intervale
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la zero
; declarații - parametri de ieșire
tablou real  $dd[n][n]$  ; diferențele divizate
pentru  $j = 0, n$ 
     $dd(0, j) = y(j)$ 
•
pentru  $i = 1, n$ 
    pentru  $j = 0, n - i$ 
         $dd(i, j) = (dd(i - 1, j) - dd(i - 1, j - 1)) / (x(j + i) - x(j))$ 
    •
•
retur

funcție Newton_evaluare( $n, x, y, dd, xcrt$ )
; evaluează polinomul de interpolare Newton în punctul  $xcrt$ 
; declarații - parametri de intrare
întreg  $n$  ; dimensiunea problemei - numărul de intervale
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la zero
tablou real  $dd[n][n]$  ; diferențele divizate
real  $xcrt$  ; punctul de evaluat
real  $ycrt$  ; valoarea funcției în punctul de evaluat
 $ycrt = dd(n, 0)$ 
pentru  $k = n - 1, 0, -1$ 
     $ycrt = dd(k, 0) + (xcrt - x_k) * ycrt$ 
•
întoarce  $ycrt$ 

```

Funcția de evaluare se bazează pe scrierea polinomului sub forma

$$g(x) = c_0 + (x - x_0) [c_1 + (x - x_1) [c_2 + \cdots [c_{n-1} + c_n(x - x_n)] \cdots]]. \quad (4.61)$$

În etapa de pregătire se efectuează $\sum_{i=1}^n 2(n - i) = 2n(n - 1)/2$ operații, această etapă având ordinul de complexitate $T_p = O(n^2)$. Deoarece etapa de evaluare are ordinul de complexitate $T_e = O(3n)$, rezultă un efort de calcul total în metoda Newton $T_N = O(n^2 + 3mn)$.

Rezumând, efortul de calcul pentru cele trei metode de interpolare globală discutate, este prezentat în tabelul 4.1.

Tabelul 4.1: Efortul de calcul pentru metodele de interpolare globală.

Metoda	Pregătire	Evaluare
Clasică	$O(2n^3/3)$	$O(2mn)$
Lagrange	$O(2n^2)$	$O(5mn)$
Newton	$O(n^2)$	$O(3mn)$

În concluzie, metodele clasică, Lagrange, Newton dau teoretic același rezultat pentru că polinomul de interpolare est unic. Dintre toate, metoda Newton este cea mai eficientă din punct de vedere al timpului de pregătire, cât și a celui de evaluare. Avantajul major este însă acela că metoda Newton permite controlul erorii de trunchiere. Evaluarea polinomului de interpolare cu controlul erorii de trunchiere este dată de procedura următoare:

```

procedură Newton_evaluare2( $n, x, y, dd, xcrt, ycrt, et$ )
; evaluează polinomul de interpolare Newton în punctul  $xcrt$ 
; declarații - parametri de intrare
întreg  $n$  ; dimensiunea problemei - numărul de intervale
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la zero
tablou real  $dd[n][n]$  ; diferențele divizate
real  $xcrt$  ; punctul de evaluat
real  $ycrt$  ; valoarea funcției în punctul de evaluat
real  $et$  ; eroarea de trunchiere
 $ycrt = dd(n - 1, 0)$ 
pentru  $k = n - 2, 0, -1$ 
     $ycrt = dd(k, 0) + (xcrt - x_k) * ycrt$ 
•
 $et = dd(n, 0)$ 
pentru  $k = 0, n - 1$ 
     $et = et * (xcrt - x_k)$ 
•
retur

```

4.3.5 Un exemplu simplu

Să ilustrăm cele trei metode prezentate până acum pentru tabelul de valori:

x	-1	2	4
y	-6	9	49

Deoarece tabelul are trei puncte, polinomul de interpolare globală va fi de gradul doi.

În metoda clasică, polinomul de interpolare este de forma

$$g(x) = c_0 + c_1x + c_2x^2, \quad (4.62)$$

unde cei trei coeficienți necunoscuți se determină prin impunerea condițiilor de interpolare

$$\begin{aligned} g(-1) &= 6, \\ g(2) &= 9, \\ g(4) &= 49, \end{aligned} \quad (4.63)$$

care conduc la rezolvarea sistemului algebric liniar

$$\begin{aligned} c_0 - c_1 + c_2 &= 6, \\ c_0 + 2c_1 + 4c_2 &= 9, \\ c_0 + 4c_1 + 16c_2 &= 49. \end{aligned} \quad (4.64)$$

Soluția acestui sistem este $c_0 = -7$, $c_1 = 2$, $c_2 = 3$, de unde rezultă expresia polinomului de interpolare $g(x) = 3x^2 + 2x - 7$.

În metoda Lagrange, se construiesc polinoamele de interpolare Lagrange:

$$l_0(x) = \frac{(x-2)(x-4)}{(-1-2)(-1-4)} = \frac{x^2 - 6x + 8}{15}, \quad (4.65)$$

$$l_1(x) = \frac{(x+1)(x-4)}{(2+1)(2-4)} = \frac{x^2 - 3x - 4}{-6}, \quad (4.66)$$

$$l_2(x) = \frac{(x+1)(x-2)}{(4+1)(4-2)} = \frac{x^2 - x - 2}{10}, \quad (4.67)$$

iar polinomul de interpolare globală este

$$g(x) = -6l_0(x) + 9l_1(x) + 49l_2(x) = 3x^2 + 2x - 7. \quad (4.68)$$

În metoda Newton se construiește mai întâi tabelul diferențelor divizate²:

-1	2	4
$f[x_0] = -6$	$f[x_1] = 9$	$f[x_2] = 49$
$f[x_0, x_1] = 5$		$f[x_1, x_2] = 20$
$f[x_0, x_1, x_2] = 3$		

² $f[x_0, x_1] = (9 + 6)/(2 + 1) = 5$, $f[x_1, x_2] = (49 - 9)/(4 - 2)$, $f[x_0, x_1, x_2] = (20 - 5)/(4 + 1)$

Polinomul de interpolare se calculează ușor ca

$$g(x) = -6 + 5(x + 1) + 3(x + 1)(x - 2) = 3x^2 + 2x - 7. \quad (4.69)$$

Dacă adăugăm un nou punct în acest tabel

x	-1	2	4	3
y	-6	9	49	10

putem calcula ușor polinomul de gradul trei. Tabelul diferențelor divizate se completează cu noile valori

-1	2	4	3
-6	9	49	10
5	20	39	
3	19		
4			

iar polinomul de interpolare este

$$h(x) = 3x^2 + 2x - 7 + 4(x + 1)(x - 2)(x - 4) = 4x^3 - 17x^2 + 10x + 25.$$

Termenul $e(x) = 4(x + 1)(x - 2)(x - 4)$ reprezintă eroarea de trunchiere a polinomului de interpolare de grad doi $g(x) = 3x^2 + 2x - 7$ pentru tabelul de valori ce conține patru puncte.

4.3.6 Interpolarea Chebyshev

Toate cele trei metode de interpolare polinomială globală descrise mai sus au un foarte mare dezavantaj dacă gridul de interpolare este uniform. Acest lucru a fost pus în evidență de Runge care a observat că la interpolarea funcției $f : [-5, 5] \rightarrow \mathbb{R}$, $f(x) = 1/(1 + x^2)$ pe o rețea echidistantă de puncte apar oscilații la capetele polinomului de interpolare, oscilații cu atât mai mari cu cât gradul polinomului este mai mare (fig. 4.7). *Efectul de oscilație a polinomului de interpolare între nodurile rețelei de interpolare se numește efect Runge.*

Efectul Runge se poate explica prin faptul că eroarea de trunchiere dată de (4.42) poate fi oricât de mare datorită faptului că $f^{(n+1)}(\xi)$ poate fi oricât de mare.

Efectul Runge se poate elimina dacă nodurile rețelei de interpolare se îndesesc către capetele domeniului. Oscilații minime ale polinomului de interpolare se obțin dacă plasarea nodurilor în interiorul domeniului de definiție se face în concordanță cu rădăcinile polinoamelor Chebyshev. În intervalul $[-1, 1]$, acestea sunt

$$x_i = \cos\left(\frac{2i + 1}{2(n + 1)}\pi\right), \quad i = 0, \dots, n, \quad (4.70)$$

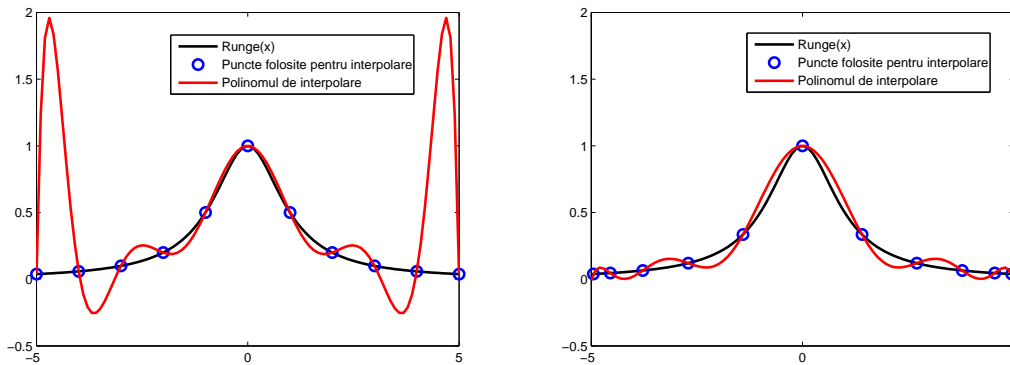


Figura 4.7: Efectul Runge: polinomul de interpolare are oscilații la capetele intervalului dacă gridul de discretizare este uniform (stânga). Efectul poate fi eliminat prin plasarea nodurilor de discretizare în conformitate cu rădăcinile polinoamelor Chebyshev (dreapta).

Printr-o transformare liniară, se poate deduce că într-un interval arbitrar $[a, b]$ nodurile Chebyshev trebuie plasate conform formulei

$$x_i = \frac{a+b}{2} + \frac{b-a}{2} \cos\left(\frac{2i+1}{2(n+1)}\pi\right). \quad (4.71)$$

În concluzie, interpolarea Chebyshev este tot interpolarea polinomială globală numai că nodurile rețelei de interpolare sunt alese în concordanță cu rădăcinile polinoamelor Chebyshev.

De cele mai multe ori însă, utilizatorul nu are libertatea de a alege punctele rețelei de interpolare. De aceea, limitarea acestei metode este aceea că se poate aplica numai funcțiilor definite prin cod și nu tabelar. În cazul funcțiilor date tabelar, pentru a obține polinoame de interpolare care să nu prezinte efect Runge, este mai eficient dacă se folosește interpolarea pe porțiuni.

4.4 Metode de interpolare polinomială pe porțiuni

Interpolarea polinomială globală poate suferi de efectul Runge, efect care poate fi eliminat în cazul funcțiilor al căror cod este cunoscut prin folosirea unei interpolări Chebyshev. În cazul funcțiilor date tabelar, interpolări satisfăcătoare se obțin doar dacă polinoamele de interpolare se caută pe porțiuni ale domeniului de definiție. Rezultatul interpolării pe porțiuni este ceea ce în matematică se numește ”funcție cu acoladă”, spre deosebire de interpolarea globală al cărei rezultat este o singură expresie compactă, valabilă pe tot domeniul de definiție.

4.4.1 Interpolarea liniară pe porțiuni

Cea mai simplă interpolare pe porțiuni este aceea în care funcțiile de bază $\varphi_k(x)$ sunt polinoame Lagrange $l_k(x)$ definite pe porțiuni. Ele satisfac de asemenea condițiile $l_k(x_k) = 1$ și $l_k(x_j) = 0$ pentru $j \neq k$, dar variația între nodurile grilei de discretizare este liniară (fig. 4.8):

$$l_0(x) = \begin{cases} \frac{x-x_1}{x_0-x_1} & \text{dacă } x \in [x_0, x_1] \\ 0 & \text{dacă } x \in (x_1, x_n] \end{cases} \quad (4.72)$$

$$l_k(x) = \begin{cases} \frac{x-x_{k-1}}{x_k-x_{k-1}} & \text{dacă } x \in [x_{k-1}, x_k] \\ \frac{x-x_{k+1}}{x_k-x_{k+1}} & \text{dacă } x \in [x_k, x_{k+1}] \\ 0 & \text{dacă } x \in [x_1, x_{k-1}] \cup (x_{k+1}, x_n] \end{cases} \quad k = 2, n-1 \quad (4.73)$$

$$l_n(x) = \begin{cases} \frac{x-x_{n-1}}{x_n-x_{n-1}} & \text{dacă } x \in [x_{n-1}, x_n] \\ 0 & \text{dacă } x \in [x_0, x_{n-1}] \end{cases} \quad (4.74)$$

Graficul polinomului de interpolare

$$g(x) = \sum_{k=0}^n c_k l_k(x) \quad (4.75)$$

este chiar linia poligonală care unește punctele rețelei de interpolare (fig. 4.9). Folosirea acestui polinom în afara domeniului de definiție se numește *extrapolare*.

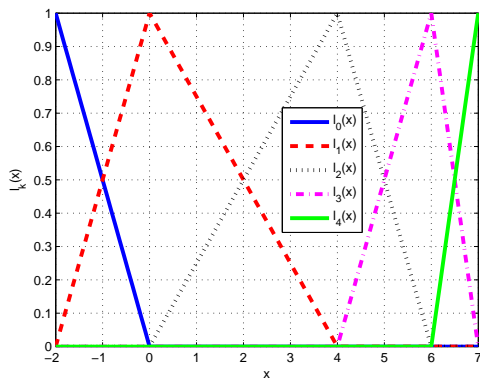


Figura 4.8: Funcții Lagrange pe porțiuni.

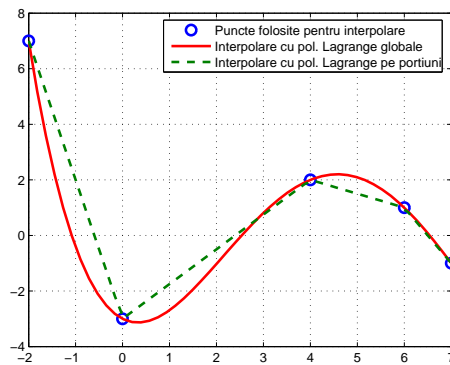


Figura 4.9: Interpolare globală și interpolare pe porțiuni.

Polinomul de interpolare are o expresie diferită pe fiecare sub-interval definit de rețeaua de discretizare

$$g(x) = \frac{y_{k+1} - y_k}{x_{k+1} - x_k} (x - x_k) + y_k, \quad \text{pentru } x \in [x_k, x_{k+1}], \quad k = 0, \dots, n-1. \quad (4.76)$$

Polinomul de interpolare este continuu, și de aceea nu are importanță în care interval (la stânga sau la dreapta) este inclus un nod al rețelei de interpolare. Procedura de evaluare a polinomului de interpolare liniară pe porțiuni se reduce în ultimă instanță la evaluarea ecuației unei drepte care trece prin două puncte date, de coordonate (x_k, y_k) și, respectiv, (x_{k+1}, y_{k+1}) . Singura dificultate este de a determina sub-intervalul din care face parte punctul de evaluat x_{crt} .

Pseudocodul următor evaluează un polinom de interpolare liniară pe porțiuni, în care căutarea intervalului în care se află punctul de evaluat se face prin *metoda căutării binare*. Ideea acestei metode este de a elimina la fiecare iterație jumătate din numărul de intervale posibile prin compararea punctului de evaluat cu punctul din rețeaua de discretizare care împarte numărul de intervale posibile în două. Pseudocodul acestei metode a fost prezentat la pagina 28.

