

Analiză numerică

Radu Tiberiu Trîmbițaș

Prefață

Lloyd N. Trefethen a propus următoarea definiție a Analizei numerice:

Analiza numerică este studiul algoritmilor pentru rezolvarea problemelor matematicii continue.

Cuvântul cheie este acela de *algoritmi*. Deși foarte multe lucrări *nu* evidențiază acest lucru, în centrul atenției Analizei numerice stă proiectarea și analiza algoritmilor de rezolvare a unei anumite clase de probleme.

Problemele sunt cele din *matematica continuă*. „Continuă” înseamnă aici faptul că variabilele ce intervin aici sunt reale sau complexe; opusul lui continuu este discret. Pe scurt, am putea spune că Analiza numerică este Algoritmica continuă, în contrast cu Algoritmica clasică, care este Algoritmica discretă.

Este clar că deoarece numerele reale și complexe nu pot fi reprezentate exact în calculator, ele trebuie să fie approximate printr-o reprezentare finită. Din acest moment intervin erorile de rotunjire și iar este clar că studiul lor este unul din obiectivele importante ale Analizei numerice. Au existat și mai există încă opinii care susțin că acesta este cel mai important obiectiv. Un argument în sprijinul acestei idei, înafară de omniprezența erorilor, este dat de metodele exacte de rezolvare a sistemelor de ecuații liniare, cum ar fi eliminarea gaussiană.

Dar, cele mai multe probleme ale matematicii continue nu pot fi rezolvate prin algoritmi așa-zisi finiți, chiar presupunând prin absurd că am lucra în aritmetică cu precizie exactă. Un prim exemplu care se poate da aici este problema rezolvării unei ecuații polinomiale. Acest lucru se evidențiază la problemele de valori și vectori proprii, dar apar în orice problemă ce presupune termeni neliniari sau derivate – determinarea zerourilor, cuadraturi, ecuații diferențiale și integrale, optimizare ș.a.m.d.

Chiar dacă erorile de rotunjire ar dispărea, Analiza numerică ar rămâne. Aproximarea numerelor, obiectivul aritmeticii în virgulă flotantă, este un subiect dificil și obo-

sitor. Un obiectiv mai profund al Analizei numerice este aproximarea necunoscutelor, nu a cantităților cunoscute. Scopul este convergența rapidă a aproximațiilor și mândria specialiștilor din acest domeniu este aceea că, pentru multe probleme s-au inventat algoritmi care converg extrem de rapid. Dezvoltarea pachetelor de calcul simbolic a micșorat importanța erorilor de rotunjire, fără a micșora importanța vitezei de convergență a algoritmilor.

Definiția de mai sus nu surprinde câteva aspecte importante : că acești algoritmi sunt implementați pe calculatoare, a căror arhitectură poate fi o parte importantă a problemei; că fiabilitatea și eficiența sunt obiective supreme; că anumiți specialiști în analiza numerică scriu programe și alții demonstrează teoreme¹; și lucrul cel mai important, că toată munca este *aplicată*, aplicată cu succes la mii de aplicații, pe milioane de computere din toată lumea. „Problemele matematicii continue” sunt problemele pe care știința și ingineria sunt construite; fără metode numerice, știința și ingineria, așa cum sunt ele practicate astăzi ar ajunge repede în impas. Ele sunt de asemenea problemele care au preocupat cei mai mulți matematicieni de la Newton până azi. La fel ca și cei ce se ocupă de matematica pură, specialiștii în Analiza numerică sunt moștenitorii marii tradiții a lui Euler, Lagrange, Gauss și a altor mari matematicieni.

Radu Tiberiu Trîmbițaș
Cluj-Napoca, iulie 2003

¹există mulți specialiști foarte buni care le fac pe amândouă

Cuprins

1. Elemente de Teoria erorilor și aritmetică în virgulă flotantă	1
1.1. Probleme numerice	2
1.2. Măsuri ale erorii	4
1.3. Eroarea propagată	5
1.4. Reprezentarea în virgulă flotantă	6
1.4.1. Parametrii reprezentării	6
1.4.2. Anularea	8
1.5. Standardizarea reprezentării în virgulă flotantă	10
1.5.1. Cantități speciale	11
1.6. Condiționarea unei probleme	11
1.7. Condiționarea unui algoritm	14
1.8. Eroarea globală	15
1.9. Probleme prost condiționate și probleme incorect puse	16
1.10. Stabilitatea	16
1.10.1. Notății asimptotice	16
1.10.2. Precizie și stabilitate	18
1.10.3. Analiza regresivă a erorilor	20
2. Rezolvarea numerică a sistemelor de ecuații algebrice liniare	23
2.1. Elemente de Analiză matricială	24
2.2. Condiționarea unui sistem liniar	30
2.3. Metode exacte	34
2.3.1. Metoda eliminării a lui Gauss	34
2.4. Metode bazate pe factorizare	39
2.4.1. Descompunerea LU	39

2.4.2.	Descompunere LUP	41
2.4.3.	Factorizarea Cholesky	43
2.4.4.	Descompunerea QR	45
2.5.	Rafinarea iterativă	52
2.6.	Algoritmul lui Strassen pentru înmulțirea matricelor	52
2.7.	Rezolvarea iterativă a sistemelor algebrice liniare	56
3.	Aproximarea funcțiilor	65
3.1.	Aproximație prin metoda celor mai mici pătrate	68
3.1.1.	Produce scalare	69
3.1.2.	Ecuatiile normale	70
3.1.3.	Eroarea în metoda celor mai mici pătrate. Convergența	73
3.2.	Exemple de sisteme ortogonale	76
3.2.1.	Exemple de polinoame ortogonale	79
3.3.	Interpolare polinomială	85
3.3.1.	Spațiul $H^n[a, b]$	85
3.3.2.	Interpolare Lagrange	87
3.3.3.	Interpolare Hermite	89
3.3.4.	Expresia erorii de interpolare	93
3.3.5.	Convergența interpolării Lagrange	98
3.4.	Calculul eficient al polinoamelor de interpolare	104
3.4.1.	Metode de tip Aitken	104
3.4.2.	Metoda diferențelor divizate	106
3.4.3.	Diferențe finite: formula lui Newton progresivă și regresivă	110
3.4.4.	Diferențe divizate cu noduri multiple	112
3.5.	Interpolare spline	114
3.5.1.	Interpolarea cu spline cubice	117
3.5.2.	Proprietăți de minimalitate ale funcțiilor spline cubice	120
4.	Aproximare uniformă	123
4.1.	Polinoamele lui Bernstein	124
4.2.	B-spline	129
4.2.1.	Noțiuni și rezultate de bază	129
4.2.2.	Algoritmul de evaluare a unui B-spline	132
4.2.3.	Aplicații în grafica pe calculator	133
4.2.4.	Exemple	136
4.3.	Funcții spline cu variație diminuată	140
4.4.	Operatori liniari și pozitivi	143
4.5.	Cea mai bună aproximare uniformă	147

5. Aproximarea funcționalelor liniare	149
5.1. Introducere	149
5.2. Derivare numerică	154
5.3. Integrare numerică	156
5.3.1. Formula trapezului și formula lui Simpson	157
5.3.2. Formule Newton-Cotes cu ponderi și formule de tip Gauss	161
5.3.3. Proprietăți ale cuadraturilor gaussiene	164
5.4. Cuadraturi adaptive	170
5.5. Cuadraturi iterate. Metoda lui Romberg	171
5.6. Cuadraturi adaptive II	174
6. Rezolvarea numerică a ecuațiilor neliniare	177
6.1. Ecuații neliniare	177
6.2. Iterații, convergență și eficiență	178
6.3. Metoda șirurilor Sturm	181
6.4. Metoda falsei poziții	183
6.5. Metoda secantei	185
6.6. Metoda lui Newton	188
6.7. Metoda aproximațiilor succesive	191
6.8. Metoda lui Newton pentru rădăcini multiple	192
6.9. Ecuații algebrice	193
6.10. Metoda lui Newton în \mathbb{R}^n	194
6.11. Metode quasi-Newton	196
6.11.1. Interpolare liniară	197
6.11.2. Metode de modificare	198
7. Vectori și valori proprii	201
7.1. Valori proprii și rădăcini ale polinoamelor	202
7.2. Terminologie și descompunere Schur	203
7.3. Iterația vectorială	206
7.4. Metoda QR – teoria	209
7.5. Metoda QR – practica	213
7.5.1. Metoda QR clasică	213
7.5.2. Deplasare spectrală	218
7.5.3. Metoda QR cu pas dublu	220
8. Rezolvarea numerică a ecuațiilor diferențiale ordinare	225
8.1. Ecuații diferențiale	226
8.2. Metode numerice	227
8.3. Descrierea locală a metodelor cu un pas	228
8.4. Exemple de metode cu un pas	229

8.4.1.	Metoda lui Euler	229
8.4.2.	Metoda dezvoltării Taylor	231
8.4.3.	Metode de tip Euler îmbunătățite	232
8.5.	Metode Runge-Kutta	233
8.6.	Descrierea globală a metodelor cu un pas	238
8.6.1.	Stabilitatea	240
8.6.2.	Convergență	243
8.6.3.	Asimptotica erorii globale	244
8.7.	Monitorizarea erorilor și controlul pasului	247
8.7.1.	Estimarea erorii globale	247
8.7.2.	Estimarea erorii de trunchiere	249
8.7.3.	Controlul pasului	252
9.	Aproximări în mai multe variabile	259
9.1.	Aproximarea funcțiilor de mai multe variabile pe un domeniu rectangular	260
9.2.	Integrarea numerică a funcțiilor de mai multe variabile	267
9.2.1.	Considerații de implementare	273
	Bibliografie	275
	Indice	279

Lista algoritmilor

2.1	Rezolvă sistemul $Ax = b$ prin metoda eliminării a lui Gauss	38
2.2	Factorizare Cholesky	45
2.3	Factorizare QR utilizând reflexii HouseHolder	48
2.4	Calculul produsului $Q^T b$	48
2.5	Calculul produsului Qx	48
2.6	Dându-se scalarii a și b , calculează elementele $c = \cos(\theta)$ și $s = \sin(\theta)$ ale unei matrice Givens	50
2.7	Factorizare $A = QR$ cu ajutorul matricei Givens; rezultatul R se scrie peste A	51
2.8	Algoritmul lui Strassen pentru înmulțirea a două matrice	53
4.1	Algoritmul Cox-deBoor pentru evaluarea unei curbe B-spline	135
4.2	Algoritmul de Casteljau pentru evaluarea unei curbe Bezier	136
5.1	Cuadratură adaptivă	170
5.2	Cuadratură adaptivă bazată pe metoda lui Simpson și extrapolare	176
6.1	Metoda secantei pentru ecuații neliniare în \mathbb{R}	188
6.2	Metoda lui Newton pentru ecuații neliniare în \mathbb{R}	191
6.3	Metoda lui Newton pentru sisteme de ecuații neliniare	196
6.4	Metoda lui Broyden pentru sisteme de ecuații neliniare	200
7.1	Transformarea RQ a unei matrice Hessenberg H , adică $H_* = RQ$ unde $H = QR$ este descompunerea QR a lui H	215
7.2	Reducere la forma Hessenberg superioară	216
7.3	Metoda QR simplă	217
7.4	QRSplit1a – metoda QR cu partiționare și tratarea cazurilor 2×2	219
7.5	Iterație QR pe o matrice Hessenberg; utilizat de algoritmul 7.4 – apel $[H_1, H_2, it] = \mathbf{QRIter}(H, t)$	219

7.6	Metoda QR cu deplasare spectrală, partiționare și tratarea valorilor proprii complexe	221
7.7	Iterație QR și partiționare	221
7.8	Metoda QR cu dublu pas, partiționare și tratarea matricelor 2×2	223
7.9	Iterație QR cu dublu pas și transformare Hessenberg	224
8.1	Metoda Runge-Kutta de ordinul 4	237
8.2	Fragment de pseudocod ce ilustrează implementarea unei metode RK cu pas variabil	257
9.1	Aproximarea unei integrale duble pe dreptunghi	274

CAPITOLUL 1

Elemente de Teoria erorilor și aritmetică în virgulă flotantă

Cuprins

1.1. Probleme numerice	2
1.2. Măsuri ale erorii	4
1.3. Eroarea propagată	5
1.4. Reprezentarea în virgulă flotantă	6
1.4.1. Parametrii reprezentării	6
1.4.2. Anularea	8
1.5. Standardizarea reprezentării în virgulă flotantă	10
1.5.1. Cantități speciale	11
1.6. Condiționarea unei probleme	11
1.7. Condiționarea unui algoritm	14
1.8. Eroarea globală	15
1.9. Probleme prost condiționate și probleme incorect puse	16
1.10. Stabilitatea	16
1.10.1. Notății asimptotice	16
1.10.2. Precizie și stabilitate	18
1.10.3. Analiza regresivă a erorilor	20

Aprecierea preciziei rezultatelor calculelor este un obiectiv important în Analiza numerică. Se disting mai multe tipuri de erori care pot limita această precizie:

1. *erori în datele de intrare;*
2. *erori de rotunjire;*
3. *erori de aproximare.*

Erorile în datele de intrare sunt în afara (dincolo de) controlului calculelor. Ele se pot datora, de exemplu, imperfecțiunilor inerente ale măsurătorilor fizice.

Erorile de rotunjire apar dacă se fac calcule cu numere a căror reprezentare se restrânge la un număr finit de cifre, așa cum se întâmplă de obicei.

Pentru al treilea tip de erori, multe metode nu dau soluția exactă a problemei date P , chiar dacă calculele se fac fără rotunjire, ci mai degrabă soluția unei alte probleme mai simple \tilde{P} , care aproximează P . De exemplu, problema P a însumării unei serii infinite:

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots$$

poate fi înlocuită cu problema mai simplă \tilde{P} de a însuma numai un număr finit de termeni ai seriei. Eroarea de aproximare astfel obținută se numește *eroare de trunchiere* (totuși, acest termen este de asemenea utilizat pentru erorile de rotunjire comise prin ștergerea ultimelor cifre ale reprezentării). Multe probleme de aproximare se obțin prin „discretizarea“ problemei originale P : integralele definite se aproximează prin sume finite, derivatele prin diferențe, etc. În astfel de cazuri, eroarea de aproximare se numește *eroare de discretizare*. Unii autori extind termenul de „eroare de trunchiere“ pentru a acoperi și eroarea de discretizare.

În acest capitol vom examina efectul general al erorilor de intrare și de rotunjire asupra rezultatelor unui calcul. Erorile de aproximare vor fi discutate în capitolele următoare când vom trata metodele numerice individual.

1.1. Probleme numerice

Combinarea dintre o problemă matematică (PM), (de natură constructivă) și specificațiile de precizie ale rezultatului (SP) se numește *problemă numerică*.

Exemplul 1.1.1. Fie $f : \mathbb{R} \rightarrow \mathbb{R}$ și $x \in \mathbb{R}$. Dorim să calculăm $y = f(x)$. În general x nu este reprezentabil în calculator; din acest motiv vom lucra cu o aproximare x^* a sa, $x^* \approx x$. De asemenea este posibil ca f să nu poată fi calculată exact; vom înlocui f cu o aproximantă a sa f_A . Valoarea calculată în calculator va fi $f_A(x^*)$. Deci problema numerică este următoarea:

PM. dându-se x și f , să se calculeze $f(x)$;

SP. $|f(x) - f_A(x^*)| < \varepsilon$, ε dat. ◇

Problemele numerice aparțin uneia din următoarele categorii:

Evaluarea unei funcționale: $l : \mathcal{F} \rightarrow \mathbb{R}$, cum ar fi, de exemplu, calcularea valorii unei funcții $f(x)$, a derivatelor $f'(x)$, $f''(x)$, ... (derivare numerică), a integralelor definite $\int_a^b f(t)dt$ (integrare numerică) și a normelor $\|f\|_p$, etc.

Rezolvarea ecuațiilor algebrice: determinarea valorilor unor necunoscute aflate în relații algebrice prin rezolvarea unor sisteme de ecuații liniare sau neliniare.

Rezolvarea unor ecuații analitice: determinarea funcțiilor (sau valorilor de funcții) soluții ale unei ecuații operatoriale, cum ar fi ecuațiile diferențiale ordinare sau cu derivate parțiale, ecuațiile integrale, ecuații funcționale, etc.

Probleme de optimizare: determinarea unor valori numerice particulare ale unor funcții, care optimizează (minimizează sau maximizează) o funcție obiectiv, cu restricții sau fără restricții. Problemele matematice constructive, din care provin problemele numerice, pot fi privite ca o aplicație abstractă $F : X \rightarrow Y$ între două spații liniare normate.

În funcție de care dintre cantitățile y , x sau F este necunoscută în ecuația

$$Fx = y,$$

avem de-a face cu o *problemă directă*, o *problemă inversă* sau o *problemă de identificare*:

	F	x	y
problemă directă	dată	dat	dorit
problemă inversă	dată	dorit	dat
problemă de identificare	dorită	dat	dat

Exemplul 1.1.2 (Integrare). Evaluarea unei integrale definite

$$Fx = \int_a^b x(t)dt$$

este o *problemă directă*. În această problemă $F : \mathcal{F} \rightarrow \mathbb{R}$ este o funcțională care aplică, de exemplu $\mathcal{F} = C[a, b]$ pe \mathbb{R} . Integrandul – în acest caz o funcție x definită pe $[a, b]$ – este dat, iar integrala $y = Fx$ trebuie determinată. ◇

Exemplul 1.1.3 (Sistem de ecuații liniare). Soluția sistemului de ecuații liniare

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad x, b \in \mathbb{R}^n$$

este o *problemă inversă* tipică. Aplicația liniară F este caracterizată de n^2 coeficienți (elementele lui A). Vectorul b (imaginea) este dat, iar vectorul x (care este transformat în b) trebuie determinat. ◇

Exemplul 1.1.4 (Analiza unei mixturi). Trebuie determinată compoziția unei mixturi de substanțe care se evaporă. Presupunând că în unitatea de timp se evaporă o cantitate fixată din fiecare substanță, se poate utiliza următorul model dependent de timp al mixturii

$$y(t) = \sum_{i=1}^m y_{i_0} e^{-k_i t}. \quad (1.1.1)$$

În afară de cantitățile inițiale y_{i_0} , $i = 1, 2, \dots, m$ și de ratele de evaporare k_1, \dots, k_m ale substanțelor necunoscute, numărul m de substanțe este de asemenea necunoscut. Dacă la momentele t_1, \dots, t_n se măsoară cantitățile y_1, \dots, y_n , trebuie determinați parametrii y_{i_0} , k_i , $i = \overline{1, m}$ și m ai aplicației

$$F : (t_1, \dots, t_n) \mapsto (y_1, \dots, y_n) = \left(\sum_{i=1}^m y_{i_0} e^{-k_i t_1}, \dots, \sum_{i=1}^m y_{i_0} e^{-k_i t_n} \right);$$

deci avem o *problemă de identificare*. Pentru orice $m \in \{1, 2, \dots, n/2\}$ (cel mult $2 \cdot n/2 = n$ necunoscute în acest caz) se poate încerca minimizarea distanței între funcția model (1.1.1) și valorile date, care poate fi definită prin

$$r_m = \sum_{j=1}^n \left(y_j - \sum_{i=1}^m y_{i_0} e^{-k_i t_j} \right)^2. \quad (1.1.2)$$

Aceasta este o problemă inversă extrem de dificilă (problemă de estimare neliniară). Minimum trebuie determinat folosind tehnici de optimizare cu restricții. Valoarea rezidului minim r_m^* face posibilă determinarea valorii lui m care poate fi utilizată pentru cea mai bună adaptare a modelului la datele de intrare. Totuși, aceasta nu înseamnă minimizarea lui r_m fără restricții asupra lui m , ci mai degrabă a nu alege m mai mare decât este necesar, în scopul de a obține modelul cel mai simplu posibil. \diamond

1.2. Măsurile ale erorii

Definiția 1.2.1 Fie X un spațiu liniar normat, $A \subseteq X$ și $x \in X$. Un element $x^* \in A$ se numește aproximantă a lui x din A (notație $x^* \approx x$).

Definiția 1.2.2 Dacă x^* este o aproximantă a lui x diferența $\Delta x = x - x^*$ se numește eroare, iar

$$\|\Delta x\| = \|x^* - x\| \quad (1.2.1)$$

se numește eroare absolută.

Definiția 1.2.3 *Expresia*

$$\delta x = \frac{\|\Delta x\|}{\|x\|}, \quad x \neq 0 \quad (1.2.2)$$

se numește eroare relativă.

Observația 1.2.4.

1. Deoarece în practică x este necunoscut, se folosește aproximarea $\delta x = \frac{\|\Delta x\|}{\|x^*\|}$.
Dacă $\|\Delta x\|$ este mic comparativ cu x^* , atunci aproximanta este bună.
2. Dacă $X = \mathbb{R}$, se lucrează cu $\delta x = \frac{\Delta x}{x}$ și $\Delta x = x^* - x$. ◇

1.3. Eroarea propagată

Fie $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $x = (x_1, \dots, x_n)$ și $x^* = (x_1^*, \dots, x_n^*)$. Dorim să evaluăm eroarea absolută și relativă Δf și respectiv δf când se aproximează $f(x)$ prin $f(x^*)$. Aceste erori se numesc *erori propagate*, deoarece ne spun cum se propagă eroarea inițială (absolută sau relativă) pe parcursul calculării lui f . Să presupunem că $x = x^* + \Delta x$, unde $\Delta x = (\Delta x_1, \dots, \Delta x_n)$. Pentru eroarea absolută avem (folosind formula lui Taylor)

$$\begin{aligned} \Delta f &= f(x_1^* + \Delta x_1, \dots, x_n^* + \Delta x_n) - f(x_1^*, \dots, x_n^*) = \\ &= \sum_{i=1}^n \Delta x_i \frac{\partial f}{\partial x_i^*}(x_1^*, \dots, x_n^*) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \Delta x_i \Delta x_j \frac{\partial^2 f}{\partial x_i^* \partial x_j^*}(\theta), \end{aligned}$$

unde $\theta \in [(x_1^*, \dots, x_n^*), (x_1^* + \Delta x_1, \dots, x_n^* + \Delta x_n)]$.

Dacă Δx_i sunt suficient de mici, atunci $\Delta x_i \Delta x_j$ sunt neglijabile comparativ cu Δx_i și obținem

$$\Delta f \approx \sum_{i=1}^n \Delta x_i \frac{\partial f}{\partial x_i^*}(x_1^*, \dots, x_n^*). \quad (1.3.1)$$

Pentru eroarea relativă avem

$$\begin{aligned} \delta f &= \frac{\Delta f}{f} \approx \sum_{i=1}^n \Delta x_i \frac{\frac{\partial f}{\partial x_i^*}(x^*)}{f(x^*)} = \sum_{i=1}^n \Delta x_i \frac{\partial}{\partial x_i^*} \ln f(x^*) = \\ &= \sum_{i=1}^n x_i^* \delta x_i \frac{\partial}{\partial x_i^*} \ln f(x^*). \end{aligned}$$

Deci

$$\delta f = \sum_{i=1}^n x_i^* \frac{\partial}{\partial x_i^*} \ln f(x^*) \delta x_i. \quad (1.3.2)$$

De o mare importanță practică este și problema inversă: cu ce precizie trebuie approximate datele pentru ca rezultatul să aibă o precizie dată? Adică, dându-se $\varepsilon > 0$, cât trebuie să fie Δx_i sau δx_i , $i = \overline{1, n}$ astfel încât Δf sau $\delta f < \varepsilon$? O metodă de rezolvare se bazează pe *principiul efectelor egale*: se presupune că toți termenii care intervin în (1.3.1) sau (1.3.2) au același efect, adică

$$\frac{\partial f}{\partial x_1^*}(x^*)\Delta x_1 = \dots = \frac{\partial f}{\partial x_n^*}(x^*)\Delta x_n.$$

Din (1.3.1) se obține

$$\Delta x_i \approx \frac{\Delta f}{n \left| \frac{\partial f}{\partial x_i^*}(x^*) \right|}. \quad (1.3.3)$$

Analog,

$$\delta x_i = \frac{\delta f}{n \left| x_i^* \frac{\partial}{\partial x_i^*} \ln f(x^*) \right|}. \quad (1.3.4)$$

1.4. Reprezentarea în virgulă flotantă

1.4.1. Parametrii reprezentării

Deși au fost propuse mai multe reprezentări pentru numerele reale, reprezentarea în virgulă flotantă este cea mai răspândită. Parametrii reprezentării în virgulă flotantă sunt baza β (întotdeauna pară), precizia p , exponentul maxim e_{\max} și cel minim e_{\min} , toate numere naturale. În general, un număr în virgulă flotantă se reprezintă sub forma

$$x = \pm d_0.d_1d_2 \dots d_{p-1} \times \beta^e, \quad 0 \leq d_i < \beta \quad (1.4.1)$$

unde $d_0.d_1d_2 \dots d_{p-1}$ se numește *semnificant* sau *mantisă*, iar e *exponent*. Valoarea lui x este

$$\pm(d_0 + d_1\beta^{-1} + d_2\beta^{-2} + \dots + d_{p-1}\beta^{-(p-1)})\beta^e. \quad (1.4.2)$$

Pentru ca reprezentarea în virgulă flotantă să fie unică, numerele flotante se *normalizează*, adică se modifică reprezentarea (nu valoarea) astfel încât $d_0 \neq 0$. Zero se reprezintă ca $1.0 \times \beta^{e_{\min}-1}$. În acest mod ordinea numerică uzuală a numerelor reale nenegative corespunde ordinii lexicografice a reprezentării lor flotante (cu exponentul în stânga semnificantului).

Termenul de *număr în virgulă flotantă* este utilizat pentru a desemna un număr real care poate fi reprezentat exact în virgulă flotantă.

Distribuția numerelor în virgulă flotantă pe axa reală apare în figura 1.1.

Fiecare interval de forma $[\beta^e, \beta^{e+1})$ din \mathbb{R} conține β^p numere în virgulă flotantă (numărul posibil de semnificanți). Intervalul $(0, \beta^{e_{\min}})$ este gol și din acest motiv se introduc *numerele denormalizate*, adică numere cu semnificantul de forma $0.d_1d_2 \dots d_{p-1}$

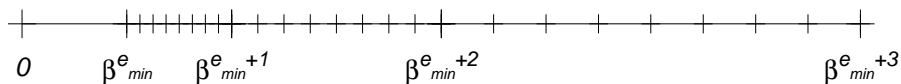


Figura 1.1: Distribuția numerelor în virgulă flotantă pe axa reală (fără denormalizare)

și exponentul $\beta^{e_{min}-1}$. Faptul că se admit numere denormalizate sau nu se consideră a fi un parametru suplimentar al reprezentării. Mulțimea numerelor în virgulă flotantă pentru un set de parametri dați ai reprezentării se va nota cu

$$\mathbb{F}(\beta, p, e_{min}, e_{max}, denorm), \quad denorm \in \{true, false\}.$$

Distribuția numerelor în virgulă flotantă când se admite denormalizarea apare în figura 1.2.

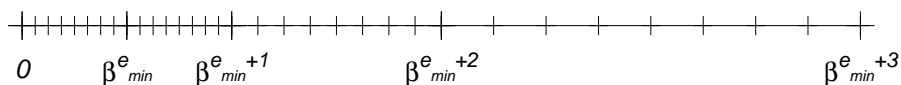


Figura 1.2: Distribuția numerelor în virgulă flotantă pe axa reală (cu denormalizare)

Această mulțime nu coincide cu \mathbb{R} din următoarele motive:

1. este o submulțime finită a lui \mathbb{Q} ;
2. pentru $x \in \mathbb{R}$ putem avea $|x| > \beta \times \beta^{e_{max}}$ (depășire superioară) sau $|x| < 1.0 \times \beta^{e_{min}}$ (depășire inferioară).

Operațiile aritmetice uzuale pe $\mathbb{F}(\beta, p, e_{min}, e_{max}, denorm)$ se notează cu \oplus , \ominus , \otimes , \oslash , iar funcțiile uzuale cu SIN, COS, EXP, LN, SQRT ș.a.m.d. ($\mathbb{F}, \oplus, \otimes$) nu este corp deoarece

$$\begin{aligned} (x \oplus y) \oplus z &\neq x \oplus (y \oplus z) & (x \otimes y) \otimes z &\neq x \otimes (y \otimes z) \\ (x \oplus y) \otimes z &\neq x \otimes z \oplus y \otimes z. \end{aligned}$$

Pentru măsurarea erorii de reprezentare, în afară de eroare relativă, se folosește *ulps* – *units in the last place* (unități în ultima poziție). Dacă numărul z se reprezintă prin $d_0.d_1d_2\dots d_{p-1} \times \beta^e$, atunci eroarea este

$$|d_0.d_1d_2\dots d_{p-1} - z/\beta^e| \beta^{p-1} \text{ulps.}$$

Eroarea relativă ce corespunde la $\frac{1}{2}$ ulps este

$$\frac{1}{2}\beta^{-p} \leq \frac{1}{2} \text{ulps} \leq \frac{\beta}{2}\beta^{-p},$$

căci eroarea absolută este $\underbrace{0.0\dots 0}_p \beta' \times \beta^e$, cu $\beta' = \frac{\beta}{2}$, iar numitorul este cuprins între β^e și $\beta \times \beta^e$. Valoarea $\text{eps} = \frac{\beta}{2}\beta^{-p}$ se numește *epsilon-ul mașinii*.

Rotunjirea implicită se face după regula cifrei pare: dacă $x = d_0.d_1\dots d_{p-1}d_p\dots$ și $d_p > \frac{\beta}{2}$ rotunjirea se face în sus, dacă $d_p < \frac{\beta}{2}$ rotunjirea se face în jos, iar dacă $d_p = \frac{\beta}{2}$ și printre cifrele eliminate există una nenulă rotunjirea se face în sus, iar în caz contrar ultima cifră păstrată este pară. Dacă notăm cu fl operația de rotunjire, operațiile aritmetice din \mathbb{F} se pot defini prin

$$x \odot y = \text{fl}(x \circ y). \quad (1.4.3)$$

Se pot alege și alte variante de rotunjire – cât mai depărtat de 0, spre $-\infty$, spre $+\infty$, spre 0 (trunchiere). În raționamentele asupra operațiilor în virgulă flotantă vom folosi următorul model

$$\forall x, y \in \mathbb{F}, \exists \delta \text{ cu } |\delta| < \text{eps} \text{ astfel încât } x \odot y = (x \circ y)(1 + \delta). \quad (1.4.4)$$

În cuvinte, orice operație în aritmetica în virgulă flotantă este exactă până la o eroare relativă de cel mult eps.

Formula (1.4.4) se numește *axioma fundamentală a aritmeticii în virgulă flotantă*.

1.4.2. Anularea

Din formulele pentru eroarea relativă (1.3.2), dacă $x \approx x(1 + \delta_x)$ și $y \approx y(1 + \delta_y)$, avem următoarele expresii pentru erorile relative ale operațiilor în virgulă flotantă:

$$\delta_{xy} = \delta_x + \delta_y \quad (1.4.5)$$

$$\delta_{x/y} = \delta_x - \delta_y \quad (1.4.6)$$

$$\delta_{x+y} = \frac{x}{x+y}\delta_x + \frac{y}{x+y}\delta_y \quad (1.4.7)$$

Singura operație critică din punct de vedere al erorii este scăderea a două cantități apropiate $x \approx y$, caz în care $\delta_{x-y} \rightarrow \infty$. Acest fenomen se numește anulare și este reprezentat grafic în figura 1.3. Aici b, b', b'' sunt cifre binare acceptabile, iar g -urile reprezintă cifre binare contaminate de eroare (gunoai - garbage digits). De notat că, gunoi - gunoi = gunoi, dar mai important, normalizarea mută prima cifră contaminată de pe poziția a 12-a pe poziția a treia.

x	=	1	0	1	1	0	0	1	0	1	b	b	g	g	g	g	e
y	=	1	0	1	1	0	0	1	0	1	b'	b'	g	g	g	g	e
x-y	=	0	0	0	0	0	0	0	0	0	b''	b''	g	g	g	g	e
	=	b''	b''	g	g	g	g	?	?	?	?	?	?	?	?	?	e-9

Figura 1.3: Anularea

Anularea este de două tipuri: *benignă*, când se scad două cantități exacte și *catastrofală*, când se scad două cantități deja rotunjite. Programatorul trebuie să fie conștient de posibilitatea apariției anulării și să încerce să o evite. Expresiile în care apare anularea trebuie rescrise, iar o anulare catastrofală trebuie întotdeauna transformată în una benignă. Vom da în continuare câteva exemple de anulări catastrofale și modul de transformare a lor .

Exemplul 1.4.1. Dacă $a \approx b$, atunci expresia $a^2 - b^2$ se transformă în $(a - b)(a + b)$. Forma inițială este de preferat în cazul când $a \gg b$ sau $b \gg a$. \diamond

Exemplul 1.4.2. Dacă anularea apare într-o expresie cu radicali, se amplifică cu conjugata:

$$\sqrt{x + \delta} - \sqrt{x} = \frac{\delta}{\sqrt{x + \delta} + \sqrt{x}}, \quad \delta \approx 0. \quad \diamond$$

Exemplul 1.4.3. Diferența valorilor unei funcții pentru argumente apropiate se transformă folosind formula lui Taylor:

$$f(x + \delta) - f(x) = \delta f'(x) + \frac{\delta^2}{2} f''(x) + \dots \quad f \in C^n[a, b]. \quad \diamond$$

Exemplul 1.4.4. La ecuația de gradul al doilea $ax^2 + bx + c = 0$, anularea poate să apară dacă $b^2 \gg 4ac$. Formulele uzuale

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad (1.4.8)$$

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} \quad (1.4.9)$$

pot să conducă la anulare astfel: pentru $b > 0$ anularea apare la calculul lui x_1 , iar pentru $b < 0$ anularea apare la calculul lui x_2 . Remediul este să amplificăm cu conjugata

$$x_1 = \frac{2c}{-b - \sqrt{b^2 - 4ac}} \quad (1.4.10)$$

$$x_2 = \frac{2c}{-b + \sqrt{b^2 - 4ac}} \quad (1.4.11)$$

și să utilizăm în primul caz formulele (1.4.10) și (1.4.9), iar în al doilea caz (1.4.8) și (1.4.11). \diamond

1.5. Standardizarea reprezentării în virgulă flotantă

Există două standarde diferite pentru calculul în virgulă flotantă: IEEE 754 care prevede $\beta = 2$ și IEEE 854 care permite $\beta = 2$ sau $\beta = 10$, dar lasă o mai mare libertate de reprezentare. Ne vom ocupa numai de primul standard.

Parametrii standardului se dau în tabela 1.1.

	Precizia			
	Simplă	Simplă extinsă	Dublă	Dublă extinsă
p	24	≥ 32	53	≥ 64
e_{\max}	+127	$\geq +1023$	+1023	$\geq +16383$
e_{\min}	-126	≤ -1022	-1022	≤ -16382
dim. exponent	8	≥ 11	11	≥ 15
dim. număr	32	≥ 43	64	≥ 79

Tabela 1.1: Parametrii reprezentării flotante

Motivele pentru formatele extinse sunt:

1. o mai bună precizie;
2. pentru conversia din binar în zecimal și invers este nevoie de 9 cifre în simplă precizie și de 17 cifre în dublă precizie.

Motivul pentru care $|e_{\min}| < e_{\max}$ este acela că $1/2^{e_{\min}}$ nu trebuie să dea depășire.

Operațiile $\oplus, \ominus, \otimes, \oslash$ trebuie să fie exact rotunjite. Precizia aceasta se asigură cu două cifre de gardă și un bit suplimentar.

Reprezentarea exponentului se numește *reprezentare cu exponent deplasat*, adică în loc de e se reprezintă $e + D$, unde D este fixat la alegerea reprezentării.

În cazul IEEE 754, $D = 127$.

1.5.1. Cantități speciale

În standardul IEEE 754 există următoarele cantități speciale:

Exponent	Semnificanț	Ce reprezintă
$e = e_{min} - 1$	$f = 0$	± 0
$e = e_{min} - 1$	$f \neq 0$	$0.f \times 2^{e_{min}}$
$e_{min} \leq e \leq e_{max}$		$1.f \times 2^e$
$e = e_{max} + 1$	$f = 0$	$\pm \infty$
$e = e_{max} + 1$	$f \neq 0$	NaN

NaN. Avem de fapt o familie de valori NaN, operațiile ilegale sau nedeterminate conduc la NaN: $\infty + (-\infty)$, $0 \times \infty$, $0/0$, ∞/∞ , $x \text{ REM } 0$, $\infty \text{ REM } y$, \sqrt{x} pentru $x < 0$. Dacă un operand este NaN rezultatul va fi tot NaN.

Infinit. $1/0 = \infty$, $-1/0 = -\infty$. Valorile infinite dau posibilitatea continuării calculului, lucru mai sigur decât abortarea sau returnarea celui mai mare număr reprezentabil.

$\frac{x}{1+x^2}$ pentru $x = \infty$ dă rezultatul 0.

Zero cu semn. Avem doi de 0: $+0$, -0 ; relațiile $+0 = -0$ și $-0 < +\infty$ sunt adevărate. Avantaje: tratarea simplă a depășirilor inferioare și discontinuităților. Se face distincție între $\log 0 = -\infty$ și $\log x = \text{NaN}$ pentru $x < 0$. Fără 0 cu semn nu s-ar putea face distincție la logaritm între un număr negativ care dă depășire superioară și 0.

1.6. Condiționarea unei probleme

Putem gândi o problemă ca o aplicație

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad y = f(x). \quad (1.6.1)$$

Ne interesează sensibilitatea aplicației într-un punct dat x la mici perturbații ale argumentului, adică cât de mare sau cât de mică este perturbația lui y comparativ cu perturbația lui x . În particular, dorim să măsurăm gradul de sensibilitate printr-un singur număr, numărul de condiționare al aplicației f în punctul x . Vom presupune că f este calculată exact, cu precizie infinită. *Condiționarea lui f este deci o proprietate inerentă a funcției f și nu depinde de nici o considerație algoritmică legată de implementarea sa.*

Aceasta nu înseamnă că determinarea condiționării unei probleme este nerelevantă pentru orice soluție algoritmică a problemei. Din contră. Motivul este acela că soluția calculată cu (1.6.1), y^* (utilizând un algoritm specific și aritmetica în virgulă flotantă) este (și acest lucru se poate demonstra) soluția unei probleme „apropiate“

$$y^* = f(x^*) \quad (1.6.2)$$

cu

$$x^* = x + \delta \quad (1.6.3)$$

și, mai mult, distanța $\|\delta\| = \|x^* - x\|$ poate fi estimată în termeni de precizie a mașinii. Deci, dacă știm cât de tare sau cât de slab reacționează aplicația la mici perturbații, cum ar fi δ în (1.6.3), putem spune ceva despre eroarea $y^* - y$ a soluției cauzată de această perturbație.

Se poate considera și condiționarea între aplicații mai generale, dar pentru implementări practice este suficient să ne limităm la cazul finit dimensional.

Fie

$$x = [x_1, \dots, x_m]^T \in \mathbb{R}^m, \quad y = [y_1, \dots, y_n]^T \in \mathbb{R}^n, \\ y_\nu = f_\nu(x_1, \dots, x_m), \quad \nu = \overline{1, n}.$$

y_ν va fi privit ca o funcție de o singură variabilă x_μ

$$\gamma_{\nu\mu} = (\text{cond}_{\nu\mu} f)(x) = \left| \frac{x_\mu \frac{\partial f_\nu}{\partial x_\mu}}{f_\nu(x)} \right|. \quad (1.6.4)$$

Aceasta ne dă o matrice de numere de condiționare

$$\Gamma(x) = \begin{pmatrix} \frac{x_1 \frac{\partial f_1}{\partial x_1}}{f_1(x)} & \cdots & \frac{x_m \frac{\partial f_1}{\partial x_m}}{f_1(x)} \\ \vdots & \ddots & \vdots \\ \frac{x_1 \frac{\partial f_n}{\partial x_1}}{f_n(x)} & \cdots & \frac{x_m \frac{\partial f_n}{\partial x_m}}{f_n(x)} \end{pmatrix} = [\gamma_{\nu\mu}(x)] \quad (1.6.5)$$

și vom lua ca *număr de condiționare*

$$(\text{cond } f)(x) = \|\Gamma(x)\|. \quad (1.6.6)$$

Altă abordare. Considerăm norma $\|\cdot\|_\infty$

$$\Delta y_\nu \approx \sum_{\mu=1}^m \frac{\partial f_\nu}{\partial x_\mu} \Delta x_\mu \quad (= f_\nu(x + \Delta x) - f_\nu(x)) \\ |\Delta y_\nu| \leq \sum_{\mu=1}^n \left| \frac{\partial f_\nu}{\partial x_\mu} \right| \Delta x_\mu \leq \max_{\mu} |\Delta x_\mu| \sum_{\mu=1}^m \left| \frac{\partial f_\nu}{\partial x_\mu} \right| \leq \\ \leq \max_{\mu} |\Delta x_\mu| \max_{\nu} \sum_{\mu=1}^m \left| \frac{\partial f_\nu}{\partial x_\mu} \right|$$

Am obținut

$$\|\Delta y\|_\infty \leq \|\Delta x\|_\infty \left\| \frac{\partial f}{\partial x} \right\|_\infty \quad (1.6.7)$$

unde

$$J(x) = \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_m} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_m} \end{bmatrix} \in \mathbb{R}^n \times \mathbb{R}^m \quad (1.6.8)$$

este matricea jacobiană a lui f

$$\frac{\|\Delta y\|_\infty}{\|y\|_\infty} \leq \frac{\|x\|_\infty \left\| \frac{\partial f}{\partial x} \right\|_\infty}{\|f(x)\|_\infty} \cdot \frac{\|\Delta x\|}{\|x\|_\infty}. \quad (1.6.9)$$

Estimarea (1.6.9) este mai grosieră decât (1.6.4).