```

funcție interpolare_lpp( $n, x, y, x_{crt}$ )
; evaluează polinomul de interpolare liniară pe porțiuni în punctul  $x_{crt}$ 
; declarații - parametri de intrare
întreg  $n$  ; dimensiunea problemei - numărul de intervale
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la zero
real  $x_{crt}$  ; punctul de evaluat
 $k = \text{cauta}(n, x, x_{crt})$ 
întoarce  $(y_{k+1} - y_k) / (x_{k+1} - x_k) * (x_{crt} - x_k) + y_k$ 

```

Interpolarea liniară pe porțiuni are dezavantajul că funcția obținută prin interpolare nu este derivabilă în nodurile rețelei de interpolare. Pentru a genera interpolări pe porțiuni care să nu fie numai continue ci și derivabile este necesară creșterea gradului polinoamelor care se folosesc pe porțiuni.

4.4.2 Interpolarea Hermite

O funcție de interpolare derivabilă s-ar putea construi dacă s-ar cunoaște valorile derivatelor funcției de interpolat în nodurile rețelei de interpolare.

Se definește astfel *problema interpolării Hermite* ca fiind problema interpolării unei funcții pe o rețea de noduri x_k în care se cunosc valorile funcției y_k și valorile derivatelor acesteia y'_k . Tabelul de valori necesar pentru o interpolare Hermite este deci de tipul

x	x_0	x_1	\cdots	x_k	\cdots	x_n
y	y_0	y_1	\cdots	y_k	\cdots	y_n
y'	y'_0	y'_1	\cdots	y'_k	\cdots	y'_n

Pe fiecare subinterval $[x_k, x_{k+1}]$ unde $k = 0, \dots, n-1$ trebuie puse patru condiții de interpolare: două pentru valorile funcției și două pentru valorile derivatei la capete:

$$\begin{cases} g(x_k) = y_k, \\ g(x_{k+1}) = y_{k+1}, \\ g'(x_k) = y'_k, \\ g'(x_{k+1}) = y'_{k+1}. \end{cases} \quad (4.77)$$

Aceste condiții asigură continuitatea și derivabilitatea polinomului de interpolare.

Deoarece se dorește ca impunerea condițiilor de interpolare să conducă la un polinom unic, rezultă că pe fiecare subinterval trebuie considerat drept polinom de interpolare unul de grad trei. Cei patru coeficienți ai polinomului de interpolare vor rezulta în mod unic din impunerea condițiilor de interpolare (4.77). Pentru simplificarea calculelor, este util ca polinomul de interpolare pe porțiuni să fie scris ca

$$g(x) = c_{0k} + c_{1k}(x - x_k) + c_{2k}(x - x_k)^2 + c_{3k}(x - x_k)^3, \quad x \in [x_k, x_{k+1}], \quad (4.78)$$

derivata lui fiind

$$g'(x) = c_{1k} + 2c_{2k}(x - x_k) + 3c_{3k}(x - x_k)^2, \quad x \in [x_k, x_{k+1}]. \quad (4.79)$$

Impunând cele patru condiții de interpolare rezultă următorul sistem de patru ecuații cu patru necunoscute

$$\begin{cases} c_{0k} = y_k, \\ c_{0k} + c_{1k}(x_{k+1} - x_k) + c_{2k}(x_{k+1} - x_k)^2 + c_{3k}(x_{k+1} - x_k)^3 = y_{k+1}, \\ c_{1k} = y'_k, \\ c_{1k} + 2c_{2k}(x_{k+1} - x_k) + 3c_{3k}(x_{k+1} - x_k)^2 = y'_{k+1}, \end{cases} \quad (4.80)$$

a cărui soluție este

$$\begin{cases} c_{0k} = y_k, \\ c_{1k} = y'_k, \\ c_{2k} = [3(y_{k+1} - y_k) - (x_{k+1} - x_k)(2y'_k + y'_{k+1})] / (x_{k+1} - x_k)^2, \\ c_{3k} = [(y'_{k+1} + y'_k) - 2(2y_{k+1} - y_k) / (x_{k+1} - x_k)] / (x_{k+1} - x_k)^2. \end{cases} \quad (4.81)$$

Problema interpolării Hermite este aceea că de cele mai multe ori nu se cunosc informații despre derivatele funcției. Evaluarea derivatelor se face numeric și această etapă este numită etapă de pregătire a interpolării Hermite. Există mai multe variante de a aproxima numeric derivatele. O primă variantă, cunoscută sub numele de *interpolare Bessel*, este de a folosi formule simple pentru aproximarea derivatelor. De exemplu, la extremități se

poate considera o aproximare liniară pe ultimul subinterval și atunci valorile se aproximează cu

$$y'_0 = \frac{y_1 - y_0}{x_1 - x_0} \quad y'_n = \frac{y_n - y_{n-1}}{x_n - x_{n-1}}, \quad (4.82)$$

iar în nodurile interioare, derivând expresia unei parabole ce interpolează trei noduri consecutive, se poate deduce o formulă aproximativă de tipul

$$y'_k = \frac{\beta^2 y_{k+1} + (\alpha^2 - \beta^2) y_k - \alpha^2 y_{k-1}}{\alpha\beta(\alpha + \beta)h}, \quad (4.83)$$

unde am notat $x_{k+1} - x_k = \alpha h$, $x_k - x_{k-1} = \beta h$. Această abordare³ nu este tocmai naturală pentru că interpolarea Hermite folosește pe porțiuni polinoame de ordinul 3 și nu de ordinul 1 sau 2 cum sunt presupuse în formulele (4.82) și (4.83).

Metoda cea mai avantajoasă de evaluare a derivatelor este cea în care ele se deduc din impunerea unor condiții de continuitate suplimentare pentru derivata a doua a polinomului de interpolare în nodurile rețelei de interpolare. În acest caz, interpolarea Hermite se numește *interpolare spline cubică clasică*. În nodurile interioare se impune condiția:

$$g''(x_k - 0) = g''(x_k + 0), \quad k = 1, \dots, n - 1. \quad (4.84)$$

Din (4.79) rezultă expresia derivatei a doua a polinomului de interpolare

$$g''(x) = 2c_{2k} + 6c_{3k}(x - x_k), \quad x \in [x_k, x_{k+1}]. \quad (4.85)$$

Condiția (4.84) devine

$$2c_{2,k-1} + 6c_{3,k-1}(x_k - x_{k-1}) = 2c_{2k}, \quad k = 1, \dots, n - 1. \quad (4.86)$$

Simplificând relația cu 2 și înlocuind expresiile coeficienților date de (4.81), relațiile (4.86) devin

$$\begin{aligned} & \frac{1}{x_k - x_{k-1}} y'_{k-1} + 2 \left(\frac{1}{x_k - x_{k-1}} + \frac{1}{x_{k+1} - x_k} \right) y'_k + \frac{1}{x_{k+1} - x_k} y'_{k+1} = \\ & = 3 \frac{y_k - y_{k-1}}{(x_k - x_{k-1})^2} + 3 \frac{y_{k+1} - y_k}{(x_{k+1} - x_k)^2}. \end{aligned} \quad (4.87)$$

Cele $n - 1$ relații (4.87) conțin $n + 1$ necunoscute. Pentru unicitate se impun încă două condiții la capete. O variantă este să se impună valorile derivatelor y'_0 și y'_n ca la interpolarea Bessel, caz în care se spune că s-au impus *condiții forțate la capete*. O altă variantă este aceea în care la capete se impune ca derivata a doua a polinomului de interpolare să fie zero:

$$g''(x_0) = 0, \quad g''(x_n) = 0. \quad (4.88)$$

³Formulele (4.82) și (4.83) sunt formule de derivare numerică, problemă explicată în detaliu în capitolul următor.

În acest caz se spune că s-au impus *condiții naturale la capete*. După impunerea acestor condiții, sistemul (4.87) devine bine formulat. El este un sistem tridiagonal din a cărui rezolvare rezultă valorile derivatelor necesare calcului coeficienților polinomului de interpolare Hermite. Din impunerea condiției naturale la capătul din stânga rezultă următoarea relație ce trebuie să existe între derivatele funcției:

$$2y'_0 + y'_1 = 3 \frac{y_1 - y_0}{x_1 - x_0}, \quad (4.89)$$

iar din impunerea condiției la capătul din dreapta rezultă

$$y'_{n-1} + 2y'_n = 3 \frac{y_n - y_{n-1}}{x_n - x_{n-1}}. \quad (4.90)$$

Oricare ar fi modul de tratare a extremităților, interpolarea spline cubică clasică are avantajul că minimizează curbura pătratică medie a polinomului de interpolare, mai precis expresia

$$\int_a^b (g''(x))^2 dx \quad (4.91)$$

este minimă față de orice altă interpolare polinomială pe porțiuni⁴. Datorită acestei proprietăți, în această interpolare nu apar oscilații nedorite între punctele rețelei de interpolare.

Algoritmul interpolării spline cubice este următorul:

funcție pregătește_spline($n, x, y, yder$)

; calculează derivatele funcției în nodurile rețelei de discretizare

; declarații - parametri de intrare

întreg n

; dimensiunea problemei - numărul de intervale

tablou real $x[n], y[n]$

; tabelul de valori, indici de la zero

tablou real $yder[n]$

; parametri de ieșire

tablou real $p[n], q[n], r[n]$

; matricea tridiagonală asamblată

; assemblează matricea tridiagonală

$q_0 = 2$

$r_0 = 1$

$b_0 = 3(y_1 - y_0)/(x_1 - x_0)$

pentru $k = 1, n - 1$

$p_k = 1/(x_k - x_{k-1})$

$q_k = 2/(x_k - x_{k-1}) + 2/(x_{k+1} - x_k)$

$r_k = 1/(x_{k+1} - x_k)$

⁴În limba engleză *spline* înseamnă, printre altele, o fâșie de lemn, destul de elastică, folosită pentru desenarea curbelor, ca un florar.

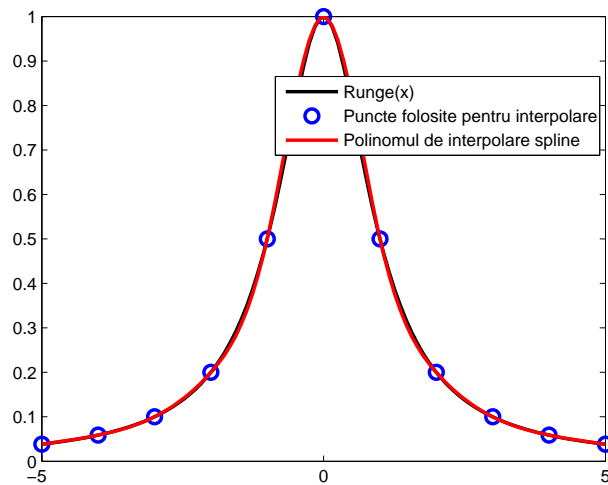


Figura 4.10: Interpolarea *spline* a funcției Runge pe o rețea uniformă de puncte.

$$b_k = 3(y_k - y_{k-1})/(x_k - x_{k-1})^2 + 3(y_{k+1} - y_k)/(x_{k+1} - x_k)^2$$

•

$$p_n = 1$$

$$q_n = 2$$

$$b_n = 3(y_n - y_{n-1})/(x_n - x_{n-1})$$

Gauss.tridiag($n + 1, p, q, r, b, yder$) ; procedură descrisă la pagina 86

funcție `aproximare_spline($n, x, y, yder, xcrt$)`

; evaluează polinomul de aproximare spline în punctul $xcrt$

; declarații - parametri de intrare

întreg n ; dimensiunea problemei - numărul de intervale

tablou real $x[n], y[n], yder[n]$; tabelul de valori, indici de la zero

real $xcrt$; punctul de evaluat

$$k = \text{cauta}(n, x, xcrt)$$

$$c0k = y_k$$

$$c1k = yder_k$$

$$h = x_{k+1} - x_k$$

$$c2k = (3(y_{k+1} - y_k) - h(2 * yder_k + yder_{k+1}))/h^2$$

$$c3k = (yder_{k+1} + yder_k - 2(y_{k+1} - y_k)/h)/h^2$$

$$hcrt = xcrt - x_k$$

$$\text{întoarce } c0k + c1k * hcrt + c2k * hcrt^2 + c3k * hcrt^3$$

Rezultatul rulării unui astfel de algoritm este prezentat în fig. 4.10.

Capitolul 5

Derivarea numerică

Necesitatea evaluării derivatelor funcțiilor este importantă în problemele de inginerie. Ele apar, de exemplu, în relații utile pentru evaluarea unor mărimi (curent, tensiune), în definirea sensibilităților unor mărimi de interes în raport cu anumiți parametri considerați variabili, sensibilități utile în proiectarea optimală a dispozitivelor. Derivarea numerică se bazează pe interpolare. Formulele de derivare numerică provin din derivarea polinomului de interpolare și ele se reduc la calcule aritmetice ce folosesc formule relativ simple.

Un alt exemplu în care apar derivate este cel în care modelul matematic al problemei conduce la ecuații sau sisteme de ecuații diferențiale, foarte puține dintre acestea putând fi rezolvate prin metode analitice. Rezolvarea ecuațiilor diferențiale folosind formule de derivare numerică pentru derivatele ce intervin este cunoscută sub numele de *metoda diferențelor finite*.

5.1 Formularea problemei derivării numerice

Vom presupune mai întâi cazul unui funcții reale ce depinde de o singură variabilă reală. Problema derivării unei astfel de funcții se formulează astfel.

Se dă funcția $f : [a, b] \rightarrow \mathbb{R}$ (cunoscută prin date sau prin cod) și **se cere** evaluarea numerică a derivatei $f'(x_0)$, unde $x_0 \in [a, b]$.

Matematic, derivata se definește ca o limită

$$f'(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0}. \quad (5.1)$$

O astfel de expresie nu poate fi calculată cu un calculator numeric deoarece evaluarea unei limite ar presupune o procedură cu un număr infinit de pași. Metodele de derivare numerică vor fi afectate în consecință și de erori de trunchiere, nu numai de erori inerente și de rotunjire. Mai mult, pe măsură ce x se apropie de x_0 și $f(x)$ se apropie de $f(x_0)$,

ceea ce înseamnă că membrul drept care trebuie evaluat tinde către cazul nedeterminat $0/0$. Este de așteptat atunci ca evaluarea unei astfel de limite cu o formulă aproximativă de tipul $(f(x) - f(x_0))/(x - x_0)$ să ridice probleme numerice atunci când x este ales prea aproape de x_0 . Pe de altă parte, evaluarea derivatei printr-o diferență divizată poate introduce erori inacceptabile atunci când x nu este suficient de apropiat de x_0 .

Problema derivării numerice este strâns legată de problema interpolării sau aproximării funcțiilor. Orice fel de derivată poate fi reprezentată ca fiind derivata funcției aproximative obținute prin interpolare sau aproximare. Dacă notăm cu $g(x)$ interpolarea sau aproximarea funcției $f(x)$ atunci derivata numerică $g'(x)$ va aproxima derivata exactă $f'(x)$:

$$g'(x) = \frac{dg}{dx} \approx f'(x) = \frac{df}{dx}. \quad (5.2)$$

Evident că precizia cu care se face interpolarea sau aproximarea va determina precizia de calcul a derivatei numerice.

5.2 Formule de derivare numerică

5.2.1 Formule de derivare cu eroare de ordinul 1

Cea mai simplă și mai robustă metodă de interpolare este metoda interpolării liniare pe porțiuni. Funcția g este derivabilă în interiorul intervalelor definite de rețeaua de interpolare (Fig.5.1). Derivata ei este o funcție constantă pe porțiuni (Fig.5.2).

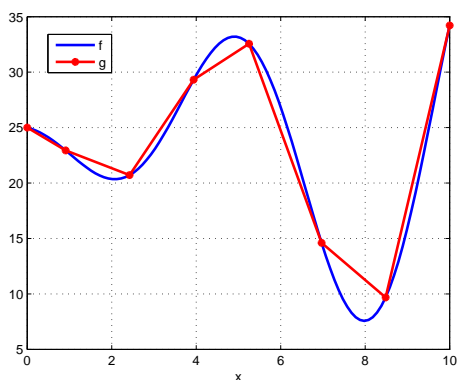


Figura 5.1: Funcția reală f și interpolarea ei pe porțiuni g .

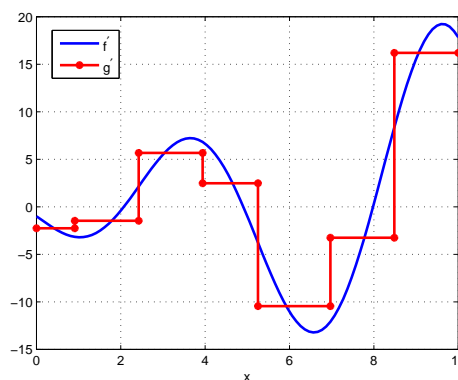


Figura 5.2: Aproximarea derivatei f' cu valoarea g' .

În nodurile de interpolare funcția nu este derivabilă. Într-un astfel de nod am putea prelungi fie valoarea derivatei din dreapta, obținând o formulă de *derivare progresivă de*

ordinul 1:

$$y'_k = \frac{y_{k+1} - y_k}{x_{k+1} - x_k}, \quad (5.3)$$

fie valoarea derivatei din stânga, obținând o formulă de *derivare regresivă de ordinul 1*:

$$y'_k = \frac{y_k - y_{k-1}}{x_k - x_{k-1}}. \quad (5.4)$$

Ordinul erorii de trunchiere specifice acestei aproximări se poate estima cu ajutorul dezvoltării în serie Taylor a funcției f în jurul punctului x_k :

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{f''(\xi)}{2!}(x - x_k)^2. \quad (5.5)$$

Evaluând această expresie în $x = x_{k+1}$ rezultă că

$$y_{k+1} = y_k + f'(x_k)(x_{k+1} - x_k) + \frac{f''(\xi)}{2!}(x_{k+1} - x_k)^2, \quad (5.6)$$

de unde

$$\frac{y_{k+1} - y_k}{x_{k+1} - x_k} - f'(x_k) = \frac{f''(\xi)}{2!}(x_{k+1} - x_k). \quad (5.7)$$

Presupunând că funcția pe care o derivăm are derivata a doua mărginită de o constantă M_2 , unde $|f''(x)| \leq M_2$, și notând cu $h = x_{k+1} - x_k$, rezultă că eroarea de trunchiere a formulei de derivare numerică progresivă folosită este

$$|e_t| \leq \frac{M_2}{2}h, \quad (5.8)$$

și se notează pe scurt $|e_t| = O(h)$. Se spune că eroarea este de ordinul 1 (1 este puterea lui h , sau 1 este gradul polinomului de interpolare folosit).

Pentru derivata regresivă (5.4) concluzia este similară, rezultatul obținându-se din evaluarea dezvoltării (5.5), care se mai scrie și sub forma

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + O((x - x_k)^2) \quad (5.9)$$

în punctul $x = x_{k-1}$. Notând $x_k - x_{k-1} = h$ rezultă

$$y_{k-1} = y_k - f'(x_k)h + O(h^2), \quad (5.10)$$

de unde

$$\frac{y_k - y_{k-1}}{h} - f'(x_k) = O(h). \quad (5.11)$$

5.2.2 Formule de derivare cu eroare de ordinul 2

Abordarea din paragraful anterior are dezavantajul că derivata funcției interpolatoare este discontinuă în nodurile rețelei de interpolare. Acest lucru poate fi corectat dacă se mărește ordinul polinomului de interpolare folosit local pentru deducerea formulelor de derivare numerică.

Să considerăm de exemplu trei puncte consecutive ale rețelei de interpolare și polinomul de interpolare de grad doi care trece prin acestea. Expresia acestui polinom (Lagrange) este

$$\begin{aligned} g(x) &= \frac{(x - x_k)(x - x_{k+1})}{(x_{k-1} - x_k)(x_{k-1} - x_{k+1})}y_{k-1} + \frac{(x - x_{k-1})(x - x_{k+1})}{(x_k - x_{k-1})(x_k - x_{k+1})}y_k + \\ &+ \frac{(x - x_{k-1})(x - x_k)}{(x_{k+1} - x_{k-1})(x_{k+1} - x_k)}y_{k+1}, \end{aligned} \quad (5.12)$$

iar derivata lui are forma

$$\begin{aligned} g'(x) &= \frac{(x - x_k) + (x - x_{k+1})}{(x_{k-1} - x_k)(x_{k-1} - x_{k+1})}y_{k-1} + \frac{(x - x_{k-1}) + (x - x_{k+1})}{(x_k - x_{k-1})(x_k - x_{k+1})}y_k + \\ &+ \frac{(x - x_{k-1}) + (x - x_k)}{(x_{k+1} - x_{k-1})(x_{k+1} - x_k)}y_{k+1}. \end{aligned} \quad (5.13)$$

Pentru a simplifica scrierea, vom nota $h_1 = x_k - x_{k-1}$ și $h_2 = x_{k+1} - x_k$. Evaluând polinomul derivat (5.13) în x_k se obține formula de *derivare centrată de ordinul 2*

$$y'_k = -\frac{h_2}{h_1(h_1 + h_2)}y_{k-1} - \frac{h_1 - h_2}{h_1h_2}y_k + \frac{h_1}{h_2(h_1 + h_2)}y_{k+1}. \quad (5.14)$$

Evaluând polinomul derivat (5.13) în punctele x_{k-1} și x_{k+1} se obține formula de *derivare progresivă de ordinul 2*

$$y'_{k-1} = -\frac{2h_1 + h_2}{h_1(h_1 + h_2)}y_{k-1} + \frac{h_1 + h_2}{h_1h_2}y_k - \frac{h_1}{h_2(h_1 + h_2)}y_{k+1}, \quad (5.15)$$

și, respectiv, de *derivare regresivă de ordinul 2*

$$y'_{k+1} = -\frac{h_2}{h_1(h_1 + h_2)}y_{k-1} - \frac{h_1 + h_2}{h_1h_2}y_k + \frac{h_1 + 2h_2}{h_2(h_1 + h_2)}y_{k+1}. \quad (5.16)$$

În cazul unui pas echidistant $h_1 = h_2 = h$, formulele de derivare de ordinul 2 devin cele din relațiile (5.17) pentru derivarea centrată, (5.18) pentru derivarea progresivă și, respectiv, (5.19) pentru derivarea regresivă:

$$y'_k = \frac{y_{k+1} - y_{k-1}}{2h}, \quad (5.17)$$

$$y'_{k-1} = \frac{-3y_{k-1} + 4y_k - y_{k+1}}{2h}, \quad (5.18)$$

$$y'_{k+1} = \frac{-y_{k-1} - 4y_k + 3y_{k+1}}{2h}. \quad (5.19)$$

Ordinul erorii de trunchiere specifice acestei aproximări se poate estima cu ajutorul dezvoltării în serie Taylor în care se rețin mai mulți termeni. De exemplu, dezvoltând funcția f în vecinătatea lui x_k și păstrând inclusiv termenul de ordin 2, eroarea de trunchiere a seriei este de ordin 3:

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{f''(x_k)}{2!}(x - x_k)^2 + O((x - x_k)^3). \quad (5.20)$$

Dacă evaluăm această dezvoltare în punctele x_{k-1} și x_{k+1} unde presupunem pas echidistant $x_{k+1} - x_k = x_k - x_{k-1} = h$, rezultă:

$$f(x_{k+1}) = f(x_k) + f'(x_k)h + \frac{f''(x_k)}{2!}h^2 + O(h^3), \quad (5.21)$$

$$f(x_{k-1}) = f(x_k) - f'(x_k)h + \frac{f''(x_k)}{2!}h^2 - O(h^3). \quad (5.22)$$

Scăzând aceste două relații rezultă

$$f(x_{k+1}) - f(x_{k-1}) = 2f'(x_k)h + O(h^3), \quad (5.23)$$

de unde rezultă că formula de derivare centrată are eroarea de trunchiere de ordinul 2:

$$\frac{y_{k+1} - y_{k-1}}{2h} - f'(x_k) = O(h^2). \quad (5.24)$$

Formulele de ordinul 2 sunt mai precise decât formulele de ordinul 1, și această proprietate poate fi ilustrată în cazul unui pas echidistant $x_{k+1} - x_k = x_k - x_{k-1} = h$ (Fig.5.3). Panta secantei care unește punctele $k-1$ și $k+1$ este mai apropiată de derivata în punctul x_k decât panta secantei care unește punctul k cu punctul $k-1$ sau cu punctul $k+1$. În cazul în care f ar fi funcție de gradul doi, derivata centrată este chiar valoarea adevărată (eroare de trunchiere este 0).

5.3 Algoritmi numerici pentru derivarea funcțiilor

Algoritmul pentru evaluarea derivatei unei funcții într-un punct depinde de modul în care este cunoscută funcția.

Dacă funcția este dată tabelar (prin date), atunci se evaluează derivatele în punctele din tabel, iar dacă se dorește evaluarea derivatei într-un punct oarecare, se poate face o interpolare liniară pe porțiuni a acestor valori. Pseudocodul acestei abordări este prezentat în paragraful 5.3.1.

Dacă funcția este dată printr-un cod, atunci avem libertatea de a alege pasul de derivare. Teoretic, cu cât pasul de derivare este mai mic, cu atât eroarea este mai mică.

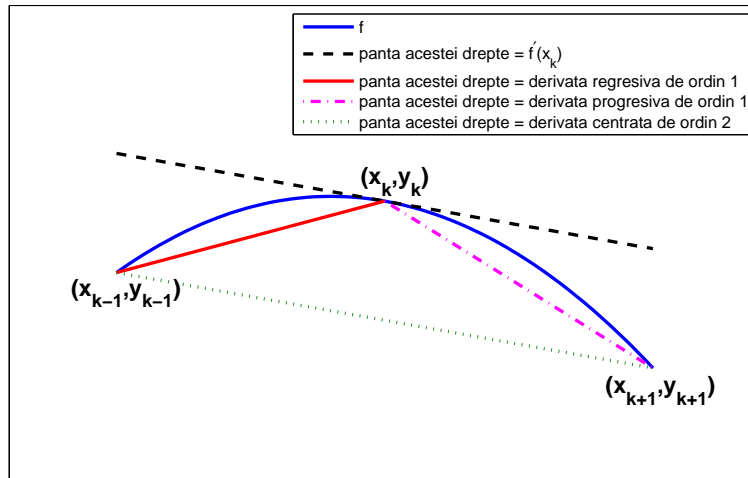


Figura 5.3: Semnificația geometrică a celor mai simple formule de derivare numerică.

Această afirmație este adevărată doar pentru eroarea de trunchiere, lucru demonstrat anterior pentru formule de derivare de ordinul 1 și 2. În calculele numerice apar și erori de rotunjire care sunt cu atât mai mari cu cât pasul de derivare este mai mic. În acest caz este necesară estimarea unui pas optim de derivare, așa cum este prezentat în paragraful 5.3.2.

5.3.1 Cazul funcțiilor date tabelar

Deoarece, în general, un polinom de interpolare pe porțiuni nu este derivabil în nodurile de interpolare, rezultând o derivată care nu este continuă (ca în Fig. 5.2), se preferă ca aproximarea derivatei unei funcții pe un domeniu să fie interpolarea (de exemplu liniară) pe porțiuni a derivatelor din nodurile de interpolare, calculate cu formule de derivare numerică. Pseudocodul următor ilustrează această abordare.