Dacă $m = n = 1$ în ambele abordări se obține

$$(\text{cond } f)(x) = \left| \frac{x f'(x)}{f(x)} \right|,$$

pentru $x \neq 0, y \neq 0$.

Dacă $x = 0 \wedge y \neq 0$ se consideră eroarea absolută pentru x și eroarea relativă pentru y

$$(\text{cond } f)(x) = \left| \frac{f'(x)}{f(x)} \right|;$$

pentru $y = 0 \wedge x \neq 0$ se ia eroarea absolută pentru y și eroarea relativă pentru x , iar pentru $x = y = 0$

$$(\text{cond } f)(x) = f'(x).$$

Exemplul 1.6.1 (Sisteme de ecuații liniare algebrice). Dându-se matricea $A \in \mathbb{R}^{m \times n}$ și vectorul $b \in \mathbb{R}^n$ să se rezolve sistemul

$$Ax = b. \quad (1.6.10)$$

Aici datele de intrare sunt elementele lui A și b , iar rezultatul este vectorul x . Pentru a simplifica lucrurile să presupunem că A este o matrice fixată care nu se schimbă și că b ar putea fi perturbat. Avem aplicația $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ dată de

$$x = f(b) := A^{-1}b,$$

care este liniară. Deci $\frac{\partial f}{\partial b} = A^{-1}$ și utilizând (1.6.9) obținem

$$\begin{aligned} (\text{cond } f)(b) &= \frac{\|b\| \|A^{-1}\|}{\|A^{-1}b\|} = \frac{\|Ax\| \|A^{-1}\|}{\|A^{-1}b\|}, \\ \max_{\substack{b \in \mathbb{R}^n \\ b \neq 0}} (\text{cond } f)(b) &= \max_{\substack{x \in \mathbb{R}^n \\ b \neq 0}} \frac{\|Ax\|}{\|x\|} \|A^{-1}\| = \|A\| \|A^{-1}\|. \quad \diamond \end{aligned} \quad (1.6.11)$$

Numărul $\|A\|\|A^{-1}\|$ se numește număr de condiționare al matricei A și se notează $\text{cond } A$.

$$\text{cond } A = \|A\|\|A^{-1}\|.$$

Vom reveni asupra acestei probleme în capitolul 2, consacrat rezolvării sistemelor de ecuații algebrice liniare.

1.7. Condiționarea unui algoritm

Fie problema

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad y = f(x). \quad (1.7.1)$$

Împreună cu f se dă un algoritm A care rezolvă problema. Adică dându-se un vector $x \in \mathbb{F}(\beta, p, e_{\min}, e_{\max}, \text{denorm})$, algoritmul A produce un vector y_A (în aritmetica în virgulă flotantă), despre care se presupune că aproximează $y = f(x)$. Astfel avem o aplicație f_A ce descrie modul în care problema f este rezolvată de algoritmul A

$$f_A : \mathbb{F}^m(\dots) \rightarrow \mathbb{F}^n(\dots), \quad y_A = f_A(x).$$

Pentru a putea analiza f_A facem următoarea ipoteză de bază

$$(IB) \quad \forall x \in \mathbb{F}^m \exists x_A \in \mathbb{F}^m : \quad f_A(x) = f(x_A). \quad (1.7.2)$$

Adică, soluția calculată corespunzând unei anumite intrări x este soluția exactă pentru o altă intrare, diferită de prima, x_A (nu neapărat număr în virgulă flotantă, și nici unic determinată) despre care sperăm că este apropiată de x . Cu cât găsim un x_A mai apropiat de x , cu atât avem mai mare încredere în algoritm.

Definim *condiționarea lui A în x* comparând eroarea relativă la intrare cu eps:

$$(\text{cond } A)(x) = \inf_{x_A} \frac{\|x_A - x\|}{\|x\|} / \text{eps}.$$

Justificare:

$$\delta_y = \frac{f_A(x) - f(x)}{f(x)} = \frac{(x_A - x)f'(\xi)}{f(x)} \approx \frac{x_A - x}{x} \cdot \frac{1}{\text{eps}} \frac{xf'(x)}{f(x)} \text{eps}$$

Infimumul se ia după toți x_A ce satisfac $y_A = f(x_A)$. În practică se poate lua orice valoare x_A și se obține o margine superioară a numărului de condiționare

$$(\text{cond } A)(x) \leq \frac{\|x_A - x\|}{\|x\|} \text{eps}. \quad (1.7.3)$$

1.8. Eroarea globală

Considerăm din nou problema:

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad y = f(x). \quad (1.8.1)$$

Aceasta este problema (matematică) idealizată, în care datele sunt numere reale, iar soluția este soluția matematică exactă. Când o rezolvăm în aritmetica în virgulă flotantă, cu precizia eps, utilizând un algoritm A , rotunjim la început toate datele și acestora nu le aplicăm f , ci f_A .

$$x^* = x \text{ rotunjit}, \quad \frac{\|x^* - x\|}{\|x\|} = \varepsilon, \quad y_A^* = f_A(x^*).$$

Aici ε este eroarea de rotunjire. Ea poate proveni și din alte surse (măsurători). Eroarea totală este

$$\frac{\|y_A^* - y\|}{\|y\|}.$$

Pe baza ipotezei (1.7.2, IB) alegând x_A optimal avem

$$\begin{aligned} f_A(x^*) &= f(x_A^*), \\ \frac{\|x_A^* - x^*\|}{\|x^*\|} &= (\text{cond } A)(x^*) \text{ eps}. \end{aligned} \quad (1.8.2)$$

Fie $y^* = f(x^*)$. Avem utilizând inegalitatea triunghiului

$$\frac{\|y_A^* - y\|}{\|y\|} \leq \frac{\|y_A^* - y^*\|}{\|y\|} + \frac{\|y^* - y\|}{\|y\|} \approx \frac{\|y_A^* - y^*\|}{\|y^*\|} + \frac{\|y^* - y\|}{\|y^*\|}$$

Am presupus că $\|y\| \approx \|y^*\|$. Din (1.8.2) rezultă pentru primul termen

$$\begin{aligned} \frac{\|y_A^* - y^*\|}{\|y^*\|} &= \frac{\|f_A(x^*) - f(x^*)\|}{\|f(x^*)\|} = \frac{\|f(x_A^*) - f(x^*)\|}{\|f(x^*)\|} \leq \\ &\leq (\text{cond } f)(x^*) \frac{\|x_A^* - x^*\|}{\|x^*\|} = (\text{cond } f)(x^*) (\text{cond } A)(x^*) \text{ eps}, \end{aligned}$$

iar pentru al doilea

$$\frac{\|y^* - y\|}{\|y\|} = \frac{\|f(x^*) - f(x)\|}{\|f(x)\|} \leq (\text{cond } f)(x) \frac{\|x^* - x\|}{\|x\|} = (\text{cond } f)(x) \varepsilon.$$

Presupunând că $(\text{cond } f)(x^*) \approx (\text{cond } f)(x)$ obținem

$$\frac{\|y_A^* - y\|}{\|y\|} \leq (\text{cond } f)(x) [\varepsilon + (\text{cond } A)(x^*) \text{ eps}]. \quad (1.8.3)$$

Interpretare: erorile în date și eps contribuie împreună la eroarea totală. Ambele sunt amplificate de condiționarea problemei, dar ultima este amplificată și de condiționarea algoritmului.

1.9. Probleme prost condiționate și probleme incorect puse

Dacă numărul de condiționare al unei probleme este mare ($(\text{cond } f)(x) \gg 1$), atunci chiar pentru erori (relative) mici trebuie să ne așteptăm la erori foarte mari în datele de ieșire. Astfel de probleme se numesc *probleme prost condiționate*. Nu este posibil să tragem o linie clară de demarcație între problemele bine condiționate și cele prost condiționate. O categorizare a unei probleme în prost condiționată sau bine condiționată depinde de specificațiile de precizie. Dacă dorim ca

$$\frac{\|y^* - y_A^*\|}{\|y\|} < \tau$$

și în (1.8.3) $(\text{cond } f)(x)\varepsilon \geq \tau$, atunci problema este sigur prost condiționată.

Este important să se aleagă o limită rezonabilă pentru eroare, căci în caz contrar, chiar dacă creștem numărul de iterații, nu vom putea crește precizia.

Dacă rezultatul unei probleme matematice depinde discontinuu de date ce variază continuu, atunci este imposibil de dat o soluție numerică a problemei în vecinătatea discontinuității. În astfel de cazuri rezultatul poate fi perturbat substanțial, chiar dacă datele de intrare sunt precise și utilizăm aritmetica de precizie multiplă. Astfel de probleme se numesc *probleme incorect puse*. O problemă incorect pusă poate să apară, de exemplu, dacă un rezultat întreg este calculat din date reale (adică date care variază continuu), de exemplu numărul de rădăcini reale ale unei funcții sau rangul unei matrice.

Exemplul 1.9.1 (Numărul de zerouri reale ale unui polinom). Ecuația

$$P_3(x, c_0) = c_0 + x - 2x^2 + x^3$$

are una, două sau trei rădăcini reale, după cum c_0 este strict pozitiv, zero sau strict negativ (vezi figura 1.4). Deci, pentru valori ale lui c_0 apropiate de zero, numărul de zerouri reale ale lui P_3 este o problemă incorect pusă. \diamond

1.10. Stabilitatea

1.10.1. Notății asimptotice

Această subsecțiune introduce notațiile asimptotice de bază și câteva abuzuri comune.

Pentru o funcție dată $g(n)$ vom nota cu $\Theta(g(n))$ mulțimea de funcții

$$\Theta(g(n)) = \{f(n) : \exists c_1, c_2, n_0 > 0 \ 0 \leq c_1 g(n) \leq f(n) \leq c_2 g(n) \ \forall n \leq n_0\}.$$

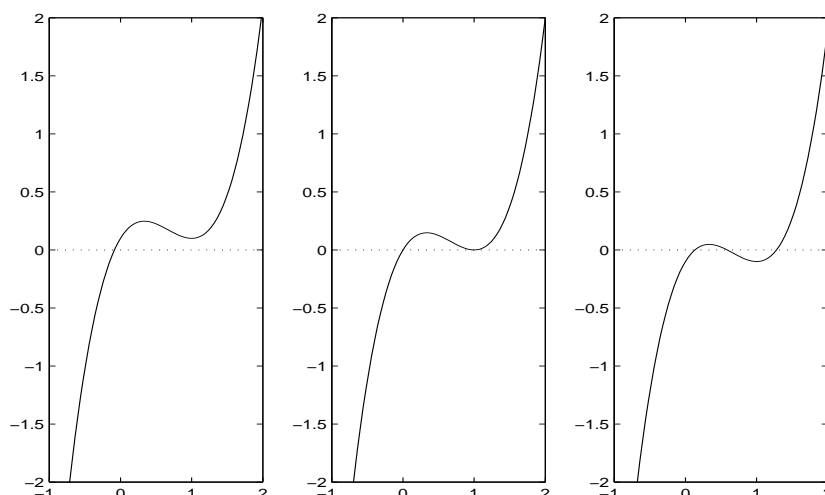


Figura 1.4: Problemă incorect pusă

Deși $\Theta(g(n))$ este o mulțime; scriem $f(n) = \Theta(g(n))$ pentru a indica faptul că $f(n) \in \Theta(g(n))$. Acest abuz de egalitate pentru a nota apartenența la o mulțime poate părea la început confuz, dar vom vedea că are anumite avantaje. Vom spune că $g(n)$ este o *margină asimptotică strânsă* (*asymptotically tight bound*) pentru $f(n)$.

Definiția mulțimii $\Theta(g(n))$ necesită ca fiecare membru al ei să fie *asimptotic nenegativ*, adică $f(n) \geq 0$ când n este suficient de mare. În consecință $g(n)$ trebuie să fie și ea asimptotic negativă, căci altfel $\Theta(g(n))$ este vidă. Din acest motiv vom presupune că fiecare funcție utilizată în interiorul notației Θ este asimptotic nenegativă. Această presupunere are loc și pentru celelalte notații asimptotice care vor fi definite în acest capitol.

Pentru o funcție dată $g(n)$ vom nota cu $O(g(n))$ mulțimea de funcții

$$O(g(n)) = \{f(n) : \exists c, n_0 \ 0 \leq f(n) \leq cg(n), \forall n \geq n_0\}.$$

Pentru a indica faptul că $f(n)$ este un membru al lui $O(g(n))$ scriem $f(n) = O(g(n))$. Observăm că $f(n) = \Theta(g(n)) \implies f(n) = O(g(n))$, deoarece notația Θ este mai tare decât notația O . Utilizând notațiile din teoria mulțimilor avem $\Theta(g(n)) \subseteq O(g(n))$. Una dintre proprietățile ciudate ale notației este aceea că $n = O(n^2)$.

Pentru o funcție dată $g(n)$ vom nota prin $\Omega(g(n))$ mulțimea de funcții

$$\Omega(g(n)) = \{f(n) : \exists c, n_0 \ 0 \leq cg(n) \leq f(n), \forall n \geq n_0\}.$$

Această notație furnizează o *margină asimptotică inferioară*. Din definițiile notațiilor asimptotice se obține imediat:

$$f(n) = \Theta(g(n)) \iff f(n) = O(g(n)) \wedge f(n) = \Omega(g(n)).$$

Spunem că funcțiile f și $g : \mathbb{N} \longrightarrow \mathbb{R}$ sunt *asimptotic echivalente*, notație \sim dacă

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1.$$

Extinderea notațiilor asimptotice la mulțimea numerelor reale este naturală. De exemplu $f(t) = O(g(t))$ înseamnă că există o constantă pozitivă C astfel încât pentru orice t suficient de apropiat de o limită subînțeleasă (de exemplu $t \rightarrow \infty$ sau $t \rightarrow 0$) avem

$$|f(t)| \leq Cg(t). \quad (1.10.1)$$

1.10.2. Precizie și stabilitate

Vom considera o *problemă* ca fiind o aplicație $f : X \longrightarrow Y$, unde X și Y sunt spații liniare normate (pentru scopurile noastre ne vom limita la cazul finit dimensional). Ne va interesa comportarea problemei într-un punct particular $x \in X$ (comportarea poate diferi de la un punct la altul). Combinația unei probleme f cu niște date prescise x se va numi *instanță a problemei*, dar se obișnuiește să se utilizeze termenul de *problemă* pentru ambele noțiuni.

Numerele complexe sunt reprezentate printr-o pereche de numere în virgulă flotantă și operațiile elementare se realizează pe această reprezentare. Rezultatul este acela că axioma (1.4.4) este valabilă și pentru numere complexe, exceptând faptul că la operațiile \otimes și \oslash eps trebuie mărit cu un factor de $2^{3/2}$ și respectiv $2^{5/2}$.

Un algoritm va fi privit ca o aplicație $f_A : X \longrightarrow Y$, unde X și Y sunt ca mai sus. Fie f o problemă, un calculator al cărui sistem de numere în virgulă flotantă satisface (1.4.4), dar nu neapărat (1.4.3), un algoritm f_A pentru f și o implementare a acestui algoritm sub formă de program pe calculator, A , toate fixate. Dându-se o dată $x \in X$, o vom rotunji la un număr în virgulă flotantă și apoi o vom furniza programului. Rezultatul este o colecție de numere în virgulă flotantă care aparțin spațiului Y (deoarece algoritmul a fost conceput să rezolve f). Vom numi acest rezultat calculat $f_A(x)$.

Exceptând cazul trivial, f_A nu poate fi continuă. Vom spune că un algoritm f_A pentru problema f este *precis*, dacă pentru orice $x \in X$, eroarea sa relativă verifică

$$\frac{\|f_A(x) - f(x)\|}{\|f(x)\|} = O(\text{eps}). \quad (1.10.2)$$

Dacă problema f este prost condiționată, obiectivul preciziei, așa cum este definit de (1.10.2) este nerezonabil de ambițios. Erorile de rotunjire în datele de intrare sunt inevitabile pe calculatoare numerice și chiar dacă toate calculele următoare se realizează perfect, această perturbație ne conduce la o modificare semnificativă a rezultatului. În loc să urmărim precizia în toate cazurile este mai rezonabil să urmărim stabilitatea.

Spunem că algoritmul f_A pentru problema f este *stabil* dacă pentru orice $x \in X$ există un \tilde{x} cu

$$\frac{\|\tilde{x} - x\|}{\|x\|} = O(\text{eps}) \quad (1.10.3)$$

astfel încât

$$\frac{\|f_A(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = O(\text{eps}). \quad (1.10.4)$$

În cuvinte, un *algoritm stabil* ne dă un răspuns aproape corect la o problemă aproape exactă.

Mulți dintre algoritmi din algebra liniară numerică satisfac o condiție care este mai puternică și mai simplă decât stabilitatea. Spunem că un algoritm f_A pentru problema f este *regresiv stabil* (backward stable) dacă

$$\forall x \in X \exists \tilde{x} \text{ cu } \frac{\|\tilde{x} - x\|}{\|x\|} = O(\text{eps}) \text{ astfel încât } f_A(x) = f(\tilde{x}). \quad (1.10.5)$$

Aceasta este o întărire a definiției stabilității în sensul că $O(\text{eps})$ din (1.10.4) a fost înlocuit cu 0. În cuvinte, un *algoritm este regresiv stabil* dacă dă răspunsul corect la o problemă aproape exactă.

Observația 1.10.1. Semnificația notației

$$\|\text{cantitate calculată}\| = O(\text{eps}) \quad (1.10.6)$$

este următoarea:

- $\|\text{cantitate calculată}\|$ reprezintă norma unui număr sau a unei colecții de numere determinate de algoritmul f_A pentru o problemă f , depinzând atât de datele de intrare $x \in X$ ale lui f cât și de eps . Un exemplu este eroarea relativă.
- Procesul implicit de trecere la limită este $\text{eps} \rightarrow 0$ (adică eps corespunde lui t din (1.10.1)).
- O se aplică uniform tuturor datelor $x \in X$. La formularea rezultatelor de stabilitate această uniformitate va fi implicită.
- Pentru orice aritmetică pe calculator particulară eps este o cantitate fixată. Când vorbim despre limita $\text{eps} \rightarrow 0$, considerăm o idealizare a unui calculator sau a unei familii de calculatoare. Ecuația (1.10.6) înseamnă că dacă rulăm algoritmul în cauză pe calculatoare ce satisfac (1.4.3) și (1.4.4) pentru un șir de eps -uri descrescătoare ce tind către zero, se garantează că $\|\text{cantitate calculată}\|$ descrește proporțional cu eps sau mai repede. Acestor calculatoare ideale li se cere să satisfacă doar (1.4.3) și (1.4.4) și nimic altceva.

- Constanta implicită din notația O poate depinde și de dimensiunile argumentelor (de exemplu pentru rezolvarea unui sistem $Ax = b$ de dimensiunile lui A și b). În general, în probleme practice, creșterea erorii datorată dimensiunii este lentă, dar pot exista situații în care să apară factori cum ar fi 2^m , care fac astfel de margini inutile în practică. \diamond

Datorită echivalenței normelor pe spații finit dimensionale, pentru problemele f și algoritmi f_A definite pe astfel de spații, proprietățile de precizie, stabilitate și stabilitate regresivă au loc sau nu independent de alegerea normelor.

1.10.3. Analiza regresivă a erorilor

Stabilitatea regresivă implică precizia în sens relativ.

Teorema 1.10.2 *Presupunem că se aplică un algoritm regresiv stabil f_A unei probleme $f: X \rightarrow Y$ cu numărul de condiționare $(\text{cond } f)(x)$ pe un calculator ce satisface (1.4.3) și (1.4.4). Atunci eroarea relativă satisface*

$$\frac{\|f_A(x) - f(x)\|}{\|f(x)\|} = O((\text{cond } f)(x) \text{ eps}). \quad (1.10.7)$$

Demonstrație. Din definiția (1.10.5) a stabilității regresive, există un $\tilde{x} \in X$ ce satisface

$$\frac{\|\tilde{x} - x\|}{\|x\|} = O(\text{eps}),$$

astfel încât $f_A(x) = f(\tilde{x})$. Din definiția (1.6.5) și (1.6.6) a lui $(\text{cond } f)(x)$ aceasta implică

$$\frac{\|f_A(x) - f(x)\|}{\|f(x)\|} \leq ((\text{cond } f)(x) + o(1)) \frac{\|\tilde{x} - x\|}{\|x\|}, \quad (1.10.8)$$

unde $o(1)$ desemnează o cantitate ce converge către zero când $\text{eps} \rightarrow 0$. Combinând aceste delimitări se obține (1.10.7). \square

Procesul urmat în demonstrația teoremei 1.10.2 este cunoscut sub numele de *analiza regresivă a erorilor* (backward error analysis). Se obține o estimare a preciziei în doi pași. Primul pas este investigarea condiționării problemei. Celălalt este investigarea stabilității propriu-zise a algoritmului. Conform teoremei 1.10.2, dacă algoritmul este regresiv stabil, atunci precizia finală reflectă acel număr de condiționare.

În afară de analiza regresivă a erorilor, există și o analiză directă sau *progresivă*. Aici se estimează erorile de rotunjire la fiecare pas și apoi modul cum se compun și în final un total (secțiunea 1.3).

Experiența a arătat că pentru cei mai mulți algoritmi ai algebrei liniare numerice analiza progresivă a erorilor este mai greu de realizat decât cea regresivă. Cei mai buni algoritmi pentru cele mai multe probleme ale algebrei liniare numerice nu fac altceva mai bun decât să calculeze soluția exactă pentru niște date ușor perturbate. Analiza regresivă este o metodă de raționament bine adaptată acestei situații.

CAPITOLUL 2

Rezolvarea numerică a sistemelor de ecuații algebrice liniare

Cuprins

2.1. Elemente de Analiză matricială	24
2.2. Condiționarea unui sistem liniar	30
2.3. Metode exacte	34
2.3.1. Metoda eliminării a lui Gauss	34
2.4. Metode bazate pe factorizare	39
2.4.1. Descompunerea LU	39
2.4.2. Descompunere LUP	41
2.4.3. Factorizarea Cholesky	43
2.4.4. Descompunerea QR	45
2.5. Rafinarea iterativă	52
2.6. Algoritmul lui Strassen pentru înmulțirea matricelor	52
2.7. Rezolvarea iterativă a sistemelor algebrice liniare	56

Există două clase de metode de rezolvare a sistemelor algebrice liniare: metode *directe* sau *exacte*, care furnizează soluția într-un număr finit de pași, în ipoteza că toate calculele se fac exact (Cramer, eliminarea gaussiană, Cholesky) și metode *iterative*, care aproximează soluția generând un șir care converge către aceasta (Jacobi, Gauss-Seidel, SOR).

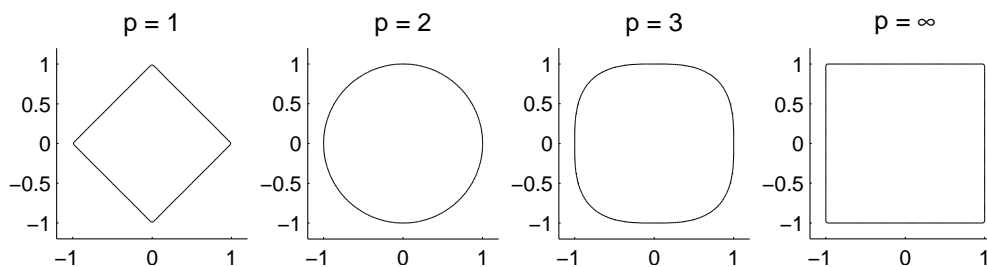


Figura 2.1: Sfera unitate pentru patru p -norme

2.1. Elemente de Analiză matricială

p -norma unui vector $x \in \mathbb{K}^n$ se definește prin

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \quad 1 \leq p < \infty.$$

Pentru $p = \infty$ norma este definită prin

$$\|x\|_\infty = \max_{i=1, n} |x_i|.$$

Norma $\|\cdot\|_2$ se numește *normă euclidiană*, $\|\cdot\|_1$ se numește *normă Minkowski*, iar $\|\cdot\|_\infty$ se numește *normă Cebîșev*. În figura 2.1 apar reprezentate sferile unitate în diverse p -norme.

Fie $A \in \mathbb{K}^{n \times n}$. Polinomul $p(\lambda) = \det(A - \lambda I)$ se numește *polinomul caracteristic* al lui A . Rădăcinile lui p se numesc *valori proprii* ale lui A , iar dacă λ este o valoare proprie a lui A vectorul $x \neq 0$ cu proprietatea că $(A - \lambda I)x = 0$ se numește *vector propriu* al lui A corespunzător valorii proprii λ .

Valoarea $\rho(A) = \max\{|\lambda| \mid \lambda \text{ valoare proprie a lui } A\}$ se numește *rază spectrală* a matricei A . Vom nota cu A^T transpusa lui A și cu A^* transpusa conjugată a lui A .

Definiția 2.1.1 O matrice A se numește:

1. normală, dacă $AA^* = A^*A$;
2. unitară, dacă $AA^* = A^*A = I$;
3. hermitiană, dacă $A = A^*$;
4. ortogonală, dacă $AA^T = A^T A = I$, A reală;
5. simetrică, dacă $A = A^T$, A reală.

Definiția 2.1.2 O normă matricială este o aplicație $\|\cdot\| : \mathbb{K}^{m \times n} \rightarrow \mathbb{R}$ care pentru orice $A, B \in \mathbb{K}^{m \times n}$ și $\alpha \in \mathbb{K}$ verifică următoarele relații

$$(NM1) \quad \|A\| \geq 0, \quad \|A\| = 0 \Leftrightarrow A = O_n;$$

$$(NM2) \quad \|\alpha A\| = |\alpha| \|A\|;$$

$$(NM3) \quad \|A + B\| \leq \|A\| + \|B\|;$$

$$(NM4) \quad \|AB\| \leq \|A\| \|B\|.$$

Primele trei proprietăți ne spun că $\|\cdot\|$ este o normă pe $\mathbb{K}^{m \times n}$, care este și spațiu vectorial de dimensiune mn , iar (NM4) este specifică normelor matriciale. Un mijloc simplu de construire a normelor matriciale este următorul: fiind dată o normă vectorială $\|\cdot\|$ pe \mathbb{C}^n , aplicația $\|\cdot\| : \mathbb{C}^{n \times n} \rightarrow \mathbb{R}$

$$\|A\| = \sup_{\substack{v \in \mathbb{C}^n \\ v \neq 0}} \frac{\|Av\|}{\|v\|} = \sup_{\substack{v \in \mathbb{C}^n \\ \|v\| \leq 1}} \|Av\| = \sup_{\substack{v \in \mathbb{C}^n \\ \|v\|=1}} \|Av\|$$

este o normă matricială numită *normă matricială subordonată* (normei vectoriale date) sau *normă naturală* (indusă de norma dată).

Observația 2.1.3. 1. Aceste norme matriciale subordonate sunt un caz particular al normei unei aplicații liniare $A : \mathbb{K}^m \rightarrow \mathbb{K}^n$.

2. O normă subordonată verifică $\|I\| = 1$.

3. Dacă matricea A este reală, marginea superioară a raportului $\|Av\|/\|v\|$ este atinsă pentru vectori reali (vezi teorema următoare). \diamond

Să calculăm acum normele subordonate normelor vectoriale $\|\cdot\|_1, \|\cdot\|_2, \|\cdot\|_\infty$.

Teorema 2.1.4 Fie $A \in \mathbb{K}^{n \times n}(\mathbb{C})$. Atunci

$$\|A\|_1 := \sup_{v \in \mathbb{C}^n \setminus \{0\}} \frac{\|Av\|_1}{\|v\|_1} = \max_j \sum_i |a_{ij}|,$$

$$\|A\|_\infty := \sup_{v \in \mathbb{C}^n \setminus \{0\}} \frac{\|Av\|_\infty}{\|v\|_\infty} = \max_i \sum_j |a_{ij}|,$$

$$\|A\|_2 := \sup_{v \in \mathbb{C}^n \setminus \{0\}} \frac{\|Av\|_2}{\|v\|_2} = \sqrt{\rho(A^*A)} = \sqrt{\rho(AA^*)} = \|A^*\|_2.$$

Norma $\|\cdot\|_2$ este invariantă la transformările unitare,

$$UU^* = I \Rightarrow \|A\|_2 = \|AU\|_2 = \|UA\|_2 = \|U^*AU\|_2.$$

Altfel spus, dacă A este normală, atunci

$$AA^* = A^*A \Rightarrow \|A\|_2 = \rho(A).$$

Demonstrație. Pentru orice vector v avem

$$\begin{aligned}\|Av\|_1 &= \sum_i \left| \sum_j a_{ij} v_j \right| \leq \sum_j |v_j| \sum_i |a_{ij}| \leq \\ &\leq \left(\max_j \sum_i |a_{ij}| \right) \|v\|_1.\end{aligned}$$

Pentru a arăta că $\max_j \sum_i |a_{ij}|$ este efectiv cel mai mic număr α pentru care are loc $\|Av\|_1 \leq \alpha \|v\|_1$, $\forall v \in \mathbb{C}^n$, să construim un vector u (care depinde de A) astfel încât

$$\|Au\|_1 = \left\{ \max_j \sum_i |a_{ij}| \right\} \|u\|_1.$$

Dacă j_0 este un indice ce verifică

$$\max_j \sum_i |a_{ij}| = \sum_i |a_{ij_0}|,$$

atunci vectorul u are componentele $u_i = 0$ pentru $i \neq j_0$, $u_{j_0} = 1$.

La fel

$$\|Av\|_\infty = \max_i \left| \sum_j a_{ij} v_j \right| \leq \left(\max_i \sum_j |a_{ij}| \right) \|v\|_\infty.$$

Fie i_0 un indice ce verifică

$$\max_i \sum_j |a_{ij}| = \sum_j |a_{i_0 j}|.$$

Vectorul u de componente $u_j = \frac{\overline{a_{i_0 j}}}{|a_{i_0 j}|}$ dacă $a_{i_0 j} \neq 0$, $u_j = 1$ dacă $a_{i_0 j} = 0$ verifică

$$\|Au\|_\infty = \left\{ \max_i \sum_j |a_{ij}| \right\} \|u\|_\infty.$$

Deoarece AA^* este hermitiană, există o descompunere proprie $AA^* = Q\Lambda Q^*$, unde Q este o matrice unitară (ale cărei coloane sunt vectori proprii) și Λ este matricea diagonală a valorilor proprii, care trebuie să fie toate reale. Dacă ar exista o valoare proprie negativă și q ar fi vectorul propriu corespunzător, am avea $0 \leq \|Aq\|_2^2 = q^T A^T A q =$

$q^T \lambda q = \lambda \|q\|_2^2$. Deci

$$\begin{aligned} \|A\|_2 &= \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \max_{x \neq 0} \frac{(x^* A^* A x)^{1/2}}{\|x\|_2} = \max_{x \neq 0} \frac{(x^* Q \Lambda Q^* x)^{1/2}}{\|x\|_2} \\ &= \max_{x \neq 0} \frac{((Q^* x)^* \Lambda Q^* x)^{1/2}}{\|Q^* x\|_2} = \max_{y \neq 0} \frac{(y^* \Lambda y)^{1/2}}{\|y\|_2} = \max_{y \neq 0} \sqrt{\frac{\sum \lambda_i y_i^2}{\sum y_i^2}} \\ &\leq \max_{y \neq 0} \sqrt{\lambda_{\max}} \sqrt{\frac{\sum y_i^2}{\sum y_i^2}}; \end{aligned}$$

egalitatea are loc dacă y este o coloană convenabil aleasă a matricei identitate.

Să arătăm că $\rho(A^*A) = \rho(AA^*)$. Dacă $\rho(A^*A) > 0$ există p astfel încât $p \neq 0$, $A^*Ap = \rho(A^*A)p$ și $Ap \neq 0$ ($\rho(A^*A) > 0$). Cum $Ap \neq 0$ și $AA^*(Ap) = \rho(A^*A)Ap$ rezultă că $0 < \rho(A^*A) \leq \rho(AA^*)$ și deci $\rho(AA^*) = \rho(A^*A)$ căci $(A^*)^* = A$. Dacă $\rho(A^*A) = 0$, avem $\rho(AA^*) = 0$. Deci în toate cazurile $\|A\|_2^2 = \rho(A^*A) = \rho(AA^*) = \|A^*\|_2^2$.

Invarianța normei $\|\cdot\|_2$ la transformări unitare nu este decât traducerea egalităților

$$\rho(A^*A) = \rho(U^*A^*AU) = \rho(A^*U^*UA) = \rho(U^*A^*UU^*AU).$$

În fine, dacă A este normală, există o matrice U astfel încât $U^*AU = \text{diag}(\lambda_i(A)) \stackrel{\text{def}}{=} \Lambda$. În aceste condiții

$$A^*A = (U\Lambda U^*)^*U\Lambda U = U\Lambda^* \Lambda U,$$

ceea ce ne arată că

$$\rho(A^*A) = \rho(\Lambda^* \Lambda) = \max_i |\lambda_i(A)|^2 = (\rho(A))^2.$$

□

Observația 2.1.5. 1) Dacă A este hermitiană sau simetrică (deci normală),

$$\|A\|_2 = \rho(A).$$

2) Dacă A este unitară sau ortogonală (deci normală),

$$\|A\|_2 = \sqrt{\rho(A^*A)} = \sqrt{\rho(I)} = 1.$$

3) Teorema 2.1.4 ne spune că matricele normale și norma $\|\cdot\|_2$ verifică

$$\|A\|_2 = \rho(A).$$

- (4) Norma $\|\cdot\|_\infty$ se mai numește norma Cebâșev sau m-normă, norma $\|\cdot\|_1$ norma lui Minkowski sau l-normă, iar norma $\|\cdot\|_2$ normă euclidiană. \diamond

Teorema 2.1.6 (1) Fie A o matrice pătratică oarecare și $\|\cdot\|$ o normă matricială oarecare (subordonată sau nu). Atunci

$$\rho(A) \leq \|A\|. \quad (2.1.1)$$

- (2) Fiind dată o matrice A și un număr $\varepsilon > 0$, există cel puțin o normă matricială subordonată astfel încât

$$\|A\| \leq \rho(A) + \varepsilon. \quad (2.1.2)$$

Demonstrație. (1) Fie p un vector ce verifică $p \neq 0$, $Ap = \lambda p$, $|\lambda| = \rho(A)$ și q un vector astfel încât $pq^T \neq 0$. Deoarece

$$\rho(A)\|pq^T\| = \|\lambda pq^T\| = \|Apq^T\| \leq \|A\|\|pq^T\|,$$

rezultă (2.1.1).

(2) Fie A o matrice dată. Există o matrice inversabilă U astfel încât $U^{-1}AU$ este triunghiulară superior (de fapt U este unitară). De exemplu

$$U^{-1}AU = \begin{pmatrix} \lambda_1 & t_{12} & t_{13} & \dots & t_{1,n} \\ & \lambda_2 & t_{23} & \dots & t_{2,n} \\ & & \ddots & & \vdots \\ & & & \lambda_{n-1} & t_{n-1,n} \\ & & & & \lambda_n \end{pmatrix}.$$

scalarii λ_i fiind valorile proprii ale matricei A . (Pentru demonstrație a se vedea teorema 7.2.6 din capitolul refcapvalpr.) Fiecărui scalar $\delta \neq 0$ îi asociem matricea

$$D_\delta = \text{diag}(1, \delta, \delta^2, \dots, \delta^{n-1}),$$

astfel ca

$$(UD_\delta)^{-1}A(UD_\delta) = \begin{pmatrix} \lambda_1 & \delta t_{12} & \delta^2 t_{13} & \dots & \delta^{n-1} t_{1n} \\ & \lambda_2 & \delta t_{23} & \dots & \delta^{n-2} t_{2n} \\ & & \ddots & & \vdots \\ & & & \lambda_{n-1} & \delta t_{n-1n} \\ & & & & \lambda_n \end{pmatrix}.$$

Fiind dat $\varepsilon > 0$, fixăm δ astfel ca

$$\sum_{j=i+1}^n |\delta^{j-i} t_{ij}| \leq \varepsilon, \quad 1 \leq i \leq n-1.$$

Atunci aplicația

$$\|\cdot\| : B \in \mathbb{K}^{n \times n} \rightarrow \|B\| = \|(UD_\delta)^{-1}B(UD_\delta)\|_\infty,$$

care depinde de A și de ε răspunde problemei. Într-adevăr, avem pe de o parte

$$\|A\| \leq \rho(A) + \varepsilon$$

conform alegerii lui δ și definiției normei $\|\cdot\|_\infty$ ($\|c_{ij}\|_\infty = \max_i \sum_j |c_{ij}|$) și pe de altă parte ea este o normă matricială subordonată normei vectoriale

$$v \in \mathbb{K}^n \rightarrow \|(UD_\delta)^{-1}v\|_\infty.$$

□

Un exemplu important de normă nesubordonată este dat de teorema următoare.

Teorema 2.1.7 *Aplicația $\|\cdot\|_E : \mathbb{K}^{n \times n} \rightarrow \mathbb{R}$ definită prin*

$$\|A\|_E = \left\{ \sum_i \sum_j |a_{ij}|^2 \right\}^{1/2} = \{tr(A^*A)\}^{1/2}$$

este o normă matricială nesubordonată, invariantă la transformările unitare

$$UU^* = I \Rightarrow \|A\|_E = \|AU\|_E = \|UA\|_E = \|U^*AU\|_E$$

și care verifică

$$\|A\|_2 \leq \|A\|_E \leq \sqrt{n}\|A\|_2, \forall A \in \mathbb{K}^{n \times n}.$$

Demonstrație. $\|\cdot\|_E$ este norma euclidiană pe $\mathbb{K}^{n \times n}$ de dimensiune n^2 . Proprietatea (NM4) se demonstrează cu inegalitatea Cauchy-Buniakowski-Schwarz

$$\begin{aligned} \|AB\|_E^2 &= \sum_{i,j} \left| \sum_k a_{jk}b_{kj} \right|^2 \leq \sum_{i,j} \left\{ \sum_k |a_{ik}|^2 \right\} \left\{ \sum_l |b_{lj}|^2 \right\} = \\ &= \left\{ \sum_{i,k} |a_{ik}|^2 \right\} \left\{ \sum_{j,l} |b_{lj}|^2 \right\} = \|A\|_E^2 \|B\|_E^2. \end{aligned}$$

Această normă nu este subordonată, deoarece $\|I\|_E = \sqrt{n}$. Dacă U este o matrice unitară

$$\begin{aligned} \|A\|_E^2 &= tr(A^*A) = tr(U^*A^*AU) = \\ &= \|AU\|_E^2 = tr(U^*A^*UA) = \|UA\|_E^2. \end{aligned}$$

În fine inegalitățile din enunț rezultă din inegalitățile

$$\rho(A^*A) \leq tr(A^*A) \leq n\rho(A^*A).$$

□

Norma $\|\cdot\|_E$ se numește *normă Frobenius*.

Teorema 2.1.8 Fie B o matrice pătratică. Următoarele afirmații sunt echivalente:

- (1) $\lim_{k \rightarrow \infty} B^k = 0$;
- (2) $\lim_{k \rightarrow \infty} B^k v = 0, \forall v \in \mathbb{K}^n$;
- (3) $\rho(B) < 1$;
- (4) Există o normă matricială subordonată astfel încât $\|B\| < 1$.

Demonstrație. (1) \Rightarrow (2)

$$\|B^k v\| \leq \|B^k\| \|v\| \Rightarrow \lim_{k \rightarrow \infty} B^k v = 0$$

(2) \Rightarrow (3) Dacă $\rho(B) \geq 1$ putem găsi p astfel încât $p \neq 0, Bp = \lambda p, |\lambda| \geq 1$. Atunci șirul de vectori $(B^k p)_{k \in \mathbb{N}}$ ar putea să nu convergă către 0.