```

procedură pregătire_derivate( $n, x, y, dy$ )
; calculează derivatele numerice în noduri
; declarații
întreg  $n$  ; numărul de puncte din tabel minus 1
tablou real  $x[n], y[n]$  ; tabelul de valori, indici de la 0
tablou real  $dy[n]$  ; valorile derivatelor
; alte declarații
...
; la capătul din stânga diferențe progresive de ordinul 2

```

$$h1 = x_1 - x_0$$

$$h2 = x_2 - x_1$$

$$dy_0 = -(2h1 + h2)/(h1(h1 + h2))dy_0 + (h1 + h2)/(h1h2)dy_1 - h1/(h2(h1 + h2))dy_2$$

; la capătul din dreapta diferențe regresive de ordinul 2

$$h1 = x_{n-1} - x_{n-2}$$

$$h2 = x_n - x_{n-1}$$

$$dy_n = -h2/(h1(h1 + h2))dy_{n-2} - (h1 + h2)/(h1h2)dy_{n-1} + (h1 + 2h2)/(h2(h1 + h2))dy_n$$

; în nodurile interioare diferențe centrate de ordinul 2

pentru $k = 1, n - 1$

$$h1 = x_k - x_{k-1}$$

$$h2 = x_{k+1} - x_k$$

$$dy_k = -h2/(h1(h1 + h2))dy_{k-1} - (h1 - h2)/(h1h2)dy_k + h1/(h2(h1 + h2))dy_{k+1}$$

•

Evaluarea derivatei funcției într-un punct oarecare $xcrt$ se face apelând procedura de interpolare liniară pe porțiuni prezentată în paragraful 4.4.1 pentru un tabel de valori format din derivatele în noduri:

$$rez = \text{interpolare_lpp}(n, x, dy, xcrt)$$

5.3.2 Cazul funcțiilor date prin cod

În cazul în care funcția este dată prin cod, pasul de discretizare folosit pentru orice formulă de derivare numerică trebuie ales de utilizator. Eroarea de trunchiere este cu atât mai mică cu cât pasul este mai mic. Dar, pe lângă eroarea de trunchiere, există și o eroare de rotunjire, care este cu atât mai mare cu cât pasul este mai mic, deoarece în formulele derivatelor se fac scăderi cu numere apropiate.

Teste numerice făcute pentru calculul numeric al derivatei funcției sinus în punctul $\pi/4$ arată că variația erorii în funcție de pasul de discretizare este cea din Fig.5.4.

Zona de instabilități numerice, care apare pentru pași mai mici decât o anumită valoare, se datorează erorilor de rotunjire care devin mai mari decât erorile de trunchiere. Pasul de la care încep să predominie erorile de trunchiere se numește *pas optim de derivare numerică*. Pasul optim nu este pasul pentru care eroarea este minimă. După cum se poate observa, erori mai mici se pot obține pentru pași aflați în zona de instabilitate. Lucrul în această zonă nu se recomandă, fiind de preferat ca h să fie în zona în care predomină erorile de trunchiere, unde erorile se pot estima.

Pasul optim de derivare se poate estima. În cele ce urmează vom estima pasul optim de derivare în cazul formulei de derivare progresivă de ordinul 1.

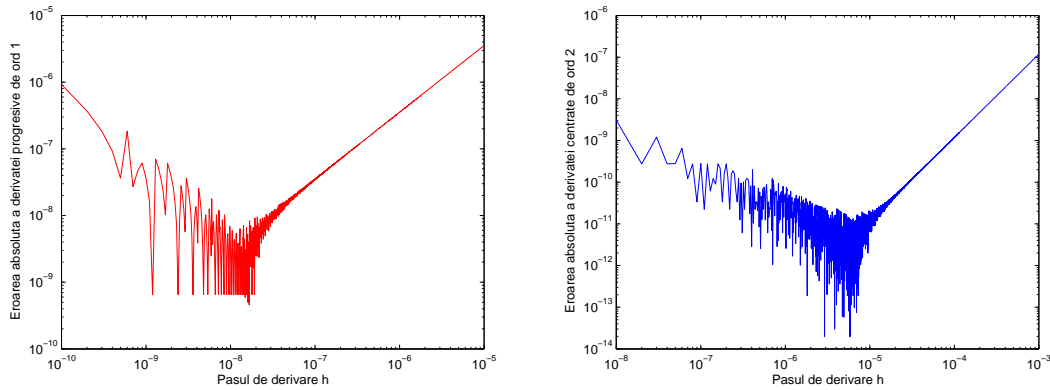


Figura 5.4: Eroarea totală în funcție de pasul de derivare pentru formula de derivare progresivă de ordinul 1 (stânga) și formula de derivare centrată de ordinul 2 (dreapta).

În paragraful 5.2.1 am arătat că eroarea de trunchiere în cazul folosirii derivatei numerice

$$y'_k = \frac{y_{k+1} - y_k}{x_{k+1} - x_k} = \frac{y_{k+1} - y_k}{h} \quad (5.25)$$

este mărginită:

$$|e_t| = |y'_k - f'(x_k)| \leq \frac{M_2}{2} h, \quad (5.26)$$

unde M_2 este o margine superioară a derivatei de ordinul 2 a funcției de derivat: $|f''(x)| \leq M_2$. Să notăm cu

$$a_t = \frac{M_2}{2} h \quad (5.27)$$

marginea erorii de trunchiere.

Să presupunem că valorile y_k sunt reprezentate cu precizia maximă admisibilă în calculator, deci ele sunt afectate doar de erori de rotunjire. În consecință $|e_{y_k}/y_k| < \text{eps}$, unde eps este zeroul mașinii. Presupunem de asemenea că pasul de discretizare h nu este afectat de erori $e_h = 0$. Aplicând formula erorii la împărțire din tabelul 1.1 rezultă că eroarea de rotunjire a relației (5.25) este majorată de

$$|e_r| \leq \frac{|e_{y_{k+1}}| + |e_{y_k}|}{h} \leq \frac{\text{eps} (|y_{k+1}| + |y_k|)}{h} \quad (5.28)$$

Dacă notăm o margine superioară M_0 pentru valorile y_{k+1} și y_k , atunci eroarea de rotunjire este majorată de

$$|e_r| \leq \frac{2M_0 \text{eps}}{h}. \quad (5.29)$$

Să notăm cu

$$a_r = \frac{2M_0 \text{eps}}{h} \quad (5.30)$$

marginea erorii de rotunjire.

Eroarea totală $e = e_t + e_r$ va fi majorată de o expresie ce depinde de pasul de derivare $|e| = |e_t + e_r| \leq |e_t| + |e_r| \leq a_t + a_r = m(h)$ unde

$$m(h) = \frac{M_2}{2}h + \frac{2M_0 \text{eps}}{h}, \quad (5.31)$$

care își va atinge valoarea minimă pentru un pas de derivare optim

$$\min_h(m(h)) = m(h_{\text{optim}}). \quad (5.32)$$

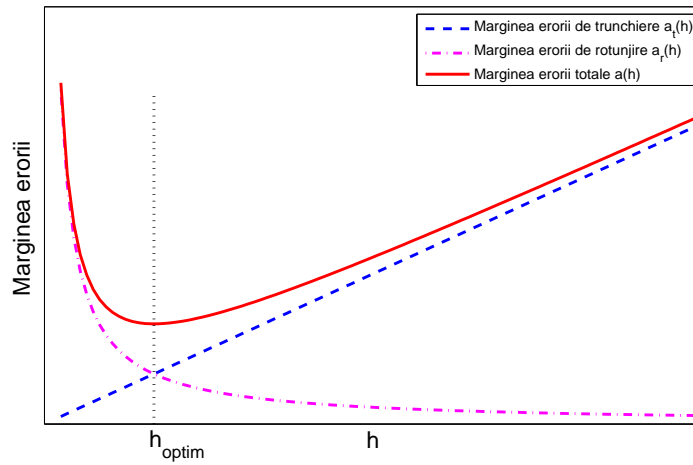


Figura 5.5: Pasul optim de derivare numerică este pasul de la care încep să predomine erorile de trunchiere.

Punând condiția de minim pentru expresia (5.31), $m'(h) = 0$, rezultă valoarea optimă a pasului de derivare

$$h_{\text{optim}} = 2\sqrt{\frac{M_0 \text{eps}}{M_2}}. \quad (5.33)$$

În concluzie, pasul optim de derivare depinde de eroarea de rotunjire, de marginea funcției și de marginea derivatei a doua. El este pasul pentru care marginea erorii totale este minimă și nu trebuie confundat cu pasul pentru care eroarea totală este minimă, pas care s-ar putea să fie mai mic decât pasul optim, dar care nu poate fi estimat apriori.

Algoritmul de mai jos calculează derivata progresivă de ordinul 1 a unei funcții date prin cod, folosind un pas de derivare optim. Pentru calculul pasului optim, trebuie estimată derivata de ordinul 2. Pentru estimarea ei se folosește o formulă de derivare în trei puncte, formulă ce va fi justificată în paragraful 5.4:

$$M_2 \approx \frac{f(x) - 2f(x+h) + f(x+2h)}{h^2}. \quad (5.34)$$

Se consideră că valoarea pasului optim este rezonabilă dacă raportul dintre două aproximații consecutive ale acestuia este cuprins în intervalul $[0.5, 2]$. S-ar putea ca algoritmul să nu fie convergent, de aceea s-a introdus un contor de iterații.

```

funcție derivata_f(x, h0, maxit)
; calculează derivata numerică a funcției f
; folosind diferențe progresive de ordinul 1
; și pas optim de derivare
real x ; punctul în care se va evalua derivata
real h0 ; pasul inițial
întreg maxit ; numărul maxim de iterații pentru calculul pasului
; alte declarații
...
eps = zeroul_mașinii () ; vezi pagina 34
f0 = f(x) ; valoarea funcției în punctul de derivare
h = h0
k = 0
repetă
    k = k + 1
    f1 = f(x + h)
    f2 = f(x + 2h)
    M2 = |f0 - 2f1 + f2|/h2 ; estimarea marginii derivatei a doua
    dacă M2 < eps ; testează dacă derivata a doua este zero
        hoptim = h ; funcția e liniară, eroarea de trunchiere e zero
    altfel
        M0 = max(|f0|, |f1|, |f2|) ; estimarea marginii funcției
        hoptim = 2√M0 * eps/M2
    r = h/hoptim ; rata de modificare a pasului
    h = hoptim
până când (k > maxit) sau (r > 0.5 și r < 2)
întoarce (f(x + hoptim) - f0)/hoptim

```

5.4 Derivate de ordin superior

Derivatele de ordin superior pot fi privite ca aplicații recursive ale derivatei de ordinul unu. O posibilitate de calcul a derivatelor de ordin superior este aceea de a folosi abordarea din paragraful 5.3.1. De exemplu, se pot calcula derivatele de ordinul doi în noduri apelând

pregătire_derivare($n, x, dy, d2y$)

iar evaluarea derivatei de ordinul doi într-un punct oarecare se poate face apelând

$rez = \text{interpolare_lpp}(n, x, d2y, xcrt)$

Experimentele numerice arată că, procedând astfel, acuratețea rezultatului se pierde pe măsură ce ordinul derivatei crește.

O altă posibilitate este de a deduce formule de derivare pornind de la polinomul de interpolare inițial. De exemplu, folosind expresia (5.13) rezultă că derivata de ordin 2 a polinomului de interpolare este

$$g''(x) = \frac{2}{(x_{k-1} - x_k)(x_{k-1} - x_{k+1})}y_{k-1} + \frac{2}{(x_k - x_{k-1})(x_k - x_{k+1})}y_k + \quad (5.35)$$

$$+ \frac{2}{(x_{k+1} - x_{k-1})(x_{k+1} - x_k)}y_{k+1}, \quad (5.36)$$

formulă care devine în cazul particular al unui pas uniform $x_{k+1} - x_k = x_k - x_{k-1} = h$

$$g''(x) = \frac{y_{k-1} - 2y_k + y_{k+1}}{h^2}. \quad (5.37)$$

Pentru evaluarea erorii unei astfel de formule se poate folosi seria Taylor cu un număr suplimentar de termeni. Astfel, adunând dezvoltările evaluate în x_{k+1} și x_{k-1} :

$$f(x_{k+1}) = f(x_k) + f'(x_k)h + \frac{f''(x_k)}{2!}h^2 + \frac{f'''(x_k)}{3!}h^3 + O(h^4), \quad (5.38)$$

$$f(x_{k-1}) = f(x_k) - f'(x_k)h + \frac{f''(x_k)}{2!}h^2 - \frac{f'''(x_k)}{3!}h^3 - O(h^4). \quad (5.39)$$

se obține

$$f(x_{k+1}) + f(x_{k-1}) = 2f(x_k) + f''(x_k)h^2 + O(h^4). \quad (5.40)$$

de unde

$$\frac{y_{k+1} - 2y_k + y_{k-1}}{h^2} - f''(x_k) = O(h^2). \quad (5.41)$$

Dacă s-ar dori estimarea unei derivate de ordinul 3, atunci derivând încă o dată (5.35) s-ar obține zero, rezultat inutil. În concluzie, polinomul de interpolare inițial trebuie să aibă un grad suficient de mare pentru ca derivata de ordin superior care interesează să fie nenulă. Acest lucru necesită considerarea a mai mult de 3 puncte consecutive din tabel. Interpolarea de grad mai mare ar putea conduce însă la apariția fenomenului Runge care va afecta și acuratețea derivatelor numerice. De aceea, calculul derivatelor de ordin superior trebuie evitat pe cât posibil în rezolvarea problemelor cu ajutorul calculatorului.

5.5 Derivate parțiale

Estimarea derivatelor parțiale ale unei funcții se face folosind aceleași abordări prezentate pentru derivatele totale. De exemplu, să considerăm o funcție reală ce depinde de două variabile, definită pe un domeniu dreptunghiular

$$f(x, y) : [a, b] \times [c, d] \rightarrow \mathbb{R}.$$

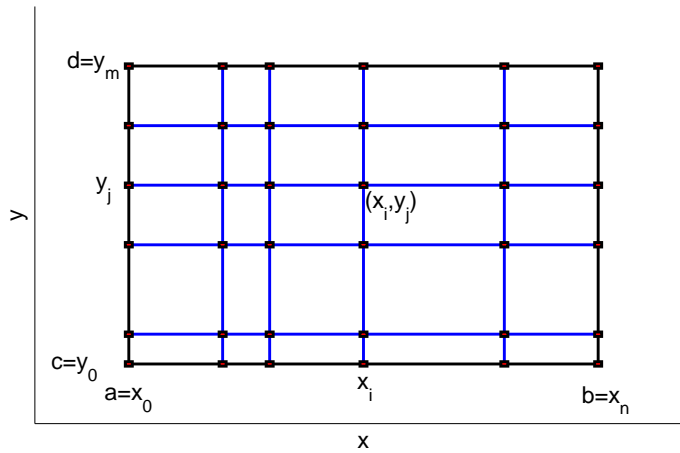


Figura 5.6: Discretizarea domeniului de definiție al funcției $f(x, y) : [a, b] \times [c, d] \rightarrow \mathbb{R}$.

Domeniul de definiție se discretizează acum cu ajutorul a două rețele de puncte, una pe axa Ox, iar cealaltă pe axa Oy:

$$a = x_0, x_1, \dots, x_i, \dots, x_n = b$$

$$c = y_0, y_1, \dots, y_j, \dots, y_m = d$$

Produsul cartezian al acestor două rețele unidimensionale dau rețeaua bidimensională de discretizare a domeniului de calcul (Fig.5.6), un nod al acestei rețele bidimensionale fiind caracterizat de coordonatele (x_i, y_j) , unde $i = 0, n, j = 0, m$. Tabelul de valori care se dă este de această dată bidimensional:

$$f(x_i, y_j) = z_{i,j}, \quad i = 0, n \quad j = 0, m. \quad (5.42)$$

Evaluarea numerică a derivatelor parțiale se reduce la evaluarea numerică a unor derivate obișnuite, deoarece

$$\left. \frac{\partial f}{\partial x} \right|_{x=x_0, y=y_0} = \lim_{x \rightarrow x_0} \frac{f(x, y_0) - f(x_0, y_0)}{x - x_0} = \lim_{x \rightarrow x_0} \frac{f_1(x) - f_1(x_0)}{x - x_0} = \frac{df_1}{dx}, \quad (5.43)$$

unde $f_1(x) = f(x, y_0)$ este o funcție care depinde de o singură variabilă reală x . Similar,

$$\left. \frac{\partial f}{\partial y} \right|_{x=x_0, y=y_0} = \lim_{y \rightarrow y_0} \frac{f(x_0, y) - f(x_0, y_0)}{y - y_0} = \lim_{y \rightarrow y_0} \frac{f_2(y) - f_2(y_0)}{y - y_0} = \frac{df_2}{dy}, \quad (5.44)$$

unde $f_2(y) = f(x_0, y)$ este o funcție care depinde de o singură variabilă reală y .

De aceea, formulele de derivare progresivă de ordinul 1 pentru calculul derivatelor parțiale sunt:

$$\frac{\partial g}{\partial x}(x_i, y_j) = \frac{z_{i+1,j} - z_{i,j}}{x_{i+1} - x_i}, \quad \frac{\partial g}{\partial y}(x_i, y_j) = \frac{z_{i,j+1} - z_{i,j}}{y_{j+1} - y_j}. \quad (5.45)$$

Un alt exemplu îl constituie discretizarea operatorului Laplace

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}. \quad (5.46)$$

În cazul în care pașii de discretizare pe cele două rețele sunt egali: $x_{i+1} - x_i = y_{j+1} - y_j = h$, pentru orice $i = 0, n-1$, $j = 0, m-1$ și folosind formula de derivare în trei puncte de ordinul 2 dată de (5.37), rezultă că o aproximare numerică a Laplaceanului funcției f într-un punct (x_i, y_j) poate fi calculată ca

$$\Delta g = \frac{z_{i-1,j} - 2z_{i,j} + z_{i+1,j}}{h^2} + \frac{z_{i,j-1} - 2z_{i,j} + z_{i,j+1}}{h^2} = \frac{z_{i-1,j} + z_{i+1,j} + z_{i,j-1} + z_{i,j+1} - 4z_{i,j}}{h^2}, \quad (5.47)$$

punctele care intervin în această formulă fiind marcate în Fig.5.7.

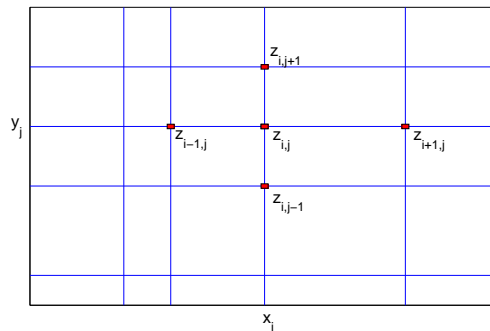


Figura 5.7: În discretizarea operatorului Laplace în punctul (x_i, y_j) apar numai valorile nodurilor marcate.

Dacă astfel de raționamente se folosesc pentru rezolvarea numerică a problemelor de inginerie formulate cu ajutorul unor ecuații cu derivate obișnuite, simple (ODE) sau parțiale (PDE), atunci se spune că pentru rezolvarea problemei se aplică *metoda diferențelor finite*.

5.6 Metoda diferențelor finite

Aplicarea metodei diferențelor finite pentru rezolvarea oricărei ecuații diferențiale presupune efectuarea următorilor pași:

- **Pasul 1:** Se alege o schemă de diferențe finite pentru aproximarea derivatelor din ecuații și se rescrie ecuația ca ecuație cu diferențe finite.
- **Pasul 2:** Se stabilește rețeaua de discretizare și se scrie ecuația discretizată pentru fiecare nod.
- **Pasul 3:** Se rezolvă sistemul de ecuații pentru determinarea valorilor necunoscute în nodurile rețelei.
- **Pasul 4:** Calculul altor valori, în afara celor din noduri, se face prin interpolare.

Se obișnuiește să se spună că pașii 1 și 2 reprezintă etapa de *preprocesare*, pasul 3 reprezintă *rezolvarea*, iar pasul 4 *postprocesarea*. În cele ce urmează vom exemplifica metoda diferențelor finite pentru ecuații diferențiale simple.

5.6.1 Rezolvarea unei ecuații diferențiale ordinare

Cel mai simplu exemplu pe care ni-l putem imagina este cel care conduce matematic la o ecuație diferențială ordinară de ordinul 1. Vom considera un condensator de capacitate $C = 4\mu\text{F}$, inițial descărcat, care se încarcă de la o sursă de tensiune continuă $E = 20\text{ mV}$ printr-un rezistor de rezistență $R = 10\ \Omega$.

Dacă notăm cu i curentul prin circuit și cu u tensiunea la bornele condensatorului, atunci teorema Kirchhoff II scrisă pe bucla circuitului este

$$Ri(t) + u(t) = E. \quad (5.48)$$

Înlocuind în această relație dependența dintre tensiunea și curentul prin condensator

$$i(t) = C \frac{du(t)}{dt}, \quad (5.49)$$

rezultă ecuația diferențială de ordinul unu, satisfăcută de u :

$$RC \frac{du(t)}{dt} + u(t) = E. \quad (5.50)$$

Soluția acestei ecuații este unică deoarece se specifică starea inițială a condensatorului:

$$u(0) = u_0 = 0. \quad (5.51)$$

Avantajul acestui exemplu este acela că are o soluție analitică

$$u(t) = (u_0 - E) \exp(-t/\tau) + E. \quad (5.52)$$

unde

$$\tau = RC, \quad (5.53)$$

este *constantă de timp* a circuitului.

Vom rezolva acum ecuația diferențială cu diferite scheme de derivare numerică. Ecuația (5.50) o rescriem ca

$$\frac{du(t)}{dt} + \frac{1}{\tau}u(t) = \frac{1}{\tau}E. \quad (5.54)$$

Vom urmări calculul numeric în intervalul de timp $[0, t_{max}]$ unde $t_{max} = 10\tau$ într-o rețea echidistantă de N puncte t_k , unde pasul de discretizare este

$$t_{k+1} - t_k = h, \quad \text{pentru } k = 1, \dots, N - 1. \quad (5.55)$$

Vom nota cu u_k valorile discrete obținute prin rezolvare numerică. Ele vor fi aproximații ale mărimii reale u :

$$u_k \approx u(t_k). \quad (5.56)$$

Varianta I - Utilizarea diferențelor finite progresive de ordinul 1

Vom discretiza ecuația (5.54) prin scrierea ei în momentul de timp t_k și folosind pentru derivată o formulă de diferențe finite progresive de ordinul 1:

$$\frac{u_{k+1} - u_k}{h} + \frac{1}{\tau}u_k = \frac{1}{\tau}E, \quad (5.57)$$

de unde

$$u_{k+1} - u_k + \frac{h}{\tau}u_k = \frac{h}{\tau}E. \quad (5.58)$$

Se observă că mărimea u_{k+1} poate fi calculată explicit cu formula

$$u_{k+1} = u_k \left(1 - \frac{h}{\tau}\right) + \frac{h}{\tau}E. \quad (5.59)$$

Implementarea acestei relații este descrisă de pseudocodul următor.

procedură mdf_ode_p1(u0,E,tau,h,N,u)

; rezolvă ecuația $\frac{du(t)}{dt} + \frac{1}{\tau}u(t) = \frac{1}{\tau}E$ cu metoda diferențelor finite

; d/dt se discretizează folosind diferențe progresive de ordinul 1

real u0 ; condiția inițială - dată

real E ; coeficient în ecuație - dată

real tau ; constantă de timp - dată

```

real h                ; pas de discretizare al intervalului de timp - dat
întreg N              ; număr de valori de timp - dat
tablou real u[N]      ; soluția discretă - rezultat
u(1) = u0
pentru k = 1,N-1
    u(k+1) = u(k)*(1-h/tau) + h*E/tau

```

•

retur

Această metodă, cunoscută și sub numele de *metoda Euler explicită*, este instabilă pentru $h > \tau$ (Fig.5.10).

Varianta a II-a - Utilizarea diferențelor finite regresive de ordinul 1

Vom discretiza ecuația (5.54) prin scrierea ei în momentul de timp t_k și folosind pentru derivată o formulă de diferențe finite regresive de ordinul 1:

$$\frac{u_k - u_{k-1}}{h} + \frac{1}{\tau}u_k = \frac{1}{\tau}E, \quad (5.60)$$

de unde

$$u_k - u_{k-1} + \frac{h}{\tau}u_k = \frac{h}{\tau}E. \quad (5.61)$$

Se observă că mărimea u_k poate fi calculată explicit cu formula

$$u_k = \left(u_{k-1} + \frac{h}{\tau}E \right) / \left(1 + \frac{h}{\tau} \right). \quad (5.62)$$

Implementarea acestei relații este descrisă de pseudocodul următor.