(3) \Rightarrow (4) $\rho(B) < 1 \Rightarrow \exists \|\cdot\|$ astfel încât $\|B\| \leq \rho(B) + \varepsilon, \forall \varepsilon > 0$ deci $\|B\| < 1$.

(4) \Rightarrow (1) Este suficient să aplicăm inegalitatea $\|B^k\| \leq \|B\|^k$. \square

2.2. Condiționarea unui sistem liniar

Fie sistemul (exemplul este datorat lui Wilson)

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix},$$

cu soluția $(1, 1, 1, 1)^T$ și considerăm sistemul perturbat, în care membrul drept este foarte puțin modificat, matricea rămânând neschimbată

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} x_1 + \delta x_1 \\ x_2 + \delta x_2 \\ x_3 + \delta x_4 \\ x_4 + \delta x_4 \end{pmatrix} = \begin{pmatrix} 32.1 \\ 22.9 \\ 33.1 \\ 30.9 \end{pmatrix},$$

cu soluția $(9.2, -12.6, 4.5, -1.1)^T$. Altfel spus, o eroare de 1/200 în date (aici componentele din membrul drept) atrage o eroare relativă de 10/1 asupra rezultatului, deci o mărire a erorii relative de ordin 2000!

Considerăm acum sistemul având de această dată matricea perturbată

$$\begin{pmatrix} 10 & 7 & 8.1 & 7.2 \\ 7.08 & 5.04 & 6 & 5 \\ 8 & 5.98 & 9.89 & 9 \\ 6.99 & 4.99 & 9 & 9.98 \end{pmatrix} \begin{pmatrix} x_1 + \Delta x_1 \\ x_2 + \Delta x_2 \\ x_3 + \Delta x_4 \\ x_4 + \Delta x_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

cu soluția $(-81, 137, -34, 22)^T$. Din nou, o variație mică în datele de intrare (aici, elementele matricei) modifică complet rezultatul (soluția sistemului liniar). Matricea are un aspect „bun“, ea este simetrică, determinantul ei este 1, iar inversa ei este

$$\begin{pmatrix} 25 & -41 & 10 & -6 \\ -41 & 68 & -17 & 10 \\ 10 & -17 & 5 & -3 \\ -6 & 10 & -3 & 2 \end{pmatrix},$$

care este de asemenea simpatcă.

Să analizăm aceste fenomene. În primul caz se dă o matrice *inversabilă* A și se compară soluțiile *exacte* x și $x + \delta x$ ale sistemelor

$$\begin{aligned} Ax &= b \\ A(x + \delta x) &= b + \delta b. \end{aligned}$$

Fie $\|\cdot\|$ o normă vectorială oarecare și $\|\cdot\|$ norma matricială subordonată. Din egalitățile $\delta x = A^{-1}\delta b$ și $b = Ax$ se deduce

$$\|\Delta x\| \leq \|A^{-1}\| \|\delta b\|, \quad \|b\| \leq \|A\| \|x\|.$$

Eroarea relativă a rezultatului $\frac{\|\delta x\|}{\|x\|}$ este majorată în funcție de eroarea relativă a datelor prin

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}.$$

În al doilea caz, când matricea variază, avem de comparat soluțiile exacte ale sistemelor

$$\begin{aligned} Ax &= b \\ (A + \Delta A)(x + \Delta x) &= b. \end{aligned}$$

Din egalitatea $\Delta x = -A^{-1}\Delta A(x + \Delta x)$ se deduce

$$\|\Delta x\| \leq \|A^{-1}\| \|\Delta A\| \|x + \Delta x\|$$

care se mai poate scrie

$$\frac{\|\Delta x\|}{\|x + \Delta x\|} \leq \|A\| \|A^{-1}\| \frac{\|\Delta A\|}{\|A\|}.$$

Dacă A este nesingulară, numărul

$$\text{cond}(A) = \|A\| \|A^{-1}\| \quad (2.2.1)$$

se numește *număr de condiționare* al matricei A .

Se poate da o estimare a numărului de condiționare în care să intervină simultan și perturbațiile lui A și ale lui b . Considerăm sistemul parametrizat, cu parametrul t

$$(A + t\Delta A)x(t) = b + t\Delta b, \quad x(0) = x.$$

Matricea A fiind nesingulară, funcția x este diferențiabilă în $t = 0$:

$$\dot{x}(0) = A^{-1}(\Delta b - \Delta A x).$$

Dezvoltarea Taylor a lui $x(t)$ este dată de

$$x(t) = x + t\dot{x}(0) + O(t^2).$$

Rezultă că eroarea absolută poate fi estimată utilizând

$$\begin{aligned} \|\Delta x(t)\| &= \|x(t) - x\| \leq |t| \|x'(0)\| + O(t^2) \\ &\leq |t| \|A^{-1}\| (\|\Delta b\| + \|\Delta A\| \|x\|) + O(t^2) \end{aligned}$$

și datorită lui $\|b\| \leq \|A\| \|x\|$ obținem pentru eroarea relativă

$$\begin{aligned} \frac{\|\Delta x(t)\|}{\|x\|} &\leq |t| \|A^{-1}\| \left(\frac{\|\Delta b\|}{\|x\|} + \|\Delta A\| \right) + O(t^2) \\ &\leq \|A\| \|A^{-1}\| |t| \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right) + O(t^2). \end{aligned} \quad (2.2.2)$$

Introducând notațiile

$$\rho_A(t) := |t| \frac{\|\Delta A\|}{\|A\|}, \quad \rho_b(t) := |t| \frac{\|\Delta b\|}{\|b\|}$$

pentru erorile relative în A și b , estimarea erorii (2.2.2) se scrie sub forma

$$\frac{\|\Delta x(t)\|}{\|x\|} \leq \text{cond}(A) (\rho_A + \rho_b) + O(t^2).$$

Exemplul 2.2.1 (Exemple de matrice prost condiționate). Matricea lui Hilbert ¹ $H_n = (h_{ij})$, cu $h_{ij} = \frac{1}{i+j-1}$, $i, j = \overline{1, n}$ are ordinul de mărime al numărului de condiționare relativ la norma euclidiană (dat de Szegő²)

$$\text{cond}_2(H_n) \sim \frac{(\sqrt{2} + 1)^{4n+4}}{2^{14/4} \sqrt{\pi n}}.$$

Pentru diverse valori ale lui n se obține

n	$\text{cond}_2(H_n)$
10	$1.6 \cdot 10^{13}$
20	$2.45 \cdot 10^{28}$
40	$7.65 \cdot 10^{58}$

Un alt exemplu este matricea Vandermonde. Dacă elementele sunt echidistante în $[-1, 1]$, atunci

$$\text{cond}_\infty(V_n) \sim \frac{1}{\pi} e^{-\frac{\pi}{4}} e^{n(\frac{\pi}{4} + \frac{1}{2} \ln 2)},$$

iar pentru $t_i = \frac{1}{i}$, $i = \overline{1, n}$ avem

$$\text{cond}_\infty(V_n) > n^{n+1}.$$

◇

1



David Hilbert (1862-1943) a fost cel mai important reprezentant al școlii matematice din Göttingen. Contribuțiile sale fundamentale în aproape toate domeniile matematicii — algebră, teoria numerelor, geometrie, ecuații integrale, calcul variațional și fundamentele matematicii — și în particular cele 23 de probleme celebre pe care le-a propus în 1900 la un congres internațional al matematicienilor de la Paris, au dat un nou impuls și o nouă direcție matematicii din secolul al XX-lea.

2



Gabor Szegő (1895-1985) Unul dintre cei mai importanți matematicieni maghiari din secolul al XX-lea. Contribuții importante în domeniul problemelor extremale și matricelor Toeplitz.

2.3. Metode exacte

2.3.1. Metoda eliminării a lui Gauss

Să considerăm sistemul liniar de n ecuații cu n necunoscute

$$Ax = b, \quad (2.3.1)$$

unde $A \in \mathbb{K}^{n \times n}$, $b \in \mathbb{K}^{n \times 1}$ sunt date, iar $x \in \mathbb{K}^{n \times 1}$ este necunoscuta, sau scris detaliat

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 & (E_1) \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 & (E_2) \\ \vdots & \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n & (E_n) \end{cases} \quad (2.3.2)$$

Metoda eliminării a lui Gauss³ are două etape:

- e1) transformarea sistemului dat într-unul echivalent, triunghiular;
- e2) rezolvarea sistemului triunghiular prin substituție inversă.

La rezolvarea sistemului (2.3.1) sau (2.3.2) sunt permise următoarele operații:

1. Ecuația E_i poate fi înmulțită cu $\lambda \in \mathbb{K}^*$. Această operație se va nota cu $(\lambda E_i) \rightarrow (E_i)$.
2. Ecuația E_j poate fi înmulțită cu $\lambda \in \mathbb{K}^*$ și adunată la ecuația E_i , iar rezultatul utilizat în locul lui E_i , notație $(E_i + \lambda E_j) \rightarrow (E_i)$.



Johann Carl Friedrich Gauss (1777-1855) a fost unul dintre cei mai mari matematicieni ai secolului al nouăsprezecelea și probabil al tuturor timpurilor. A trăit aproape toată viața în Göttingen, unde a fost directorul observatorului astronomic 40 de ani. În timpul studenției la Göttingen, Gauss a descoperit că poligonul cu 17 laturi poate fi construit cu rigla și compasul, rezolvând astfel o problemă deschisă a antichității. În dizertația sa a dat prima demonstrație a teoremei fundamentale a algebrei. A avut contribuții fundamentale în teoria numerelor, geometrie diferențială și neeuclidiană, funcții eliptice și hipergeometrice, mecanică cerească și geodezie și diverse ramuri ale fizicii, în special magnetism și optică. Eforturile sale de calcul în mecanica cerească și geodezie, bazate pe principiul celor mai mici pătrate, au necesitat rezolvarea manuală a unor sisteme de ecuații liniare mari, la care a utilizat metodele cunoscute astăzi sub numele de eliminare gaussiană și metoda relaxării. Lucrările lui Gauss în domeniul cuadraturilor numerice continuă munca predecesorilor săi Newton și Cotes.

3. Ecuațiile E_i și E_j pot fi interschimbate, notație $(E_i) \longleftrightarrow (E_j)$.

Pentru a exprima convenabil operațiile necesare pentru transformarea sistemului în unul triunghiular vom lucra cu matricea extinsă:

$$\tilde{A} = [A, b] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & a_{1,n+1} \\ a_{21} & a_{22} & \dots & a_{2n} & a_{2,n+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} & a_{n,n+1} \end{array} \right]$$

cu $a_{i,n+1} = b_i$.

Presupunând că $a_{11} \neq 0$, vom elimina coeficienții lui x_1 din E_j , pentru $j = \overline{2, n}$ prin operația $(E_j - (a_{j1}/a_{11})E_1) \rightarrow (E_j)$. Vom proceda apoi la fel cu coeficienții lui x_i , pentru $i = \overline{2, n-1}$, $j = i+1, n$. Aceasta este posibil dacă $a_{ii} \neq 0$.

Procedura poate fi descrisă astfel: se formează o secvență de matrice extinse $\tilde{A}^{(1)}$, $\tilde{A}^{(2)}$, ..., $\tilde{A}^{(n)}$, unde $\tilde{A}^{(1)} = A$ și $\tilde{A}^{(k)}$ are elementele $a_{ij}^{(k)}$ date de

$$\left(E_i - \frac{a_{i,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}} E_{k-1} \right) \longrightarrow (E_i)$$

sau desfășurat

$$a_{i,j}^{(k)} = \begin{cases} a_{ij}^{(k-1)}, & \text{pentru } i = \overline{1, k-1}, j = \overline{1, n+1} \\ 0, & \text{pentru } i = \overline{k, n}, j = \overline{1, k-1} \\ a_{ij}^{(k-1)} - \frac{a_{i,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}} a_{k-1,j}^{(k-1)} & \text{pentru } i = \overline{k, n}, j = \overline{k, n+1} \end{cases}.$$

Observația 2.3.1. Notația $a_{ij}^{(p)}$ semnifică valoarea elementului a_{ij} la pasul p . ◇

Astfel

$$\tilde{A}^{(k)} = \left[\begin{array}{cccc|cccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \dots & a_{1,k-1}^{(1)} & a_{1k}^{(1)} & \dots & a_{1n}^{(1)} & a_{1,n+1}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \dots & a_{2,k-1}^{(2)} & a_{2,k}^{(2)} & \dots & a_{2n}^{(2)} & a_{2,n+1}^{(2)} \\ \vdots & & \ddots & \dots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & & & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & & & & a_{k-1,k-1}^{(k-1)} & a_{k-1,k}^{(k-1)} & \dots & a_{k-1,n}^{(k-1)} & a_{k-1,n+1}^{(k-1)} \\ \vdots & & & & 0 & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} & a_{k,n+1}^{(k)} \\ \vdots & & & & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & 0 & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} & a_{n,n+1}^{(k)} \end{array} \right]$$

reprezintă sistemul liniar echivalent în care variabila x_{k-1} a fost eliminată din ecuațiile E_k, E_{k+1}, \dots, E_n . Sistemul corespunzător lui $\tilde{A}^{(n)}$ este un sistem triunghiular echivalent cu sistemul inițial

$$\begin{cases} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \dots + a_{1n}^{(1)}x_n = a_{1,n+1}^{(1)} \\ a_{22}^{(2)}x_2 + \dots + a_{2n}^{(2)}x_n = a_{2,n+1}^{(2)} \\ \vdots \\ a_{nn}^{(n)}x_n = a_{n,n+1}^{(n)} \end{cases}.$$

Se obține

$$x_n = \frac{a_{n,n+1}^{(n)}}{a_{n,n}^{(n)}}$$

și în general

$$x_i = \frac{1}{a_{ii}^{(i)}} \left(a_{i,n+1}^{(i)} - \sum_{j=i+1}^n a_{ij}^{(i)} x_j \right), \quad i = \overline{1, n-1}.$$

Pentru ca procedura să fie aplicabilă trebuie ca $a_{ii}^{(i)} \neq 0, i = \overline{1, n}$. Elementul $a_{ii}^{(i)}$ se numește *element pivot*. Dacă în timpul algoritmului de eliminare gaussiană la pasul k se obține pivotul $a_{kk}^{(k)} = 0$, se poate face interschimbarea de linii $(E_k) \leftrightarrow (E_p)$, unde $k+1 \leq p \leq n$ este cel mai mic întreg cu proprietatea $a_{pk}^{(k)} \neq 0$. În practică sunt necesare astfel de operații chiar și în cazul când pivotul este nenul. Aceasta din cauză că un pivot mic în comparație cu elementele care urmează după el în aceeași coloană duce la erori de rotunjire substanțiale. Aceasta se poate remedia alegând ca pivot elementul din aceeași coloană care este situat sub diagonală și care are cea mai mare valoare absolută, adică determinând p astfel încât

$$|a_{pk}^{(k)}| = \max_{k \leq i \leq n} |a_{ik}^{(k)}|$$

și apoi făcând $(E_k) \leftrightarrow (E_p)$. Această tehnică se numește *pivotare maximală pe coloană* sau *pivotare parțială*.

O altă tehnică care reduce erorile și preîntâmpină anulările în aritmetica flotantă este *pivotarea scalată pe coloană*. La primul pas se definește un factor de scalare pentru fiecare linie

$$s_i = \max_{j=\overline{1, n}} |a_{ij}| \text{ sau } s_i = \sum_{j=1}^n |a_{ij}|.$$

Dacă există i astfel încât $s_i = 0$, matricea este singulară.

La pașii următori se determină ce interschimbări se vor realiza. La pasul i se determină cel mai mic întreg $p, i \leq p \leq n$, pentru care

$$\frac{|a_{pi}|}{s_p} = \max_{1 \leq j \leq n} \frac{|a_{ji}|}{s_j}$$

și apoi $(E_i) \leftrightarrow (E_p)$. Efectul scalării este de a ne asigura că elementul cel mai mare din fiecare coloană are mărimea relativă 1 înainte de a realiza comparațiile pentru interschimbarea liniilor. Scalarea se realizează doar în scop de comparație, așa că împărțirea cu factorul de scalare nu produce erori de rotunjire. O a treia metodă este cea a pivotării totale (sau maximale). La această metodă la pasul k se determină

$$\max\{|a_{ij}|, i = \overline{k, n}, j = \overline{k, n}\}$$

și se realizează și interschimbări de linii și de coloane.

Pivotarea a fost introdusă de Goldstine și von Neumann în 1947 [23].

Dacă matricea A este singulară și are rangul $p - 1$ atunci se obține

$$\tilde{A}^{(p)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1,p-1}^{(1)} & a_{1p}^{(1)} & \dots & a_{1n}^{(1)} & a_{1,n+1}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2,p-1}^{(2)} & a_{2p}^{(2)} & \dots & a_{2n}^{(2)} & a_{2,n+1}^{(2)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & a_{p-1,p-1}^{(p-1)} & a_{p-1,p}^{(p-1)} & \dots & a_{p-1,n}^{(p-1)} & a_{p-1,n+1}^{(p-1)} \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 & a_{p,n+1}^{(p)} \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 & a_{n,n+1}^{(n)} \end{bmatrix}.$$

Dacă $a_{i,n+1}^{(p)} = b_i^{(p)} = 0$, $i = \overline{p, n}$, atunci sistemul este compatibil nedeterminat, iar dacă există $q \in \{p, \dots, n\}$ astfel încât $a_{q,n+1}^{(p)} = b_q^{(p)} \neq 0$ sistemul este incompatibil.

Deci metodele de eliminare de tip Gauss permit atât rezolvarea cât și discuția sistemelor de ecuații liniare.

Observația 2.3.2. Dăm câteva sugestii care pot duce la îmbunătățirea timpului de execuție.

1. La metodele de pivotare nu este nevoie să se realizeze fizic interschimbarea de linii și/sau coloane, ci se pot păstra unul (sau doi vectori) de permutări $p(q)$ cu semnificația $p[i](q[i])$ este linia (coloana) care a fost interschimbată cu linia (coloana) i .
2. Elementele de sub diagonală (care devin 0) pot să nu fie calculate.
3. Matricea A se poate inversa rezolvând sistemul $Ax = e_k$, $k = \overline{1, n}$, unde e_k sunt vectorii bazei canonice din \mathbb{K}^n . Metoda se numește metoda ecuațiilor simultane. \diamond

Să analizăm metoda eliminării a lui Gauss. Descrierea ei apare în algoritmul 2.1.

Ca măsură a complexității vom considera numărul de operații aritmetice flotante, desemnate prescurtat prin *flops*.

În corpul ciclului interior, liniile 10–11 avem $2n - 2i + 3$ flops, deci pentru întreg ciclul $(n - i)(2n - 2i + 3)$ flops. Pentru ciclul exterior avem un total general de

$$\sum_{i=1}^{n-1} (n - i)(2n - 2i + 3) \sim \frac{2n^3}{3}.$$

Pentru substituția inversă avem $\Theta(n^2)$ flops.

Total general, $\Theta(n^3)$.

Algoritmul 2.1 Rezolvă sistemul $Ax = b$ prin metoda eliminării a lui Gauss

Intrare: Matricea extinsă $A = (a_{ij})$, $i = \overline{1, n}$, $j = \overline{1, n + 1}$

Ieșire: Soluțiile x_1, \dots, x_n sau un mesaj de eroare

```

{Eliminare}
1: for  $i := 1$  to  $n - 1$  do
2:   Fie  $p$  cel mai mic întreg  $i \leq p \leq n$  și  $a_{pi} \neq 0$ 
3:   if  $\nexists p$  then
4:     mesaj (' $\nexists$  soluție unică'); STOP
5:   end if
6:   if  $p \neq i$  then
7:      $(E_p) \leftrightarrow (E_i)$ 
8:   end if
9:   for  $j := i + 1$  to  $n$  do
10:     $m_{ji} := a_{ji}/a_{ii}$ ;
11:     $(E_j - m_{ji}E_i) \rightarrow (E_j)$ ;
12:   end for
13: end for
14: if  $a_{nn} = 0$  then
15:   mesaj (' $\nexists$  soluție unică'); STOP
16: end if
{Substituție inversă}
17:  $x_n := a_{n,n+1}/a_{nn}$ ;
18: for  $i := n - 1$  downto  $1$  do
19:    $x_i = \left[ a_{i,n+1} - \sum_{j=i+1}^n a_{ij}x_j \right] / a_{ii}$ ;
20: end for
21: Extrage  $(x_1, \dots, x_n)$  {succes} STOP.
```

2.4. Metode bazate pe factorizare

2.4.1. Descompunerea LU

Teorema 2.4.1 Dacă eliminarea gaussiană pentru sistemul $Ax = b$ se poate realiza fără interschimbări de linii, atunci A se poate factoriza în $A = LU$, L - triunghiulară inferior, U - triunghiulară superior (Perechea (L, U) se numește descompunerea LU a matricei A).

Avantaje. $Ax = b \Leftrightarrow LUX = b \Leftrightarrow Ly = b \wedge Ux = y$.

Dacă avem de rezolvat mai multe sisteme $Ax = b_i$, $i = \overline{1, m}$, fiecare rezolvare durează $\Theta(n^3)$; dacă se factorizează la început rezolvarea unui sistem durează $\Theta(n^2)$, factorizarea durează $\Theta(n^3)$.

Observația 2.4.2. U este matricea triunghiulară superior obținută în urma eliminării gaussiene, iar L este matricea multiplicatorilor m_{ij} . \diamond

Dacă eliminarea gaussiană se face cu interschimbări avem de asemenea $A = LU$, dar L nu este triunghiulară inferior.

Metoda obținută se numește factorizare LU .

Situații când eliminarea gaussiană se face fără interschimbări:

- A este diagonal dominantă pe linii, adică

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = \overline{1, n}$$

- A este pozitiv definită ($\forall x \neq 0 \ x^*Ax > 0$).

Demonstrația teoremei 2.4.1. (schiță) Pentru $n > 1$ partiționăm A astfel

$$A = \left[\begin{array}{c|ccc} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{array} \right] = \left[\begin{array}{cc} a_{11} & w^* \\ v & A' \end{array} \right],$$

unde: v - vector coloană de dimensiune $n - 1$, w^* - vector linie de dimensiune $n - 1$.
Putem factoriza A prin

$$A = \left[\begin{array}{cc} a_{11} & w^* \\ v & A' \end{array} \right] = \left[\begin{array}{cc} 1 & 0 \\ v/a_{11} & I_{n-1} \end{array} \right] \left[\begin{array}{cc} a_{11} & w^* \\ 0 & A' - vw^*/a_{11} \end{array} \right].$$

Matricea $A' - vw^*/a_{11}$ se numește *complement Schur* al lui A în raport cu a_{11} . Se continuă apoi cu descompunerea recursivă a complementului Schur:

$$A' - vw^*/a_{11} = L'U'.$$

$$\begin{aligned}
A &= \begin{bmatrix} 1 & 0 \\ v/a_{11} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & A' - vw^*/a_{11} \end{bmatrix} = \\
&= \begin{bmatrix} 1 & 0 \\ v/a_{11} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & L'U' \end{bmatrix} = \\
&= \begin{bmatrix} 1 & 0 \\ v/a_{11} & L' \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & U' \end{bmatrix}.
\end{aligned}$$

□

Exemplul 2.4.3. Să se calculeze descompunerea LU a matricii

$$A = \left[\begin{array}{c|ccc} 2 & 3 & 1 & 5 \\ \hline 6 & 13 & 5 & 19 \\ 2 & 19 & 10 & 23 \\ 4 & 10 & 11 & 31 \end{array} \right]$$

Matricea inițială este

$$\begin{array}{c|ccc} 2 & 3 & 1 & 5 \\ \hline 3 & 4 & 2 & 4 \\ 1 & 16 & 9 & 18 \\ 2 & 4 & 9 & 21 \end{array}$$

iar primul complement Schur

$$\begin{aligned}
A' - vw^*/a_{11} &= \begin{pmatrix} 13 & 5 & 19 \\ 19 & 10 & 23 \\ 10 & 11 & 31 \end{pmatrix} - \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix} (3 \ 1 \ 5) = \\
&= \begin{pmatrix} 13 & 5 & 19 \\ 19 & 20 & 23 \\ 10 & 11 & 31 \end{pmatrix} - \begin{pmatrix} 9 & 3 & 15 \\ 3 & 1 & 5 \\ 6 & 2 & 10 \end{pmatrix} = \begin{pmatrix} 4 & 2 & 4 \\ 16 & 9 & 18 \\ 4 & 9 & 21 \end{pmatrix}.
\end{aligned}$$

Continuăm cu descompunerea recursivă a complementelor Schur de ordinul 2 și 1:

$$\begin{array}{c|cc} 2 & 3 & 1 & 5 \\ \hline 3 & 4 & 2 & 4 \\ \hline 1 & 4 & 1 & 2 \\ 2 & 1 & 7 & 17 \end{array}$$

$$\begin{pmatrix} 9 & 18 \\ 9 & 21 \end{pmatrix} - \begin{pmatrix} 4 \\ 1 \end{pmatrix} (2, 4) = \begin{pmatrix} 1 & 2 \\ 7 & 17 \end{pmatrix}$$

$$\begin{array}{ccc|c} 2 & 3 & 1 & 5 \\ 3 & 4 & 2 & 4 \\ 1 & 4 & 1 & 2 \\ \hline 2 & 1 & 7 & 3 \end{array}$$

$$A' - vw^*/a_{11} = 17 - 7 \cdot 2 = 3$$

Verificare:

$$\begin{pmatrix} 2 & 3 & 1 & 5 \\ 6 & 13 & 5 & 19 \\ 2 & 19 & 10 & 23 \\ 4 & 10 & 11 & 31 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 3 & 1 & 0 & 0 \\ 1 & 4 & 1 & 0 \\ 2 & 1 & 7 & 1 \end{pmatrix} \begin{pmatrix} 2 & 3 & 1 & 5 \\ 0 & 4 & 2 & 4 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 3 \end{pmatrix} \quad \diamond$$

Avem mai multe posibilități de alegere pentru u_{ii} și ℓ_{ii} , $i = \overline{1, n}$. De exemplu dacă $\ell_{ii} = 1$ avem *factorizare Doolittle*, dacă $u_{ii} = 1$ avem *factorizare Crout*.

2.4.2. Descompunere LUP

Ideea din spatele descompunerii LUP este de a găsi 3 matrice pătratice L, U și P unde L - triunghiulară inferior, U - triunghiulară superior, P matrice de permutare, astfel încât $PA = LU$.

Tripletul (L, U, P) se va numi *descompunerea LUP* a matricei A .

Rezolvarea sistemului $Ax = b$ este echivalentă cu rezolvarea a două sisteme triunghiulare, deoarece

$$Ax = b \Leftrightarrow LUx = Pb \Leftrightarrow Ly = Pb \wedge Ux = y$$

și

$$Ax = P^{-1}LUx = P^{-1}Ly = P^{-1}Pb = b.$$

Vom alege ca pivot în locul lui a_{11} elementul a_{k1} . Efectul este înmulțirea cu o matrice de permutare Q :

$$QA = \begin{bmatrix} a_{k1} & w^* \\ v & A' \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ v/a_{k1} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & A' - vw^*/a_{k1} \end{bmatrix}.$$

Determinăm mai departe descompunerea LUP a complementului Schur.

$$P'(A' - vw^*/a_{k1}) = L'U'.$$

Definim

$$P = \begin{bmatrix} 1 & 0 \\ 0 & P' \end{bmatrix} Q,$$

care este tot o matrice de permutare. Avem acum

$$\begin{aligned}
 PA &= \begin{bmatrix} 1 & 0 \\ 0 & P' \end{bmatrix} QA = \\
 &= \begin{bmatrix} 1 & 0 \\ 0 & P' \end{bmatrix} \begin{bmatrix} 1 & 0 \\ v/a_{k1} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & A' - vw^*/a_{k1} \end{bmatrix} = \\
 &= \begin{bmatrix} 1 & 0 \\ P'v/a_{k1} & P' \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & A' - vw^*/a_{k1} \end{bmatrix} = \\
 &= \begin{bmatrix} 1 & 0 \\ P'v/a_{k1} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & P'(A' - vw^*/a_{k1}) \end{bmatrix} = \\
 &= \begin{bmatrix} 1 & 0 \\ P'v/a_{k1} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & L'U' \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ P'v/a_{k1} & L' \end{bmatrix} \begin{bmatrix} a_{k1} & w^* \\ 0 & U' \end{bmatrix}.
 \end{aligned}$$

De notat că în acest raționament, atât vectorul coloană cât și complementul Schur se înmulțesc cu matricea de permutare P' .

Exemplul 2.4.4. Să se calculeze descompunerea LUP a matricei

$$\begin{pmatrix} 2 & 0 & 2 & 0.6 \\ 3 & 3 & 4 & -2 \\ 5 & 5 & 4 & 2 \\ -1 & -2 & 3.4 & -1 \end{pmatrix}$$

$$\begin{array}{c|cccc} 1 & 2 & 0 & 2 & 0.6 \\ \hline 2 & 3 & 3 & 4 & -2 \\ 3 & 5 & 5 & 4 & 2 \\ 4 & -1 & -2 & 3.4 & -1 \end{array} \qquad \begin{array}{c|cccc} 3 & 5 & 5 & 4 & 2 \\ \hline 2 & 3 & 3 & 4 & -2 \\ 1 & 2 & 0 & 2 & 0.6 \\ 4 & -1 & -2 & 3.4 & -1 \end{array}$$

$$\begin{pmatrix} 3 & 4 & -2 \\ 0 & 2 & 0.6 \\ -2 & 3.4 & -1 \end{pmatrix} - \begin{pmatrix} 0.6 \\ 0.4 \\ 0.2 \end{pmatrix} (5, 4, 2) = \begin{pmatrix} 0 & 1.6 & -3.2 \\ -2 & 0.4 & -0.2 \\ -1 & 4.2 & -0.6 \end{pmatrix}$$

$$\begin{array}{ccccc} 3 & 5 & 5 & 4 & 2 \\ 2 & 0.6 & 0 & 1.6 & -3.2 \\ 1 & 0.4 & -2 & 0.4 & -0.2 \\ 4 & 0.2 & -1 & 4.2 & -0.6 \end{array} \qquad \begin{array}{ccccc} 3 & 5 & 5 & 4 & 2 \\ 2 & 0.6 & 0 & 1.6 & -3.2 \\ 1 & 0.4 & -2 & 0.4 & -0.2 \\ 4 & -0.2 & -1 & 4.2 & -0.6 \end{array}$$

$$\begin{array}{ccccc} 3 & 5 & 5 & 4 & 2 \\ 1 & 0.4 & -2 & 0.4 & -0.2 \\ 2 & 0.6 & 0 & 1.6 & -3.2 \\ 4 & -0.2 & -1 & 4.2 & 0.6 \end{array} \qquad \begin{array}{ccccc} 3 & 5 & 5 & 4 & 2 \\ 1 & 0.4 & -2 & 0.4 & -0.2 \\ 2 & 0.6 & 0 & 1.6 & -3.2 \\ 4 & -0.2 & -0.5 & 4 & -0.5 \end{array}$$

$$\begin{array}{ccccc}
 3 & 5 & 5 & 4 & 2 \\
 1 & 0.4 & -2 & 0.4 & -0.2 \\
 2 & 0.6 & 0 & 1.6 & -3.2 \\
 4 & -0.2 & 0.5 & 4 & -0.5
 \end{array}
 \quad
 \begin{array}{ccccc}
 3 & 5 & 5 & 4 & 2 \\
 1 & 0.4 & -2 & 0.4 & -0.2 \\
 4 & -0.2 & 0.5 & 4 & -0.5 \\
 2 & 0.6 & 0 & 1.6 & -3.2
 \end{array}$$

$$\begin{array}{ccccc}
 & & 3 & 5 & 5 & 4 & 2 \\
 & & 1 & 0.4 & -2 & 0.4 & -0.2 \\
 & & 4 & -0.2 & 0.5 & 4 & -0.5 \\
 & & 2 & 0.6 & 0 & 0.4 & -3
 \end{array}$$

$$\begin{aligned}
 & \begin{pmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 & 0 & 2 & 0.6 \\ 3 & 3 & 4 & -2 \\ 5 & 5 & 4 & 2 \\ -1 & -2 & 3.4 & -1 \end{pmatrix} = \\
 & = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.4 & 1 & 0 & 0 \\ -0.2 & 0.5 & 1 & 0 \\ 0.6 & 0 & 0.4 & 1 \end{pmatrix} \begin{pmatrix} 5 & 5 & 4 & 2 \\ 0 & -2 & 0.4 & -0.2 \\ 0 & 0 & 4 & -0.5 \\ 0 & 0 & 0 & -3 \end{pmatrix} \quad \diamond
 \end{aligned}$$

2.4.3. Factorizarea Cholesky

Matricele hermitiene pozitiv definite pot fi descompuse în factori triunghiulari de două ori mai repede decât matricele generale. Algoritmul standard pentru aceasta, factorizarea Cholesky⁴, este o variantă a eliminării gaussiene ce operează atât la stânga cât și la dreapta matricei, păstrând și exploatarea simetriei.

Sistemele cu matrice hermitiene pozitiv definite joacă un rol important în algebra liniară numerică și în aplicații. Datorită legilor fundamentale ale fizicii, multe matrice care intervin în probleme practice sunt de acest tip.

Reamintim câteva proprietăți ale matricelor hermitiene. Dacă A este o matrice $m \times m$ hermitiană pozitiv definită și X este o matrice de tip $m \times n$ de rang maxim, cu $m \geq n$, atunci matricea X^*AX este de asemenea hermitiană pozitiv definită deoarece $(X^*AX)^* = X^*A^*X = X^*AX$ și pentru orice vector $x \neq 0$ avem $Xx \neq 0$ și astfel $x^*(X^*AX)x = (Xx)^*A(Xx) > 0$. Alegând pe post de X o matrice $m \times n$ cu un 1 în fiecare coloană și zero în rest, putem scrie orice submatrice principală $n \times n$ a lui A

⁴Andre-Louis Cholesky (1875-1918) ofițer francez, specialist în topografie și geodezie, a activat în Creta și Africa de nord înainte începerii primului război mondial. A dezvoltat metoda care îi poartă numele și a aplicat-o la calculul soluțiilor ecuațiilor normale pentru probleme de aproximare în sensul celor mai mici pătrate care apar în geodezie. Lucrarea sa a fost publicată postum în 1924, de către camaradul său Benoît, în Bulletin Géodésique. Se pare că a luat parte la misiunea militară franceză din România în timpul primului război mondial.

sub forma X^*AX . De aceea, orice submatrice principală a lui A trebuie să fie pozitiv definită. În particular, orice element diagonal al lui A este un număr real pozitiv.

Valorile proprii ale unei matrice hermitiene pozitiv definite sunt de asemenea numere reale pozitive, Dacă $Ax = \lambda x$ pentru $x \neq 0$, avem $x^*Ax = \lambda x^*x > 0$ și deci $\lambda > 0$ și reciproc, dacă toate valorile proprii sunt definite, atunci A este pozitiv definită. Vectorii proprii ce corespund valorilor proprii distincte ale unei matrice hermitiene sunt ortogonali. Presupunem că $Ax_1 = \lambda_1 x_1$ și $Ax_2 = \lambda_2 x_2$ cu $\lambda_1 \neq \lambda_2$. Atunci

$$\lambda_2 x_1^* x_2 = x_1^* A x_2 = \overline{x_2^* A x_1} = \overline{\lambda_1 x_2^* x_1} = \lambda_1 x_1^* x_2,$$

așă că $(\lambda_1 - \lambda_2)x_1^* x_2 = 0$. Deoarece $\lambda_1 \neq \lambda_2$, rezultă că $x_1^* x_2 = 0$. Matricele hermitiene sunt de asemenea normale ($AA^* = A^*A = A^2$).

O factorizare Cholesky a unei matrice A este o descompunere de forma

$$A = R^*R, \quad r_{jj} > 0, \quad (2.4.1)$$

unde R este o matrice triunghiulară superior.

Teorema 2.4.5 Orice matrice hermitiană pozitiv definită $A \in \mathbb{C}^{m \times m}$ are o factorizare Cholesky (2.4.1) unică.

Demonstrație. (Existența) Deoarece A este hermitiană și pozitiv definită $a_{11} > 0$ și putem pune $\alpha = \sqrt{a_{11}}$. Observăm că

$$\begin{aligned} A &= \begin{bmatrix} a_{11} & w^* \\ w & K \end{bmatrix} \\ &= \begin{bmatrix} \alpha & 0 \\ w/\alpha & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & K - ww^*/a_{11} \end{bmatrix} \begin{bmatrix} \alpha & w^*/\alpha \\ 0 & I \end{bmatrix} = R_1^* A_1 R_1. \end{aligned} \quad (2.4.2)$$

Acesta este pasul de bază care este repetat în factorizarea Cholesky. Submatricea $K - ww^*/a_{11}$ fiind o submatrice principală de tip $(m-1) \times (m-1)$ a matricei pozitiv definite $R_1^* A R_1^{-1}$ este pozitiv definită și deci elementul ei situat în colțul din stânga sus este pozitiv. Se arată prin inducție că toate submatricele care apar în cursul factorizării sunt pozitiv definite și procesul nu poate eșua. Continuăm cu factorizarea lui $A_1 = R_2^* A_2 R_2$ și astfel $A = R_1^* R_2^* A_2 R_2 R_1$; procesul poate continua până se ajunge la colțul din dreapta jos, obținându-se

$$A = \underbrace{R_1^* R_2^* \dots R_m^*}_{R^*} \underbrace{R_m \dots R_2 R_1}_R,$$

care are chiar forma dorită.

(Unicitatea) De fapt procedeul de mai sus stabilește și unicitatea. La fiecare pas (2.4.2), valoare $\alpha = \sqrt{a_{11}}$ este determinată din forma factorizării R^*R și odată ce α este determinat, prima linie a lui R_1^* este de asemenea determinată. Deoarece cantitățile analoge sunt determinate la fiecare pas al reducerii, întreaga factorizare este unică. \square

Când se implementează factorizarea Cholesky, este nevoie să se reprezinte explicit doar jumătate din matricea asupra căreia se operează. Această simplificare permite evitarea a jumătate din operațiile aritmetice. Se dă mai jos una din multele posibilități de prezentare formală a algoritmului (algoritmul 2.2). Matricea A de la intrare conține diagonala principală și jumătatea de deasupra diagonalei principale a matricei hermitiene și pozitiv definite de tip $m \times m$ ce urmează a fi factorizată. În implementări practice se pot utiliza scheme de memorare comprimată ce evită irosirea spațiului pentru jumătate din matricea pătratică. Matricea de ieșire reprezintă factorul triunghiular superior din factorizarea $A = R^*R$. Fiecare iterație externă corespunde unei singure factorizări elementare: partea triunghiulară superior a submatricei $R_{k:m,k:m}^*$ reprezintă partea supra-diagonală a matricei hermitiene ce trebuie factorizată la pasul k .

Algoritmul 2.2 Factorizare Cholesky

```

 $R := A;$ 
for  $k := 1$  to  $m$  do
  for  $j := k + 1$  to  $m$  do
     $R_{j,j:m} := R_{j,j:m} - R_{k,j:m}R_{k,j}/R_{k,k}$ 
  end for
   $R_{k,k:m} := R_{k,k:m} / \sqrt{R_{k,k}}$ 
end for

```

Operațiile aritmetice în factorizarea Cholesky (algoritmul 2.2) sunt dominate de ciclul interior. O singură execuție a liniei

$$R_{j,j:m} := R_{j,j:m} - R_{k,j:m}R_{k,j}/R_{k,k}$$

necesită o împărțire, $m - j + 1$ înmulțiri și $m - j + 1$ scăderi, deci un total de $\sim 2(m - j)$ flops. Acest calcul este repetat pentru fiecare j de la $k + 1$ la m și acest ciclu este repetat pentru fiecare k de la 1 la m . Suma se evaluează direct

$$\sum_{k=1}^m \sum_{j=k+1}^m 2(m - j) \sim 2 \sum_{k=1}^m \sum_{j=1}^k j \sim \sum_{k=1}^m k^2 \sim \frac{1}{3}m^3 \text{ flops,}$$

deci jumătate din volumul de calcule necesar pentru eliminarea gaussiană.

2.4.4. Descompunerea QR

Teorema 2.4.6 Fie $A \in \mathbb{R}^{m \times n}$, cu $m \geq n$. Atunci există o matrice ortogonală unică Q de tip $m \times n$ și o matrice triunghiulară superior unică, de tip $n \times n$ R cu diagonala pozitivă ($r_{ii} > 0$) astfel încât $A = QR$.