```

procedură mdf_ode_r1(u0,E,tau,h,N,u)
; rezolvă ecuația  $\frac{du(t)}{dt} + \frac{1}{\tau}u(t) = \frac{1}{\tau}E$  cu metoda diferențelor finite
; d/dt se discretizează folosind diferențe regresive de ordinul 1
real u0                ; condiția inițială - dată
real E                 ; coeficient în ecuație - dată
real tau               ; constantă de timp - dată
real h                 ; pas de discretizare al intervalului de timp - dat
întreg N              ; număr de valori de timp - dat
tablou real u[N]      ; soluția discretă - rezultat
u(1) = u0
pentru k = 2,N
    u(k) = (h*E/tau + u(k-1))/(1 + h/tau)

```

•

retur

Metoda nu mai este instabilă pentru $h > \tau$ (Fig.5.10). Această metodă este cunoscută și sub numele de *metoda Euler implicită*¹ pentru rezolvarea ecuațiilor diferențiale ordinare.

Varianta a III-a - Utilizarea diferențelor finite centrate de ordinul 2

Vom discretiza ecuația (5.54) prin scrierea ei în momentul de timp t_k și folosind pentru derivată o formulă de diferențe finite centrate de ordinul 2 dată de relația (5.17):

$$\frac{u_{k+1} - u_{k-1}}{2h} + \frac{1}{\tau}u_k = \frac{1}{\tau}E, \quad (5.63)$$

de unde

$$-u_{k-1} + \frac{2h}{\tau}u_k + u_{k+1} = \frac{2h}{\tau}E. \quad (5.64)$$

Se observă că mărimea u_k nu mai poate fi calculată explicit. Relația (5.64) poate fi scrisă pentru $k = 2, \dots, N - 1$ și reprezintă $N - 2$ ecuații cu N necunoscute. Pentru a obține un sistem bine formulat, trebuie să adăugăm încă două relații. Acestea sunt condițiile la capete. La $t = 0$ vom impune condiția inițială:

$$u(1) = u_0, \quad (5.65)$$

iar pentru ultimul punct $k = N$ vom scrie ecuația discretizată dar în care vom folosi pentru derivată o formulă de diferențe regresive de ordinul 2 dată de relația (5.19):

$$\frac{u_{N-2} - 4u_{N-1} + 3u_N}{2h} + \frac{1}{\tau}u_N = \frac{1}{\tau}E, \quad (5.66)$$

de unde

$$u_{N-2} - 4u_{N-1} + \left(3 + \frac{2h}{\tau}\right)u_N = \frac{2h}{\tau}E. \quad (5.67)$$

Implementarea acestei metode este descrisă de pseudocodul următor.

procedură mdf_ode_c2(u0,E,tau,h,N,u)

; rezolvă ecuația $\frac{du(t)}{dt} + \frac{1}{\tau}u(t) = \frac{1}{\tau}E$ cu metoda diferențelor finite

; d/dt se discretizează folosind diferențe centrate de ordinul 2

real u0 ; condiția inițială - dată

real E ; coeficient în ecuație - dată

real tau ; constantă de timp - dată

real h ; pas de discretizare al intervalului de timp - dat

întreg N ; număr de valori de timp - dat

tablou real u[N] ; soluția discretă - rezultat

tablou real A[N,N] ; matricea coeficienților sistemului asamblat - stocata rar

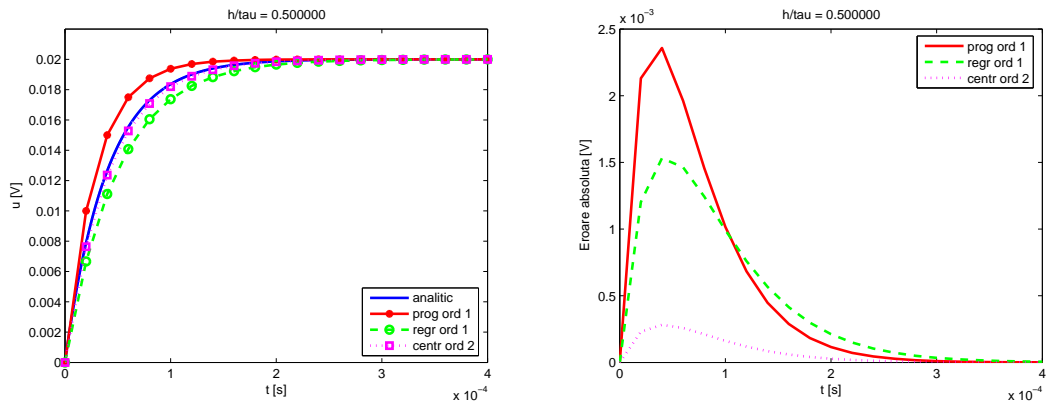
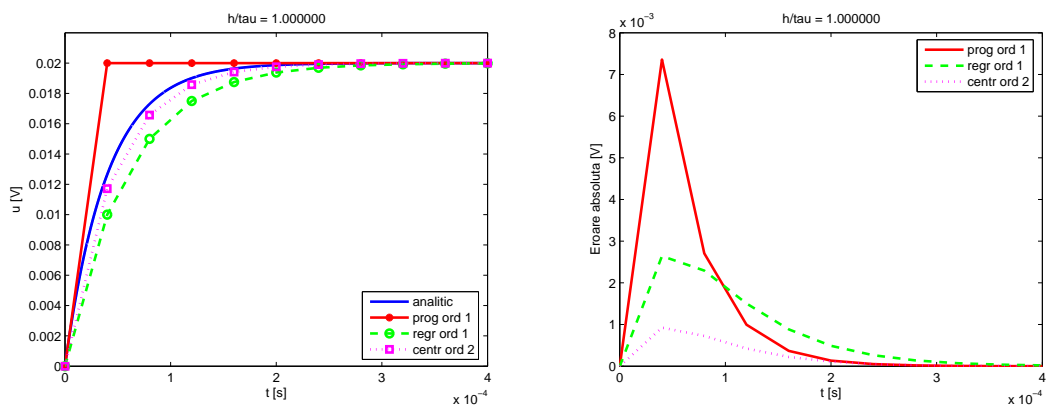
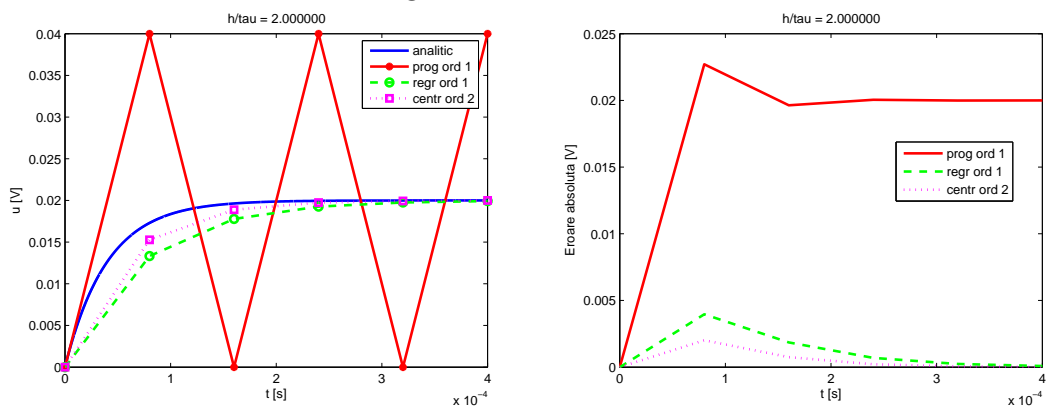
¹Ecuția generală este $dy/dt = f(t, y)$ unde f este o funcție în general neliniară. Folosirea derivatei regresive pentru f conduce la o relație implicită pentru calculul lui u_k , relație ce reprezintă o ecuație neliniară. În cazul particular studiat, f este liniară și de aceea soluția se poate calcula explicit.