Demonstrație. Va rezulta din algoritmul 2.3, care va fi dat în această secțiune. \square

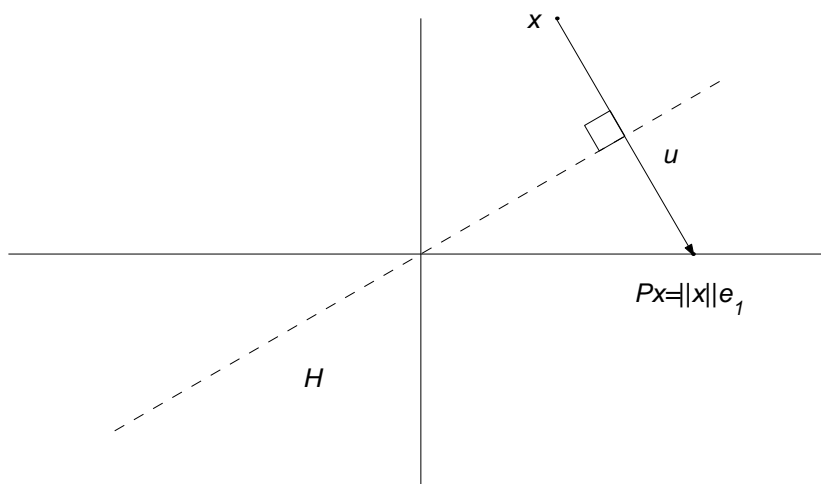


Figura 2.2: Un reflector Householder

Transformări Householder

O transformare Householder⁵ (sau reflexie) este o matrice de forma $P = I - 2uu^T$, unde $\|u\|_2 = 1$. Se verifică ușor că $P = P^T$ și că $PP^T = (I - 2uu^T)(I - 2uu^T) = I - 4uu^T + 4uu^Tuu^T = I$, deci P este o matrice simetrică și ortogonală. Ea se numește reflexie deoarece Px este reflexia lui x față de hiperplanul H ce trece prin origine și este ortogonal pe u (figura 2.2).

Dându-se un vector x , este ușor de găsit reflexia Householder care anulează toate componentele lui x , exceptând prima: $Px = [c, 0, \dots, 0]^T = ce_1$. Vom face aceasta după cum urmează. Scriem $Px = x - 2u(u^T x) = ce_1$, deci $u = \frac{1}{2(u^T x)}(x - ce_1)$, adică u este o combinație liniară a lui x și e_1 . Deoarece $\|x\|_2 = \|Px\|_2 = |c|$, u trebuie să fie paralel cu vectorul $\tilde{u} = x \pm \|x\|_2 e_1$, deci $u = \tilde{u} / \|\tilde{u}\|_2$. Se poate verifica că orice alegere de semn ne conduce la un u ce satisface $Px = ce_1$, atât timp cât $\tilde{u} \neq 0$. Vom utiliza $\tilde{u} = x + \text{sign}(x_1)\|x\|_2 e_1$, deoarece astfel nu va apare nici o anulare flotantă la calculul

Alston S. Householder (1904-1993), matematician american. Contribuții importante în biologia matematică și mai ales în ⁵algebra liniară numerică. Cartea sa cea mai importantă, „The theory of matrices in numerical analysis” a avut un impact deosebit asupra dezvoltării analizei numerice și informaticii.



componentelor lui \tilde{u} . Dacă $x_1 = 0$, vom conveni să luăm $\text{sign}(x_1) = 1$. Rezumând, vom lua

$$\tilde{u} = \begin{bmatrix} x_1 + \text{sign}(x_1)\|x\|_2 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \text{ cu } u = \frac{\tilde{u}}{\|\tilde{u}\|_2}.$$

Vom scrie aceasta sub forma $u = \text{House}(x)$. În practică, putem memora \tilde{u} în locul lui u pentru a reduce efortul de calcul al lui u și să utilizăm formula $P = I - \frac{2}{\|\tilde{u}\|_2^2} \tilde{u}\tilde{u}^T$ în loc de $P = I - 2uu^T$.

Exemplul 2.4.7. Vom arăta cum se poate calcula descompunerea QR a unei matrice 5×4 utilizând transformări Householder. Acest exemplu poate face procedeul mai ușor de înțeles în cazul general. În calculele de mai jos P_i sunt matrice ortogonale 5×5 , x este un element generic nenul, iar o o poziție nulă.

1. Alegem P_1 astfel încât $A_1 \equiv P_1 A = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & x & x & x \\ o & x & x & x \\ o & x & x & x \end{bmatrix}$.
2. Alegem $P_2 = \begin{bmatrix} I_1 & 0 \\ 0 & P'_2 \end{bmatrix}$ astfel încât $A_2 \equiv P_2 A_1 = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & x & x \\ o & o & x & x \end{bmatrix}$.
3. Alegem $P_3 = \begin{bmatrix} I_2 & 0 \\ 0 & P'_3 \end{bmatrix}$ astfel încât $A_3 \equiv P_3 A_2 = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & o & x \\ o & o & o & x \end{bmatrix}$.
4. Alegem $P_4 = \begin{bmatrix} I_3 & 0 \\ 0 & P'_4 \end{bmatrix}$ astfel încât $A_4 \equiv P_4 A_3 = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & o & x \\ o & o & o & o \end{bmatrix}$.

Aici am ales matricele Householder P'_i pentru a anula elementele subdiagonale din coloana i ; aceasta nu distruge zerourile deja introduse în coloanele precedente. Să notăm matricea finală 5×4 cu $\tilde{R} = A_4$. Atunci $A = P_1^T P_2^T P_3^T P_4^T \tilde{R} = QR$, unde Q este

formată din primele patru coloane a lui $P_1^T P_2^T P_3^T P_4^T = P_1 P_2 P_3 P_4$ (deoarece P_i sunt simetrice iar R este formată din primele patru linii ale lui R). \diamond

Algoritmul 2.3 descrie procesul de calcul pentru descompunerea QR bazată pe transformări HouseHolder.

Algoritmul 2.3 Factorizare QR utilizând reflexii HouseHolder

```

1: for  $i := 1$  to  $\min(m - 1, n)$  do
2:    $u_i := \text{House}(A_{i:m,i});$ 
3:    $P'_i := I - 2u_i u_i^T;$ 
4:    $A_{i:m,i:n} := P'_i A_{i:m,i:n};$ 
5: end for

```

Algoritmul 2.4 Calculul produsului $Q^T b$

```

1: for  $i := 1$  to  $n$  do
2:    $\gamma := -2u_i^T b_{i:m};$ 
3:    $b_{i:m} := b_{i:m} + \gamma u_i;$ 
4: end for

```

Algoritmul 2.5 Calculul produsului Qx

```

1: for  $k := n$  downto  $1$  do
2:    $x_{k:m} := x_{k:m} - 2u_k(u_k^* x_{k:m});$ 
3: end for

```

Iată câteva detalii de implementare. Nu avem niciodată nevoie să construim explicit P_i ci doar să efectuăm înmulțirea

$$(I - 2u_i u_i^T) A_{i:m,i:n} = A_{i:m,i:n} - 2u_i (u_i^T A_{i:m,i:n}),$$

care este mai puțin costisitoare. Pentru a memora P_i avem nevoie doar de u_i sau de \tilde{u}_i și $\|\tilde{u}_i\|_2$. Acestea pot fi memorate în coloana i al lui A ; de fapt nu este nevoie ca ea să fie schimbată! Astfel descompunerea QR poate fi scrisă peste A , unde Q este memorată în forma factorizată $P_1 \dots P_{n-1}$, iar P_i este memorată prin vectorul de reflexie \tilde{u}_i , în coloana i a lui A , sub diagonală. (Avem nevoie de un tablou suplimentar este memorată în forma factorizată $P_1 \dots P_{n-1}$ iar P_i este memorată ca \tilde{u}_i în partea subdiagonală a coloanei i a lui A . Mai este nevoie de un tablou suplimentar de lungime n pentru elementul de sus al lui \tilde{u}_i , deoarece diagonală este deja ocupată de elementele R_{ii} .)

Pornind de la observația

$$Ax = b \Leftrightarrow QRx = b \Leftrightarrow Rx = Q^T b,$$

putem alege următoarea strategie pentru rezolvarea sistemului de ecuații liniare $Ax = b$:

Exemplul 2.4.8. Vom ilustra doi pași intermediari din calculul descompunerii QR a unei matrice 5×4 utilizând rotații Givens. Pentru a trece de la

$$\begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & x & x \\ o & o & x & x \end{bmatrix} \quad \text{la} \quad \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & o & x \\ o & o & o & x \end{bmatrix}$$

efectuăm înmulțirile

$$\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & c & s \\ & & & -s & c \end{bmatrix} \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & x & x \\ o & o & x & x \end{bmatrix} = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & x & x \\ o & o & o & x \end{bmatrix}$$

și

$$\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & c' & s' & \\ & & -s' & c' & \\ & & & & 1 \end{bmatrix} \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & x & x \\ o & o & o & x \end{bmatrix} = \begin{bmatrix} x & x & x & x \\ o & x & x & x \\ o & o & x & x \\ o & o & o & x \\ o & o & o & x \end{bmatrix} .$$

◇

Costul descompunerii QR cu rotații Givens este de două ori costul descompunerii cu reflexii Householder. Ele sunt utilizate în alte aplicații (de exemplu valori proprii).

Iată câteva detalii de implementare. Calculul sinusului și cosinusului din matricea Givens este dat în algoritmul 2.6. El necesită 5 flops și un radical. Nu necesită funcții tri-

Algoritmul 2.6 Dându-se scalarii a și b , calculează elementele $c = \cos(\theta)$ și $s = \sin(\theta)$ ale unei matrice Givens

```

function [c, s] =givens(a, b)
if b = 0 then
    c := 1;    s := 0;
else
    if |b| > |a| then
        τ := -a/b; s := 1/√(1 + τ²); c := sτ;
    else
        τ := -b/a; c := 1/√(1 + τ²); s := cτ;
    end if
end if

```

gonometrice inverse și nu dă depășire superioară. Descompunerea QR bazată pe rotații

Algoritmul 2.7 Factorizare $A = QR$ cu ajutorul matricei Givens; rezultatul R se scrie peste A .

Intrare: $A \in \mathbb{R}^{m \times n}$

Ieșire: $R = Q^T A$, unde R este triunghiulară superior.

```

for  $j := 1$  to  $n$  do
  for  $i := m - 1$  downto  $j + 1$  do
     $[c, s] = \text{givens}(A_{i-1,j}, A_{i,j});$ 
     $A_{i-1:i,j:n} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T A_{i-1:i,j:n};$ 
  end for
end for

```

Givens este descrisă de algoritmul 2.7 Acest algoritm necesită $3n^2(m - n/3)$ flops. De notat că se pot utiliza și alte secvențe de rotații; de exemplu primele două **for**-uri din algoritmul 2.7 se pot înlocui cu

```

for  $i := m$  downto  $2$ 
  for  $j := 1$  to  $\min\{i - 1, n\}$ 

```

Astfel, zerourile vor fi introduse linie cu linie. O altă schimbare posibilă se referă la planele de rotație: în loc să rotim liniile $i - 1$ și i ca în algoritmul 2.7, am putea roti liniile j și i :

```

for  $j := 1$  to  $n$  do
  for  $i := m - 1$  downto  $j + 1$  do
     $[c, s] = \text{givens}(A_{j,j}, A_{i,j});$ 
     $A_{[j\ i],j:n} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T A_{[j\ i],j:n};$ 
  end for
end for

```

Deoarece o rotație Givens anulează exact un element, trebuie să memorăm informații despre rotație pe poziția acelu element. Vom face asta după cum urmează. Fie $s = \sin \theta$ și $c = \cos \theta$. Dacă $|s| < |c|$ se memorează $s \cdot \text{sign}(c)$; altfel se memorează $\frac{\text{sign}(s)}{c}$. Pentru a recupera s și c din valoarea memorată (să o numim p) vom proceda astfel: dacă $|p| < 1$, atunci $s := p$ și $c = \sqrt{1 - s^2}$; altfel $c := 1/p$ și $s := \sqrt{1 - c^2}$. Motivul pentru care nu memorăm s și nu calculăm $c = \sqrt{1 - s^2}$ este acela că dacă s este apropiat de 1, c va fi imprecis (datorită anulării flotante). Se poate recupera fie s și c , fie $-s$ și $-c$; acest lucru este convenabil în practică.

Secvența de rotații Givens se poate aplica în așa fel ca să se aplice mai puține flops decât în variantele prezentate aici. Se ajunge astfel la rotații Givens rapide (vezi [24, capitolul 5]).

2.5. Rafinarea iterativă

Dacă metoda de rezolvare pentru $Ax = b$ este nestabilă, atunci $A\bar{x} \neq b$, unde \bar{x} este valoarea calculată. Vom calcula corecția Δx astfel încât

$$A(\bar{x} + \Delta x_1) = b \Rightarrow A\Delta x_1 = b - A\bar{x}$$

Se rezolvă sistemul și se obține un nou \bar{x} , $\bar{x} := x + \Delta x_1$. Dacă din nou $Ax \neq b$, se calculează o nouă corecție până când

$$\|\Delta x_i - \Delta x_{i-1}\| < \varepsilon \quad \text{sau} \quad \|Ax - b\| < \varepsilon.$$

Calculul cantității $r = b - Ax$, numită *reziduu*, se va efectua în dublă precizie.

2.6. Algoritmul lui Strassen pentru înmulțirea matricelor

Fie $A, B \in \mathbb{K}^{n \times n}$. Dorim să calculăm $C = AB$. Presupunem că $n = 2^k$. Partiționăm A și B

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \quad B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \quad C = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix}$$

În algoritmul clasic avem 8 înmulțiri și 4 adunări pentru un pas, iar timpul de execuție $T(n) = \Theta(n^3)$, deoarece $T(n) = 8T(n/2) + \Theta(n^2)$.

Interesul este de a reduce numărul de înmulțiri. Volker Strassen a descoperit o metodă de a reduce numărul de înmulțiri pe pas la 7. Se calculează următoarele cantități:

$$\begin{aligned} p_1 &= (a_{11} + a_{22})(b_{11} + b_{22}), \\ p_2 &= (a_{21} + a_{22})b_{11}, \\ p_3 &= a_{11}(b_{12} - b_{22}), \\ p_4 &= a_{22}(b_{21} - b_{11}), \\ p_5 &= (a_{11} + a_{12})b_{22}, \\ p_6 &= (a_{21} - a_{11})(b_{11} + b_{12}), \\ p_7 &= (a_{12} - a_{22})(b_{21} + b_{22}), \end{aligned}$$

$$\begin{aligned} c_{11} &= p_1 + p_4 - p_5 + p_7, \\ c_{12} &= p_3 + p_5, \\ c_{21} &= p_2 + p_4, \\ c_{22} &= p_1 + p_3 - p_2 + p_6. \end{aligned}$$

Procedura este descrisă detaliat de algoritmul 2.8. Deoarece avem 7 înmulțiri și 18

Algoritmul 2.8 Algoritmul lui Strassen pentru înmulțirea a două matrice

Presupunem că $m = 2^q$ și că $A, B \in \mathbb{R}^{m \times m}$. Dacă $m_{\min} = 2^d$ cu $d \leq q$, atunci acest algoritm calculează $C = AB$ aplicând recursiv procedura lui Strassen de $q - d$ ori.

function $C = \mathbf{strass}(A, B, m, m_{\min})$

if $m \leq m_{\min}$ **then**

$C := AB;$

else

$n := m/2; u := 1 : n; v := n + 1 : m;$

$P_1 := \mathbf{strass}(A(u, u) + A(v, v), B(u, u) + B(v, v), n, m_{\min});$

$P_2 := \mathbf{strass}(A(v, u) + A(v, v), B(u, u), n, m_{\min});$

$P_3 := \mathbf{strass}(A(u, u), B(u, v) - B(v, v), n, m_{\min});$

$P_4 := \mathbf{strass}(A(v, v), B(v, u) - B(u, u), n, m_{\min});$

$P_5 := \mathbf{strass}(A(u, u) + A(u, v), B(v, v), n, m_{\min});$

$P_6 := \mathbf{strass}(A(v, u) - A(u, u), B(u, u) + B(u, v), n, m_{\min});$

$P_7 := \mathbf{strass}(A(u, v) - A(v, v), B(v, u) + B(v, v), n, m_{\min});$

$C(u, u) := P_1 + P_4 - P_5 + P_7;$

$C(u, v) := P_3 + P_5;$

$C(v, u) := P_2 + P_4;$

$C(v, v) := P_1 + P_3 - P_2 + P_6;$

end if

adunări sau scăderi pe pas, timpul de execuție verifică relația de recurență

$$T(n) = 7T(n/2) + \Theta(n^2),$$

cu soluția

$$T(n) = \Theta(n^{\log_2 7}) \sim 28n^{\log_2 7}.$$

Algoritmul se poate generaliza pentru matrice cu dimensiunea $n = m \cdot 2^k$.

Dacă n este impar ultima coloană a lui se calculează prin metode standard și se execută algoritmul lui Strassen pentru matrice de dimensiune $n - 1$.

$$m \cdot 2^{k+1} \rightarrow m \cdot 2^k$$

p -urile se pot calcula în paralel; la fel și c -urile.

Accelerarea teoretică obținută pentru înmulțirea matricelor se traduce printr-o accelerare a inversării matricelor, deci și a rezolvării sistemelor de ecuații liniare. Dacă notăm cu $M(n)$ timpul de înmulțire a două matrice pătratice de ordinul n și cu $I(n)$ timpul de inversare a unei matrice de același ordin, atunci $M(n) = \Theta(n)$. Vom demonstra aceasta în două etape: arătăm că $M(n) = O(I(n))$ și apoi că $I(n) = O(M(n))$.

Teorema 2.6.1 (Înmulțirea nu este mai grea decât inversarea) *Dacă putem inversa o matrice $n \times n$ în timp $I(n)$, unde $I(n) = \Omega(n^2)$ satisface condiția de regularitate $I(3n) = O(I(n))$, atunci putem înmulți două matrice de ordinul n în timp $O(I(n))$.*

Demonstrație. Fie A și B două matrice $n \times n$. Vrem să calculăm $C = AB$. Definim matricea D de ordinul $3n \times 3n$ astfel:

$$D = \begin{pmatrix} I_n & A & 0 \\ 0 & I_n & B \\ 0 & 0 & I_n \end{pmatrix}.$$

Inversa lui D este

$$D^{-1} = \begin{pmatrix} I_n & -A & AB \\ 0 & I_n & -B \\ 0 & 0 & I_n \end{pmatrix},$$

și, astfel, putem calcula produsul AB luând submatricea de ordinul $n \times n$ din colțul din dreapta sus al matricei D^{-1} .

Putem calcula matricea D într-un timp de ordinul $\Theta(n^2) = O(I(n))$ și conform condiției de regularitate o putem inversa într-un timp $O(I(3n)) = O(I(n))$. Deci

$$M(n) = O(I(n)).$$

□

Să observăm că $I(n)$ satisface condiția de regularitate doar dacă $I(n)$ nu are salturi mari în valoare. De exemplu, dacă $I(n) = \Theta(n^c \log^d n)$, pentru orice constante $c > 0$, $d \geq 0$, atunci $I(n)$ satisface condiția de regularitate.

Demonstrația că inversarea matricelor nu este mai grea decât înmulțirea lor se bazează pe proprietățile matricelor simetrice, pozitiv definite.

Teorema 2.6.2 (Inversarea nu este mai grea decât înmulțirea) *Dacă putem înmulți două matrice reale $n \times n$ în timp $M(n)$ unde $M(n) = \Omega(n^2)$ și $M(n)$ satisface condițiile de regularitate $M(n) = O(M(n+k))$ pentru orice k , $0 \leq k \leq n$ și $M(n/2) \leq cM(n)$, pentru orice constantă $c < 1/2$, atunci putem calcula inversa unei matrice reale nesingulare de ordinul $n \times n$, în timp $O(M(n))$.*

Demonstrație. Putem presupune că n este multiplu de 2, deoarece

$$\begin{pmatrix} A & 0 \\ 0 & I_k \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & 0 \\ 0 & I_k \end{pmatrix},$$

pentru orice $k \in \mathbb{N}^*$. Așadar, alegând pe k astfel încât $n+k$ să fie o putere a lui 2, extindem matricea la o dimensiune care este puterea următoare a lui 2 și obținem

răspunsul dorit din răspunsul problemei pentru matricea extinsă în acest mod. Condițiile de regularitate ne asigură că extinderea nu cauzează creșterea timpului decât cu cel mult un factor constant.

Pentru moment, să presupunem că A este o matrice $n \times n$ simetrică și pozitiv-definită; o vom partiționa în patru submatrice de ordinul $n/2 \times n/2$:

$$A = \begin{pmatrix} B & C^T \\ C & D \end{pmatrix}. \quad (2.6.1)$$

Dacă

$$S = D - CB^{-1}C^T \quad (2.6.2)$$

este complementul Schur al lui A în raport cu B , atunci

$$A^{-1} = \begin{pmatrix} B^{-1} + B^{-1}C^T S^{-1}CB^{-1} & -B^{-1}C^T S^{-1} \\ -S^{-1}CB^{-1} & S^{-1} \end{pmatrix}, \quad (2.6.3)$$

relație care se poate verifica prin înmulțire. Matricele B^{-1} și S^{-1} există, deoarece A este simetrică și pozitiv definită, deoarece atât B cât și S sunt simetrice și pozitiv definite. În plus, $B^{-1}C^T = (CB^{-1})^T$ și $B^{-1}C^T S^{-1} = (S^{-1}CB^{-1})^T$. Putem folosi ecuațiile (2.6.2) și (2.6.3) pentru a specifica un algoritm recursiv în care intervin patru înmulțiri de matrice de ordinul $n/2 \times n/2$:

$$\begin{aligned} & C \cdot B^{-1}, \\ & (C \cdot B^{-1}) \cdot C^T \\ & S^{-1} \cdot (CB^{-1}) \\ & (CB^{-1})^T \cdot (S^{-1}CB^{-1}). \end{aligned}$$

Întrucât matricele de ordinul $n/2 \times n/2$ se înmulțesc folosind un algoritm pentru matrice de ordinul $n \times n$, inversarea matricelor simetrice pozitiv-definite se poate rezolva în timpul

$$I(n) \leq 2I(n/2) + 4M(n) + O(n^2) = 2I(n/2) + O(M(n)) = O(M(n)).$$

Rămâne de studiat cazul când A este inversabilă, dar nu este simetrică și pozitiv definită. Deoarece $A^T A$ este simetrică și pozitiv definită, problema inversării lui A se reduce la problema inversării lui $A^T A$. Reducerea se bazează pe observația că, atunci când A este o matrice nesingulară de ordinul $n \times n$, avem

$$A^{-1} = (A^T A)^{-1} A^T,$$

deoarece $((A^T A)^{-1} A^T) = (A^T A)^{-1} (A^T A) = I_n$, și inversa unei matrice este unică. Prin urmare, putem calcula A^{-1} , înmulțind întâi pe A^T cu A pentru a obține $A^T A$,

inversând apoi matricea simetrică și pozitiv definită $A^T A$ prin algoritmul divide et impera prezentat, și, în final înmulțind rezultatul cu A^T . Fiecare din acești pași necesită un timp $O(M(n))$. Astfel, orice matrice nesingulară cu elemente reale poate fi inversată într-un timp de ordinul $O(M(n))$. \square

Demonstrația teoremei 2.6.2 sugerează un mod de a rezolva ecuații de forma $Ax = b$, cu A nesingulară, fără pivotare, rezolvând ecuația echivalentă $(A^T A)x = b$ prin factorizare Cholesky. Dar în practică, descompunerea LUP funcționează mai bine, iar transformarea propusă mărește numărul de condiționare.

2.7. Rezolvarea iterativă a sistemelor algebrice liniare

Dorim să calculăm soluția sistemului liniar

$$Ax = b, \quad (2.7.1)$$

când A este inversabilă. Presupunem că am găsit o matrice T și un vector c astfel încât $I - T$ să fie inversabilă și astfel încât punctul fix unic al ecuației

$$x = Tx + c \quad (2.7.2)$$

să fie egal cu soluția sistemului $Ax = b$. Fie x^* soluția lui (2.7.1) sau, echivalent, a lui (2.7.2).

Iterațiile: se dă un $x^{(0)}$ arbitrar; se definește șirul $(x^{(k)})$ prin

$$x^{(k+1)} = Tx^{(k)} + c, \quad k \in \mathbb{N} \quad (2.7.3)$$

Lema 2.7.1 Dacă $\rho(X) < 1$, există $(I - X)^{-1}$ și

$$(I - X)^{-1} = I + X + X^2 + \dots + X^k + \dots$$

Demonstrație. Fie

$$S_k = I + X + \dots + X^k$$

$$(I - X)S_k = I - X^{k+1}$$

$$\lim_{k \rightarrow \infty} (I - X)S_k = I \Rightarrow \lim_{k \rightarrow \infty} S_k = (I - X)^{-1}$$

căci $X^{k+1} \rightarrow 0 \Leftrightarrow \rho(X) < 1$. \square

Teorema 2.7.2 Propozițiile următoare sunt echivalente

- (1) metoda (2.7.3) este convergentă
 (2) $\rho(T) < 1$
 (3) $\|T\| < 1$ pentru cel puțin o normă matricială.

Demonstrație.

$$\begin{aligned} x^{(k)} &= Tx^{(k-1)} + c = T(Tx^{(k-2)} + c) + c = \dots = \\ &= T^k x^{(0)} + (I + T + \dots + T^{k-1})c. \end{aligned}$$

(2.7.3) convergentă $\Leftrightarrow I - T$ inversabilă $\Leftrightarrow \rho(T) < 1 \Leftrightarrow \exists \|\cdot\|$ astfel încât $\|T\| < 1$ (din teorema 2.1.8). \square

Aplicând teorema de punct fix a lui Banach se obține:

Teorema 2.7.3 Dacă există $\|\cdot\|$ astfel încât $\|T\| < 1$, șirul $(x^{(k)})$ definit de (2.7.3) este convergent pentru orice $x^{(0)} \in \mathbb{R}^n$ și are loc

$$\|x^* - x^{(k)}\| \leq \frac{\|T\|^k}{1 - \|T\|} \|x^{(1)} - x^{(0)}\| \leq \frac{\|T\|}{1 - \|T\|} \|x^{(1)} - x^{(0)}\|. \quad (2.7.4)$$

O metodă iterativă de rezolvare a unui sistem algebric liniar $Ax = b$ pornește de la o aproximație inițială $x^{(0)} \in \mathbb{R}^n(\mathbb{C}^n)$ și generează un șir de vectori $\{x^{(k)}\}$ care converge către soluția x^* a sistemului. Aceste tehnici transformă sistemul inițial într-un sistem echivalent de forma $x = Tx + c$, $T \in \mathbb{K}^{n \times n}$, $c \in \mathbb{K}^n$. Se generează un șir de forma $x^{(k)} = Tx^{(k-1)} + c$. Criteriul de oprire este

$$\|x^{(k)} - x^{(k-1)}\| \leq \frac{1 - \|T\|}{\|T\|} \varepsilon. \quad (2.7.5)$$

El are la bază rezultatul următor:

Propoziția 2.7.4 Dacă x^* este soluția sistemului (2.7.2), cu $\|T\| < 1$, atunci

$$\|x^* - x^{(k)}\| \leq \frac{\|T\|}{1 - \|T\|} \|x^{(k)} - x^{(k-1)}\|. \quad (2.7.6)$$

Demonstrație. Fie $p \in \mathbb{N}^*$. Avem

$$\|x^{(k+p)} - x^{(k)}\| \leq \|x^{(k+1)} - x^{(k)}\| + \dots + \|x^{(k+p)} - x^{(k+p-1)}\|. \quad (2.7.7)$$

Pe de altă parte, din (2.7.3), rezultă că

$$\|x^{(m+1)} - x^{(m)}\| \leq \|T\| \|x^{(m)} - x^{(m-1)}\|$$

sau, pentru $k < m$

$$\|x^{(m+1)} - x^{(m)}\| \leq \|T\|^{m-k+1} \|x^{(k)} - x^{(k-1)}\|.$$

Aplicând aceste inegalități, pentru $m = \overline{k, k+p-1}$ relația (2.7.6) devine

$$\begin{aligned} \|x^{(k+p)} - x^{(k)}\| &\leq (\|T\| + \dots + \|T\|^p) \|x^{(k)} - x^{(k-1)}\| \\ &\leq (\|T\| + \dots + \|T\|^p + \dots) \|x^{(k)} - x^{(k-1)}\|. \end{aligned}$$

Deoarece $\|T\| < 1$, avem

$$\|x^{(k+p)} - x^{(k)}\| \leq \frac{\|T\|}{1 - \|T\|} \|x^{(k)} - x^{(k-1)}\|,$$

din care prin trecere la limită în raport cu p se obține (2.7.6). \square

Dacă $\|T\| \leq 1/2$, inegalitatea (2.7.6) devine

$$\|x^* - x^{(k)}\| \leq \|x^k - x^{(k-1)}\|,$$

iar criteriul de oprire

$$\|x^k - x^{(k-1)}\| \leq \varepsilon.$$

Tehnicile iterative sunt rar utilizate pentru a rezolva sisteme de mici dimensiuni, deoarece timpul necesar pentru a obține precizia dorită depășește timpul necesar pentru eliminarea gaussiană. Pentru sisteme rare (a căror matrice are multe zerouri), de mari dimensiuni, metodele iterative sunt eficiente și din punct de vedere al spațiului și din punct de vedere al timpului.

Fie sistemul $Ax = b$. Presupunem că putem scrie matricea inversabilă A sub forma $A = M - N$. Dacă M este ușor de inversat (diagonală, triunghiulară, etc.) este mai convenabil să procedăm astfel:

$$Ax = b \Leftrightarrow Mx = Nx + b \Leftrightarrow x = M^{-1}Nx + M^{-1}b.$$

Ultima ecuație are forma $x = Tx + c$, cu $T = M^{-1}N = I - M^{-1}A$. Se obține șirul

$$x^{(k+1)} = M^{-1}Nx^{(k)} + M^{-1}b, \quad k \in \mathbb{N},$$

$x^{(0)}$ vector arbitrar.

Prima descompunere pe care o considerăm este $A = D - L - U$, unde

$$(D)_{ij} = a_{ij}\delta_{ij}, \quad (-L)_{ij} = \begin{cases} a_{ij}, & i > j \\ 0, & \text{altfel} \end{cases}$$

$$(-U)_{ij} = \begin{cases} a_{ij}, & i < j \\ 0, & \text{altfel} \end{cases}$$

Se ia $M = D$, $N = L + U$. Se obține succesiv

$$Ax = b \Leftrightarrow Dx = (L + U)x + b \Leftrightarrow x = D^{-1}(L + U)x + D^{-1}b$$

Deci $T = T_J = D^{-1}(L + U)$, $c = c_J = D^{-1}b$. Metoda se numește *metoda lui Jacobi* ⁶ (D inversabilă, de ce?)

Altă descompunere este $A = D - L - U$, $M = D - L$, $N = U$. Se obține $T_{GS} = (D - L)^{-1}U$, $c_{GS} = (D - L)^{-1}b$ – numită *metoda Gauss-Seidel* ($D - L$ inversabilă, de ce?)

$$\text{Să examinăm iterațiile Jacobi } x_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k-1)} \right).$$

La calculul lui $x_i^{(k)}$ se utilizează componentele lui $x^{(k-1)}$ (substituția simultană). Deoarece pentru $i > 1$, $x_1^{(k)}, \dots, x_{i-1}^{(k)}$ au fost deja calculați și se presupune că sunt aproximații mai bune ale componentelor soluției decât $x_1^{(k-1)}, \dots, x_{i-1}^{(k-1)}$ pare rezonabil să calculăm $x_i^{(k)}$ utilizând valorile cele mai recente, adică

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

Se pot da condiții necesare și suficiente pentru convergența metodei lui Jacobi și a metodei Gauss-Seidel

$$\rho(T_J) < 1$$

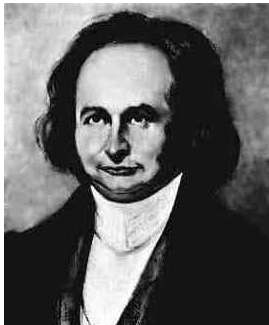
$$\rho(T_{GS}) < 1$$

și condiții suficiente: pentru o normă dată avem

$$\|T_J\| < 1$$

$$\|T_{GS}\| < 1.$$

6



Carl Gustav Jacob Jacobi (1804-1851) a fost contemporan al lui Gauss, și unul dintre cei mai importanți matematicieni germani din secolul al XIX-lea. Numele său este legat de funcții eliptice, ecuațiile cu derivate parțiale ale dinamicii, mecanică cerească. Matricea derivatelor parțiale poartă de asemenea numele său. Este inventatorul metodei iterative de rezolvare a sistemelor algebrice liniare numită metoda lui Jacobi și pe care a aplicat-o în mecanica cerească.

Putem îmbunătăți metoda Gauss-Seidel introducând un parametru ω și alegând

$$M = \frac{D}{\omega} - L.$$

Avem

$$A = \left(\frac{D}{\omega} - L \right) - \left(\frac{1-\omega}{\omega} D + U \right),$$

iar iterația obținută este

$$\left(\frac{D}{\omega} - L \right) x^{(k+1)} = \left(\frac{1-\omega}{\omega} D + U \right) x^{(k)} + b$$

Se obține matricea

$$\begin{aligned} T &= T_\omega = \left(\frac{D}{\omega} - L \right)^{-1} \left(\frac{1-\omega}{\omega} D + U \right) \\ &= (D - \omega L)^{-1} ((1-\omega)D + \omega U). \end{aligned}$$

Metoda se numește *metoda relaxării*. Avem următoarele variante:

- $\omega > 1$ suprelaxare (SOR - Successive Over Relaxation);
- $\omega < 1$ subrelaxare;
- $\omega = 1$ Gauss-Seidel.

Următoarele două teoreme se referă la convergența metodei relaxării. Scriem matricea metodei sub forma

$$T_\omega = (I - \omega \bar{L})^{-1} [(1 - \omega \bar{L})I + \omega \bar{U}],$$

unde $\bar{L} = D^{-1}L$ și $\bar{U} = D^{-1}U$.

Teorema 2.7.5 (Kahan) Are loc $\rho(T_\omega) \geq |\omega - 1|$. De aici rezultă condiția necesară $0 < \omega < 2$.

Demonstrație. Scriem polinomul caracteristic al lui T_ω sub forma $p(\lambda) = \det(\lambda I - T_\omega) = \det((I - \omega \bar{L})(\lambda I - T_\omega)) = \det((\lambda + \omega - 1)I - \omega \lambda \bar{L} - \omega \bar{U})$; deci avem

$$p(0) = \pm \prod_{i=1}^n \lambda_i(T_\omega) = \pm \det((\omega - 1)I) = \pm(\omega - 1)^n,$$

ceea ce implică $\max_i |\lambda_i(T_\omega)| \geq |\omega - 1|$. \square

Teorema 2.7.6 (Ostrowski-Reich) *Dacă A este o matrice pozitiv definită și $0 < \omega < 2$, SOR converge pentru orice alegere a aproximației inițiale $x^{(0)}$.*

Demonstrație. Demonstrația se face în doi pași. Fie $M = \omega^{-1}(D - \omega L)$. Atunci:

- (1) definim $Q = A^{-1}(2M - A)$ și arătăm ca $\operatorname{Re}\lambda_i(Q) > 0$, pentru orice i .
- (2) arătăm că $T_\omega = (Q - I)(Q + I)$, de unde rezultă că $|\lambda_i(T_\omega)| < 1$, pentru orice i .

Pentru primul pas, observăm că $Qx = \lambda x$ implică $(2M - A)x = \lambda Ax$. Adunând această ecuație cu transpusa ei conjugată se obține $x^*(M + M^* - A)x = \operatorname{Re}\lambda(x^*Ax)$. Astfel $\operatorname{Re}\lambda = x^*(M + M^* - A)x/x^*Ax = x^*(\frac{2}{\omega} - 1)Dx/x^*Ax > 0$, căci A și $(\frac{2}{\omega} - 1)D$ sunt pozitiv definite.

Pentru a demonstra (2), observăm că

$$(Q - I)(Q + I)^{-1} = (2A^{-1}M - 2I)(2A^{-1}M)^{-1} = I - M^{-1}A = T_\omega,$$

deci

$$|\lambda(T_\omega)| = \left| \frac{\lambda(Q) - 1}{\lambda(Q) + 1} \right| = \left| \frac{(\operatorname{Re}\lambda(Q) - 1)^2 + (\operatorname{Im}\lambda(Q))^2}{(\operatorname{Re}\lambda(Q) + 1)^2 + (\operatorname{Im}\lambda(Q))^2} \right|^{1/2}.$$

□

Teoremele 2.7.5 și 2.7.6 ne dau o condiție necesară și suficientă pentru convergența metodei relaxării dacă matricea sistemului este simetrică și pozitiv definită: $0 < \omega < 2$.

Observația 2.7.7. Pentru metoda Jacobi (și Gauss-Seidel) o condiție suficientă de convergență este

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad (A \text{ diagonal dominantă pe linii}).$$

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ji}|, \quad (A \text{ diagonal dominantă pe coloane}). \quad \diamond$$

Valoarea optimală pentru ω este

$$\omega_O = \frac{2}{1 + \sqrt{1 - (\rho(T_J))^2}}.$$

Ea este valabilă doar pentru anumite tipuri de matrice (de exemplu tridiagonale pe blocuri).

Teorema 2.7.8 Fie B o matrice pătratică și $\|\cdot\|$ o normă matricială oarecare. Atunci

$$\lim_{k \rightarrow \infty} \|B^k\|^{1/k} = \rho(B).$$

Demonstrație. Cum $\rho(B) \leq \|B\|$ (teorema 2.1.6 și cum $\rho(B) = (\rho(B^k))^{1/k}$ avem

$$\rho(B) \leq \|B^k\|^{1/k}, \forall k \in \mathbb{N}^*.$$

Se va stabili că, pentru orice $\varepsilon > 0$ există $\ell = \ell(\varepsilon) \in \mathbb{N}$ astfel încât

$$k \geq \ell \Rightarrow \|B^k\|^{1/k} \leq \rho(B) + \varepsilon,$$

ceea ce demonstrează relația. Fie deci $\varepsilon > 0$ dat. Deoarece matricea

$$B_\varepsilon = \frac{B}{\rho(B) + \varepsilon}$$

verifică $\rho(B_\varepsilon) < 1$ se deduce din teorema 2.7.2 că $\lim_{k \rightarrow \infty} B_\varepsilon^k = 0$. Prin urmare, există un întreg $\ell = \ell_\varepsilon$ astfel încât

$$k \geq \ell \Rightarrow \|B_\varepsilon^k\| = \frac{\|B\|^k}{(\rho(B) + \varepsilon)^k} \leq 1$$

ceea ce este chiar relația căutată. \square

Dăm o teoremă referitoare la matrice de forma $(I + B)$.

Teorema 2.7.9 (1) Fie $\|\cdot\|$ o normă matricială subordonată și B o matrice ce verifică

$$\|B\| < 1.$$

Atunci matricea $I + B$ este inversabilă și $\|I + B\|^{-1} \leq \frac{1}{1 - \|B\|}$.

(2) Dacă o matrice de forma $(I + B)$ este singulară, atunci în mod necesar

$$\|B\| \geq 1$$

pentru orice normă matricială subordonată sau nu.

Demonstrație. (1) Deoarece

$$(I + B)u = 0 \Rightarrow \|u\| = \|Bu\|$$

$$\|B\| < 1 \wedge u \neq 0 \Rightarrow \|Bu\| < \|u\|,$$

pentru norma vectorială corespondentă, se deduce că

$$(I + B)u = 0 \Rightarrow u = 0.$$

Matricea $I + B$ fiind inversabilă putem scrie

$$(I + B)^{-1} = I - B(I + B)^{-1}$$

de unde

$$\|(I + B)^{-1}\| \leq 1 + \|B\| \|(I + B)^{-1}\|,$$

ceea ce conduce la inegalitatea căutată.

(2) A spune că matricea $(I + B)$ este singulară revine la a spune că -1 este valoare proprie a lui B . În aceste condiții aplicarea teoremei 2.1.6 ne arată că $\|B\| \geq \rho(B) \geq 1$. \square

Cum se alege între mai multe metode iterative convergente pentru rezolvarea aceluiași sistem liniar $Au = b$? Pentru a fixa ideile, presupunem că B este normală. Atunci

$$\|B^k e_0\|_2 \leq \|B^k\|_2 \|e_0\|_2 = (\rho(B))^k \|e_0\|_2$$

și această inegalitate este cea mai bună posibilă, în sensul că, pentru orice întreg $k \geq 0$, există un vector $e_0(k) \neq 0$ pentru care ea devine egalitate.

În cazul matricelor normale, metoda este cu atât mai rapidă, cu cât $\rho(B)$ este mai mic, căci

$$\sup_{\|e_0\|_2=1} \|B^k e_0\|_2^{1/k} = \rho(B) \text{ pentru orice } k \geq 0.$$

În cazul general (matricea B oarecare, normă vectorială oarecare), concluzia este identică: asimptotic, vectorul de eroare $e_k = B^k e_0$ se comportă ca și $(\rho(B))^k$, cum precizează rezultatul următor.

Teorema 2.7.10 (1) Fie $\|\cdot\|$ o normă vectorială oarecare și u astfel încât

$$u = Bu + c.$$

Se consideră metoda iterativă

$$u_{k+1} = Bu_k + c, \quad k \geq 0.$$

Atunci

$$\lim_{k \rightarrow \infty} \left\{ \sup_{\|u_0 - u\|=1} \|u_k - u\|^{1/k} \right\} = \rho(B).$$

(2) Fie $\|\cdot\|$ o normă vectorială oarecare și fie u astfel încât

$$u = Bu + c = \tilde{B}u + \tilde{c}.$$

Se consideră metodele iterative

$$\tilde{u}_{k+1} = \tilde{B}\tilde{u}_k + \tilde{c}, \quad k \geq 0, \quad u_{k+1} = Bu_k + c, \quad k \geq 0$$

cu $\rho(B) < \rho(\tilde{B})$, $u_0 = \tilde{u}_0$. Atunci, oricare ar fi numărul $\varepsilon > 0$ există un întreg $\ell(\varepsilon)$ astfel încât

$$k \geq \ell \Rightarrow \sup_{\|u_0 - u\|=1} \left\{ \frac{\|\tilde{u}_k - u\|}{\|u_k - u\|} \right\}^{1/k} \geq \frac{\rho(\tilde{B})}{\rho(B) + \varepsilon}.$$

Demonstrație. Fie $\|\cdot\|$ norma matricială subordonată. Pentru orice k natural putem scrie

$$(\rho(B))^k = \rho(B^k) \leq \|B^k\| = \sup_{\|e_0\|=1} \|B^k e_0\|$$

astfel ca

$$\rho(B) \leq \sup_{\|e_0\|=1} \|B^k e_0\|^{1/k} = \|B^k\|^{1/k}$$

și aserțiunea (1) decurge din teorema 2.7.8. Conform aceleiași teoreme, fiind dat un $\varepsilon > 0$ există un întreg $\ell(\varepsilon)$ astfel încât

$$k \geq \ell \Rightarrow \sup_{\|e_0\|=1} \|B^k e_0\|^{1/k} \leq (\rho(B) + \varepsilon).$$

Prin urmare, pentru orice $k \geq \ell$, există un vector $e_0 = e_0(k)$ astfel încât

$$\|e_0\| = 1 \quad \text{și} \quad \|\tilde{B}^k e_0\|^{1/k} = \|\tilde{B}^k\|^{1/k} \geq \rho(\tilde{B})$$

și (2) este demonstrată. \square

Deci studiul metodelor iterative răspunde la două probleme.

(1) Fiind dată o metodă iterativă cu matricea B , să se determine dacă metoda este convergentă, adică dacă $\rho(B) < 1$, sau echivalent, dacă există o normă matricială astfel ca $\|B\| < 1$.

(2) Fiind date două metode iterative convergente, metoda iterativă cea mai rapidă este cea a cărei matrice are cea mai mică rază spectrală.

CAPITOLUL 3

Aproximarea funcțiilor

Cuprins

3.1. Aproximație prin metoda celor mai mici pătrate	68
3.1.1. Produse scalare	69
3.1.2. Ecuațiile normale	70
3.1.3. Eroarea în metoda celor mai mici pătrate. Convergența	73
3.2. Exemple de sisteme ortogonale	76
3.2.1. Exemple de polinoame ortogonale	79
3.3. Interpolare polinomială	85
3.3.1. Spațiul $H^n[a, b]$	85
3.3.2. Interpolare Lagrange	87
3.3.3. Interpolare Hermite	89
3.3.4. Expresia erorii de interpolare	93
3.3.5. Convergența interpolării Lagrange	98
3.4. Calculul eficient al polinoamelor de interpolare	104
3.4.1. Metode de tip Aitken	104
3.4.2. Metoda diferențelor divizate	106
3.4.3. Diferențe finite: formula lui Newton progresivă și regresivă	110
3.4.4. Diferențe divizate cu noduri multiple	112
3.5. Interpolare spline	114

3.5.1. Interpolarea cu spline cubice	117
3.5.2. Proprietăți de minimalitate ale funcțiilor spline cubice	120

Capitolul prezent este legat de aproximarea funcțiilor. Funcțiile în cauză pot să fie definite pe un continuu - de regulă un interval - sau pe o mulțime finită de puncte. Prima situație apare în contextul funcțiilor speciale (elementare sau transcendente) pe care dorim să le evaluăm ca parte a unei subrutine. Deoarece o astfel de evaluare trebuie să se reducă la un număr finit de operații aritmetice, trebuie în ultimă instanță să aproximăm funcțiile prin intermediul polinoamelor sau funcțiilor raționale. A doua situație este întâlnită în științele fizice, când măsurătorile unor cantități fizice se fac în funcție de alte cantități (cum ar fi timpul). În ambele cazuri dorim să aproximăm o funcție dată, cât mai bine posibil, în termeni de funcții mai simple.

În general o schemă de aproximare poate fi descrisă după cum urmează.

Se dă o funcție $f \in X$ ce urmează a fi aproximată, împreună cu o clasă Φ de aproximante și o normă $\|\cdot\|$ ce măsoară mărimea funcțiilor. Căutăm o aproximare $\hat{\varphi} \in \Phi$ a lui f astfel încât

$$\|f - \hat{\varphi}\| \leq \|f - \varphi\| \text{ pentru orice } \varphi \in \Phi. \quad (3.0.1)$$

Această problemă se numește *problemă de cea mai bună aproximare* a lui f cu elemente din Φ , iar funcția $\hat{\varphi}$ se numește (*element de*) *cea mai bună aproximare* a lui f relativ la norma $\|\cdot\|$.

Cunoscându-se o bază $\{\pi_j\}_{j=1}^n$ a lui Φ putem scrie

$$\Phi = \Phi_n = \left\{ \varphi : \varphi(t) = \sum_{j=1}^n c_j \pi_j(t), c_j \in \mathbb{R} \right\}. \quad (3.0.2)$$

Φ este un spațiu liniar finit dimensional sau o submulțime a acestuia.

Exemplul 3.0.11. $\Phi = \mathbb{P}_m$ - mulțimea polinoamelor de grad cel mult m . O bază a sa este $e_j(t) = t^j$, $j = 0, 1, \dots, m$. Deci $\dim \mathbb{P}_m = m + 1$. Polinoamele sunt cele mai utilizate aproximante pentru funcții pe domenii mărginite (intervale sau mulțimi finite). Motivul - teorema lui Weierstrass - orice funcție din $C[a, b]$ poate fi aproximată oricât de bine printr-un polinom de grad suficient de mare. \diamond

Exemplul 3.0.12. $\Phi = \mathbb{S}_m^k(\Delta)$ spațiul funcțiilor spline polinomiale de grad m și cu clasa de netezime k pe subdiviziunea

$$\Delta : a = t_1 < t_2 < t_3 < \dots < t_{N-1} < t_N = b$$

a intervalului $[a, b]$. Acestea sunt funcții polinomiale pe porțiuni de grad $\leq m$, racordate în t_1, \dots, t_{N-1} , astfel încât toate derivatele până la ordinul k să fie continue pe $[a, b]$. Presupunem $0 \leq k < m$. Pentru $k = m$ se obține \mathbb{P}_m . Dacă $k = -1$ permitem discontinuități în punctele de joncțiune. \diamond

Exemplul 3.0.13. $\Phi = \mathbb{T}_m[0, 2\pi]$ spațiul polinoamelor trigonometrice de grad cel mult m pe $[0, 2\pi]$. Acestea sunt combinații liniare ale funcțiilor

$$\begin{aligned}\pi_k(t) &= \cos(k-1)t & k = \overline{1, m+1}, \\ \pi_{m+1-k}(t) &= \sin kt & k = \overline{1, m}.\end{aligned}$$

Dimensiunea spațiului este $n = 2m + 1$. Astfel de aproximante sunt alegeri naturale dacă funcția de aproximat este periodică de perioadă 2π . (Dacă f are perioada p se face schimbarea de variabilă $t \rightarrow tp/2\pi$.) \diamond

De notat că mulțimea funcțiilor raționale

$$\Phi = \mathbb{R}_{r,s} = \{\varphi : \varphi = p/q, p \in \mathbb{P}_r, q \in \mathbb{P}_s\},$$

nu este spațiu liniar.

Câteva alegeri posibile ale normei, atât pentru funcții continue, cât și pentru cele discrete apar în tabelul 3.1. Cazul continuu presupune un interval $[a, b]$ și o funcție pondere $w(t)$ (posibil și $w(t) \equiv 1$) definită pe intervalul $[a, b]$ și pozitivă, exceptând zerourile izolate. Cazul discret presupune o mulțime de N puncte distincte t_1, t_2, \dots, t_N împreună cu ponderile w_1, w_2, \dots, w_N (posibil $w_i = 1, i = \overline{1, N}$). Intervalul $[a, b]$ poate fi nemărginit, dacă funcția pondere w este astfel încât integrala pe $[a, b]$ care definește norma să aibă sens.

normă continuă	tip	normă discretă
$\ u\ _\infty = \max_{a \leq t \leq b} u(t) $	L^∞	$\ u\ _\infty = \max_{1 \leq i \leq N} u(t_i) $
$\ u\ _1 = \int_a^b u(t) dt$	L^1	$\ u\ _1 = \sum_{i=1}^N u(t_i) $
$\ u\ _{1,w} = \int_a^b u(t) w(t) dt$	L_w^1	$\ u\ _{1,w} = \sum_{i=1}^n w_i u(t_i) $
$\ u\ _{2,w} = \left(\int_a^b u(t) ^2 w(t) dt \right)^{1/2}$	L_w^2	$\ u\ _{2,w} = \left(\sum_{i=1}^N w_i u(t_i) ^2 \right)^{1/2}$

Tabela 3.1: Tipuri de aproximare și normele asociate

Deci combinând normele din tabelă cu spațiile liniare din exemple se obține o problemă de cea mai bună aproximare (3.0.1) cu sens. În cazul continuu, funcția dată f și funcțiile φ din clasa Φ trebuie definite pe $[a, b]$ și norma $\|f - \varphi\|$ să aibă sens. La fel, f și φ trebuie definite în punctele t_i în cazul discret.

De notat că dacă cea mai bună aproximantă $\widehat{\varphi}$ în cazul discret este astfel încât $\|f - \widehat{\varphi}\| = 0$, atunci $\widehat{\varphi}(t_i) = f(t_i)$, pentru $i = 1, 2, \dots, N$. Spunem că $\widehat{\varphi}$ interpolatează f în punctele t_i și numim această problemă de aproximare *problemă de interpolare*.

Cele mai simple probleme de aproximare sunt problema celor mai mici pătrate (vezi secțiunea 3.1) și problema de interpolare, iar spațiul cel mai simplu este cel al polinoamelor.

Înainte de a începe cu problema celor mai mici pătrate, introducem un instrument notațional care ne permite să tratăm cazul continuu și cel discret simultan. Definim în cazul continuu

$$\lambda(t) = \begin{cases} 0, & \text{dacă } t < a \text{ (când } -\infty < a), \\ \int_a^t w(\tau)d\tau, & \text{dacă } a \leq t \leq b, \\ \int_a^b w(\tau)d\tau, & \text{dacă } t > b \text{ (când } b < \infty). \end{cases} \quad (3.0.3)$$

Astfel putem scrie, pentru orice funcție continuă u

$$\int_{\mathbb{R}} u(t)d\lambda(t) = \int_a^b u(t)w(t)dt, \quad (3.0.4)$$

căci $d\lambda(t) \equiv 0$ în afara lui $[a, b]$ și $d\lambda(t) = w(t)dt$ în interiorul lui. Vom numi $d\lambda$ *măsură* (pozitivă) *continuă*. *Măsura discretă* (numită și „măsura Dirac“) asociată mulțimii de puncte $\{t_1, t_2, \dots, t_N\}$ este o măsură $d\lambda$ care este nenulă numai în punctele t_i și are aici valoarea w_i . Astfel în acest caz

$$\int_{\mathbb{R}} u(t)d\lambda(t) = \sum_{i=1}^N w_i u(t_i). \quad (3.0.5)$$

O definiție mai precisă se poate da cu ajutorul integralei Stieltjes, dacă definim $\lambda(t)$ ca fiind o funcție în scară cu saltul în t_i egal cu w_i . În particular, definim norma lui L^2 prin

$$\|u\|_{2,d\lambda} = \left(\int_{\mathbb{R}} |u(t)|^2 d\lambda(t) \right)^{\frac{1}{2}} \quad (3.0.6)$$

și obținem norma continuă sau discretă după cum λ este ca în (3.0.3) sau o funcție în scară ca în (3.0.5).

Vom numi *supportul* lui $d\lambda$ – notat cu $\text{suppd}\lambda$ – intervalul $[a, b]$ în cazul continuu (presupunem că w este pozitivă pe $[a, b]$ exceptând zerourile izolate) și mulțimea $\{t_1, t_2, \dots, t_N\}$ în cazul discret. Spunem că mulțimea de funcții π_j din (3.0.2) este liniar independentă pe $\text{suppd}\lambda$ dacă

$$\forall t \in \text{suppd}\lambda \quad \sum_{j=1}^n c_j \pi_j(t) \equiv 0 \Rightarrow c_1 = c_2 = \dots = c_k = 0. \quad (3.0.7)$$

3.1. Aproximație prin metoda celor mai mici pătrate

Vom particulariza problema (3.0.1) luând ca normă norma din L^2

$$\|u\|_{2,d\lambda} = \left(\int_{\mathbb{R}} |u(t)|^2 d\lambda(t) \right)^{\frac{1}{2}}, \quad (3.1.1)$$

unde $d\lambda$ este fie o măsură continuă (conform (3.0.3)) sau discretă (conform (3.0.5)) și utilizând aproximanta φ dintr-un spațiu linear n -dimensional

$$\Phi = \Phi_n = \left\{ \varphi : \varphi(t) = \sum_{j=1}^n c_j \pi_j(t), c_j \in \mathbb{R} \right\}. \quad (3.1.2)$$

Presupunem că funcțiile π_j sunt linear independente pe $\text{supp}d\lambda$ și că integrala din (3.1.1) are sens pentru $u = \pi_j$ sau $u = f$.

Problema astfel obținută se numește *problemă de aproximare în sensul celor mai mici pătrate* sau *problemă de aproximare în medie pătratică*. Soluția ei a fost dată la începutul secolului al XIX-lea de către Gauss și Legendre ¹.

Presupunem că funcțiile de bază π_j sunt linear independente pe $\text{supp}d\lambda$. Vom presupune că integrala din (3.1.1) are sens pentru $u = \pi_j, j = 1, \dots, n$ și $u = f$.

Soluția problemei de cea mai bună aproximare se exprimă mai ușor cu ajutorul produselor scalare.

3.1.1. Produse scalare

Dându-se o măsură continuă sau discretă $d\lambda$ și două funcții u și v având norma finită putem defini produsul scalar

$$(u, v) = \int_{\mathbb{R}} u(t)v(t)d\lambda(t). \quad (3.1.3)$$

Inegalitatea Cauchy-Buniakovski-Schwarz

$$\|(u, v)\| \leq \|u\|_{2,d\lambda} \|v\|_{2,d\lambda}$$

ne spune că integrala în (3.1.3) este bine definită.

Produsul scalar real are următoarele proprietăți utile:

- (i) simetria $(u, v) = (v, u)$;
- (ii) omogenitatea $(\alpha u, v) = \alpha(u, v), \alpha \in \mathbb{R}$;

Adrien Marie Legendre (1752-1833) matematician francez, bine cunoscut pentru tratatul său asupra integralelor eliptice, dar și pentru lucrările sale de teoria numerelor și geometrie. Este considerat, alături de Gauss, inițiator (în 1805) al metodei celor mai mici pătrate, deși Gauss a utilizat metoda încă din 1794, dar a publicat-o doar în 1809.



(iii) aditivitatea $(u + v, w) = (u, w) + (v, w)$;

(iv) pozitiv definirea $(u, u) \geq 0$ și $(u, u) = 0 \Leftrightarrow u \equiv 0$ pe $\text{supp}d\lambda$.

Omogenitatea și aditivitatea ne dau liniaritatea

$$(\alpha_1 u_1 + \alpha_2 u_2, v) = \alpha_1 (u_1, v) + \alpha_2 (u_2, v) \quad (3.1.4)$$

Relația (3.1.4) se poate extinde la combinații liniare finite. De asemenea

$$\|u\|_{2, d\lambda}^2 = (u, u). \quad (3.1.5)$$

Spunem că u și v sunt *ortogonale* dacă

$$(u, v) = 0. \quad (3.1.6)$$

Mai general, putem considera sisteme ortogonale $\{u_k\}_{k=1}^n$:

$$(u_i, u_j) = 0 \text{ dacă } i \neq j, \quad u_k \neq 0 \text{ pe } \text{supp}d\lambda; \quad i, j = \overline{1, n}, \quad k = \overline{1, n}. \quad (3.1.7)$$

Pentru un astfel de sistem are loc teorema generalizată a lui Pitagora

$$\left\| \sum_{k=1}^n \alpha_k u_k \right\|^2 = \sum_{k=1}^n |\alpha_k|^2 \|u_k\|^2. \quad (3.1.8)$$

O consecință importantă a lui (3.1.8) este aceea că orice sistem ortogonal este liniar independent pe $\text{supp}d\lambda$. Într-adevăr, dacă membrul stâng al lui (3.1.8) se anulează, atunci și membrul drept se anulează și deoarece $\|u_k\|^2 > 0$, din ipoteză rezultă $\alpha_1 = \alpha_2 = \dots = \alpha_n = 0$.

3.1.2. Ecuațiile normale

Suntem acum în măsură să rezolvăm problema de aproximare în sensul celor mai mici pătrate. Din (3.1.5) putem scrie pătratul erorii din L^2 sub forma

$$E^2[\varphi] := \|\varphi - f\|^2 = (\varphi - f, \varphi - f) = (\varphi, \varphi) - 2(\varphi, f) + (f, f).$$

Înlocuind pe φ cu expresia sa din (3.1.2) se obține

$$\begin{aligned} E^2[\varphi] &= \int_{\mathbb{R}} \left(\sum_{j=1}^n c_j \pi_j(t) \right)^2 d\lambda(t) - 2 \int_{\mathbb{R}} \left(\sum_{j=1}^n c_j \pi_j(t) \right) f(t) d\lambda(t) + \\ &+ \int_{\mathbb{R}} f^2(t) d\lambda(t). \end{aligned} \quad (3.1.9)$$

Pătratul erorii din L^2 este o funcție quadratică de coeficienții c_1, \dots, c_n ai lui φ . Problema celei mai bune aproximații în L^2 revine la a minimiza această funcție pătratică; ea se rezolvă anulând derivatele parțiale. Se obține

$$\frac{\partial}{\partial c_i} E^2[\varphi] = 2 \int_{\mathbb{R}} \left(\sum_{j=1}^n c_j \pi_j(t) \right) \pi_i(t) d\lambda(t) - 2 \int_{\mathbb{R}} \pi_i(t) f(t) d\lambda(t) = 0$$

adică

$$\sum_{j=1}^n (\pi_i, \pi_j) c_j = (\pi_i, f), \quad i = 1, 2, \dots, n. \quad (3.1.10)$$

Aceste ecuații se numesc *ecuații normale* pentru problema celor mai mici pătrate. Ele formează un sistem linear de forma

$$Ac = b \quad (3.1.11)$$

unde matricea A și vectorul b au elementele

$$A = [a_{ij}], \quad a_{ij} = (\pi_i, \pi_j), \quad b = [b_i], \quad b_i = (\pi_i, f). \quad (3.1.12)$$

Datorită simetriei produsului scalar, A este o matrice simetrică. Mai mult, A este pozitiv definită, adică

$$x^T Ax = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j > 0 \text{ dacă } x \neq [0, 0, \dots, 0]^T. \quad (3.1.13)$$

Funcția (3.1.13) se numește *formă pătratică* (deoarece este omogenă de grad 2). Pozitiv definirea lui A ne spune că forma pătratică ai cărei coeficienți sunt elementele lui A este întotdeauna nenegativă și zero numai dacă variabilele x_i se anulează.

Pentru a demonstra (3.1.13) să inserăm definiția lui a_{ij} și să utilizăm proprietățile (i)-(iv) ale produsului scalar

$$x^T Ax = \sum_{i=1}^n \sum_{j=1}^n x_i x_j (\pi_i, \pi_j) = \sum_{i=1}^n \sum_{j=1}^n (x_i \pi_i, x_j \pi_j) = \left\| \sum_{i=1}^n x_i \pi_i \right\|^2.$$

Aceasta este evident nenegativă. Ea este zero numai dacă $\sum_{i=1}^n x_i \pi_i \equiv 0$ pe $\text{supp} d\lambda$, care pe baza linear independenței lui π_i implică $x_1 = x_2 = \dots = x_n = 0$.

Este un rezultat cunoscut din algebra liniară că o matrice A simetrică pozitiv definită este nesingulară. Într-adevăr, determinantul său, precum și minorii principali sunt strict pozitivi. Rezultă că sistemul de ecuații normale (3.1.10) are soluție unică. Corespunde această soluție minimumului lui $E[\varphi]$ în (3.1.9)? Matricea hessiană $H = [\partial^2 E^2 / \partial c_i \partial c_j]$ trebuie să fie pozitiv definită. Dar $H = 2A$, deoarece E^2 este o funcție quadratică. De aceea, H , ca și A , este într-adevăr pozitiv definită și soluția ecuațiilor normale ne dă

minimul dorit. Problema de aproximare în sensul celor mai mici pătrate are o soluție unică, dată de

$$\widehat{\varphi}(t) = \sum_{j=1}^n \widehat{c}_j \pi_j(t) \quad (3.1.14)$$

unde $\widehat{c} = [\widehat{c}_1, \widehat{c}_2, \dots, \widehat{c}_n]^T$ este vectorul soluție al ecuațiilor normale (3.1.10). Aceasta rezolvă problema de aproximare în sensul celor mai mici pătrate complet în teorie, dar nu și în practică. Referitor la o mulțime generală de funcții de bază liniar independente, pot apărea următoarele dificultăți:

(1) Sistemul de ecuații normale (3.1.10) poate fi prost condiționat. Un exemplu simplu este următorul: $\text{supp} d\lambda = [0, 1]$, $d\lambda(t) = dt$ pe $[0, 1]$ și $\pi_j(t) = t^{j-1}$, $j = 1, 2, \dots, n$. Atunci

$$(\pi_i, \pi_j) = \int_0^1 t^{i+j-2} dt = \frac{1}{i+j-1}, \quad i, j = 1, 2, \dots, n,$$

adică matricea A este matricea Hilbert. Prost condiționarea ecuațiilor normale se datorează alegerii neinspirate a funcțiilor de bază. Acestea devin aproape liniar dependente când exponentul crește. O altă sursă de degradare provine din elementele membrului drept $b_j = \int_0^1 \pi_j(t) f(t) dt$. Când j este mare $\pi_j(t) = t^{j-1}$ se comportă pe $[0, 1]$ ca o funcție discontinuă. Un polinom care oscilează mai rapid pe $[0, 1]$ ar fi de preferat, căci ar angaja mai viguros funcția f .

(2) Al doilea dezavantaj este faptul că toți coeficienții \widehat{c}_j din (3.1.14) depind de n , adică $\widehat{c}_j = \widehat{c}_j^{(n)}$, $j = 1, 2, \dots, n$. Mărirea lui n ne dă un nou sistem de ecuații mai mare și cu o soluție complet diferită. Acest fenomen se numește *nepermanența coeficienților* \widehat{c}_j .

Amândouă neajunsurile (1) și (2) pot fi eliminate (sau măcar atenuate) alegând ca funcții de bază un sistem ortogonal,

$$(\pi_i, \pi_j) = 0 \text{ dacă } i \neq j \quad (\pi_i, \pi_j) = \|\pi_j\|^2 > 0 \quad (3.1.15)$$

Atunci sistemul de ecuații normale devine diagonal și poate fi rezolvat imediat cu formula

$$\widehat{c}_j = \frac{(\pi_j, f)}{(\pi_j, \pi_j)}, \quad j = 1, 2, \dots, n. \quad (3.1.16)$$

Evident, acești coeficienți \widehat{c}_j sunt independenți de n și odată calculați rămân la fel pentru orice n mai mare. Avem acum proprietatea de permanență a coeficienților. De asemenea nu trebuie să rezolvăm sistemul de ecuații normale, ci putem aplica direct (3.1.16). Aceasta nu înseamnă că nu sunt probleme numerice în (3.1.16). Într-adevăr, numitorul $\|\pi_j\|^2$ crește odată cu j , în timp ce integrala de la numărător (sau termenii individuali în cazul discret) au același ordin de mărime ca și f . De aceea ne așteptăm ca valorile \widehat{c}_j ale coeficienților să descrească rapid. Din acest motiv pot să apară erori de anulare atunci

când se calculează produsul scalar de la numărător. Problemele de anulare pot fi ocolite într-o oarecare măsură calculând \hat{c}_j în forma alternativă

$$\hat{c}_j = \frac{1}{(\pi_j, \pi_j)} \left(f - \sum_{k=1}^{j-1} c_k \pi_k, \pi_j \right), \quad j = 1, 2, \dots, n, \quad (3.1.17)$$

unde suma vidă (când $j = 1$) se ia egală cu zero. Evident din ortogonalitatea lui π_j , (3.1.17) este echivalentă cu (3.1.16) din punct de vedere matematic, dar nu neapărat și numeric.

Dăm un algoritm care calculează \hat{c}_j cu (3.1.17), dar în același timp și $\hat{\varphi}(t)$:

```

 $s_0 := 0;$ 
for  $j = 1, 2, \dots, n$  do
   $\hat{c}_j := \frac{1}{|\pi_j|^2} (f - s_{j-1}, \pi_j);$ 
   $s_j := s_{j-1} + \hat{c}_j \pi_j(t).$ 
end for

```

Orice sistem $\{\hat{\pi}_j\}$ care este liniar independent pe $\text{supp}d\lambda$ poate fi ortogonalizat (în raport cu măsura $d\lambda$) prin *procedeeul Gram-Schmidt*. Se ia

$$\pi_1 = \hat{\pi}_1$$

și apoi, pentru $j = 2, 3, \dots$ se calculează recursiv

$$\pi_j = \hat{\pi}_j - \sum_{k=1}^{j-1} c_k \pi_k, \quad c_k = \frac{(\hat{\pi}_j, \pi_k)}{(\pi_k, \pi_k)}, \quad k = \overline{1, j-1}.$$

Atunci fiecare π_j astfel determinat este ortogonal pe toate funcțiile precedente.

3.1.3. Eroarea în metoda celor mai mici pătrate. Convergența

Am văzut că dacă $\Phi = \Phi_n$ constă din n funcții π_j , $j = 1, 2, \dots, n$ care sunt liniar independente pe $\text{supp}d\lambda$, atunci problema de aproximare în sensul celor mai mici pătrate pentru $d\lambda$

$$\min_{\varphi \in \Phi_n} \|f - \varphi\|_{2, d\lambda} = \|f - \hat{\varphi}\|_{2, d\lambda} \quad (3.1.18)$$

are o soluție unică $\hat{\varphi} = \hat{\varphi}_n$, dată de (3.1.14). Există multe moduri de a selecta baza $\{\pi_j\}$ a lui Φ_n și de aceea mai multe moduri de a reprezenta soluția, care conduc totuși la aceeași funcție. Eroarea în sensul celor mai mici pătrate – cantitatea din dreapta relației (3.1.18) – este independentă de alegerea funcțiilor de bază (deși calculul soluției, așa cum s-a menționat anterior, nu este). În studiul acestor erori, putem presupune fără a

restrânge generalitatea că baza π_j este un sistem ortogonal (fiecare sistem liniar independent poate fi ortogonalizat prin procedeul Gram-Schmidt). Avem conform (3.1.16)

$$\widehat{\varphi}_n(t) = \sum_{j=1}^n \widehat{c}_j \pi_j(t), \quad \widehat{c}_j = \frac{(\pi_j, f)}{(\pi_j, \pi_j)}. \quad (3.1.19)$$

Observăm întâi că eroarea $f - \widehat{\varphi}_n$ este ortogonală pe Φ_n , adică

$$(f - \widehat{\varphi}_n, \varphi) = 0, \quad \forall \varphi \in \Phi_n \quad (3.1.20)$$

unde produsul scalar este cel din (3.1.3). Deoarece φ este o combinație liniară de π_k , este suficient să arătăm (3.1.20) pentru fiecare $\varphi = \pi_k$, $k = 1, 2, \dots, n$. Înlocuind $\widehat{\varphi}_n$ cu expresia sa din (3.1.19) în (3.1.20), găsim

$$(f - \widehat{\varphi}_n, \pi_k) = \left(f - \sum_{j=1}^n \widehat{c}_j \pi_j, \pi_k \right) = (f, \pi_k) - \widehat{c}_k (\pi_k, \pi_k) = 0,$$

ultima ecuație rezultând din formula pentru \widehat{c}_k din (3.1.19). Rezultatul din (3.1.20) are o interpretare geometrică simplă. Dacă reprezentăm funcțiile ca vectori și spațiul Φ_n ca un plan, atunci pentru orice funcție f care înțeapă planul Φ_n , aproximanta în sensul celor mai mici pătrate $\widehat{\varphi}_n$ este *proiecția ortogonală* a lui f pe Φ_n , vezi figura 3.1.

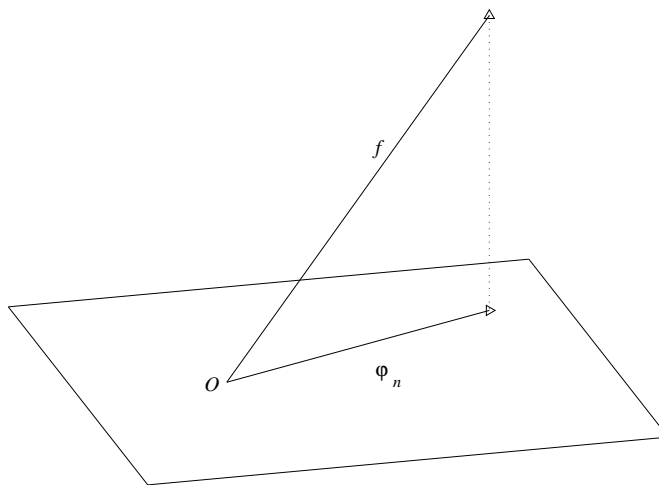


Figura 3.1: Aproximația în sensul celor mai mici pătrate ca proiecție ortogonală

În particular, alegând $\varphi = \widehat{\varphi}_n$ în (3.1.20) obținem

$$(f - \widehat{\varphi}_n, \widehat{\varphi}_n) = 0$$

și de aceea, deoarece $f = (f - \widehat{\varphi}) + \widehat{\varphi}$, conform teoremei lui Pitagora și generalizării sale (3.1.8)

$$\begin{aligned} \|f\|^2 &= \|f - \widehat{\varphi}\|^2 + \|\widehat{\varphi}\|^2 = \|f - \widehat{\varphi}_n\|^2 + \left\| \sum_{j=1}^n \widehat{c}_j \pi_j \right\|^2 \\ &= \|f - \widehat{\varphi}_n\|^2 + \sum_{j=1}^n |\widehat{c}_j|^2 \|\pi_j\|^2. \end{aligned}$$

Exprimând primul termen din dreapta obținem

$$\|f - \widehat{\varphi}_n\| = \left\{ \|f\|^2 - \sum_{j=1}^n |\widehat{c}_j|^2 \|\pi_j\|^2 \right\}^{1/2}, \quad \widehat{c}_j = \frac{(\pi_j, f)}{(\pi_j, \pi_j)}. \quad (3.1.21)$$

De notat că expresia dintre acolade trebuie să fie nenegativă.

Formula (3.1.21) este de interes teoretic, dar de utilitate practică limitată. De notat, într-adevăr, că pe măsură ce eroarea se apropie de nivelul eps al preciziei mașinii, calculul erorii din membrul drept al lui (3.1.21) nu poate produce ceva mai mic decât $\sqrt{\text{eps}}$ datorită erorilor comise în timpul scăderilor sub radical (astfel se poate obține chiar un rezultat negativ sub radical). Utilizând în loc definiția

$$\|f - \widehat{\varphi}_n\| = \left\{ \int_{\mathbb{R}} [f(t) - \widehat{\varphi}_n(t)]^2 d\lambda(t) \right\}^{\frac{1}{2}},$$

împreună, probabil, cu o regulă de cuadratură (pozitivă) potrivită se garantează că se produce un rezultat nenegativ, potențial la fel de mic ca $O(\text{eps})$.

Dacă se dă acum un șir de spații liniare Φ_n , $n = 1, 2, 3, \dots$, avem evident

$$\|f - \widehat{\varphi}_1\| \geq \|f - \widehat{\varphi}_2\| \geq \|f - \widehat{\varphi}_3\| \geq \dots,$$

care rezultă nu numai din (3.1.21), dar mai direct din faptul că

$$\Phi_1 \subset \Phi_2 \subset \Phi_3 \subset \dots$$

Deoarece există o infinitate de astfel de spații, atunci secvența de erori din L^2 , fiind monoton descrescătoare, trebuie să convergă la o limită. Este limita 0? Dacă este așa, spunem că aproximarea prin metoda celor mai mici pătrate converge în medie când $n \rightarrow \infty$. Este evident din (3.1.21) că o condiție necesară și suficientă pentru aceasta este

$$\sum_{j=1}^{\infty} |\widehat{c}_j|^2 \|\pi_j\|^2 = \|f\|^2. \quad (3.1.22)$$

Un mod echivalent de a formula convergența este următorul: dându-se f cu $\|f\| < \infty$, adică $\forall f \in L^2_{d\lambda}$ și dându-se un $\varepsilon > 0$ arbitrar de mic, există un întreg $n = n_\varepsilon$ și o funcție $\varphi^* \in \Phi_n$ astfel încât $\|f - \varphi^*\| \leq \varepsilon$. O clasă de spații Φ_n având această proprietate se numește completă în raport cu norma $\|\cdot\| = \|\cdot\|_{2,d\lambda}$. Vom numi relația (3.1.22) relația de completitudine. Pentru un interval finit $[a, b]$ putem defini completitudinea lui $\{\varphi_n\}$ de asemenea pentru norma uniformă $\|\cdot\|_\infty$ pe $[a, b]$. Se poate presupune că $f \in C[a, b]$ și $\pi_j \in C[a, b]$ și numim $\{\varphi_n\}$ completă în norma $\|\cdot\|_\infty$ dacă pentru orice $f \in C[a, b]$ și orice $\varepsilon > 0$ există $n = n_\varepsilon$ și un $\varphi^* \in \Phi_n$ astfel încât $\|f - \varphi^*\|_\infty \leq \varepsilon$. Este ușor de văzut că completitudinea lui $\{\varphi_n\}$ în norma $\|\cdot\|_\infty$ pe (a, b) implică completitudinea lui $\{\varphi_n\}$ în norma din L^2 $\|\cdot\|_{2,d\lambda}$, unde $\text{supp}d\lambda = [a, b]$ și deci convergența procesului de aproximare prin metoda celor mai mici pătrate. Într-adevăr, fie $\varepsilon > 0$ arbitrar și fie n și $\varphi^* \in \Phi_n$ astfel încât

$$\|f - \varphi^*\|_\infty \leq \frac{\varepsilon}{\left(\int_{\mathbb{R}} d\lambda(t)\right)^{1/2}}.$$

Aceasta este posibil din ipoteză. Atunci

$$\begin{aligned} \|f - \varphi^*\|_{2,d\lambda} &= \left(\int_{\mathbb{R}} [f(t) - \varphi^*(t)]^2 d\lambda(t) \right)^{1/2} \leq \\ &\leq \|f - \varphi^*\|_\infty \left(\int_{\mathbb{R}} d\lambda(t) \right)^{1/2} \leq \\ &\leq \frac{\varepsilon}{\left(\int_{\mathbb{R}} d\lambda(t)\right)^{1/2}} \left(\int_{\mathbb{R}} d\lambda(t) \right)^{1/2} = \varepsilon, \end{aligned}$$

așa cum s-a afirmat.

Exemplul 3.1.1. $\Phi_n = \mathbb{P}_{n-1}$. Aici completitudinea lui $\{\varphi_n\}$ în norma $\|\cdot\|_\infty$ (pe un interval finit $[a, b]$) este o consecință a teoremei de aproximare a lui Weierstrass. Astfel aproximația polinomială în sensul celor mai mici pătrate pe un interval finit converge întotdeauna (în medie). \diamond

3.2. Exemple de sisteme ortogonale

Unul dintre cele mai utilizate sisteme este sistemul trigonometric cunoscut din analiza Fourier. Un alt sistem larg utilizat este cel al polinoamelor ortogonale.

(1) **Sistemul trigonometric** este format din funcțiile:

$$1, \cos t, \cos 2t, \cos 3t, \dots, \sin t, \sin 2t, \sin 3t, \dots$$

El este ortogonal pe $[0, 2\pi]$ în raport cu măsura

$$d\lambda(t) = \begin{cases} dt & \text{pe } [0, 2\pi] \\ 0 & \text{în rest} \end{cases}$$

Avem

$$\int_0^{2\pi} \sin kt \sin \ell t dt = \begin{cases} 0, & \text{pentru } k \neq \ell \\ \pi, & \text{pentru } k = \ell \end{cases} \quad k, \ell = 1, 2, 3, \dots$$

$$\int_0^{2\pi} \cos kt \cos \ell t dt = \begin{cases} 0, & k \neq \ell \\ 2\pi, & k = \ell = 0 \\ \pi, & k = \ell > 0 \end{cases} \quad k, \ell = 0, 1, 2, \dots$$

$$\int_0^{2\pi} \sin kt \cos \ell t dt = 0, \quad k = 1, 2, 3, \dots, \quad \ell = 0, 1, 2, \dots$$

Aproximarea are forma

$$f(t) = \frac{a_0}{2} + \sum_{k=1}^{\infty} (a_k \cos kt + b_k \sin kt). \quad (3.2.1)$$

Utilizând (3.1.16) obținem

$$a_k = \frac{1}{\pi} \int_0^{2\pi} f(t) \cos kt dt, \quad k = 1, 2, \dots$$

$$b_k = \frac{1}{\pi} \int_0^{2\pi} f(t) \sin kt dt, \quad k = 1, 2, \dots \quad (3.2.2)$$

numiți *coeficienții Fourier* ai lui f . Ei sunt coeficienții (3.1.16) pentru sistemul trigonometric. Prin extensie coeficienții (3.1.16) pentru orice sistem ortogonal (π_j) se vor numi coeficienții Fourier ai lui f relativ la acest sistem. În particular, recunoaștem în seria Fourier trunchiată pentru $k = m$ cea mai bună aproximare a lui f din clasa polinoamelor trigonometrice de grad $\leq n$ relativ la norma

$$\|u\|_2 = \left(\int_0^{2\pi} |u(t)|^2 dt \right)^{1/2}.$$

(2) **Polinoame ortogonale.** Dându-se o măsură $d\lambda$, știm că orice sistem finit de puteri $1, t, t^2, \dots$ este liniar independent pe $[a, b]$, dacă $\text{supp} d\lambda = [a, b]$, iar $1, t, \dots, t^{n-1}$ sunt liniar independente pe $\text{supp} d\lambda = \{t_1, t_2, \dots, t_N\}$. Deoarece o mulțime de vectori liniar independenți a unui spațiu liniar poate fi ortogonalizată prin procedeul Gram-Schmidt, orice măsură $d\lambda$ de tipul considerat generează o mulțime unică de polinoame monice $\pi_j(t, d\lambda)$, $j = 0, 1, 2, \dots$ ce satisfac

$$\text{grad} \pi_j = j, \quad j = 0, 1, 2, \dots$$

$$\int_{\mathbb{R}} \pi_k(t) \pi_\ell(t) d\lambda(t) = 0, \quad \text{dacă } k \neq \ell. \quad (3.2.3)$$

Aceste polinoame se numesc *polinoame ortogonale* relativ la măsura $d\lambda$. Vom permite indicilor să meargă de la 0. Mulțimea $\{\pi_j\}$ este infinită dacă $\text{suppd}\lambda = [a, b]$ și constă din exact N polinoame $\pi_0, \pi_1, \dots, \pi_{N-1}$ dacă $\text{suppd}\lambda = \{t_1, \dots, t_N\}$. În ultimul caz polinoamele se numesc *polinoame ortogonale discrete*.