```

tablou real tl[N]           ; vectorul termenilor liberi ai sistemului asamblat
A = 0                       ; pseudocod simplificat
tl = 0
A(1,1) = 1                   ; conditia initiala tl(1) = u0
pentru k = 2,N-1           ; noduri interioare
    A(k,k-1) = -1
    A(k,k) = 2*h/tau
    A(k,k+1) = 1
    tl(k) = 2*h*E/tau
•
A(N,N-2) = 1                 ; conditia la tmax = N*h
A(N,N-1) = -4
A(N,N) = (3 + 2*h/tau)
tl(N) = 2*h*E/tau
u = "A-1tl"             ; rezolva sistemul algebric liniar asamblat
retur

```

Metoda nu este instabilă și este mai precisă decât implementările anterioare (Fig.5.8, 5.9, 5.10). Acest lucru era de așteptat deoarece formulele de derivare numerică de ordinul 2 sunt mai precise decât cele de ordinul 1. Creșterea acurateții se face însă pe seama creșterii complexității algoritmului.

Figura 5.8: Cazul $h = \tau/2$.Figura 5.9: Cazul $h = \tau$.Figura 5.10: Cazul $h = 2\tau$.

5.6.2 Rezolvarea unei ecuații cu derivate parțiale de tip eliptic

Exemplul din acest paragraf va conduce la o ecuație diferențială cu derivate parțiale de ordinul 2, de tip eliptic.

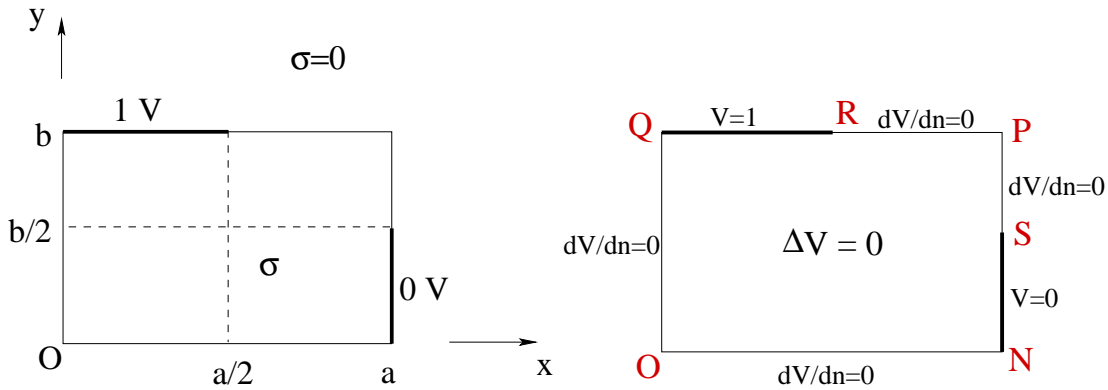


Figura 5.11: Problema 2D de regim electrocinetic: domeniul de calcul (stânga), condiții de frontieră (dreapta).

Vom considera un conductor omogen, de conductivitate σ , situat într-un mediu perfect izolant. Conductorul are o dimensiune după axa Oz mult mai lungă decât celelalte două, astfel încât problema poate fi tratată ca bidimensională. Fig.5.11 reprezintă o secțiune perpendiculară pe direcția dimensiunii foarte mari. Această secțiune este un dreptunghi, de dimensiuni a, b . Conductorul are două borne supraconductoare, una aflată la potențial $V_0 = 1V$, iar cealaltă aflată la potențial nul. Bornele sunt plasate ca în figură. Dorim să reprezentăm liniile echipotențiale și spectrul liniilor de curent. Ar fi interesant să calculăm și mărimi globale cum ar fi puterea disipată în domeniu precum și rezistența pe unitatea de lungime a acestui rezistor. La aceste aspecte vom reveni după ce vom discuta problema integrării numerice.

Înainte de orice, o problemă de câmp electromagnetic trebuie corect formulată, în conformitate cu teorema de unicitate demonstrată pentru regimul studiat.

Problema fiind de regim electrocinetic staționar, formularea se va face în V – potențial electrocinetic, unde $\mathbf{E} = -\text{grad } V$, \mathbf{E} fiind intensitatea câmpului electric.

Ecuația de ordinul doi satisfăcută de potențialul electrocinetic este:

$$-\text{div}(\sigma \text{ grad } V) = 0 \quad (5.68)$$

unde $V : \mathcal{D} \rightarrow \mathbb{R}$ este potențialul definit pe domeniul bidimensional $\mathcal{D} = \text{ONPQ}$. Ecuația (5.68) este o ecuație eliptică, de tip Laplace generalizată.

Deoarece domeniul este omogen (conductivitatea σ are aceeași valoare în orice punct

al domeniului), ecuația se simplifică la o ecuație de tip Laplace:

$$\Delta V = 0, \quad (5.69)$$

unde $\Delta = \text{div grad}$ este operatorul Laplace, care în coordonate carteziene (cazul 2D) are expresia

$$\Delta V = \frac{\partial^2 V}{\partial x^2} + \frac{\partial^2 V}{\partial y^2}. \quad (5.70)$$

În consecință, ecuația diferențială ce trebuie rezolvată este o ecuație cu derivate parțiale:

$$\frac{\partial^2 V}{\partial x^2} + \frac{\partial^2 V}{\partial y^2} = 0, \quad (5.71)$$

unde $V = V(x, y) : \mathcal{D} \rightarrow \mathbb{R}$.

Pentru buna formulare a problemei, trebuie impuse pentru potențial condițiile de frontieră. Frontierele supraconductoare (SN și QR) au potențial constant, și în consecință pe ele trebuie impuse condiții Dirichlet în concordanță cu valorile potențialului. Pe frontierele de lângă izolat (NOQ și RPS), trebuie impuse condiții Neumann nule².

Condițiile de frontieră impuse potențialului sunt:

$$V = V_0 \quad \text{pe QR (Dirichlet)} \quad (5.72)$$

$$V = 0 \quad \text{pe SN (Dirichlet)} \quad (5.73)$$

$$\frac{\partial V}{\partial n} = 0 \quad \text{pe NOQ} \cup \text{RPS (Neumann)} \quad (5.74)$$

Și pentru o astfel de problemă se pot găsi rezolvări analitice bazate, de exemplu, pe metoda separării variabilelor sau pe aproximarea liniilor de câmp.

Față de cazul unidimensional discutat în exemplul anterior, aici rețeaua de discretizare este una bidimensională, iar în loc de aproximarea unei derivate simple de ordinul 1, acum avem de aproximat derivate parțiale de ordinul 2.

Rețeaua de discretizare bidimensională poate fi privită ca fiind obținută din produsul cartezian dintre o rețea de discretizare pe axa Ox și o rețea de discretizare pe axa Oy:

$$0 = x_1 < x_2 < \dots < x_{n_x} = a,$$

$$0 = y_1 < y_2 < \dots < y_{n_y} = b,$$

astfel încât, un nod al rețelei este identificat prin poziția proiecțiilor sale pe cele două axe:

$$(x_i, y_j) \quad i = 1, \dots, n_x, \quad j = 1, \dots, n_y.$$

²În izolat nu există curent de conducție, în consecință $\mathbf{J} \cdot \mathbf{n} = 0$ la frontiera cu izolantul, deci $-\sigma \text{grad } V \cdot \mathbf{n} = 0$, echivalent cu $-\sigma \frac{\partial V}{\partial n} = 0$ pe acele frontiere.

Metoda diferențelor finite va determina valorile aproximative ale potențialului în nodurile acestui grid:

$$V(x_i, y_j) \stackrel{not}{=} V_{i,j} \quad i = 1, n_x, \quad j = 1, n_y. \quad (5.75)$$

Forma discretizată (aproximativă) a ecuației potențialului este obținută aproximând prin diferențe finite ecuația (5.71). Forma specifică depinde de tipul nodului: interior sau pe frontiera Neumann. Nodurile de pe frontiera Dirichlet nu ridică probleme, ecuația asociată unui astfel de nod constând în simpla atribuire a valorii potențialului nodului respectiv.

Ecuația asociată unui nod interior:

Ecuația aproximativă pentru un nod interior se deduce înlocuind derivatele parțiale după x și după y cu formule de tipul (5.37).

Pentru a simplifica scrierea formulelor, vom presupune că cele două griduri pe Ox și, respectiv, Oy sunt uniforme, cu pași h_x și respectiv h_y . Derivatele parțiale se vor înlocui cu:

$$\frac{\partial^2 V}{\partial x^2}(x_i, y_j) = \frac{V_{i+1,j} - 2V_{i,j} + V_{i-1,j}}{h_x^2}, \quad (5.76)$$

$$\frac{\partial^2 V}{\partial y^2}(x_i, y_j) = \frac{V_{i,j+1} - 2V_{i,j} + V_{i,j-1}}{h_y^2}. \quad (5.77)$$

Ecuația discretizată asociată unui nod interior devine

$$\frac{V_{i+1,j} - 2V_{i,j} + V_{i-1,j}}{h_x^2} + \frac{V_{i,j+1} - 2V_{i,j} + V_{i,j-1}}{h_y^2} = 0, \quad (5.78)$$

sau

$$2 \left(\frac{1}{h_x^2} + \frac{1}{h_y^2} \right) V_{i,j} - \frac{1}{h_x^2} V_{i+1,j} - \frac{1}{h_x^2} V_{i-1,j} - \frac{1}{h_y^2} V_{i,j+1} - \frac{1}{h_y^2} V_{i,j-1} = 0. \quad (5.79)$$

Relația (5.79) arată că potențialul într-un nod interior este o combinație liniară a potențialelor nodurilor învecinate. Dacă pașii de discretizare pe cele două axe ar fi egali ($h_x = h_y$), atunci potențialul într-un nod interior este media aritmetică a potențialelor celor patru noduri vecine.

În cazul unor rețele neuniforme, folosind notațiile din Fig.5.12-dreapta, ecuația discretizată asociată unui nod interior este

$$\begin{aligned} V_O \left(\frac{1}{h_B h_D} + \frac{1}{h_A h_C} \right) - V_A \frac{1}{h_A (h_A + h_C)} - V_B \frac{1}{h_B (h_B + h_D)} - \\ - V_C \frac{1}{h_C (h_A + h_C)} - V_D \frac{1}{h_D (h_B + h_D)} = 0, \end{aligned} \quad (5.80)$$

unde $h_A = \|OA\|$, $h_B = \|OB\|$, $h_C = \|OC\|$, $h_D = \|OD\|$.

Ecuația asociată unui nod pe frontieră dreaptă:

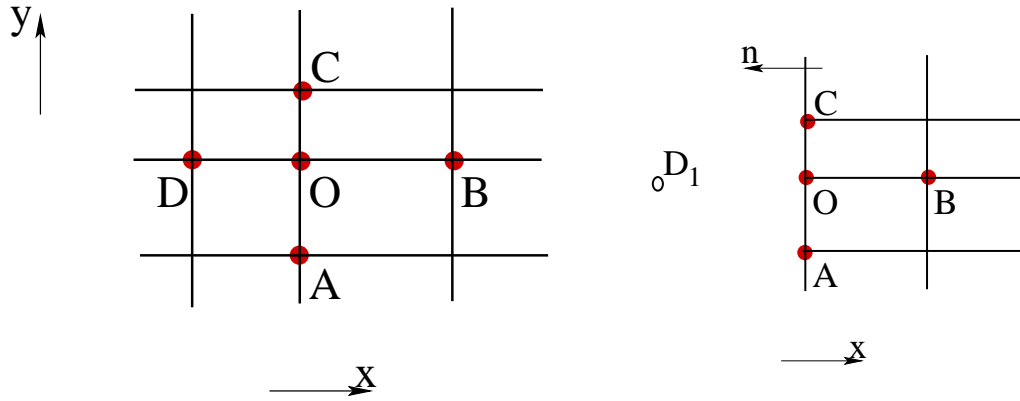


Figura 5.12: Nodurile marcate intervin în scrierea ecuațiilor: stânga - cazul unui nod interior; dreapta - cazul unui nod pe frontieră Neumann dreaptă.

În cazul unui punct pe frontieră O (fig. 5.12- dreapta), să notăm cu $g = -\sigma \frac{\partial V}{\partial n}$ condiția de frontieră Neumann și cu g_O valoarea ei medie în punctul O:

$$g_O = \frac{g_{OA}h_A + g_{OC}h_C}{h_A + h_C} \quad (5.81)$$

unde g_{OA} este condiția de frontieră asociată segmentului OA, iar g_{OC} este condiția de frontieră asociată segmentului OC.

Pentru deducerea ecuației aproximative pentru punctul O, vom considera un nod “fantomă” D_1 , situat în exterior, la o distanță $h_D = \|OD_1\|$ de punctul O. Pentru simplitate putem considera $h_D = h_B$.

Din condiția de frontieră Neumann, deducem valoarea potențialului fantomă în funcție de valoarea acestei condiții. Pentru cazul din figură (normala exterioară în sens contrar axei Ox) rezultă:

$$\frac{\partial V}{\partial x} = \frac{g}{\sigma}. \quad (5.82)$$

Scriind relația aproximativă (5.17) pentru condiția de frontieră Neumann, rezultă că:

$$V_{D_1} = V_B - 2h_B \frac{g_O}{\sigma}. \quad (5.83)$$

Pentru nodul O se scrie acum o ecuație de tipul (5.80), ca pentru un nod interior și, înlocuind expresia (5.83), rezultă ecuația finală:

$$V_O \left(\frac{1}{h_B^2} + \frac{1}{h_A h_C} \right) - V_A \frac{1}{h_A(h_A + h_C)} - V_B \frac{1}{h_B^2} - V_C \frac{1}{h_C(h_A + h_C)} = -\frac{g_O}{\sigma h_B}. \quad (5.84)$$

În cazul unei condiții Neumann nule, relația devine:

$$V_O \left(\frac{1}{h_B^2} + \frac{1}{h_A h_C} \right) - V_A \frac{1}{h_A(h_A + h_C)} - V_B \frac{1}{h_B^2} - V_C \frac{1}{h_C(h_A + h_C)} = 0. \quad (5.85)$$

Ecuția asociată unui nod pe frontieră - colț exterior:

Nodurile situate în colțuri reprezintă o problemă pentru metoda diferențelor finite. Aceasta deoarece pentru un colț nu poate fi definită direcția normalei³. Dacă notăm $g = -\sigma \frac{\partial V}{\partial n}$, și notăm cu g_{OA} valoarea lui g pe segmentul OA (fără capătul O) și cu g_{OB} valoarea lui g pe segmentul OB (fără capătul O), atunci vom presupune că în nodul O (fig. 5.13) este cunoscută derivata după o direcție care rezultă din prelungirea în O a valorilor funcției g .

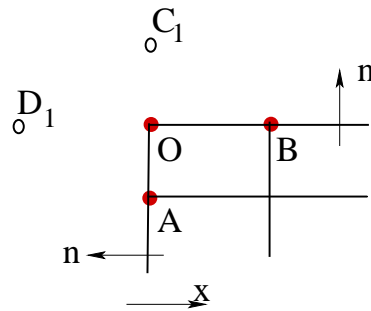


Figura 5.13: Nodurile marcate intervin în scrierea ecuației unui nod pe frontieră colț exterior.

Pentru deducerea ecuației vom considera două noduri “fantomă” C_1 și D_1 , situate în exterior, la distanțele $h_D = \|OD_1\|$ și $h_C = \|OC_1\|$ de punctul O. Pentru simplitate putem considera $h_D = h_B$ și $h_C = h_A$.

Din condiția de frontieră Neumann pe segmentul OB, scrisă în nodul O, putem deduce valoarea potențialului fantomă C_1 :

$$V_{C_1} = V_A - 2h_A \frac{g_{OB}}{\sigma}. \quad (5.86)$$

Din condiția de frontieră Neumann pe segmentul OA, scrisă în nodul O, putem deduce valoarea potențialului fantomă D_1 :

$$V_{D_1} = V_B - 2h_B \frac{g_{OA}}{\sigma}. \quad (5.87)$$

Pentru nodul O se scrie acum o ecuație de tipul (5.80), ca pentru un nod interior și, înlocuind expresiile (5.86) și (5.87), rezultă ecuația finală:

$$V_O \left(\frac{1}{h_A^2} + \frac{1}{h_B^2} \right) - V_A \frac{1}{h_A^2} - V_B \frac{1}{h_B^2} = -\frac{g_{OB}}{\sigma h_A} - \frac{g_{OA}}{\sigma h_B}. \quad (5.88)$$

În cazul unor condiții Neumann nule, relația devine:

$$V_O \left(\frac{1}{h_A^2} + \frac{1}{h_B^2} \right) - V_A \frac{1}{h_A^2} - V_B \frac{1}{h_B^2} = 0. \quad (5.89)$$

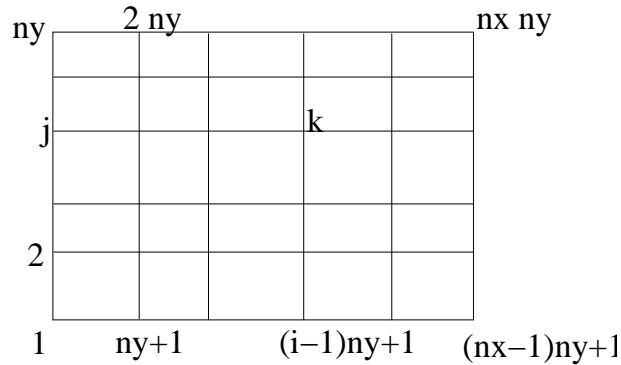
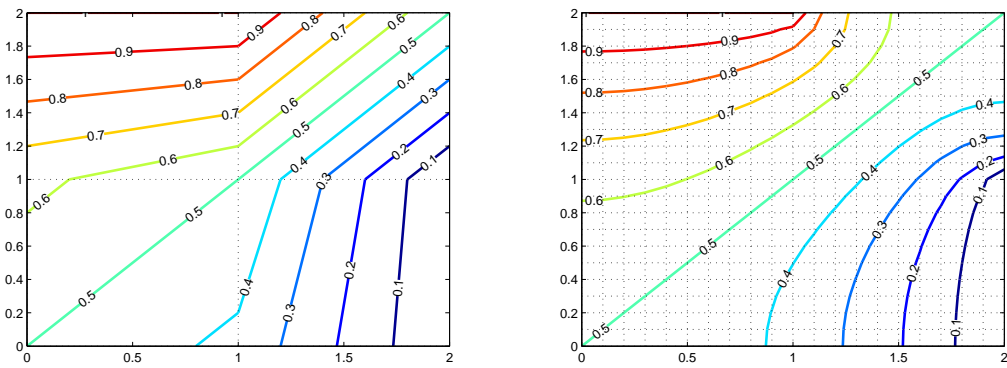


Figura 5.14: Numerotarea nodurilor.

În concluzie, discretizarea problemei conduce la un sistem de ecuații algebrice liniar, prin a cărui rezolvare vom obține valorile potențialelor în nodurile rețelei de discretizare. Pentru asamblarea acestui sistem cu ajutorul unui algoritm, este util ca nodurile să fie numerotate astfel încât necunoscutele să reprezinte un vector, nu o matrice cum sugerează notațiile de până acum. Vom numerota nodurile de jos în sus și de la stânga la dreapta, ca în Fig.5.14. Se poate verifica ușor că relația între numărul nodului k și pozițiile proiecțiilor sale pe cele două axe i și j este

$$k = (i - 1)n_y + j, \quad i = 1, \dots, n_x, \quad j = 1, \dots, n_y. \quad (5.90)$$

unde n_x și n_y reprezintă numărul de puncte de discretizare pe axele Ox, respectiv Oy.

Figura 5.15: Linii echipotențiale pentru o rețea cu $n_x = n_y = 3$ (stânga) și o rețea cu $n_x = n_y = 21$ (dreapta).

³În realitate, colțul reprezintă o modelare brutală a realității. Trecerea de la o suprafață la alta se face prin racordări.

5.6.3 Rezolvarea unei ecuații cu derivate parțiale de tip hiperbolic

În acest paragraf vom rezolva cu metoda diferențelor finite cel mai simplu exemplu de ecuație hiperbolică, și anume *ecuația undei scalare*:

$$\frac{\partial^2 u}{\partial t^2} = v^2 \frac{\partial^2 u}{\partial x^2}, \quad (5.91)$$

unde mărimea scalară necunoscută $u(x, t) : [0, a] \times [0, T] \rightarrow \mathbb{R}$ depinde de o singură coordonată spațială și de timp, iar v este o constantă cunoscută. Soluția acestei ecuații este o *undă* care se propagă cu viteză v .

Buna formulare a acestei probleme necesită impunerea

- *condiției inițiale* $u(x, 0) = h_0(x)$,
- *condițiilor la capete* - relații în care intervin mărimile $u(0, t) = h_1(t)$ și $u(a, t) = h_2(t)$.

Se poate demonstra că soluția ecuației undei scalare (5.91) este

$$u(x, t) = ud(x, t) + ui(x, t) + U_0, \quad (5.92)$$

unde $ud(x, t)$ se numește *undă directă* și este o funcție care se poate scrie sub forma

$$ud(x, t) = f(x - vt), \quad (5.93)$$

iar $ui(x, t)$ se numește *undă inversă* și este o funcție care se poate scrie sub forma

$$ui(x, t) = g(x + vt). \quad (5.94)$$

Funcțiile f și g rezultă în mod univoc, din impunerea condițiilor inițiale și condițiilor la capete. Mărimea $f(x - vt)$ este o undă directă deoarece o anumită valoare a acestei funcții, într-un anumit punct x și într-un anumit moment de timp t , se va regăsi după un interval de timp Δt în punctul $x + \Delta x$, unde $\Delta x = v\Delta t$. Unda directă se propagă deci în sensul pozitiv al axei Ox. Un raționament similar se poate face pentru unda inversă.

În conceperea algoritmului bazat pe metoda diferențelor finite avem nevoie de o discretizare spațială și de una temporală. Pentru a simplifica prezentarea, vom presupune că ambele discretizări sunt uniforme și vom nota cu Δx pasul discretizării spațiale și cu Δt pasul discretizării temporale. Vom nota cu N numărul de puncte de discretizare spațiale și cu M numărul de pași de timp simulați. În consecință, timpul maxim de simulare este $T = M\Delta t$.

Ca exemplu, vom considera condiții inițiale nule și condiții Dirichlet la ambele capete⁴:

- *condiția inițială* nulă $u(x, 0) = h_0(x) = 0$,
- *condiția la capătul din stânga* - excitația cu un impuls Gauss: $u(0, t) = h_1(t) = \exp(-(t - T/10)^2 / (2 * (T/50)^2))$
- *condiția la capătul din dreapta* $u(a, t) = h_2(t) = 0$.

Mărimea $u(x, t)$ este discretizată atât în spațiu cât și în timp. Prin rezolvarea cu metoda diferențelor finite, vom obține aproximații pentru valorile reale $u(x_k, t_j)$.

$$u(x_k, t_j) \approx u_k^{(j)}, \quad k = 1, \dots, N, \quad j = 1, \dots, M. \quad (5.95)$$

Alegând pentru derivatele ce intervin în (5.91) formule de derivare ce provin din polinomul de interpolare de ordin doi, rezultă următoarea relație discretizată:

$$\frac{u_k^{(j-1)} - 2u_k^{(j)} + u_k^{(j+1)}}{(\Delta t)^2} = v^2 \frac{u_{k-1}^{(j)} - 2u_k^{(j)} + u_{k+1}^{(j)}}{(\Delta x)^2}, \quad (5.96)$$

de unde rezultă că mărimea u , într-un anumit punct k și la un anumit moment de timp $j + 1$, depinde explicit de valoarea în acel punct și în cele învecinate în momentul de timp anterior j și de valoarea în acel punct în momentul $j - 1$:

$$u_k^{(j+1)} = \left(\frac{v\Delta t}{\Delta x} \right)^2 \left(u_{k-1}^{(j)} - 2u_k^{(j)} + u_{k+1}^{(j)} \right) + 2u_k^{(j)} - u_k^{(j-1)}, \quad k = 2, \dots, N-1, \quad j = 0, \dots, M-1. \quad (5.97)$$

Momentul -1 îl vom considera identic cu momentul inițial 0.

Rezultatul simulării cu metoda diferențelor finite este prezentat în Fig.5.16. În momentul reflexiei apare o undă inversă deoarece condiția Dirichlet la capătul din dreapta impune ca mărimea să fie zero. În consecință, unda inversă la capătul din dreapta este exact opusul undei directe, și ea se propagă în sens contrar axei Ox.

Atunci când o problemă necesită atât o discretizare spațială cât și una temporală, pentru ca soluția numerică să fie stabilă este necesar să fie îndeplinită *condiția lui Courant*

$$|v|\Delta t \leq \Delta x. \quad (5.98)$$

⁴În problemele în care apare propagare, este util uneori ca frontiera domeniului spațial să fie modelată ca o frontieră absorbantă, "invizibilă" din punct de vedere al propagării mărimilor în domeniul spațial modelat. Condiția de frontieră absorbantă pentru capătul din dreapta al problemei studiate este

$$\frac{\partial u}{\partial t} + v \frac{\partial u}{\partial x} = 0.$$

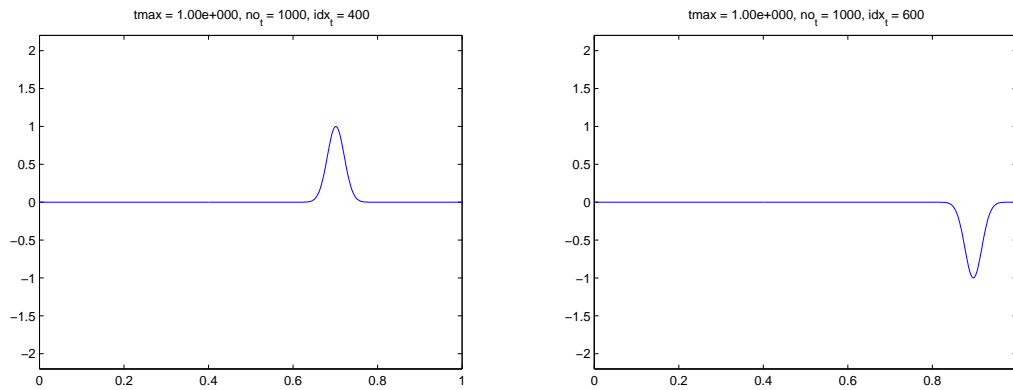


Figura 5.16: Propagarea unui impuls Gauss (stânga) și rezultatul reflexiei după atingerea unei frontiere pe care s-a impus condiție Dirichlet nulă (dreapta).

În concluzie, în rezolvarea ecuațiilor diferențiale cu metoda diferențelor finite, trebuie avute în vedere aspectul stabilității algoritmului și aspectul instabilității numerice datorate intrării în zona în care predomină erorile de rotunjire. Stabilitatea se analizează diferit dacă ecuația este cu derivate ordinare (ODE) sau cu derivate parțiale (PDE), dar această analiză depășește scopul acestui capitol introductiv despre problema derivării numerice. Formulată pe scurt, trebuie ca pașii de discretizare trebuie să fie suficient de mici pentru ca algoritmul să fie stabil, dar nu prea mici ca să nu se intre în zona erorilor de rotunjire. De asemenea, un algoritm eficient ar trebui să poată să folosească pași de discretizare neuniformi, adaptați tipului de variație a soluției.

Anexa A

Fișiere utile pentru implementarea în C