Între trei polinoame ortogonale monice² consecutive există o relație liniară. Mai exact, există constantele reale $\alpha_k = \alpha_k(d\lambda)$ și $\beta_k = \beta_k(d\lambda) > 0$ (depinzând de măsura $d\lambda$) astfel încât

$$\begin{aligned}\pi_{k+1}(t) &= (t - \alpha_k)\pi_k(t) - \beta_k\pi_{k-1}(t), \quad k = 0, 1, 2, \dots \\ \pi_{-1}(t) &= 0, \quad \pi_0(t) = 1.\end{aligned}\tag{3.2.4}$$

(Se subînțelege că (3.2.4) are loc pentru orice $k \in \mathbb{N}$ dacă $\text{suppd}\lambda = [a, b]$ și numai pentru $k = \overline{0, N-2}$ dacă $\text{suppd}\lambda = \{t_1, t_2, \dots, t_N\}$).

Pentru a demonstra (3.2.4) și a obține expresiile coeficienților să observăm că

$$\pi_{k+1}(t) - t\pi_k(t)$$

este un polinom de grad $\leq k$, și deci poate fi exprimat ca o combinație liniară a lui $\pi_0, \pi_1, \dots, \pi_k$. Scriem această combinație sub forma

$$\pi_{k+1} - t\pi_k(t) = -\alpha_k\pi_k(t) - \beta_k\pi_{k-1}(t) + \sum_{j=0}^{k-2} \gamma_{k,j}\pi_j(t)\tag{3.2.5}$$

(sumele vide se consideră nule). Înmulțim scalar ambii membri ai relației anterioare cu π_k și obținem

$$(-t\pi_k, \pi_k) = -\alpha_k(\pi_k, \pi_k)$$

adică

$$\alpha_k = \frac{(t\pi_k, \pi_k)}{(\pi_k, \pi_k)}, \quad k = 0, 1, 2, \dots\tag{3.2.6}$$

La fel, înmulțind scalar cu π_{k-1} obținem

$$(-t\pi_k, \pi_{k-1}) = -\beta_k(\pi_{k-1}, \pi_{k-1}).$$

Deoarece $(t\pi_k, \pi_{k-1}) = (\pi_k, t\pi_{k-1})$ și $t\pi_{k-1}$ diferă de π_k printr-un polinom de grad $< k$ se obține prin ortogonalitate $(t\pi_k, \pi_{k-1}) = (\pi_k, \pi_k)$, deci

$$\beta_k = \frac{(\pi_k, \pi_k)}{(\pi_{k-1}, \pi_{k-1})}, \quad k = 1, 2, \dots\tag{3.2.7}$$

Înmulțind (3.2.5) cu π_ℓ , $\ell < k-1$, se obține

$$\gamma_{k,\ell} = 0, \quad \ell = 0, 1, \dots, k-1.\tag{3.2.8}$$

²Un polinom se numește *monic* dacă coeficientul său dominant este 1.

Formula de recurență (3.2.4) ne dă o modalitate practică de determinare a polinoamelor ortogonale. Deoarece $\pi_0 = 1$, putem calcula α_0 cu (3.2.6) pentru $k = 0$. Se continuă apoi cu π_1 , utilizând (3.2.4) pentru $k = 0$. Mai departe, utilizând alternativ (3.2.6), (3.2.7) și (3.2.4) putem calcula câte polinoame ortogonale dorim. Procedeu – numit procedura lui Stieltjes ³ – este foarte potrivit pentru polinoame ortogonale discrete, căci în acest caz produsul scalar se exprimă prin sume finite. În cazul continuu, calculul produsului scalar necesită calcul de integrale, ceea ce complică lucrurile. Din fericire, pentru multe măsuri speciale importante, coeficienții se cunosc explicit. Cazul special când măsura este simetrică (adică $d\lambda(t) = w(t)$ cu $w(-t) = w(t)$ și $\text{supp}d\lambda$ simetrică față de origine) merită o atenție specială deoarece în acest caz $\alpha_k = 0, \forall k \in \mathbb{N}$, conform lui (3.2.1) căci

$$(\pi_k, \pi_k) = \int_{\mathbb{R}} w(t)t\pi_k^2(t)dt = \int_a^b w(t)t\pi_k^2(t)dt = 0,$$

deoarece avem o integrală dintr-o funcție impară pe un domeniu simetric față de origine.

3.2.1. Exemple de polinoame ortogonale

Polinoamele lui Legendre

Se definesc prin așa-numita formulă a lui Rodrigues

$$\pi_k(t) = \frac{k!}{(2k)!} \frac{d^k}{dt^k} (t^2 - 1)^k. \quad (3.2.9)$$

Verificăm întâi ortogonalitatea pe $[-1, 1]$ în raport cu măsura $d\lambda(t) = dt$. Pentru orice $0 \leq \ell < k$, prin integrare repetată prin părți se obține:

$$\int_{-1}^1 \frac{d^k}{dt^k} (t^2 - 1)^k = \sum_{m=0}^{\ell} \ell(\ell-1) \dots (\ell-m+1) t^{\ell-m} \frac{d^{k-m-1}}{dt^{k-m-1}} (t^2 - 1)^k \Big|_{-1}^1 = 0, \quad (3.2.10)$$

3



Thomas Jan Stieltjes (1856-1894), matematician olandez, a studiat la Institutul Tehnic din Delft, dar nu și-a luat niciodată licența datorită aversiunii pe care o avea pentru examene. A lucrat la Observatorul astronomic din Leyda, în calitate de „calculator asistent pentru calcule astronomice”. Lucrările sale timpurii au atras atenția lui Hermite, care i-a asigurat un post universitar la Toulouse. Prietenia lor a fost exemplară. Stieltjes este cunoscut pentru lucrările sale despre fracții continue și problema momentelor, în care printre altele a introdus integrala care îi poartă numele. A murit de tuberculoză la 38 de ani.

ultima relație având loc deoarece $0 \leq k - m - 1 < k$. Deci,

$$(\pi_k, p) = 0, \quad \forall p \in \mathbb{P}_{k-1},$$

demonstrându-se astfel ortogonalitatea. Datorită simetriei, putem scrie

$$\pi_k(t) = t^k + \mu_k t^{k-2} + \dots, \quad k \geq 2$$

și observând (din nou datorită simetriei) că relația de recurență are forma

$$\pi_{k+1}(t) = t\pi_k(t) - \beta_k \pi_{k-1}(t),$$

obținem

$$\beta_k = \frac{t\pi_k(t) - \pi_{k+1}(t)}{\pi_{k-1}(t)},$$

care este valabilă pentru orice t . Făcând $t \rightarrow \infty$,

$$\beta_k = \lim_{t \rightarrow \infty} \frac{t\pi_k(t) - \pi_{k+1}(t)}{\pi_{k-1}(t)} = \lim_{t \rightarrow \infty} \frac{(\mu_k - \mu_{k+1})t^{k-1} + \dots}{t^{k-1} + \dots} = \mu_k - \mu_{k+1}.$$

(Dacă $k = 1$, punem $\mu_1 = 0$.) Din formula lui Rodrigues rezultă

$$\begin{aligned} \pi_k(t) &= \frac{k!}{(2k)!} \frac{d^k}{dt^k} (t^{2k} - kt^{2k-2} + \dots) \\ &= \frac{k!}{(2k)!} (2k(2k-1) \dots (k+1)t^k - k(2k-2)(2k-3) \dots (k-1)t^{k-1} + \dots) \\ &= t^k - \frac{k(k-1)}{2(2k-1)} t^{k-2} + \dots, \end{aligned}$$

așa că

$$\mu_k = \frac{k(k-1)}{2(2k-1)}, \quad k \geq 2.$$

Deci,

$$\beta_k = \mu_k - \mu_{k+1} = \frac{k^2}{(2k-1)(2k+1)}$$

și deoarece $\mu_1 = 0$,

$$\beta_k = \frac{1}{4 - k^{-2}}, \quad k \geq 1. \quad (3.2.11)$$

Polinoamele Cebîșev de speța I

Polinoamele lui Cebîșev ⁴ de speța I se pot defini prin relația

$$T_n(x) = \cos(n \arccos x), \quad n \in \mathbb{N}. \quad (3.2.12)$$

Din identitatea trigonometrică

$$\cos(k+1)\theta + \cos(k-1)\theta = 2 \cos \theta \cos k\theta$$

și din (3.2.12), punând $\theta = \arccos x$ se obține

$$\begin{aligned} T_{k+1}(x) &= 2xT_k(x) - T_{k-1}(x) \quad k = 1, 2, 3, \dots \\ T_0(x) &= 1, \quad T_1(x) = x. \end{aligned} \quad (3.2.13)$$

De exemplu

$$\begin{aligned} T_2(x) &= 2x^2 - 1, \\ T_3(x) &= 4x^3 - 3x, \\ T_4(x) &= 8x^4 - 8x^2 + 1 \end{aligned}$$

ș.a.m.d.

Din relația (3.2.13) se obține pentru coeficientul dominant al lui T_n valoarea 2^{n-1} (dacă $n \geq 1$), deci polinomul Cebîșev de speța I monic este

$$\overset{\circ}{T}_n(x) = \frac{1}{2^{n-1}} T_n(x), \quad n \geq 0, \quad \overset{\circ}{T}_0 = T_0. \quad (3.2.14)$$

Din (3.2.12) se pot obține rădăcinile lui T_n

$$x_k^{(n)} = \cos \theta_k^{(n)}, \quad \theta_k^{(n)} = \frac{2k-1}{2n} \pi, \quad k = \overline{1, n}. \quad (3.2.15)$$

Ele sunt proiecțiile pe axa reală ale punctelor de pe cercul unitate de argument $\theta_k^{(n)}$; figura 3.2 ilustrează acest lucru pentru $n=4$.

Pe intervalul $[-1, 1]$ T_n oscilează de la +1 la -1, atingând aceste valori extreme în punctele

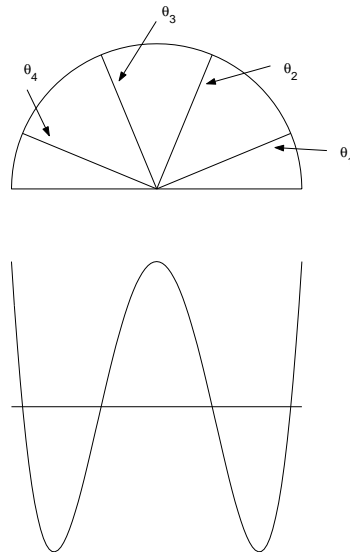
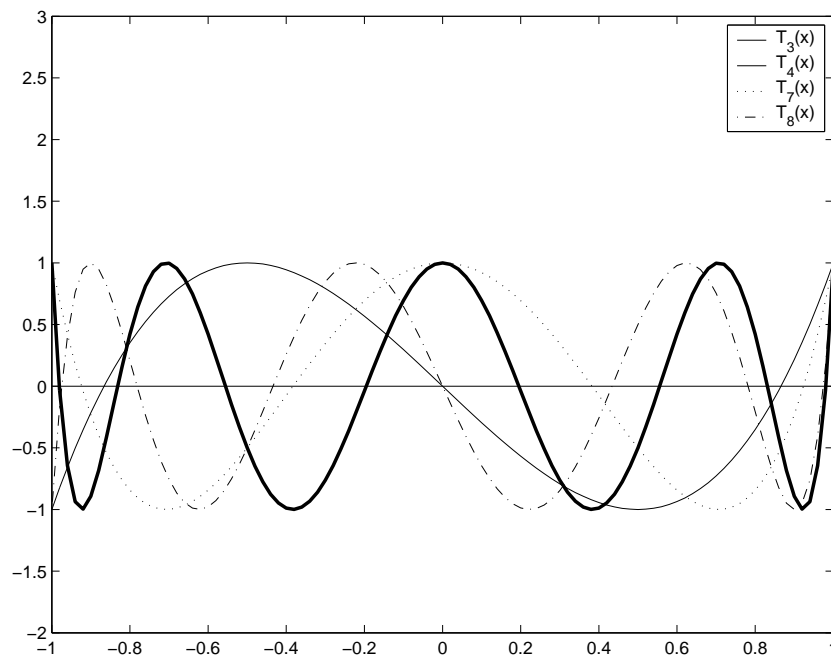
$$y_k^{(n)} = \cos \eta_k^{(n)}, \quad \eta_k^{(n)} = \frac{k\pi}{n}, \quad k = \overline{0, n}.$$

În figura 3.3 apar graficele unor polinoame Cebîșev de speța I.

4



Pafnuti Levovici Cebîșev (1821-1894), matematician rus, cel mai important reprezentant al școlii matematice din Sankt Petersburg. A avut contribuții de pionierat în domeniul teoriei numerelor, calculului probabilităților și teoriei aproximării. Este considerat fondatorul teoriei constructive a funcțiilor, dar a lucrat și în mecanică și în balistică.

Figura 3.2: Polinomul Cebîșev T_4 și rădăcinile saleFigura 3.3: Polinoamele Cebîșev T_3, T_4, T_7, T_8 pe $[-1, 1]$

Polinoamele Cebîșev de speța I sunt ortogonale în raport cu măsura

$$d\lambda(x) = \frac{dx}{\sqrt{1-x^2}}, \quad \text{pe } [-1, 1].$$

Se verifică ușor din (3.2.12) că

$$\begin{aligned} \int_{-1}^1 T_k(x)T_\ell(x) \frac{dx}{\sqrt{1-x^2}} &= \int_0^\pi T_k(\cos \theta)T_\ell(\cos \theta)d\theta \\ &= \int_0^\pi \cos k\theta \cos \ell\theta d\theta = \begin{cases} 0, & \text{dacă } k \neq \ell, \\ \pi, & \text{dacă } k = \ell = 0, \\ \pi/2, & \text{dacă } k = \ell \neq 0. \end{cases} \end{aligned} \quad (3.2.16)$$

Dezvoltarea în serie Fourier de polinoame Cebîșev este dată de

$$f(x) = \sum_{j=0}^{\infty} c_j T_j(x) := \frac{1}{2}c_0 + \sum_{j=1}^{\infty} c_j T_j(x), \quad (3.2.17)$$

unde

$$c_j = \frac{2}{\pi} \int_{-1}^1 f(x)T_j(x) \frac{dx}{\sqrt{1-x^2}}, \quad j \in \mathbb{N}.$$

Păstrând în (3.2.17) numai termenii de grad cel mult n se obține o aproximare polinomială utilă de grad n

$$\tau_n(x) = \sum_{j=0}^n c_j T_j(x), \quad (3.2.18)$$

având eroarea

$$f(x) - \tau_n(x) = \sum_{j=n+1}^{\infty} c_j T_j(x) \approx c_{n+1} T_{n+1}(x). \quad (3.2.19)$$

Aproximanta din (3.2.18) este cu atât mai bună cu cât coeficienții din extremitatea dreaptă tind mai repede către zero. Eroarea (3.2.19) oscilează în esență între $+c_{n+1}$ și $-c_{n+1}$ și este deci de mărime „uniformă”. Acest lucru contrastează puternic cu dezvoltarea Taylor în jurul lui $x = 0$, unde polinomul de grad n are eroarea proporțională cu x^{n+1} pe $[-1, 1]$.

Dintre toate polinoamele monice de grad n $\overset{\circ}{T}_n$ are norma uniformă cea mai mică.

Teorema 3.2.1 (Cebîșev) Pentru orice polinom monic $\overset{\circ}{p}_n$ de grad n are loc

$$\max_{-1 \leq x \leq 1} \left| \overset{\circ}{p}_n(x) \right| \geq \max_{-1 \leq x \leq 1} \left| \overset{\circ}{T}_n(x) \right| = \frac{1}{2^{n-1}}, \quad n \geq 1, \quad (3.2.20)$$

unde $\overset{\circ}{T}_n(x)$ este dat de (3.2.14).

Demonstrație. Se face prin reducere la absurd. Presupunem că

$$\max_{-1 \leq x \leq 1} |\overset{\circ}{p}_n(x)| < \frac{1}{2^{n-1}}. \quad (3.2.21)$$

Atunci polinomul $d_n(x) = \overset{\circ}{T}_n(x) - \overset{\circ}{p}_n(x)$ (de grad $\leq n - 1$) satisface

$$d_n(y_0^{(n)}) > 0, d_n(y_1^{(n)}) < 0, d_n(y_2^{(n)}) > 0, \dots, (-1)^n d_n(y_n^{(n)}) > 0. \quad (3.2.22)$$

Deoarece d_n are n schimbări de semn, el este identic nul; aceasta contrazice (3.2.22) și astfel (3.2.21) nu poate fi adevărată. \square

Rezultatul (3.2.20) se poate interpreta în modul următor: cea mai bună aproximare uniformă din \mathbb{P}_{n-1} pe $[-1, 1]$ a lui $f(x) = x^n$ este dată de $x^n - \overset{\circ}{T}_n(x)$, adică, de agregarea termenilor până la gradul $n - 1$ din $\overset{\circ}{T}_n$ luați cu semnul minus. Din teoria aproximațiilor uniforme se știe că cea mai bună aproximare polinomială uniformă este unică. Deci, egalitatea în (3.2.20) poate avea loc numai dacă $\overset{\circ}{p}_n(x) = \overset{\circ}{T}_n(x)$.

Polinoamele lui Cebîșev de speța a doua

Se definesc prin

$$Q_n(t) = \frac{\sin[(n+1) \arccos t]}{\sqrt{1-t^2}}, \quad t \in [-1, 1].$$

Ele sunt ortogonale pe $[-1, 1]$ în raport cu măsura $d\lambda(t) = w(t)dt$, $w(t) = \sqrt{1-t^2}$. Relația de recurență este

$$Q_{n+1}(t) = 2tQ_n(t) - Q_{n-1}(t), \quad Q_0(t) = 1, \quad Q_1(t) = 2t.$$

Polinoamele lui Laguerre

Polinoamele lui Laguerre ⁵ sunt ortogonale pe $[0, \infty)$ în raport cu ponderea $w(t) = t^\alpha e^{-t}$. Se definesc prin

$$l_n^\alpha(t) = \frac{e^t t^{-\alpha}}{n!} \frac{d^n}{dt^n} (t^{n+\alpha} e^{-t}) \text{ pentru } \alpha > 1.$$

5



Edmond Laguerre (1834-1886) matematician francez, activ în Paris, cu contribuții esențiale în geometrie, algebră și analiză.

Relația de recurență este

$$nl_n^\alpha(t) - (2n - 1 + \alpha - t)l_{n-1}^\alpha(t) + (n - 1 - \alpha)l_{n-2}^\alpha(t) = 0.$$

Polinoamele lui Hermite

Polinoamele lui Hermite se definesc prin

$$H_n(t) = (-1)^n e^{t^2} \frac{d^n}{dt^n} (e^{-t^2}).$$

Ele sunt ortogonale pe $(-\infty, \infty)$ în raport cu ponderea $w(t) = e^{-t^2}$ și verifică relația de recurență

$$H_{n+1}(t) = 2tH_n(t) - 2nH_{n-1}(t),$$

$$H_0(t) = 1, \quad H_1(t) = 2t.$$

Polinoamele lui Jacobi

Sunt ortogonale pe $[-1, 1]$ în raport cu ponderea

$$w(t) = (1 - t)^\alpha (1 + t)^\beta.$$

Pentru detalii privind aspectele algoritmice legate de polinoame ortogonale a se vedea [21].

3.3. Interpolare polinomială

3.3.1. Spațiul $H^n[a, b]$

Pentru $n \in \mathbb{N}^*$, definim

$$H^n[a, b] = \{f : [a, b] \rightarrow \mathbb{R} : f \in C^{n-1}[a, b], f^{(n-1)} \text{ absolut continuă pe } [a, b]\}. \quad (3.3.1)$$

Orice funcție $f \in H^n[a, b]$ admite o reprezentare de tip Taylor cu restul sub formă integrală

$$f(x) = \sum_{k=0}^{n-1} \frac{(x-a)^k}{k!} f^{(k)}(a) + \int_a^x \frac{(x-t)^{n-1}}{(n-1)!} f^{(n)}(t) dt. \quad (3.3.2)$$

$H^n[a, b]$ este un spațiu liniar.

Observația 3.3.1. Funcția $f : I \rightarrow \mathbb{R}$, I interval, se numește absolut continuă pe I dacă $\forall \varepsilon > 0 \exists \delta > 0$ astfel încât oricare ar fi un sistem finit de subintervale disjuncte ale lui I $\{(a_k, b_k)\}_{k=1, n}$ cu proprietatea $\sum_{k=1}^n (b_k - a_k) < \delta$ să avem

$$\sum_{k=1}^n |f(b_k) - f(a_k)| < \varepsilon. \quad \diamond$$

Teorema următoare, datorată lui Peano ⁶, de o importanță deosebită în analiza numerică, este o teoremă de reprezentare a funcționalelor liniare reale, definite pe $H^n[a, b]$.

Teorema 3.3.2 (Peano) Fie L o funcțională reală, continuă, definită pe $H^n[a, b]$. Dacă $\text{Ker} L = \mathbb{P}_{n-1}$ atunci

$$Lf = \int_a^b K(t) f^{(n)}(t) dt, \quad (3.3.3)$$

unde

$$K(t) = \frac{1}{(n-1)!} L[(\cdot - t)_+^{n-1}] \quad (\text{nucleul lui Peano}). \quad (3.3.4)$$

Observația 3.3.3. Funcția

$$z_+ = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

se numește *parte pozitivă*, iar z_+^n se numește *putere trunchiată*. \(\diamond\)

Demonstrație. f admite o reprezentare de tip Taylor, cu restul în formă integrală

$$f(x) = T_{n-1}(x) + R_{n-1}(x)$$

unde

$$R_{n-1}(x) = \int_a^x \frac{(x-t)^{n-1}}{(n-1)!} f^{(n)}(t) dt = \frac{1}{(n-1)!} \int_a^b (x-t)_+^{n-1} f^{(n)}(t) dt.$$



Giuseppe Peano (1858-1932), matematician italian activ la Torino, cu contribuții fundamentale în logica matematică, teoria mulțimilor, și fundamentele matematicii. Teoremele generale de existență din domeniul teoriei ecuațiilor diferențiale îi poartă numele.

Aplicând L obținem

$$Lf = \underbrace{LT_{n-1}}_0 + LR_{n-1} \Rightarrow Lf = \frac{1}{(n-1)!} L \left(\int_a^b (\cdot - t)_+^{n-1} f^{(n)}(t) dt \right) =$$

$$\stackrel{cont}{=} \frac{1}{(n-1)!} \int_a^b L(\cdot - t)_+^{n-1} f^{(n)}(t) dt.$$

□

Observația 3.3.4. Concluzia teoremei rămâne valabilă și dacă L nu este continuă, ci are forma

$$Lf = \sum_{i=0}^{n-1} \int_a^b f^{(i)}(x) d\mu_i(x), \quad \mu_i \in BV[a, b],$$

unde $BV[a, b]$ este mulțimea funcțiilor cu variație mărginită pe $[a, b]$. ◇

Corolarul 3.3.5 Dacă K păstrează semn constant pe $[a, b]$ și $f^{(n)}$ este continuă pe $[a, b]$, atunci există $\xi \in [a, b]$ astfel încât

$$Lf = \frac{1}{n!} f^{(n)}(\xi) Le_n, \quad (3.3.5)$$

unde $e_k(x) = x^k$, $k \in \mathbb{N}$.

Demonstrație. Deoarece K păstrează u{a} semn constant putem aplica în (3.3.3) teorema de medie

$$Lf = f^{(n)}(\xi) \int_a^b K_n(t) dt, \quad \xi \in [a, b].$$

Luând $f = e_n$ se obține chiar (3.3.5). □

3.3.2. Interpolare Lagrange

Fie intervalul închis $[a, b] \subset \mathbb{R}$, $f : [a, b] \rightarrow \mathbb{R}$ și o mulțime de $m+1$ puncte distincte $\{x_0, x_1, \dots, x_m\} \subset [a, b]$.

Teorema 3.3.6 Există un polinom și numai unul $L_m f \in \mathbb{P}_m$ astfel încât

$$\forall i = 0, 1, \dots, m, \quad (L_m f)(x_i) = f(x_i); \quad (3.3.6)$$

acest polinom se scrie sub forma

$$(L_m f)(x) = \sum_{i=0}^m f(x_i) \ell_i(x), \quad (3.3.7)$$

unde

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^m \frac{x - x_j}{x_i - x_j}. \quad (3.3.8)$$

Definiția 3.3.7 Polinomul $L_m f$ definit astfel se numește polinom de interpolare Lagrange ⁷ a lui f relativ la punctele x_0, x_1, \dots, x_m , iar funcțiile $l_i(x)$, $i = \overline{0, m}$, se numesc polinoame de bază (fundamentale) Lagrange asociate acelor puncte.

Demonstrație. Se verifică imediat că $l_i \in \mathbb{P}_m$ și că $l_i(x_j) = \delta_{ij}$ (simbolul lui Kronecker); rezultă că polinomul $L_m f$ definit de (3.3.6) este de grad cel mult m și verifică (3.3.7). Presupunem că există un alt polinom $p_m^* \in \mathbb{P}_m$ care verifică (3.3.7) și punem $q_m = L_m - p_m^*$; avem $q_m \in \mathbb{P}_m$ și $\forall i = 0, 1, \dots, m$, $q_m(x_i) = 0$; deci q_m având $(m + 1)$ rădăcini distincte este identic nul, de unde unicitatea lui L_m . \square

Observația 3.3.8. Polinomul fundamental l_i este deci unicul polinom care verifică

$$l_i \in \mathbb{P}_m \text{ și } \forall j = 0, 1, \dots, m, \quad l_i(x_j) = \delta_{ij}.$$

Punând

$$u(x) = \prod_{j=0}^m (x - x_j),$$

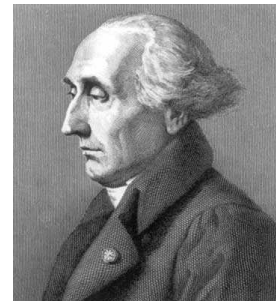
din (3.3.8) se deduce că $\forall x \neq x_i, \quad l_i(x) = \frac{u(x)}{(x - x_i)u'(x_i)}$. \diamond

Demonstrând teorema 3.3.6 am demonstrat de fapt existența și unicitatea soluției problemei generale de interpolare Lagrange:

(PGIL) Fiind date $b_0, b_1, \dots, b_m \in \mathbb{R}$, să se determine

$$p_m \in \mathbb{P}_m \text{ astfel încât } \forall i = 0, 1, \dots, m, \quad p_m(x_i) = b_i. \quad (3.3.9)$$

Joseph Louis Lagrange (1736-1813), protejat al lui Euler. Clai-raut scria despre tânărul Lagrange: „...un tânăr nu mai puțin remarcabil prin talent decât prin modestie; temperamentul său este blând și melancolic; nu cunoaște altă plăcere decât studiul.” Lagrange a avut contribuții fundamentale în calculul variațional, teoria numerelor și analiză matematică. Este cunoscut și pentru reprezentarea pe care a dat-o restului din formula lui Taylor. A dat formula de interpolare în 1794. Lucrarea sa *Mécanique Analytique*, publicată în 1788, l-a făcut unul din fondatorii mecanicii analitice.



Problema (3.3.9) conduce la un sistem liniar de $(m + 1)$ ecuații cu $(m + 1)$ necunoscute (coeficienții lui p_m).

Din teoria sistemelor liniare se știe că

$$\{\text{existența unei soluții } \forall b_0, b_1, \dots, b_m\} \Leftrightarrow \{\text{unicitatea soluției}\} \Leftrightarrow$$

$$\{(b_0 = b_1 = \dots = b_m = 0) \Rightarrow p_m \equiv 0\}$$

Punem $p_m = a_0 + a_1x + \dots + a_mx^m$

$$a = (a_0, a_1, \dots, a_m)^T, \quad b = (b_0, b_1, \dots, b_m)^T$$

și notăm cu $V = (v_{ij})$ matricea pătratică de ordin $m + 1$ cu elementele $v_{ij} = x_i^j$. Ecuația (3.3.9) se scrie sub forma

$$Va = b.$$

Matricea V este inversabilă (determinantul ei este Vandermonde); se arată ușor că $V^{-1} = U^T$, unde $U = (u_{ij})$ cu $\ell_i(x) = \sum_{k=0}^m u_{ik}x^k$; se obține în acest mod un procedeu puțin costisitor de inversare a matricei Vandermonde și prin urmare și de rezolvare a sistemului (3.3.9).

Exemplul 3.3.9. Polinomul de interpolare Lagrange corespunzător unei funcții f și nodurilor x_0 și x_1 este

$$(L_1f)(x) = \frac{x - x_1}{x_0 - x_1}f(x_0) + \frac{x - x_0}{x_1 - x_0}f(x_1),$$

adică dreapta care trece prin punctele $(x_0, f(x_0))$ și $(x_1, f(x_1))$. Analog, polinomul de interpolare Lagrange corespunzător unei funcții f și nodurilor x_0, x_1 și x_2 este

$$(L_2f)(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)}f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)}f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}f(x_2),$$

adică parabola care trece prin punctele $(x_0, f(x_0))$, $(x_1, f(x_1))$ și $(x_2, f(x_2))$. Interpretarea lor geometrică apare în figura 3.4. \diamond

3.3.3. Interpolare Hermite

În loc să facem să coincidă f și polinomul de interpolare în punctele x_i din $[a, b]$, am putea face ca f și polinomul de interpolare să coincidă împreună cu derivatele lor până la ordinul r_i în punctele x_i . Se obține:

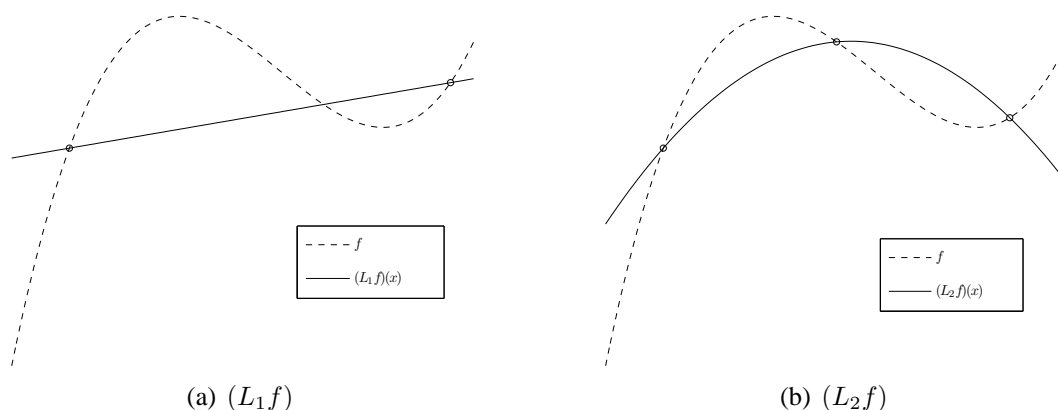


Figura 3.4: Interpretarea geometrică a lui L_1f (stânga) și L_2f

Teorema 3.3.10 Fiind date $(m+1)$ puncte distincte x_0, x_1, \dots, x_m din $[a, b]$ și $(m+1)$ numere naturale r_0, r_1, \dots, r_m , punem $n = m + r_0 + r_1 + \dots + r_m$. Atunci, fiind dată o funcție f , definită pe $[a, b]$ și admițând derivate de ordin r_i în punctele x_i există un singur polinom și numai unul $H_n f$ de grad $\leq n$ astfel încât

$$\forall (i, \ell), 0 \leq i \leq m, 0 \leq \ell \leq r_i \quad (H_n f)^{(\ell)}(x_i) = f^{(\ell)}(x_i), \quad (3.3.10)$$

unde $f^{(\ell)}(x_i)$ este derivata de ordinul ℓ a lui f în x_i .

Definiția 3.3.11 Polinomul definit în acest mod se numește polinom de interpolare al lui Hermite ⁸ al funcției f relativ la punctele x_0, x_1, \dots, x_m și la întregii r_0, r_1, \dots, r_m .

Demonstrație. Ecuația (3.3.10) conduce la un sistem liniar de $(n+1)$ ecuații cu $(n+1)$ necunoscute (coeficienții lui $H_n f$), deci este suficient să arătăm că sistemul omogen corespunzător admite doar soluția nulă, adică relațiile

$$H_n f \in \mathbb{P}_n \text{ și } \forall (i, \ell), 0 \leq i \leq k, 0 \leq \ell \leq r_i, (H_n f)^{(\ell)}(x_i) = 0$$



Charles Hermite (1822-1901) matematician francez de frunte, membru al Academiei Franceze, cunoscut pentru lucrările sale în domeniul teoriei numerelor, algebră și analiză. A devenit faimos după ce a dat, în 1873, demonstrația transcendenței numărului e .

ne asigură că pentru orice $i = 0, 1, \dots, m$ x_i este rădăcină de ordinul $r_i + 1$ a lui $H_n f$; prin urmare $H_n f$ are forma

$$(H_n f)(x) = q(x) \prod_{i=0}^m (x - x_i)^{r_i+1},$$

unde q este un polinom. Cum $\sum_{i=0}^m (r_i + 1) = n + 1$, acest lucru nu este compatibil cu apartenența lui H_n la \mathbb{P}_n , decât dacă $q \equiv 0$ și deci $H_n \equiv 0$. \square

Observația 3.3.12. 1) Dându-se numerele reale $b_{i\ell}$ pentru orice pereche (i, ℓ) astfel încât $0 \leq i \leq k$ și $0 \leq \ell \leq r_i$, am arătat că problema generală de interpolare Hermite

$$\begin{aligned} &\text{să se determine } p_n \in \mathbb{P}_n \text{ a.î. } \forall (i, \ell) \text{ cu } 0 \leq i \leq m \text{ și} \\ &0 \leq \ell \leq r_i, p_n^{(\ell)}(x_i) = b_{i\ell} \end{aligned} \quad (3.3.11)$$

admite o soluție și numai una. În particular, dacă alegem pentru o pereche (i, ℓ) dată $b_{i\ell} = 1$ și $b_{jn} = 0$, $\forall (j, m) \neq (i, \ell)$ se obține un polinom de bază (fundamental) de interpolare Hermite relativ la punctele x_0, x_1, \dots, x_m și la întregii r_0, r_1, \dots, r_m . Polinomul de interpolare Hermite definit prin (3.3.10) se obține cu ajutorul polinoamelor de bază (fundamentale) cu formula

$$(H_n f)(x) = \sum_{i=0}^m \sum_{\ell=0}^{r_i} f^{(\ell)}(x) h_{i\ell}(x). \quad (3.3.12)$$

Punând

$$q_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^k \left(\frac{x - x_j}{x_i - x_j} \right)^{r_j+1},$$

se verifică că polinoamele de bază $h_{i\ell}$ sunt definite prin relațiile de recurență

$$h_{i r_i}(x) = \frac{(x - x_i)^{r_i}}{r_i!} q_i(x)$$

și pentru $\ell = r_{i-1}, r_{i-2}, \dots, 1, 0$

$$h_{i\ell}(x) = \frac{(x - x_i)^\ell}{\ell!} q_i(x) - \sum_{j=l+1}^{r_i} \binom{j}{\ell} q_i^{(j-\ell)}(x_i) h_{ij}(x).$$

2) Matricea V asociată sistemului liniar (3.3.11) se numește matrice Vandermonde generalizată; ea este inversabilă, iar elementele matricei ei inverse sunt coeficienții polinoamelor $h_{i\ell}$.

- 3) Interpolarea Lagrange este un caz particular al interpolării Hermite (pentru $r_i = 0$, $i = 0, 1, \dots, m$); polinomul Taylor este un caz particular pentru $m = 0$ și $r_0 = n$. \diamond

Vom prezenta o expresie mai convenabilă a polinoamelor fundamentale Hermite, obținută de Dimitrie D. Stancu în 1957. Ele verifică relațiile

$$\begin{aligned} h_{kj}^{(p)}(x_\nu) &= 0, \quad \nu \neq k, \quad p = \overline{0, r_\nu} \\ h_{kj}^{(p)}(x_k) &= \delta_{jp}, \quad p = \overline{0, r_k} \end{aligned} \quad (3.3.13)$$

pentru $j = \overline{0, r_k}$ și $\nu, k = \overline{0, m}$. Introducând notațiile

$$u(x) = \prod_{k=0}^m (x - x_k)^{r_k+1}$$

și

$$u_k(x) = \frac{u(x)}{(x - x_k)^{r_k+1}},$$

din (3.3.13) rezultă că h_{kj} are forma

$$h_{kj}(x) = u_k(x)(x - x_k)^j g_{kj}(x), \quad g_{kj} \in \mathbb{P}_{r_k-j}. \quad (3.3.14)$$

Dezvoltând g_{kj} cu formula lui Taylor, avem

$$g_{kj}(x) = \sum_{\nu=0}^{r_k-j} \frac{(x - x_k)^\nu}{\nu!} g_{kj}^{(\nu)}(x_k); \quad (3.3.15)$$

mai rămân de determinat valorile lui $g_{kj}^{(\nu)}(x_k)$, $\nu = \overline{0, r_k - j}$. Scriind (3.3.14) sub forma

$$(x - x_k)^j g_{kj}(x) = h_{kj}(x) \frac{1}{u_k(x)},$$

și aplicând formula lui Leibniz pentru derivata de ordinul $j + \nu$ a produsului se obține

$$\sum_{s=0}^{j+\nu} \binom{j+\nu}{s} [(x - x_k)^j]^{(j+\nu-s)} g_{kj}^{(s)}(x) = \sum_{s=0}^{j+\nu} \binom{j+\nu}{s} h_{kj}^{(j+\nu-s)}(x) \left[\frac{1}{u_k(x)} \right]^{(s)}.$$

Lând $x = x_k$, toți termenii din ambii membri se vor anula, cu excepția celor corespunzătorii lui $s = \nu$. Avem deci

$$\binom{j+\nu}{\nu} j! g_{kj}^{(\nu)}(x_k) = \binom{j+\nu}{\nu} \left[\frac{1}{u_k(x)} \right]_{x=x_k}^{(\nu)}, \quad \nu = \overline{0, r_k - j}.$$

Am obținut

$$g_{kj}^{(\nu)}(x_k) = \frac{1}{j!} \left[\frac{1}{u_k(x)} \right]_{x=x_k}^{(\nu)},$$

iar din (3.3.15) și (3.3.14) avem în final

$$h_{kj}(x) = \frac{(x - x_k)^j}{j!} u_k(x) \sum_{\nu=0}^{r_k-j} \frac{(x - x_k)^\nu}{\nu!} \left[\frac{1}{u_k(x)} \right]_{x=x_k}^{(\nu)}.$$

Propoziția 3.3.13 *Operatorul H_n este proiector, adică*

- *este liniar* ($H_n(\alpha f + \beta g) = \alpha H_n f + \beta H_n g$);
- *este idempotent* ($H_n \circ H_n = H_n$).

Demonstrație. Liniaritatea rezultă imediat din formula (3.3.12). Datorită unicității polinomului de interpolare Hermite $H_n(H_n f) - H_n f$ este identic nul, deci $H_n(H_n f) = H_n f$ și am arătat idempotența. \square

Exemplul 3.3.14. Polinomul de interpolare Hermite corespunzător unei funcții f și nodurilor duble 0 și 1 are expresia

$$(H_3 f)(x) = h_{00}(x)f(0) + h_{10}(x)f(1) + h_{01}(x)f'(0) + h_{11}(x)f'(1),$$

unde

$$\begin{aligned} h_{00}(x) &= (x-1)^2(2x+1), \\ h_{01}(x) &= x(x-1)^2, \\ h_{10}(x) &= x^2(3-2x), \\ h_{11}(x) &= x^2(x-1). \end{aligned}$$

Dacă se adaugă nodul $x = \frac{1}{2}$, calitatea aproximării crește (vezi figura 3.5). \diamond

3.3.4. Expresia erorii de interpolare

Reamintim că norma unui operator liniar P_n se poate defini prin

$$\|P_n\| = \max_{f \in C[a,b]} \frac{\|P_n f\|}{\|f\|}, \quad (3.3.16)$$

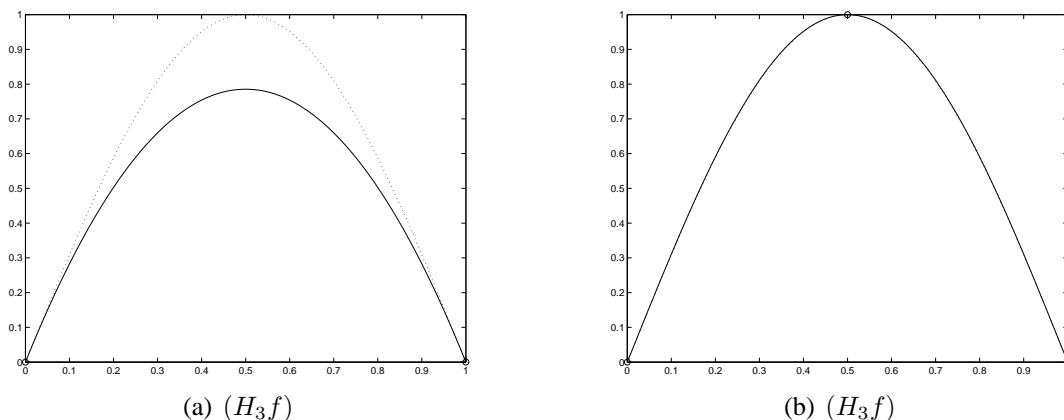


Figura 3.5: Polinoamele de interpolare Hermite (H_3f) (—) corespunzător funcției $f : [0, 1] \rightarrow \mathbb{R}$, $f(x) = \sin \pi x$ și nodurilor duble $x_0 = 0$ și $x_1 = 1$ (\cdots)(stânga) și (H_5f) (—) corespunzător funcției $f : [0, 1] \rightarrow \mathbb{R}$, $f(x) = \sin \pi x$ (\cdots) și nodurilor duble $x_0 = 0$, $x_1 = \frac{1}{2}$ și $x_2 = 1$

unde în membrul drept se ia o normă convenabilă pentru funcții. Luând norma L^∞ , din formula lui Lagrange se obține

$$\begin{aligned} \|(L_m f)(\cdot)\|_\infty &= \max_{a \leq x \leq b} \left| \sum_{i=0}^m f(x_i) \ell_i(x) \right| \\ &\leq \|f\|_\infty \max_{a \leq x \leq b} \sum_{i=0}^m |\ell_i(x)|. \end{aligned} \quad (3.3.17)$$

Fie $\|\lambda_m\|_\infty = \lambda_m(x_\infty)$. Egalitatea are loc pentru o funcție $\varphi \in C[a, b]$, liniară pe porțiuni și care verifică $\varphi(x_i) = \text{sgn} \ell_i(x_\infty)$, $i = \overline{0, m}$. Deci,

$$\|P_n\|_\infty = \Lambda_m, \quad (3.3.18)$$

unde

$$\Lambda_m = \|\lambda_m\|_\infty, \quad \lambda_m(x) = \sum_{i=0}^m |\ell_i(x)|. \quad (3.3.19)$$

Funcția $\lambda_m(x)$ și maximul său Λ_m se numesc *funcția lui Lebesgue*⁹ și respectiv *constanta lui Lebesgue* pentru interpolarea Lagrange. Ele furnizează o primă estimare a erorii de interpolare: fie $\mathcal{E}_m(f)$ eroarea în cea mai bună aproximare a lui f prin polinoame de grad $\leq m$,

$$\mathcal{E}_m(f) = \min_{p \in \mathbb{P}_m} \|f - p\|_\infty = \|f - \hat{p}_m\|_\infty, \quad (3.3.20)$$

unde \hat{p}_m este polinomul de grad m de cea mai bună aproximare a lui f . Utilizând faptul că operatorul L_m este proiector și formulele (3.3.17) și (3.3.19), se găsește

$$\begin{aligned} \|f - L_m f\| &= \|f - \hat{p}_m - L_m(f - \hat{p}_m)\|_\infty \\ &\leq \|f - \hat{p}_m\|_\infty + \Lambda_m \|f - \hat{p}_m\|_\infty; \end{aligned}$$

adică,

$$\|f - L_m f\|_\infty \leq (1 + \Lambda_m) \mathcal{E}_m(f). \quad (3.3.21)$$

Astfel, cu cât f poate fi aproximată mai bine prin polinoame de grad $\leq m$, cu atât este mai mică eroarea de interpolare. Din păcate, Λ_m nu este uniform mărginită: indiferent de cum se aleg nodurile $x_i = x_i^{(m)}$, $i = \overline{0, m}$ se poate arăta că $\Lambda_m > O(\log m)$ când $m \rightarrow \infty$. Totuși, nu este posibil să tragem, pe baza teoremei de aproximare a lui Weierstrass (adică din $\mathcal{E}_m \rightarrow 0$, $m \rightarrow \infty$), concluzia că interpolarea Lagrange converge uniform pentru orice funcție f , nici chiar pentru noduri judicios alese; se știe că, de fapt convergența nu are loc.

Dacă dorim să utilizăm polinomul de interpolare Lagrange sau Hermite pentru a aproxima funcția f într-un punct $x \in [a, b]$, distinct de nodurile de interpolare (x_0, \dots, x_m) , trebuie să estimăm eroarea comisă $(R_n f)(x) = f(x) - (H_n f)(x)$. Dacă nu posedăm nici o informație referitoare la f în afara punctelor x_i , este clar că nu putem spune nimic despre $(R_n f)(x)$; într-adevăr este posibil să schimbăm f în afara punctelor x_i fără a modifica $(H_n f)(x)$. Trebuie deci să facem ipoteze suplimentare, care vor fi ipoteze de regularitate asupra lui f . Să notăm cu $C^m[a, b]$ spațiul funcțiilor reale de m ori continuu diferențiabile pe $[a, b]$. Avem următoarea teoremă referitoare la estimarea erorii în interpolarea Hermite.

9



Henry Lebesgue (1875-1941), matematician francez, cunoscut pentru lucrările sale fundamentale în domeniul teoriei funcțiilor reale, și în special pentru introducerea măsurii și integralei care îi poartă numele.

Teorema 3.3.15 *Presupunem că $f \in C^n[\alpha, \beta]$ și există $f^{(n+1)}$ pe (α, β) , unde $\alpha = \min\{x, x_0, \dots, x_m\}$ și $\beta = \max\{x, x_0, \dots, x_m\}$; atunci, pentru orice $x \in [\alpha, \beta]$, există un $\xi_x \in (\alpha, \beta)$ astfel încât*

$$(R_n f)(x) = \frac{1}{(n+1)!} u_n(x) f^{(n+1)}(\xi_x), \quad (3.3.22)$$

unde

$$u_n(x) = \prod_{i=0}^m (x - x_i)^{r_i+1}.$$

Demonstrație. Dacă $x = x_i$, $(R_n f)(x) = 0$ și (3.3.22) se verifică trivial. Presupunem că x este distinct de x_i și considerăm, pentru x fixat, funcția auxiliară

$$F(z) = \begin{vmatrix} u_n(z) & (R_n f)(z) \\ u_n(x) & (R_n f)(x) \end{vmatrix}.$$

Se observă că $F \in C^n[\alpha, \beta]$, $\exists F^{(n+1)}$ pe (α, β) , $F(x) = 0$ și $F^{(j)}(x_k) = 0$ pentru $k = \overline{0, m}$, $j = \overline{0, r_k}$. Deci, F are $(n+2)$ zerouri, luând în considerare și ordinele de multiplicitate. Aplicând succesiv teorema lui Rolle generalizată, rezultă că există cel puțin un $\xi \in (\alpha, \beta)$ astfel încât $F^{(n+1)}(\xi) = 0$, adică

$$F^{(n+1)}(\xi) = \begin{vmatrix} (n+1)! & f^{(n+1)}(\xi) \\ u_n(x) & (R_n f)(x) \end{vmatrix} = 0, \quad (3.3.23)$$

unde s-a ținut cont că $(R_n f)^{(n+1)} = f^{(n+1)} - (H_n f)^{(n+1)} = f^{(n+1)}$. Exprimând $(R_n f)(x)$ din (3.3.23) se obține (3.3.22). \square

Corolarul 3.3.16 *Punem $M_{n+1} = \max_{x \in [a, b]} |f^{(n+1)}(x)|$; o margine superioară a erorii de interpolare $(R_n f)(x) = f(x) - (H_n f)(x)$ este dată prin*

$$|(R_n f)(x)| \leq \frac{M_{n+1}}{(n+1)!} |u_n(x)|.$$

Deoarece H_n este proiector, rezultă că R_n este de asemenea proiector; în plus $\text{Ker } R_n = \mathbb{P}_n$, deoarece $R_n f = f - H_n f = f - f = 0$, $\forall f \in \mathbb{P}_n$. Deci, putem aplica lui R_n teorema lui Peano.

Teorema 3.3.17 *Dacă $f \in C^{n+1}[a, b]$, atunci*

$$(R_n f)(x) = \int_a^b K_n(x; t) f^{(n+1)}(t) dt, \quad (3.3.24)$$

unde

$$K_n(x; t) = \frac{1}{n!} \left\{ (x-t)_+^n - \sum_{k=0}^m \sum_{j=0}^{r_k} h_{kj}(x) [(x_k - t)_+^n]^{(j)} \right\}. \quad (3.3.25)$$

Demonstrație. Aplicând teorema lui Peano, avem

$$(R_n f)(x) = \int_a^b K_n(x; t) f^{(n+1)}(t) dt$$

și ținând cont că

$$K_n(x; t) = R_n \left[\frac{(x-t)_+^n}{n!} \right] = \frac{(x-t)_+^n}{n!} - H_n \left[\frac{(x-t)_+^n}{n!} \right],$$

teorema rezultă imediat. \square

Deoarece interpolarea Lagrange este un caz particular al interpolării Hermite pentru $r_i = 0, i = 0, 1, \dots, m$ din teorema 3.3.15 se obține:

Corolarul 3.3.18 *Presupunem că $f \in C^m[\alpha, \beta]$ și există $f^{(m+1)}$ pe (α, β) , unde $\alpha = \min\{x, x_0, \dots, x_m\}$ și $\beta = \max\{x, x_0, \dots, x_m\}$; atunci, pentru orice $x \in [\alpha, \beta]$, există un $\xi_x \in (\alpha, \beta)$ astfel încât*

$$(R_m f)(x) = \frac{1}{(m+1)!} u_m(x) f^{(m+1)}(\xi_x), \quad (3.3.26)$$

unde

$$u_m(x) = \prod_{i=0}^m (x - x_i).$$

De asemenea, din teorema 3.3.17 avem:

Corolarul 3.3.19 *Dacă $f \in C^{m+1}[a, b]$, atunci*

$$(R_m f)(x) = \int_a^b K_m(x; t) f^{(m+1)}(t) dt \quad (3.3.27)$$

unde

$$K_m(x; t) = \frac{1}{m!} \left[(x-t)_+^m - \sum_{k=0}^m \ell_k(x) (x_k - t)_+^m \right]. \quad (3.3.28)$$

Exemplul 3.3.20. Pentru polinoamele de interpolare din exemplul 3.3.9 resturile corespunzătoare sunt

$$(R_1 f)(x) = \frac{(x-x_0)(x-x_1)}{2} f''(\xi)$$

și respectiv

$$(R_2 f)(x) = \frac{(x-x_0)(x-x_1)(x-x_2)}{6} f'''(\xi). \quad \diamond$$

Exemplul 3.3.21. Restul din formula de interpolare Hermite cu nodurile duble 0 și 1 pentru $f \in C^4[\alpha, \beta]$ este

$$(R_3 f)(x) = \frac{x^2(x-1)^2}{6!} f^{(4)}(\xi). \quad \diamond$$

Exemplul 3.3.22. Luăm $f(x) = e^x$. Avem pentru $x \in [a, b]$, $M_{n+1} = e^b$ și oricum am alege punctele x_i , $|u_n(x)| \leq (b-a)^{n+1}$, de unde

$$\max_{x \in [a, b]} |(R_n f)(x)| \leq \frac{(b-a)^{n+1}}{(n+1)!} e^b.$$

Se deduce că

$$\lim_{n \rightarrow \infty} \left\{ \max_{x \in [a, b]} |(R_n f)(x)| \right\} = \lim_{n \rightarrow \infty} \|R_n f\| = 0,$$

adică $H_n f$ converge uniform către f pe $[a, b]$ când n tinde la ∞ . De fapt se poate demonstra un rezultat analog pentru orice funcție dezvoltabilă în serie întreagă în jurul punctului $x = \frac{a+b}{2}$ cu raza de convergență $r > \frac{3}{2}(b-a)$. \diamond

3.3.5. Convergența interpolării Lagrange

Să definim ce înțelegem prin convergență. Presupunem că se dă un tablou triunghiular de noduri de interpolare $x_i = x_i^{(m)}$, având exact $m+1$ noduri distincte pentru orice $m = 0, 1, 2, \dots$

$$\begin{array}{cccc} x_0^{(0)} & & & \\ x_0^{(1)} & x_1^{(1)} & & \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \\ \vdots & \vdots & \vdots & \ddots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_m^{(m)} \\ \vdots & \vdots & \vdots & & \vdots \end{array} \quad (3.3.29)$$

Presupunem că toate nodurile sunt conținute într-un interval finit $[a, b]$. Atunci pentru orice m definim

$$P_m(x) = L_m(f; x_0^{(m)}, x_1^{(m)}, \dots, x_m^{(m)}; x), \quad x \in [a, b]. \quad (3.3.30)$$

Spunem că interpolarea Lagrange bazată pe tabelul de noduri (3.3.29) converge dacă

$$p_m(x) \rightrightarrows f(x), \quad \text{când } n \rightarrow \infty \text{ pe } [a, b]. \quad (3.3.31)$$

Convergența depinde evident de comportarea derivatei de ordinul k a lui f , $f^{(k)}$, când $k \rightarrow \infty$. Presupunem că $f \in C^\infty[a, b]$ și că

$$|f^{(k)}(x)| \leq M_k \text{ pentru } a \leq x \leq b, k = 0, 1, 2, \dots \quad (3.3.32)$$

Deoarece $|x_i - x_i^{(m)}| \leq b - a$ când $x \in [a, b]$ și $x_i^{(n)} \in [a, b]$ avem

$$\left| (x - x_0^{(m)}) \dots (x - x_m^{(m)}) \right| < (b - a)^{m+1}, \quad (3.3.33)$$

deci

$$|f(x) - (L_m f)(x)| \leq (b - a)^{m+1} \frac{M_{m+1}}{(m + 1)!}, \quad x \in [a, b]. \quad (3.3.34)$$

Deci avem convergență dacă

$$\lim_{k \rightarrow \infty} \frac{(b - a)^k}{k!} M_k = 0. \quad (3.3.35)$$

Să arătăm că (3.3.35) este adevărată dacă f este analitică într-o vecinătate suficient de mare din \mathbb{C} ce conține intervalul $[a, b]$. Mai concret fie C_r discul circular (închis) cu centrul în mijlocul intervalului $[a, b]$ și de rază r și presupunem că $r > \frac{1}{2}(b - a)$, astfel că $[a, b] \subset C_r$. Presupunem că f este analitică în C_r . Atunci putem estima derivata în (3.3.32) cu formula lui Cauchy

$$f^{(k)}(x) = \frac{k!}{2\pi i} \int_{C_r} \frac{f(z)}{(z - x)^{k+1}} dz, \quad x \in [a, b]. \quad (3.3.36)$$

Observând că $|z - x| \geq r - \frac{1}{2}(b - a)$ (vezi figura 3.6) obținem

$$|f^{(k)}(x)| \leq \frac{k!}{2\pi} \frac{\max_{z \in \partial C_r} |f(z)|}{\left[r - \frac{1}{2}(b - a)\right]^{k+1}} \cdot 2\pi r;$$

putem lua pentru M_k în (3.3.32)

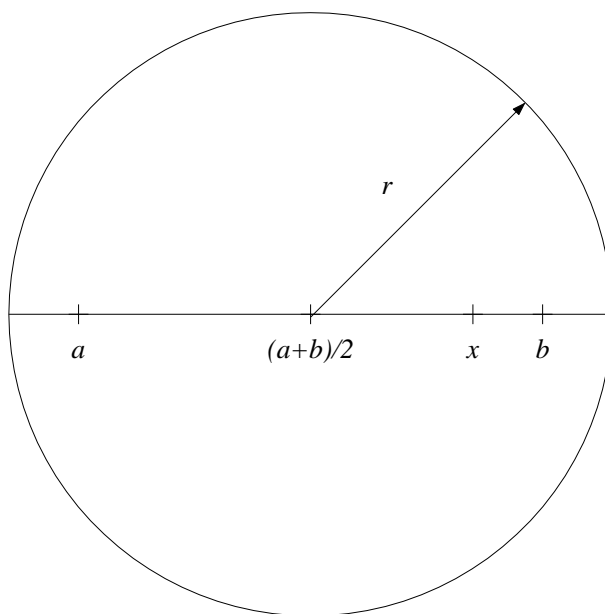
$$M_k = \frac{r}{r - \frac{1}{2}(b - a)} \max_{z \in \partial C_r} |f(z)| \frac{k!}{\left[r - \frac{1}{2}(b - a)\right]^k} \quad (3.3.37)$$

și (3.3.35) are loc dacă

$$\left(\frac{b - a}{r - \frac{1}{2}(b - a)} \right)^k \rightarrow 0 \text{ când } k \rightarrow \infty,$$

adică, dacă $b - a < r - \frac{1}{2}(b - a)$ sau echivalent

$$r > \frac{3}{2}(b - a). \quad (3.3.38)$$

Figura 3.6: Discul circular C_r

Am arătat că interpolarea Lagrange converge (uniform pe $[a, b]$) pentru o mulțime arbitrară de noduri (3.3.29) (toate conținute în $[a, b]$) dacă f este analitică în discul circular C_r centrat în $(a + b)/2$ și având raza suficient de mare astfel ca (3.3.29) să aibă loc.

Deoarece acest rezultat utilizează o estimare grosieră (în particular (3.3.33), domeniul cerut de analiticitate pentru f nu este prea îngust. Utilizând metode mai rafinate, se poate arăta următorul lucru. Fie $d\mu(t)$ distribuția limită a nodurilor de interpolare, adică

$$\int_a^x d\mu(t), \quad a < x \leq b,$$

raportul dintre numărul de noduri $x_i^{(m)}$ din $[a, x]$ și numărul total de noduri, asimptotic când $n \rightarrow \infty$. (Când nodurile sunt uniform distribuite pe intervalul $[a, b]$, atunci $d\mu(t) = \frac{dt}{b-a}$). O curbă cu potențial logaritm constant este locul geometric al punctelor $z \in \mathbb{C}$ cu proprietatea

$$u(z) = \int_a^b \ln \frac{1}{|z - t|} d\mu(t) = \gamma,$$

unde γ este o constantă. Pentru γ negativ foarte mare, aceste curbe arată ca niște cercuri cu raza foarte mare și centrul în $(a + b)/2$. Pe măsură ce γ crește, curbele se „comprimă”

spre intervalul $[a, b]$. Domeniul important (ce înlocuiește C_r) este domeniul

$$C_\Gamma = \{z \in \mathbb{C} : u(z) \geq \Gamma\},$$

în sensul că dacă f este analitică în orice domeniu C ce conține C_r în interiorul său (nu contează cât de strâns C acoperă pe C_r), atunci

$$|f(z) - (L_m f)(z)| \rightarrow 0, \text{ când } n \rightarrow \infty, \quad (3.3.39)$$

uniform pentru $z \in C_\Gamma$.

Exemplul 3.3.23. Noduri echidistante: $d\mu(t) = dt/(b-a)$, $a \leq t \leq b$. În acest caz C_Γ este un domeniu în formă de lentilă, așa cum se arată în figura 3.7. Astfel, avem convergența uniformă în C_Γ (nu doar pe $[a, b]$ ca mai sus) cu condiția ca f să fie analitică într-o regiune puțin mai largă decât C_Γ . Fie $\Gamma = \sup \gamma$, unde supremumul se ia pentru toate curbele $u(z) = \gamma$ ce conțin pe $[a, b]$ în interior.

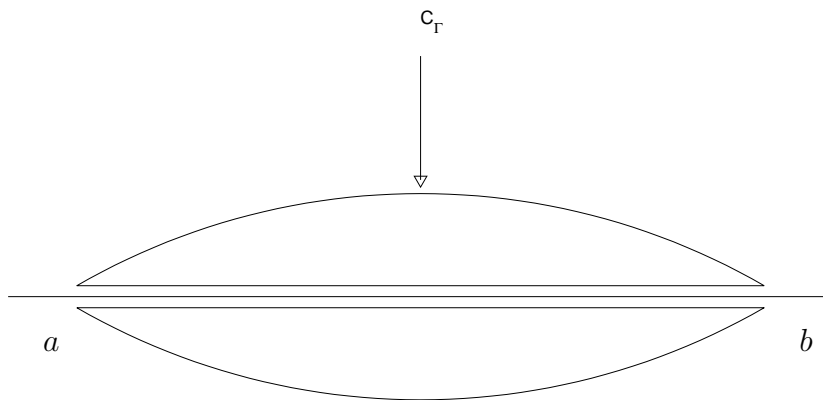


Figura 3.7: Domeniul C_Γ pentru noduri uniform distribuite

Exemplul 3.3.24. Distribuția arcsin pe $[-1, 1]$:

$$d\mu(t) = \frac{1}{\pi} \frac{dt}{\sqrt{1-t^2}}.$$

Nodurile sunt în acest caz rădăcinile polinomului Cebîșev de speța I. Ele sunt mai dens distribuite în apropierea capetelor intervalului $[-1, 1]$. În acest caz $C_\Gamma = [-1, 1]$, așa că interpolarea Lagrange converge uniform pe $[-1, 1]$ dacă f este analitică pe $[-1, 1]$, adică în orice regiune al cărei interior conține intervalul $[-1, 1]$. \diamond

Care este semnificația polinoamelor Cebîșev pentru interpolare?

Reamintim că eroarea de interpolare (pe $[-1, 1]$ pentru o funcție de clasă $C^{m+1}[-1, 1]$) este dată de

$$f(x) - (L_m f)(x) = \frac{f^{(m+1)}(\xi(x))}{(m+1)!} \prod_{i=0}^m (x - x_i). \quad (3.3.40)$$

Primul factor este independent de alegerea nodurilor x_i . Punctul intermediar $\xi(x)$ depinde de x_i , dar de obicei majorăm $|f^{(m+1)}|$ prin $\|f^{(m+1)}\|_\infty$, ceea ce înlătură această dependență. Pe de altă parte, produsul din al doilea factor, inclusiv norma sa

$$\left\| \prod_{i=0}^m (\cdot - x_i) \right\|_\infty, \quad (3.3.41)$$

depinde puternic de x_i . Are sens, deci, să încercăm să minimizăm (3.3.41) după toți $x_i \in [-1, 1]$. Deoarece produsul din (3.3.41) este un polinom monic de grad $m+1$, din teorema 3.2.1 rezultă că nodurile optimale $x_i = \widehat{x}_i^{(m)}$ din (3.3.40) sunt rădăcinile lui T_{m+1} , adică

$$\widehat{x}_i^{(m)} = \cos \frac{2i+1}{2m+2} \pi, \quad i = \overline{0, m}. \quad (3.3.42)$$

Pentru aceste noduri, avem conform lui (3.2.20)

$$\|f(\cdot) - (L_m f)(\cdot)\|_\infty \leq \frac{\|f^{(m+1)}\|_\infty}{(m+1)!} \cdot \frac{1}{2^m}. \quad (3.3.43)$$

Se cuvine a compara acest factor cu marginea mai brută dată în (3.3.34) care în acest caz pe intervalul $[-1, 1]$ este $2^{m+1}/(m+1)!$.

Deoarece conform (3.3.40) curba de eroare $y = f - L_m f$ pentru punctele Cebîșev (3.3.42) este în esență *echilibrată* (modulo variația factorului $f^{(m+1)}$) și astfel liberă de oscilațiile violente pe care le-am văzut pentru puncte echidistante ne vom aștepta la proprietăți mai favorabile pentru tablouri triunghiulare (3.3.29) formate din noduri Cebîșev. Se poate arăta că dacă $f \in C^1[-1, 1]$, atunci

$$(L_m f) \left(x; \widehat{x}_0^{(m)}, \widehat{x}_1^{(m)}, \dots, \widehat{x}_m^{(m)} \right) \rightrightarrows f(x), \quad n \rightarrow \infty, \quad (3.3.44)$$

pe $[-1, 1]$. Astfel, nu avem nevoie de analiticitatea lui f pentru ca (3.3.44) să aibă loc.

Exemplul 3.3.25 (Exemplul lui Runge). Considerăm funcția

$$f(x) = \frac{1}{1+x^2}, \quad x \in [-5, 5]$$

și nodurile

$$x_k^{(m)} = -5 + 10 \frac{k}{m}, \quad k = \overline{0, m}. \quad (3.3.45)$$

Nodurile sunt echidistante pe $[-5, 5]$, deci asimptotic uniform distribuite. Observăm că f are doi poli în $z = \pm i$. Acești poli sunt așezați în interiorul regiunii C_Γ din figura 3.7 pentru intervalul $[-5, 5]$, deci f nu este analitică în C_Γ . Din acest motiv nu ne așteptăm să avem convergență pe întreg intervalul $[-5, 5]$. Se poate demonstra că

$$\lim_{m \rightarrow \infty} |f(x) - p_m(f; x)| = \begin{cases} 0 & \text{dacă } |x| < 3.633 \dots \\ \infty & \text{dacă } |x| > 3.633 \dots \end{cases} \quad (3.3.46)$$

Având în minte figura 3.7 acest rezultat nu este surprinzător. Graficul pentru $m = 10, 13, 16$ apare în figura 3.8. \diamond

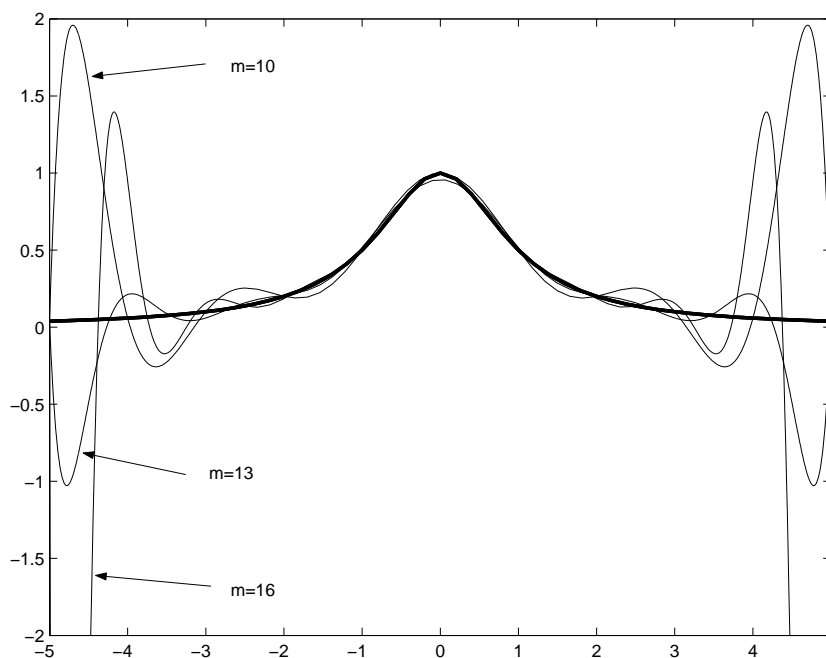


Figura 3.8: O ilustrare grafică a contraexemplului lui Runge

Exemplul 3.3.26 (Exemplul lui Bernstein). Luăm funcția

$$f(x) = |x|, \quad x \in [-1, 1]$$

și nodurile echidistante.

$$x_k^{(m)} = -1 + \frac{2k}{m}, \quad k = 0, 1, 2, \dots, m. \quad (3.3.47)$$

Problema analiticității nu se pune, deoarece f nu este derivabilă în $x = 0$. Se obține că

$$\lim_{m \rightarrow \infty} |f(x) - L_m(f; x)| = \infty \quad \forall x \in [-1, 1]$$

exceptând punctele $x = -1$, $x = 0$ și $x = 1$. Vezi figura 3.9(a), pentru $m = 20$. Convergența în $x = \pm 1$ este trivială deoarece acestea sunt noduri de interpolare și deci eroarea în aceste puncte este 0. Același lucru este adevărat pentru $x = 0$, când n este impar, dar nu și când n este par. Eșecul convergenței pentru aceste noduri se explică doar parțial prin insuficiența regularității a lui f . Un alt motiv este distribuția uniformă a nodurilor. Există exemple mai bune de distribuții ale nodurilor, cum ar fi distribuția *arcsin* din exemplul 3.3.24. În figura 3.9(b) se dă graficul pentru $m = 17$. \diamond

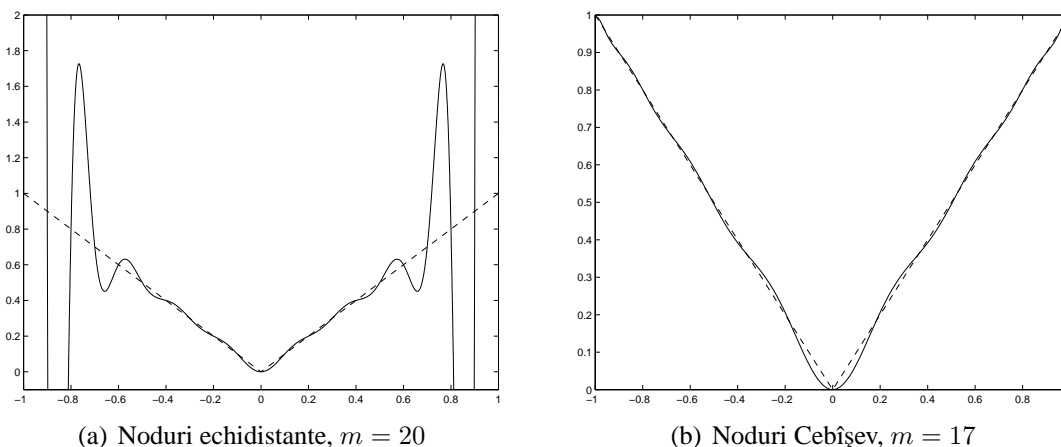


Figura 3.9: Comportarea interpolării Lagrange pentru $f : [-1, 1] \rightarrow \mathbb{R}$, $f(x) = |x|$.

3.4. Calculul eficient al polinoamelor de interpolare

3.4.1. Metode de tip Aitken

În multe situații gradul necesar pentru a atinge precizia dorită în interpolarea polinomială este necunoscut. El se poate determina din expresia restului, dar pentru aceasta este necesar să cunoaștem $\|f^{(m+1)}\|_\infty$. Vom nota cu P_{m_1, m_2, \dots, m_k} polinomul de interpolare Lagrange având nodurile x_{m_1}, \dots, x_{m_k} .

Propoziția 3.4.1 Dacă f este definită în x_0, \dots, x_k , $x_j \neq x_i$, $0 \leq i, j \leq k$, atunci

$$\begin{aligned} P_{0,1,\dots,k} &= \frac{(x-x_j)P_{0,1,\dots,j-1,j+1,\dots,k}(x) - (x-x_i)P_{0,1,\dots,i-1,i+1,\dots,k}(x)}{x_i - x_j} = \\ &= \frac{1}{x_i - x_j} \begin{vmatrix} x - x_j & P_{0,1,\dots,i-1,i+1,\dots,k}(x) \\ x - x_i & P_{0,1,\dots,j-1,j+1,\dots,k}(x) \end{vmatrix} \end{aligned} \quad (3.4.1)$$

Demonstrație. $Q = P_{0,1,\dots,i-1,i+1,\dots,k}$, $\widehat{Q} = P_{0,1,\dots,j-1,j+1,k}$

$$P(x) = \frac{(x-x_j)\widehat{Q}(x) - (x-x_i)Q(x)}{x_i - x_j}$$

$$P(x_r) = \frac{(x_r - x_j)\widehat{Q}(x_r) - (x_r - x_i)Q(x_r)}{x_i - x_j} = \frac{x_i - x_j}{x_i - x_j} f(x_r) = f(x_r)$$

pentru $r \neq i \wedge r \neq j$, căci $Q(x_r) = \widehat{Q}(x_r) = f(x_r)$. Dar

$$P(x_i) = \frac{(x_i - x_j)\widehat{Q}(x_i) - (x_i - x_j)Q(x_i)}{x_i - x_j} = f(x_i)$$

și

$$P(x_j) = \frac{(x_j - x_i)\widehat{Q}(x_j) - (x_j - x_i)Q(x_j)}{x_i - x_j} = f(x_j),$$

deci $P = P_{0,1,\dots,k}$. \square

În acest mod am stabilit o relație de recurență între un polinom de interpolare Lagrange de gradul k și două polinoame de interpolare Lagrange de gradul $k-1$. Calculele pot fi așezate în formă tabelară

$$\begin{array}{cccccc} x_0 & P_0 & & & & \\ x_1 & P_1 & P_{0,1} & & & \\ x_2 & P_2 & P_{1,2} & P_{0,1,2} & & \\ x_3 & P_3 & P_{2,3} & P_{1,2,3} & P_{0,1,2,3} & \\ x_4 & P_4 & P_{3,4} & P_{2,3,4} & P_{1,2,3,4} & P_{0,1,2,3,4} \end{array}$$

Să presupunem că în acest moment $P_{0,1,2,3,4}$ nu ne asigură precizia dorită. Se poate selecta un nou nod și adăuga o nouă linie tabelii

$$x_5 \quad P_5 \quad P_{4,5} \quad P_{3,4,5} \quad P_{2,3,4,5} \quad P_{1,2,3,4,5} \quad P_{0,1,2,3,4,5}$$

iar elementele vecine de pe linie, coloană sau diagonală se pot compara pentru a vedea dacă s-a obținut precizia dorită.

Metoda de mai sus se numește *metoda lui Neville*.

Notățiile pot fi simplificate

$$Q_{i,j} := P_{i-j, i-j+1, \dots, i-1, i},$$

$$Q_{i,j-1} = P_{i-j+1, \dots, i-1, i},$$

$$Q_{i-1, j-1} := P_{i-j, i-j+1, \dots, i-1}.$$

Din (3.4.1) rezultă

$$Q_{i,j} = \frac{(x - x_{i-j})Q_{i,j-1} - (x - x_i)Q_{i-1, j-1}}{x_i - x_{i-j}},$$

pentru $j = 1, 2, 3, \dots, i = j + 1, j + 2, \dots$

În plus, $Q_{i,0} = f(x_i)$. Obținem tabelul

$$\begin{array}{cccccc} x_0 & Q_{0,0} & & & & \\ x_1 & Q_{1,0} & Q_{1,1} & & & \\ x_2 & Q_{2,0} & Q_{2,1} & Q_{2,2} & & \\ x_3 & Q_{3,0} & Q_{3,1} & Q_{3,2} & Q_{3,3} & \end{array}$$

Dacă procedeul de interpolare converge, atunci șirul $Q_{i,i}$ converge și el și s-ar putea lua drept criteriu de oprire

$$|Q_{i,i} - Q_{i-1, i-1}| < \varepsilon.$$

Pentru a rapidiza algoritmul nodurile se vor ordona crescător după valorile $|x_i - x|$. *Metoda lui Aitken* este similară cu metoda lui Neville. Ea construiește tabelul

$$\begin{array}{cccccc} x_0 & P_0 & & & & \\ x_1 & P_1 & P_{0,1} & & & \\ x_2 & P_2 & P_{0,2} & P_{0,1,2} & & \\ x_3 & P_3 & P_{0,3} & P_{0,1,3} & P_{0,1,2,3} & \\ x_4 & P_4 & P_{0,4} & P_{0,1,4} & P_{0,1,2,4} & P_{0,1,2,3,4} \end{array}$$

Pentru a calcula o nouă valoare se utilizează valoarea din vârful coloanei precedente și valoarea din aceeași linie, coloana precedentă.

3.4.2. Metoda diferențelor divizate

Vom nota cu $L_k f$ polinomul de interpolare Lagrange cu nodurile x_0, x_1, \dots, x_k pentru $k = 0, 1, \dots, n$. Vom construi L_m prin recurență. Avem

$$(L_0 f)(x) = f(x_0)$$

pentru $k \geq 1$, polinomul $L_k - L_{k-1}$ este de grad k , se anulează în punctele x_0, x_1, \dots, x_k și deci este de forma:

$$(L_k f)(x) = (L_{k-1} f)(x) + f[x_0, x_1, \dots, x_k](x - x_0)(x - x_1) \dots (x - x_{k-1}), \quad (3.4.2)$$

unde $f[x_0, x_1, \dots, x_k]$ desemnează coeficientul lui x^k din $(L_k f)(x)$. Se deduce expresia polinomului de interpolare $L_m f$ cu nodurile x_0, x_1, \dots, x_n

$$(L_m f)(x) = f(x_0) + \sum_{k=1}^m f[x_0, x_1, \dots, x_k](x - x_0)(x - x_1) \dots (x - x_{k-1}), \quad (3.4.3)$$

numită forma Newton ¹⁰ a polinomului de interpolare Lagrange.

Formula (3.4.3) reduce calculul prin recurență al lui $L_m f$ la cel al coeficienților $f[x_0, x_1, \dots, x_k]$, $k = \overline{0, m}$.

Are loc

Lema 3.4.2

$$\forall k \geq 1 \quad f[x_0, x_1, \dots, x_k] = \frac{f[x_1, x_2, \dots, x_k] - f[x_0, x_1, \dots, x_{k-1}]}{x_k - x_0} \quad (3.4.4)$$

și

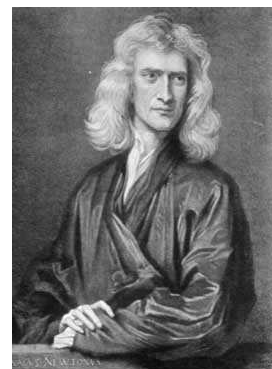
$$f[x_i] = f(x_i), \quad i = 0, 1, \dots, k.$$

Demonstrație. Notăm, pentru $k \geq 1$ cu $L_{k-1}^* f$ polinomul de interpolare pentru f de grad $k - 1$ și cu nodurile x_1, x_2, \dots, x_k ; coeficientul lui x^{k-1} este $f[x_1, x_2, \dots, x_k]$. Polinomul q_k de grad k definit prin

$$q_k(x) = \frac{(x - x_0)(L_{k-1}^* f)(x) - (x - x_k)(L_{k-1} f)(x)}{x_k - x_0}$$

coincide cu f în punctele x_0, x_1, \dots, x_k și deci $q_k(x) \equiv (L_k f)(x)$. Formula (3.4.4) se obține identificând coeficientul lui x^k din cei doi membri. \square

Sir Isaac Newton (1643 - 1727) a fost una dintre cele mai remarcabile figuri ale matematicii și fizicii din vremea sa. Nu numai că dat legile fundamentale ale fizicii moderne, dar a fost și unul dintre inventatorii calculului diferențial și integral (alături de Leibniz, cu care a intrat într-o polemică de o viață privind prioritatea). Lucrarea sa care a avut cea mai mare influență a fost *Principia*, care conține ideile sale asupra interpolării și utilizării ei la integrare.



Definiția 3.4.3 Cantitatea $f[x_0, x_1, \dots, x_k]$ se numește diferență divizată de ordinul k a lui f în punctele x_0, x_1, \dots, x_k .

Altă notație utilizată este $[x_0, \dots, x_k; f]$.

Din definiție rezultă că $f[x_0, x_1, \dots, x_k]$ este independentă de ordinea punctelor x_i și ea poate fi calculată în funcție de $f(x_0), \dots, f(x_m)$. Într-adevăr polinomul de interpolare Lagrange de grad $\leq m$ relativ la punctele x_0, \dots, x_m se scrie

$$(L_m f)(x) = \sum_{i=0}^m l_i f(x_i)$$

și coeficientul lui x^m este

$$f[x_0, \dots, x_m] = \sum_{i=0}^m \frac{f(x_i)}{\prod_{\substack{j=0 \\ j \neq i}}^m (x_i - x_j)}. \quad (3.4.5)$$

Diferențele divizate se pot obține prin algoritmul tabelar următor, bazat pe formula (3.4.4), care este mai flexibil și mai puțin costisitor decât aplicarea formulei (3.4.5)

$$\begin{array}{ccccccc}
 x_0 & f[x_0] & \rightarrow & f[x_0, x_1] & \rightarrow & f[x_0, x_1, x_2] & \rightarrow & f[x_0, x_1, x_2, x_3] \\
 & & \nearrow & & \nearrow & & \nearrow & \\
 x_1 & f[x_1] & \rightarrow & f[x_1, x_2] & \rightarrow & f[x_1, x_2, x_3] & & \\
 & & \nearrow & & \nearrow & & & \\
 x_2 & f[x_2] & \rightarrow & f[x_2, x_3] & & & & \\
 & & \nearrow & & & & & \\
 x_3 & f[x_3] & & & & & & \\
 \vdots & & & & & & &
 \end{array}$$

Prima coloană conține valorile lui f , a doua valorile diferențelor divizate de ordinul I ș.a.m.d; se trece de la o coloană la următoarea utilizând formula (3.4.4): fiecare element este diferența dintre elementul situat în stânga și dedesubtul lui și elementul situat imediat în stânga lui, împărțită la diferența dintre valoarea lui x determinată mergând pe diagonală în jos și valoarea lui x situată la stânga pe orizontală. Diferențele divizate care apar în formula lui Newton (3.4.3) sunt cele $m + 1$ elemente de pe prima linie a tabelului diferențelor divizate. Calculul lor necesită $n(n + 1)$ adunări și $\frac{1}{2}n(n + 1)$ împărțiri. Adăugarea unui nou punct $(x_{m+1}, f[x_{m+1}])$ necesită generarea diagonalei următoare. $L_{m+1}f$ poate fi obținut din $L_m f$ adăugând termenul $f[x_0, \dots, x_{m+1}](x - x_0) \dots (x - x_{m+1})$.