```
// fisier nrutil.h
#include<stdio.h>
#include<stdlib.h>
#include<math.h>
#include<malloc.h>
#include<time.h>

typedef float **MATRIX;
typedef float *VECTOR;
typedef int **IMATRIX;
typedef int *IVECTOR;

void nrerror (char error_text[]);
VECTOR vector (int nl, int nh);
IVECTOR ivector (int nl, int nh);
MATRIX matrix (int nrl, int nrh, int ncl, int nch);
IMATRIX imatrix (int nrl, int nrh, int ncl, int nch);
void free_vector (VECTOR v, int nl, int nh);
void free_ivector (IVECTOR v, int nl, int nh);
void free_matrix (MATRIX m, int nrl, int nrh, int ncl, int nch);
void free_imatrix (IMATRIX m, int nrl, int nrh, int ncl, int nch);

// fisier nrutil.c
#include "nrutil.h"

void
```

```
nrrerror (char error_text[])
```

```
{
    fprintf (stderr, "Run-time error...\n");
    fprintf (stderr, "%s\n", error_text);
    fprintf (stderr, "...now exiting to system...\n");
    exit (1);
}
```

```
VECTOR
```

```
vector (int nl, int nh)
```

```
{
    VECTOR v;

    v = (float *) malloc ((unsigned) (nh - nl + 1) * sizeof (float));
    if (!v)
        nrrerror ("allocation failure in vector()");
    return v - nl;
}
```

```
IVECTOR
```

```
ivector (int nl, int nh)
```

```
{
    IVECTOR v;

    v = (int *) malloc ((unsigned) (nh - nl + 1) * sizeof (int));
    if (!v)
        nrrerror ("allocation failure in ivector()");
    return v - nl;
}
```

```
MATRIX
```

```
matrix (int nrl, int nrh, int ncl, int nch)
```

```
{
    int i;
    MATRIX m;

    m = (float **) malloc ((unsigned) (nrh - nrl + 1) * sizeof (float *));
    if (!m)
        nrrerror ("allocation failure 1 in matrix()");
    m -= nrl;
}
```

```
for (i = nrl; i <= nrh; i++)
{
    m[i] = (float *) malloc ((unsigned) (nch - ncl + 1) * sizeof (float));
    if (!m[i])
nrerror ("allocation failure 2 in matrix()");
    m[i] -= ncl;
}
return m;
}
```

IMATRIX

```
imatix (int nrl, int nrh, int ncl, int nch)
```

```
{
    int i;
    IMATRIX m;

    m = (int **) malloc ((unsigned) (nrh - nrl + 1) * sizeof (int *));
    if (!m)
        nrerror ("allocation failure 1 in imatrix()");
    m -= nrl;

    for (i = nrl; i <= nrh; i++)
    {
        m[i] = (int *) malloc ((unsigned) (nch - ncl + 1) * sizeof (int));
        if (!m[i])
nrerror ("allocation failure 2 in imatrix()");
        m[i] -= ncl;
    }
    return m;
}
```

void

```
free_vector (VECTOR v, int nl, int nh)
```

```
{
    free ((char *) (v + nl));
}
```

void

```
free_ivector (IVECTOR v, int nl, int nh)
```

```
{
    free ((char *) (v + n1));
}
```

```
void
free_matrix (MATRIX m, int nrl, int nrh, int ncl, int nch)
{
    int i;

    for (i = nrh; i >= nrl; i--)
        free ((char *) (m[i] + ncl));
    free ((char *) (m + nrl));
}
```

```
void
free_imatrix (IMATRIX m, int nrl, int nrh, int ncl, int nch)
{
    int i;

    for (i = nrh; i >= nrl; i--)
        free ((char *) (m[i] + ncl));
    free ((char *) (m + nrl));
}
```

Referințe

- [1] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine și H. Van der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods, 2nd Edition*. SIAM, Philadelphia, PA, 1994. http://www.netlib.org/linalg/html_templates/Templates.html.
- [2] T. A. Davis. *Direct Methods for Sparse Linear Systems*. SIAM, Philadelphia, 2006. Part of the SIAM Book Series on the Fundamentals of Algorithms, <http://www.cise.ufl.edu/research/sparse/CSparse/>.
- [3] Gene Golub și Charles van Loan. *Matrix Computations*. The Johns Hopkins University Press, 1996.
- [4] D. Ioan, I. Munteanu, B. Ionescu, M. Popescu, R. Popa, M. Lăzărescu și G. Ciuprina. *Metode numerice în ingineria electrică*. MatrixROM, București, 1998.
- [5] Valeriu Iorga și Boris Jora. *Metode numerice*. Editura Albastră, 2004. Cluj-Napoca.
- [6] B. Jora, C. Popeea și S. Barbulea. *Metode de calcul numeric în automatică*. Ed. Enciclopedică, 1996. București.
- [7] Cleve Moler. *Numerical Computing with MATLAB*. SIAM, 2004. <http://www.mathworks.com/moler/>.
- [8] A.M. Morega. *Modelarea numerică pentru probleme la limită în inginerie*. MatrixROM, București, 1998.
- [9] William Press, Saul Teukolsky, William Vetterling și Brian Flannery. *Numerical Recipes in C*. Cambridge University Press, Cambridge, UK, a doua ed., 1992. <http://www.nrbook.com/a/bookcpdf.php>.
- [10] T.H. Cormen C.E. Leiserson R.R. Rivest. *Introduction to algorithms*. MIT Press and McGraw-Hill, 1990. <http://www.cs.dartmouth.edu/thc/CLRS3e/>.

-
- [11] Yousef Saad. *Iterative Methods for Sparse Linear Systems*. PWS Publishing, 1996. <http://www-users.cs.umn.edu/~saad/books.html>.
- [12] J.R. Shewchuk. An introduction to the conjugate gradient method without the agonizing pain. Rap. tehn., School of Computer Science, Carnegie Mellon University, Pittsburgh, 1994. <http://www.cs.cmu.edu/quake-papers/painless-conjugate-gradient.pdf>.
- [13] G.W. Stewart. Building an old-fashioned sparse solver. Rap. tehn., University of Maryland, Institute of Advanced Computer Studies, Department of Computer Science, 2003. <http://http://www.cs.umd.edu/stewart/>.
- [14] Anca Tomescu, I.B.L. Tomescu și F.M.G. Tomescu. *Modelarea numerică a câmpului electromagnetic. Preliminarii. Câmpuri statice și staționare*. MatrixRom, 2003. București.
- [15] F.M.G. Tomescu. *Fundamentals of Electrical Engineering, Electric Circuits*. MatrixRom, 2011. București.
- [16] Lloyd Trefethen și David Bau. *Numerical Linear Algebra*. SIAM - Society for Industrial and Applied Philadelphia, 1997.
- [17] Henk A. van der Vorst. *Iterative Krylov Methods for Large Linear Systems*. Cambridge University Press, 2003.

Lucrarea de față este o primă parte a unei prezentări articulate a metodelor fundamentale de calcul numeric de interes pentru ingineri, cu aplicații directe în ingineria electrică. Acest prim volum este structurat urmând o alternare și întrepătrundere a unei abordări formale a subiectului cu exemplificări antrenante și ilustrative, într-o manieră clară, ușor accesibilă, fără a renunța însă la rigoarea adecvată materiei prezentate. Lucrarea include extinderi către subiecte particulare de interes curent, și se constituie într-o construcție completă și riguroasă, în care domeniul abordat este prezentat clar, cu ilustrări pertinente, acompaniate de o prezentare grafică excelentă. Ca abordare nouă, deosebită, de o reală valoare pedagogică, a unei materii de mare utilitate în activitatea inginerului modern, lucrarea va fi bine-venită pentru cei interesați în aplicarea metodelor de calcul numeric la rezolvarea problemelor științifice și ingineresti.

București, 22 septembrie 2013

Prof.dr.ing F.M.G. TOMESCU,
Facultatea de Inginerie Electrică,
Universitatea Politehnica – București

Aveți în mână o lucrare referitoare la Metodele numerice, scrisă de o ingineră și adresată inginerilor de formație electrică. Asta nu o face mai puțin riguroasă, dar precis o face mai orientată spre aplicații, deci mai utilă pentru cei care vor să depășească bariera înțelegerii conceptuale și vor să intre în domeniul dezvoltării de noi proceduri informatice sau măcar în cel al (re)folosirii lor profesionale. Autoarea continuă și consolidează abordarea tradițională dezvoltată în Laboratorul de Metode Numerice (LMN) din Facultatea de Inginerie Electrică. Aceasta se bazează pe observația că gândirea algoritmică este o deprindere clar diferită, atât de gândirea matematică – deductivă, cât și de cea științifică - inductivă. Înțelegerea algoritmilor și concepția acestora presupune un mod specific de gândire, iar dacă problemele care trebuie rezolvate nu sunt doar de tipul așezării reginelor pe tabla de șah, ci sunt unele numerice, atunci soluția trebuie bazată pe fundamente matematice solide. Toate acestea ar rămâne pur speculative, dacă nu se au în vedere și aplicațiile ingineresti. Aceasta este scopul lucrării, autoarea fiind nu numai profesor, ci și un talentat cercetător.

Lucrarea, este rezultatul unei îndelungi experiențe de lucru cu studenții și masteranzii, dar este marcată și de activitatea de cercetare și dezvoltare. Așa se explică faptul că în această carte sunt prezentate doar aspecte fundamentale, simple, dar profunde. Sunt evitate cele complicate sau foarte avansate, pentru ca țelul lucrării este de a transmite un mod de gândire și mai puțin o cantitate mare de informații. După ce acesta este însușit, lucrurile avansate pot fi găsite pe net sau în monografiile. Dar fără acest prim pas, n-ai ști nici ce și nici cum să cauți. Este un model de curs universitar, pentru primul ciclu universitar, care te ajută să poți învăța ulterior singur. Recomand cu căldură, celor care îți doresc să devină adevărați profesioniști în ramuri ale ingineriei electrice, în cel mai larg sens să parcurgă aceasta lucrare, în toate detaliile ei.

București, 9 octombrie 2013

Prof.dr.ing Daniel IOAN,
Seful Laboratorului de Metode Numerice
Facultatea de Inginerie Electrică,
Universitatea Politehnica – București



9786062500085