Observația 3.4.4. Eroarea de interpolare este dată de

$$f(x) - (L_m f)(x) = u_m(x)f[x_0, x_1, \dots, x_m, x]. \quad (3.4.6)$$

Într-adevăr, este suficient să observăm că

$$(L_m f)(t) + u_m(t)f[x_0, \dots, x_m; x]$$

este conform lui (3.4.3) polinomul de interpolare (în t) al lui f în punctele x_0, x_1, \dots, x_m, x . Se deduce din teorema referitoare la restul formulei de interpolare Lagrange (3.3.22) că există $\xi \in (a, b)$ astfel încât

$$f[x_0, x_1, \dots, x_m] = \frac{1}{m!} f^{(m)}(\xi) \quad (3.4.7)$$

(formula de medie pentru diferențe divizate). \diamond

Diferența divizată se poate scrie sub forma unui cât a doi determinanți.

Teorema 3.4.5 *Are loc*

$$f[x_0, \dots, x_m] = \frac{(Wf)(x_0, \dots, x_m)}{V(x_0, \dots, x_m)} \quad (3.4.8)$$

unde

$$(Wf)(x_0, \dots, x_m) = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^{m-1} & f(x_0) \\ 1 & x_1 & x_1^2 & \dots & x_1^{m-1} & f(x_1) \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^{m-1} & f(x_m) \end{vmatrix}, \quad (3.4.9)$$

iar $V(x_0, \dots, x_m)$ este determinantul Vandermonde.

Demonstrație. Se dezvoltă $(Wf)(x_0, \dots, x_m)$ după elementele ultimei coloane și ținând cont că fiecare complement algebric este un determinant Vandermonde, se obține

$$\begin{aligned} f[x_0, \dots, x_m] &= \frac{1}{V(x_0, \dots, x_m)} \sum_{i=0}^m V(x_0, \dots, x_{i-1}, x_{i+1}, \dots, x_m) f(x_i) = \\ &= \sum_{i=0}^m (-1)^{m-i} \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_n - x_i)}, \end{aligned}$$

din care după schimbarea semnelor ultimilor $m - i$ termeni rezultă (3.4.5). \square

3.4.3. Diferențe finite: formula lui Newton progresivă și regresivă

În cazul când punctele de interpolare x_i nu sunt echidistante se utilizează algoritmul lui Newton descris anterior; în cazul când punctele sunt echidistante se poate construi un algoritm mai simplu și mai puțin costisitor. Istoric acești algoritmi au avut o mare importanță pentru interpolarea funcțiilor ale căror valori erau tabelate; apariția calculatoarelor moderne a diminuat acest interes, dar noile (co)procesoare flotante i-au readus în actualitate.

Presupunem că funcția f este cunoscută în punctele x_i cu pasul h

$$x_i = x_0 + ih, \quad i = 0, 1, \dots$$

Se definesc *diferențele progresive* prin

$$\Delta^0 f_i = f(x_i) = f_i,$$

și $\forall k \geq 0$

$$\Delta^{k+1} f_i = \Delta(\Delta^k f_i) = \Delta^k f_{i+1} - \Delta^k f_i.$$

Cu ajutorul formulei (3.4.5) se verifică imediat că

$$\Delta^k f_i = f[x_i, x_{i+1}, \dots, x_{i+k}] k! h^k.$$

Forma polinomului Newton $L_m f$ de gradul m a lui f în punctele x_0, x_1, \dots, x_m se simplifică. Dacă $s = (x - x_0)/h$,

$$(x - x_0)(x - x_1) \dots (x - x_{k-1}) f[x_0, x_1, \dots, x_k] = \frac{s(s-1) \dots (s-k+1)}{k!} \Delta^k f_0$$

și deci

$$\begin{aligned} (L_m f)(x) &= f_0 + \frac{s}{1!} \Delta f_0 + \frac{s(s-1)}{2!} \Delta^2 f_0 + \dots \\ &\quad + \frac{s(s-1) \dots (s-m+1)}{m!} \Delta^m f_0. \end{aligned}$$

Utilizând pentru $s \in \mathbb{R}$, $k \in \mathbb{N}$ notația

$$\binom{s}{k} = \frac{s(s-1) \dots (s-k+1)}{k!}$$

(coeficient binomial generalizat), se obține *formula lui Newton progresivă*

$$(L_m f)(x) = f_0 + \binom{s}{1} \Delta f_0 + \binom{s}{2} \Delta^2 f_0 + \dots + \binom{s}{n} \Delta^n f_0, \quad (3.4.10)$$

unde $s = (x - x_0)/h$.

În practică această formulă servește la calculul lui $(L_m f)(x)$ în puncte apropiate de începutul tabeli. Diferențele finite succesive se obțin din tabloul triunghiular

x_0	f_0				
		↘			
			Δf_0		
		↗		↘	
x_1	f_1			$\Delta^2 f_0$	
		↘			
			Δf_1		$\Delta^3 f_0$
		↗		↘	
x_2	f_2			$\Delta^2 f_1$	\vdots
		↘			\vdots
			Δf_2	\vdots	\vdots
		↗		\vdots	\vdots
x_3	f_3			\vdots	\vdots
				\vdots	\vdots
\vdots	\vdots			\vdots	\vdots

unde $\begin{matrix} a & & \\ & \searrow & \\ & & c \text{ înseamnă } c = b - a. \\ & \nearrow & \\ b & & \end{matrix}$

Eroarea comisă după m pași, se scrie pentru $x = x_0 + sh$

$$f(x) - (L_m f)(x) = h^{m+1} \binom{s}{m+1} f^{(m+1)}(\xi_x), \quad (3.4.11)$$

unde ξ_x aparține celui mai mic interval ce conține x_0, x_m și x .

Analog se introduce operatorul de *diferență finită regresivă* ∇ prin

$$\begin{aligned} \nabla^0 f_i &= f_i \\ \nabla^1 f_i &= f_i - f_{i-1} \\ \nabla^{k+1} f_i &= \nabla(\nabla^k f_i) = \nabla^k f_i - \nabla^k f_{i-1} \end{aligned}$$

Efectuând schimbarea de variabilă $s = \frac{x - x_n}{m}$ se obține

$$\begin{aligned} (L_m f)(x) &= f_m + \frac{s}{1!} \nabla f_m + \frac{s(s+1)}{2!} \nabla^2 f_m + \dots \\ &+ \frac{s(s+1) \dots (s+n-1)}{m!} \nabla^m f_m, \end{aligned}$$

care se mai poate scrie

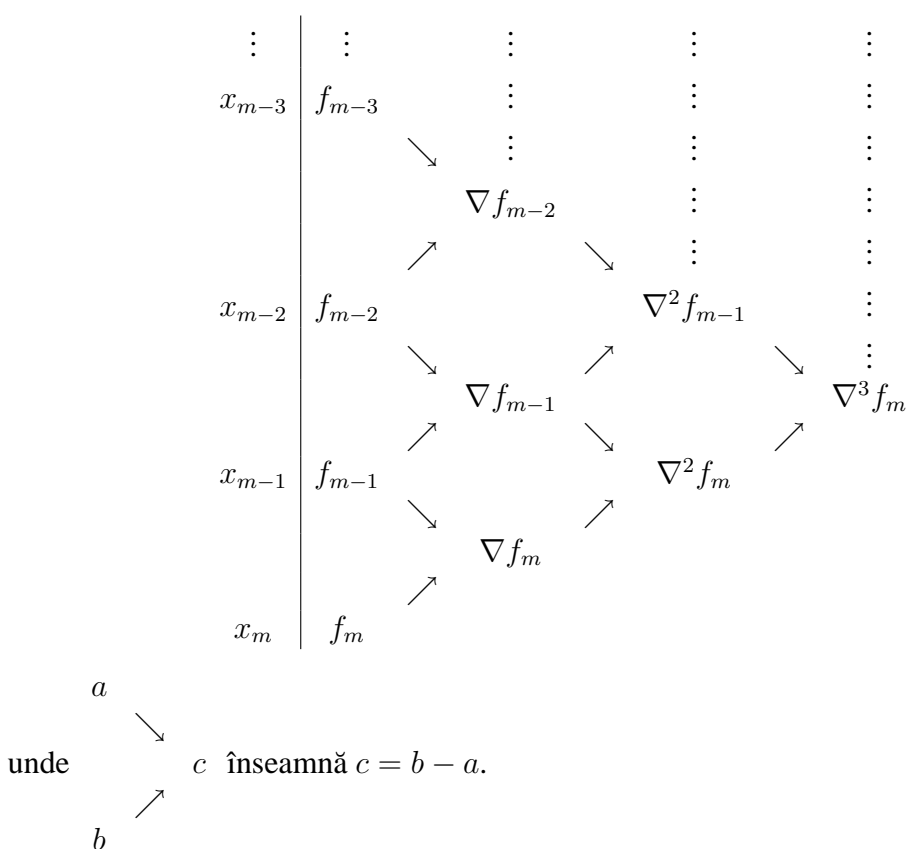
$$(L_m f)(x) = f_m + \binom{s}{1} \nabla f_m + \binom{s+1}{2} \nabla^2 f_m + \cdots + \binom{s+m-1}{m} \nabla^m f_m, \quad (3.4.12)$$

numită *formula lui Newton regresivă*.

Eroarea de interpolare se scrie sub forma

$$f(x) - (L_m f)(x) = h^{m+1} \binom{s+m}{m+1} f^{(m+1)}(\xi_x).$$

unde ξ_x aparține celui mai mic interval ce conține x_0, x_m și x . Formula (3.4.12) servește pentru calculul polinomului de interpolare în valori apropiate de capătul tablei. Diferențele regresive se pot calcula cu tabela



3.4.4. Diferențe divizate cu noduri multiple

Formula (3.4.8) servește ca bază pentru introducerea diferenței divizate cu noduri multiple: dacă $f \in C^m[a, b]$ și $\alpha \in [a, b]$, atunci

$$\lim_{x_0, \dots, x_n \rightarrow \alpha} [x_0, \dots, x_n; f] = \lim_{\xi \rightarrow \alpha} \frac{f^m(\xi)}{m!} = \frac{f^{(m)}(\alpha)}{m!}$$

Aceasta justifică relația

$$[\underbrace{\alpha, \dots, \alpha}_{m+1}; f] = \frac{1}{m!} f^{(m)}(\alpha).$$

Reprezentând aceasta ca pe un câț de doi determinanți se obține

$$(Wf) \left(\underbrace{\alpha, \dots, \alpha}_{m+1} \right) = \begin{vmatrix} 1 & \alpha & \alpha^2 & \dots & \alpha^{m-1} & f(\alpha) \\ 0 & 1 & 2\alpha & \dots & (m-1)\alpha^{m-2} & f'(\alpha) \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & (m-1)! & f^{(m-1)}(\alpha) \end{vmatrix}$$

și

$$V \left(\underbrace{\alpha, \dots, \alpha}_{m+1} \right) = \begin{vmatrix} 1 & \alpha & \alpha^2 & \dots & \alpha^m \\ 0 & 1 & 2\alpha & \dots & m\alpha^{m-1} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & m! \end{vmatrix},$$

adică cei doi determinanți sunt constituiți din linia relativă la nodul α și derivatele succesive ale acesteia până la ordinul m în raport cu α .

Generalizarea pentru mai multe noduri este următoarea:

Definiția 3.4.6 Fie $r_k \in \mathbb{N}$, $k = \overline{0, m}$, $n = r_0 + \dots + r_m$. Presupunem că există $f^{(j)}(x_k)$, $k = \overline{0, m}$, $j = \overline{0, r_k - 1}$. *Mărima*

$$[\underbrace{x_0, \dots, x_0}_{r_0}, \underbrace{x_1, \dots, x_1}_{r_1}, \dots, \underbrace{x_m, \dots, x_m}_{r_m}; f] = \frac{(Wf)(x_0, \dots, x_0, \dots, x_m, \dots, x_m)}{V(x_0, \dots, x_0, \dots, x_m, \dots, x_m)},$$

unde

$$(Wf)(x_0, \dots, x_0, \dots, x_m, \dots, x_m) = \begin{vmatrix} 1 & x_0 & \dots & x_0^{r_0-1} & \dots & x_0^{n-1} & f(x_0) \\ 0 & 1 & \dots & (r_0-1)x_0^{r_0-2} & \dots & (n-1)x_0^{n-2} & f'(x_0) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & (r_0-1)! & \dots & \prod_{p=1}^{r_0-1} (n-p)x_0^{n-r_0} & f^{(r_0-1)}(x_0) \\ 1 & x_m & \dots & x_m^{r_m-1} & \dots & x_m^{n-1} & f(x_m) \\ 0 & 1 & \dots & (r_m-1)x_m^{r_m-2} & \dots & (n-1)x_m^{n-2} & f'(x_m) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & (r_m-1)! & \dots & \prod_{p=1}^{r_m-1} (n-p)x_m^{n-r_m} & f^{(r_m-1)}(x_m) \end{vmatrix},$$

iar $V(x_0, \dots, x_0, \dots, x_m, \dots, x_m)$ este ca mai sus, exceptând ultima coloană care este

$$(x_0^n, nx_0^{n-1}, \dots, \prod_{p=0}^{r_0-2} (n-p)x_0^{n-r_0+1}, \dots, x_m^n, nx_m^{n-1}, \dots, \prod_{p=0}^{r_m-2} x_m^{n-r_m+1})^T$$

$z_0 = x_0$	$f[z_0]$	$f[z_0, z_1] = f'(x_0)$	
$z_1 = x_0$	$f[z_1]$	$f[z_1, z_2] = \frac{f(z_2) - f(z_1)}{z_2 - z_1}$	$f[z_0, z_1, z_2] = \frac{f[z_1, z_2] - f[z_0, z_1]}{z_2 - z_0}$
$z_2 = x_1$	$f[z_2]$	$f[z_2, z_3] = f'(x_1)$	$f[x_1, z_2, z_3] = \frac{f[z_3, z_2] - f[z_2, z_1]}{z_3 - z_1}$
$z_3 = x_1$	$f[z_3]$	$f[z_3, z_4] = \frac{f(z_4) - f(z_3)}{z_4 - z_3}$	$f[z_2, z_3, z_4] = \frac{f[z_4, z_3] - f[z_3, z_2]}{z_4 - z_2}$
$z_4 = x_2$	$f[z_4]$	$f[z_4, z_5] = f'(x_2)$	$f[z_3, z_4, z_5] = \frac{f[z_5, z_4] - f[z_4, z_3]}{z_5 - z_3}$
$z_5 = x_2$	$f[z_5]$		

Tabela 3.2: Tabelă de diferențe divizate pentru noduri duble

se numește diferență divizată cu nodurile multiple x_k , $k = \overline{0, m}$ și ordinele de multiplimitate r_k , $k = \overline{0, m}$.

Generalizând forma Newton a polinomului de interpolare Lagrange se obține o metodă bazată pe diferențele divizate cu noduri multiple pentru polinomul de interpolare Hermite.

Presupunem că se dau nodurile x_i , $i = \overline{0, m}$ și valorile $f(x_i)$, $f'(x_i)$. Definim secvența de noduri $z_0, z_1, \dots, z_{2n+1}$ prin $z_{2i} = z_{2i+1} = x_i$, $i = \overline{0, m}$. Construim acum tabela diferențelor divizate utilizând nodurile z_i , $i = \overline{0, 2m+1}$. Deoarece $z_{2i} = z_{2i+1} = x_i$ pentru orice i , $f[z_{2i}, z_{2i+1}]$ este o diferență divizată cu nod dublu și este egală cu $f'(x_i)$, deci vom utiliza $f'(x_0), f'(x_1), \dots, f'(x_m)$ în locul diferențelor divizate de ordinul I

$$f[z_0, z_1], f[z_2, z_3], \dots, f[z_{2m}, z_{2m+1}].$$

Restul diferențelor se obțin în manieră obișnuită, așa cum se arată în tabelul 3.2. Ideea poate fi extinsă și pentru alte interpolări Hermite. Se pare că metoda este datorată lui Powell.

3.5. Interpolare spline

Fie Δ o diviziune a lui $[a, b]$

$$\Delta : a = x_1 < x_2 < \dots < x_{n-1} < x_n = b \quad (3.5.1)$$

Vom utiliza un polinom de grad mic pe subintervalul $[x_i, x_{i+1}]$, $i = \overline{1, n-1}$. Motivul este acela că pe intervale suficient de mici funcțiile pot fi approximate arbitrar de bine prin polinoame de grad mic, chiar 0 sau 1.

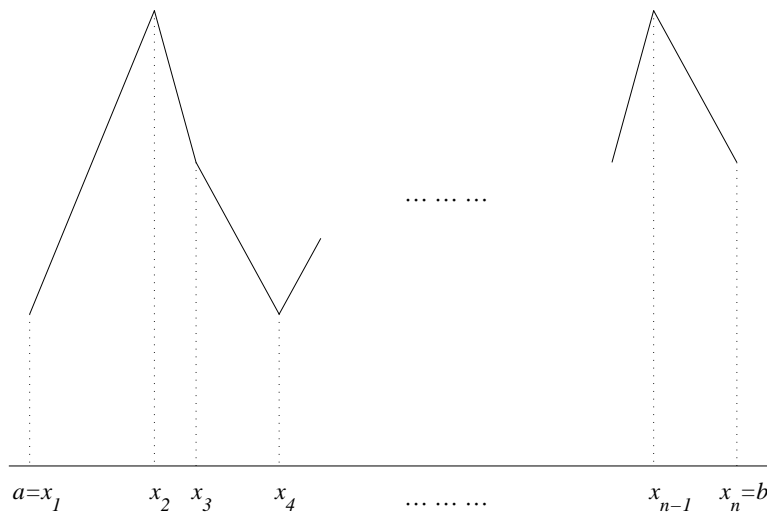


Figura 3.10: Spline liniare

Am introdus în exemplul 3.0.12 spațiul

$$\mathbb{S}_m^k(\Delta) = \{s : s \in C^k[a, b], s|_{[x_i, x_{i+1}]} \in \mathbb{P}_m, i = 1, 2, \dots, n-1\} \quad (3.5.2)$$

$m \geq 0, k \in \mathbb{N} \cup \{-1\}$, numit spațiul funcțiilor spline polinomiale de grad m și clasă de netezime k . Dacă $k = m$, atunci funcțiile $s \in \mathbb{S}_m^m(\Delta)$ sunt polinoame.

Pentru $m = 1$ și $k = 0$ se obțin *spline liniare*.

Dorim să găsim $s \in \mathbb{S}_1^0(\Delta)$ astfel încât

$$s(x_i) = f_i, \text{ unde } f_i = f(x_i), \quad i = 1, 2, \dots, n.$$

Soluția este trivială, vezi figura 3.10. Pe intervalul $[x_i, x_{i+1}]$

$$s(f; x) = f_i + (x - x_i)f[x_i, x_{i+1}], \quad (3.5.3)$$

iar

$$|f(x) - s(f(x))| \leq \frac{(\Delta x_i)^2}{8} \max_{x \in [x_i, x_{i+1}]} |f''(x)|. \quad (3.5.4)$$

Rezultă că

$$\|f(\cdot) - s(f, \cdot)\|_\infty \leq \frac{1}{8} |\Delta|^2 \|f''\|_\infty. \quad (3.5.5)$$

Dimensiunea lui $\mathbb{S}_1^0(\Delta)$ se calculează astfel: deoarece avem $n - 1$ porțiuni și pe fiecare 2 coeficienți (2 grade de libertate) și fiecare condiție reduce numărul de grade de

libertate cu 1, avem în final

$$\dim \mathbb{S}_1^0(\Delta) = 2n - 2 - (n - 2) = n.$$

O bază a spațiului este dată de așa-numitele funcții *B-spline*:

Punem $x_0 = x_1$, $x_{n+1} = x_n$, pentru $i = \overline{1, n}$

$$B_i(x) = \begin{cases} \frac{x - x_{i-1}}{x_i - x_{i-1}}, & \text{pentru } x_{i-1} \leq x \leq x_i \\ \frac{x_{i+1} - x}{x_{i+1} - x_i}, & \text{pentru } x_i \leq x \leq x_{i+1} \\ 0, & \text{în rest} \end{cases} \quad (3.5.6)$$

Pentru $i = 1$ prima și pentru $i = n$ a doua ecuație se ignoră.

Funcția B_i se numește *pălărie chinezească*. Graficul funcțiilor B_i apare în figura 3.11.

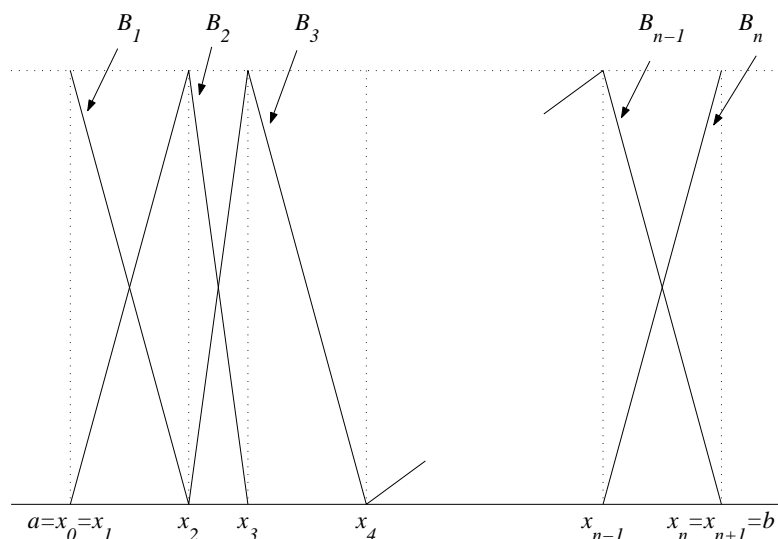


Figura 3.11: Funcții B-spline de grad 1

Ele au proprietatea

$$B_i(x_j) = \delta_{ij},$$

sunt liniar independente, deoarece

$$s(x) = \sum_{i=1}^n c_i B_i(x) = 0 \wedge x \neq x_j \Rightarrow c_j = 0.$$

și

$$\langle B_i \rangle_{i=\overline{1,n}} = S_1^0(\Delta),$$

B_i joacă același rol ca polinoamele fundamentale Lagrange l_i .

3.5.1. Interpolarea cu spline cubice

Funcțiile spline cubice sunt cele mai utilizate.

Vom discuta întâi problema interpolării pentru $s \in \mathbb{S}_3^1(\Delta)$. Continuitatea derivatei de ordinul I pentru $s_3(f; \cdot)$ se poate realiza impunând valorile primei derivate în fiecare punct x_i , $i = 1, 2, \dots, n$. Astfel fie m_1, m_2, \dots, m_n numere arbitrare date și notăm

$$s_3(f; \cdot)|_{[x_i, x_{i+1}]} = p_i(x), \quad i = 1, 2, \dots, n-1 \quad (3.5.7)$$

Realizăm $s_3(f; x_i) = m_i$, $i = \overline{1, n}$, luând fiecare bucată ca soluție unică a problemei de interpolare Hermite, și anume

$$\begin{aligned} p_i(x_i) &= f_i, & p_i(x_{i+1}) &= f_{i+1}, & i &= \overline{1, n-1}, \\ p_i'(x_i) &= m_i, & p_i'(x_{i+1}) &= m_{i+1} \end{aligned} \quad (3.5.8)$$

Vom rezolva problema folosind interpolarea Newton. Diferențele divizate sunt

$$\begin{array}{cccc} x_i & f_i & m_i & \frac{f[x_i, x_{i+1}] - m_i}{\Delta x_i} & \frac{m_{i+1} + m_i - 2f[x_i, x_{i+1}]}{(\Delta x_i)^2} \\ x_i & f_i & f[x_i, x_{i+1}] & \frac{m_{i+1} - f[x_i, x_{i+1}]}{\Delta x_i} & \\ x_{i+1} & f_{i+1} & m_{i+1} & & \\ x_{i+1} & f_{i+1} & & & \end{array}$$

și deci forma Newton a polinomului de interpolare Hermite este

$$\begin{aligned} p_i(x) &= f_i + (x - x_i)m_i + (x - x_i)^2 \frac{f[x_i, x_{i+1}] - m_i}{\Delta x_i} + \\ &+ (x - x_i)^2 (x - x_{i+1}) \frac{m_{i+1} + m_i - 2f[x_i, x_{i+1}]}{(\Delta x_i)^2}. \end{aligned}$$

Forma Taylor a lui p_i pentru $x_i \leq x \leq x_{i+1}$ este

$$p_i(x) = c_{i,0} + c_{i,1}(x - x_i) + c_{i,2}(x - x_i)^2 + c_{i,3}(x - x_i)^3 \quad (3.5.9)$$

și deoarece $x - x_{i+1} = x - x_i - \Delta x_i$, prin identificare avem

$$\begin{aligned} c_{i,0} &= f_i \\ c_{i,1} &= m_i \\ c_{i,2} &= \frac{f[x_i, x_{i+1}] - m_i}{\Delta x_i} - c_{i,3} \Delta x_i \\ c_{i,3} &= \frac{m_{i+1} + m_i - 2f[x_i, x_{i+1}]}{(\Delta x_i)^2} \end{aligned} \quad (3.5.10)$$

Deci, pentru a calcula $s_3(f; x)$ într-un punct care nu este nod, trebuie în prealabil să localizăm intervalul $[x_i, x_{i+1}] \ni x$.

Să calculăm coeficienții cu (3.5.10) și să evaluăm spline-ul cu (3.5.9). Vom discuta câteva alegeri posibile pentru m_1, m_2, \dots, m_n .

Interpolare Hermite cubică pe porțiuni

Se alege $m_i = f'(x_i)$ (presupunând că aceste derivate sunt cunoscute). Se ajunge la o schemă strict locală, în care fiecare bucată poate fi determinată independent de cealaltă. Mai mult, eroarea este

$$|f(x) - p_i(x)| \leq \left(\frac{1}{2}\Delta x_i\right)^4 \max_{x \in [x_i, x_{i+1}]} \frac{|f^{(4)}(x)|}{4!}, \quad x_i \leq x \leq x_{i+1}. \quad (3.5.11)$$

Deci

$$\|f(\cdot) - s_3(f; \cdot)\|_\infty \leq \frac{1}{384} |\Delta|^4 \|f^{(4)}\|_\infty. \quad (3.5.12)$$

Pentru puncte echidistante

$$|\Delta| = (b - a)/(n - 1)$$

și deci

$$\|f(\cdot) - s_3(f; \cdot)\|_\infty = O(n^{-4}), \quad n \rightarrow \infty. \quad (3.5.13)$$

Interpolare cu spline cubice

Cerem ca $s_3(f; \cdot) \in \mathbb{S}_3^2(\Delta)$, adică continuitatea derivatelor de ordinul al II-lea. Aceasta înseamnă cu notația (3.5.7)

$$p''_{i-1}(x_i) = p''_i(x_i), \quad i = \overline{2, n-1}, \quad (3.5.14)$$

care convertită în coeficienți Taylor (3.5.9) dă

$$2c_{i-1,2} + 6c_{i-1,3}\Delta x_{i-1} = 2c_{i,2}, \quad i = \overline{2, n-1}.$$

Înlocuind cu valorile explicite (3.5.10) pentru coeficienți se ajunge la sistemul liniar

$$\Delta x_i m_{i-1} + 2(\Delta x_{i-1} + \Delta x_i) m_i + (\Delta x_{i-1}) m_{i+1} = b_i, \quad i = \overline{2, n-1} \quad (3.5.15)$$

unde

$$b_i = 3\{\Delta x_i f[x_{i-1}, x_i] + \Delta x_{i-1} f[x_i, x_{i+1}]\} \quad (3.5.16)$$

Avem un sistem de $n - 2$ ecuații liniare cu n necunoscute m_1, m_2, \dots, m_n . Odată alese m_1 și m_n , sistemul devine tridiagonal și se poate rezolva eficient prin eliminare gaussiană, prin factorizare sau cu o metodă iterativă.

Se dau în continuare câteva alegeri posibile pentru m_1 și m_n .

Spline complete(racordate, limitate). Luăm $m_1 = f'(a)$, $m_n = f'(b)$. Se știe că pentru acest tip de spline, dacă $f \in C^4[a, b]$

$$\|f^{(r)}(\cdot) - s^{(r)}(f; \cdot)\|_\infty \leq c_r |\Delta|^{4-r} \|f^{(n)}\|_\infty, \quad r = 0, 1, 2, 3 \quad (3.5.17)$$

unde $c_0 = \frac{5}{384}$, $c_1 = \frac{1}{24}$, $c_2 = \frac{3}{8}$, iar c_3 depinde de raportul $\frac{|\Delta|}{\min_i \Delta x_i}$.

Spline care utilizează derivatele secunde. Impunem condițiile $s_3''(f; a) = f''(a)$; $s_3''(f; b) = f''(b)$. Aceste condiții conduc la două ecuații suplimentare

$$\begin{aligned} 2m_1 + m_2 &= 3f[x_1, x_2] - \frac{1}{2}f''(a)\Delta x_1 \\ m_{n-1} + 2m_n &= 3f[x_{n-1}, x_n] + \frac{1}{2}f''(b)\Delta x_{n-1} \end{aligned} \quad (3.5.18)$$

Prima ecuație se pune la începutul sistemului (3.5.15), iar a doua la sfârșitul lui, păstrându-se astfel structura tridiagonală a sistemului.

Spline cubice naturale. Impunând $s''(f; a) = s''(f; b) = 0$, se obțin două ecuații noi din (3.5.18) luând $f''(a) = f''(b) = 0$.

Avantajul – este nevoie numai de valori ale lui f , nu și ale derivatelor, dar prețul plătit este degradarea preciziei la $O(|\Delta|^2)$ în vecinătatea capetelor (în afară de cazul când $f''(a) = f''(b) = 0$).

”Not-a-knot spline”. (C. deBoor). Cerem ca $p_1(x) \equiv p_2(x)$ și $p_{n-2}(x) \equiv p_{n-1}(x)$; adică primele două părți și respectiv ultimele două trebuie să coincidă. Într-adevăr, aceasta înseamnă că primul punct interior x_2 și ultimul x_{n-1} sunt ambele inactive. Se obțin încă două ecuații suplimentare exprimând continuitatea lui $s_3'''(f; x)$ în $x = x_2$ și $x = x_{n-1}$. Condiția de continuitate a lui $s_3(f, \cdot)$ în x_2 și x_{n-1} revine la egalitatea coeficienților dominanți $c_{1,3} = c_{2,3}$ și $c_{n-2,3} = c_{n-1,3}$. De aici se obțin ecuațiile

$$\begin{aligned} (\Delta x_2)^2 m_1 + [(\Delta x_2)^2 - (\Delta x_1)^2] m_2 - (\Delta x_1)^2 m_3 &= \beta_1 \\ (\Delta x_2)^2 m_{n-2} + [(\Delta x_2)^2 - (\Delta x_1)^2] m_{n-1} - (\Delta x_1)^2 m_n &= \beta_2, \end{aligned}$$

unde

$$\begin{aligned} \beta_1 &= 2\{(\Delta x_2)^2 f[x_1, x_2] - (\Delta x_1)^2 f[x_2, x_3]\} \\ \beta_2 &= 2\{(\Delta x_{n-1})^2 f[x_{n-2}, x_{n-1}] - (\Delta x_{n-2})^2 f[x_{n-1}, x_n]\}. \end{aligned}$$

Prima ecuație se adaugă pe prima poziție iar a doua pe ultima poziție a sistemului format din cele $n - 2$ ecuații date de (3.5.15) și (3.5.16). Sistemul obținut nu mai este tridiagonal, dar el se poate transforma în unul tridiagonal combinând ecuațiile 1 cu 2 și $n - 1$ cu n . După aceste transformări prima și ultima ecuație devin

$$\Delta x_2 m_1 + (\Delta x_2 + \Delta x_1) m_2 = \gamma_1 \quad (3.5.19)$$

$$(\Delta x_{n-1} + \Delta x_{n-2}) m_{n-1} + \Delta x_{n-2} m_n = \gamma_2, \quad (3.5.20)$$

unde

$$\begin{aligned}\gamma_1 &= \frac{1}{\Delta x_2 + \Delta x_1} \{f[x_1, x_2] \Delta x_2 [\Delta x_1 + 2(\Delta x_1 + \Delta x_2)] + (\Delta x_1)^2 f[x_2, x_3]\} \\ \gamma_2 &= \frac{1}{\Delta x_{n-1} + \Delta x_{n-2}} \{ \Delta x_{n-1}^2 f[x_{n-2}, x_{n-1}] + \\ &\quad [2(\Delta x_{n-1} + \Delta x_{n-2}) + \Delta x_{n-1}] \Delta x_{n-2} f[x_{n-1}, x_n] \}.\end{aligned}$$

3.5.2. Proprietăți de minimalitate ale funcțiilor spline cubice

Funcțiile spline cubice complete și naturale au proprietăți interesante de optimalitate. Pentru a le formula, este convenabil să considerăm nu numai subdiviziunea Δ ci și

$$\Delta' : a = x_0 = x_1 < x_2 < x_3 < \dots < x_{n-1} < x_n = x_{n+1} = b, \quad (3.5.21)$$

în care capetele sunt noduri duble. Aceasta înseamnă că ori de câte ori interpolăm pe Δ' , interpolăm valorile funcției pe punctele interioare, iar la capete valorile funcției și ale derivatei. Prima teoremă se referă la funcții spline cubice complete $s_{\text{compl}}(f; \cdot)$.

Teorema 3.5.1 Pentru orice funcție $g \in C^2[a, b]$ care interpoalează f pe Δ' , are loc

$$\int_a^b [g''(x)]^2 dx \geq \int_a^b [s''_{\text{compl}}(f; x)]^2 dx, \quad (3.5.22)$$

cu egalitate dacă și numai dacă $g(\cdot) = s_{\text{compl}}(f; \cdot)$.

Observația 3.5.2. $s_{\text{compl}}(f; \cdot)$ din teorema 3.5.1 interpoalează f pe Δ' și dintre toți interpolanții de acest tip, derivata sa de ordinul II are norma minimă. \diamond

Demonstrație. Folosim notația prescurtată $s_{\text{compl}} = s$. Teorema rezultă imediat, dacă arătăm că

$$\int_a^b [g''(x)]^2 dx = \int_a^b [g''(x) - s''(x)]^2 dx + \int_a^b [s''(x)]^2 dx. \quad (3.5.23)$$

Aceasta implică imediat (3.5.22) și faptul că egalitatea în (3.5.22) are loc dacă și numai dacă $g''(x) - s''(x) \equiv 0$, din care integrând de două ori de la a la x și utilizând proprietățile de interpolare ale lui s și g în $x = a$ se obține $g(x) = s(x)$. Relația (3.5.23) este echivalentă cu

$$\int_a^b s''(x)[g''(x) - s''(x)] dx = 0. \quad (3.5.24)$$

Integrând prin părți și ținând cont că $s'(b) = g'(b) = f'(b)$ și $s'(a) = g'(a) = f'(a)$ se obține

$$\int_a^b s''(x)[g''(x) - s''(x)] dx = \quad (3.5.25)$$

$$\begin{aligned}
&= s''(x)[g'(x) - s'(x)] \Big|_a^b - \int_a^b s'''(x)[g'(x) - s'(x)]dx = \\
&= - \int_a^b s'''(x)[g'(x) - s'(x)]dx.
\end{aligned}$$

Deoarece s''' este constantă pe porțiuni

$$\begin{aligned}
\int_a^b s'''(x)[g'(x) - s'(x)]dx &= \sum_{\nu=1}^{n-1} s'''(x_\nu + 0) \int_{x_\nu}^{x_{\nu+1}} [g'(x) - s'(x)]dx = \\
&= \sum_{\nu=1}^{n-1} s'''(x_{\nu+0}) [g(x_{\nu+1}) - s(x_{\nu+1}) - (g(x_\nu) - s(x_\nu))] = 0
\end{aligned}$$

căci atât s cât și g interpolează f pe Δ . Aceasta demonstrează (3.5.24) și deci și teorema. \square

Pentru interpolarea pe Δ , calitatea de a fi optimal revine funcțiilor spline naturale de interpolare $s_{nat}(f; \cdot)$.

Teorema 3.5.3 Pentru orice funcție $g \in C^2[a, b]$ ce interpolează f pe Δ , are loc

$$\int_a^b [g''(x)]^2 dx \geq \int_a^b [s''_{nat}(f; x)]^2 dx \quad (3.5.26)$$

cu egalitate dacă și numai dacă $g(\cdot) = s_{nat}(f; \cdot)$.

Demonstrația este analoagă cu a teoremei 3.5.1, deoarece (3.5.23) are loc din nou căci $s''(b) = s''(a) = 0$.

Punând $g(\cdot) = s_{compl}(f; \cdot)$ în teorema 3.5.3 se obține

$$\int_a^b [s''_{compl}(f; x)]^2 dx \geq \int_a^b [s''_{nat}(f; x)]^2 dx. \quad (3.5.27)$$

Deci, într-un anumit sens, spline-ul cubic natural este cel mai neted interpolant.

Proprietatea exprimată în teorema 3.5.3 stă la originea numelui de spline. Un spline este o vergea flexibilă folosită pentru a desena curbe. Dacă forma sa este dată de ecuația $y = g(x)$, $x \in [a, b]$ și dacă spline-ul trebuie să treacă prin punctele (x_i, g_i) , atunci se presupune că spline-ul are o formă ce minimizează energia potențială

$$\int_a^b \frac{[g''(x)]^2 dx}{(1 + [g'(x)]^2)^3},$$

pentru toate funcțiile g supuse aceluiași restricții. Pentru variații lente ale lui g ($\|g'\|_\infty \ll 1$) aceasta aproximează bine proprietatea de minim din teorema 3.5.3.

CAPITOLUL 4

Aproximare uniformă

Cuprins

4.1. Polinoamele lui Bernstein	124
4.2. B-spline	129
4.2.1. Noțiuni și rezultate de bază	129
4.2.2. Algoritmul de evaluare a unui B-spline	132
4.2.3. Aplicații în grafica pe calculator	133
4.2.4. Exemple	136
4.3. Funcții spline cu variație diminuată	140
4.4. Operatori liniari și pozitivi	143
4.5. Cea mai bună aproximare uniformă	147

Aproximarea uniformă a funcțiilor are două aspecte fundamentale:

- aproximarea unei funcții f prin funcții ce converg uniform către f pe domeniul considerat;
- cea mai bună aproximare uniformă a lui f prin funcții dintr-o mulțime dată.

Vom pune accentul pe studiul primei probleme.

Fie \mathcal{B} o clasă de funcții definite pe un interval $[a, b] \subset \mathbb{R}$ și $\mathcal{A} \subset \mathcal{B}$. Pentru $f \in \mathcal{B}$ și $\varepsilon \in \mathbb{R}_+$ date, se pune problema determinării unei funcții $g \in \mathcal{A}$ astfel încât

$$|f(x) - g(x)| < \varepsilon, \quad x \in [a, b].$$

Această problemă își are originea în binecunoscuta teoremă de aproximare uniformă a lui Weierstrass ¹.

Teorema 4.0.4 (Weierstrass 1885) Fie $(C[a, b], \|\cdot\|)$. Subspațiul $\mathbb{P} \leq C[a, b]$ al polinoamelor de grad arbitrar este dens în $C[a, b]$ ($\overline{\mathbb{P}} = C[a, b]$), adică $\forall f \in C[a, b] \forall \varepsilon \in \mathbb{R}_+ \exists p \in \mathbb{P}$ astfel încât $|f(x) - p(x)| < \varepsilon, \forall x \in [a, b]$. (Echivalent $\exists (p_n) \subset \mathbb{P}$ astfel încât $p_n \rightrightarrows f$.)

O importanță deosebită în studiul aproximării uniforme a funcțiilor a avut problema pusă de E. Borel în 1905. Deoarece șirul polinoamelor de interpolare Lagrange $(L_m f)$ nu converge către f , el nu poate fi folosit la demonstrația teoremei lui Weierstrass. Din acest motiv Borel a propus căutarea unor procedee de interpolare mai generale, care să permită contruirea unor șiruri de polinoame $(P_n f)$ ce converg uniform către f dacă $f \in C[a, b]$. Astfel de procedee ar permite să se dea demonstrații constructive ale teoremei lui Weierstrass. Una dintre primele soluții pentru problema lui Borel a fost dată de Sergi N. Bernstein în 1912 prin introducerea polinoamelor care în poartă numele.

4.1. Polinoamele lui Bernstein

Definiția 4.1.1 Fie $f : [0, 1] \rightarrow \mathbb{R}$. Operatorul B_m definit prin relația

$$(B_m f)(x) = \sum_{k=0}^m p_{mk}(x) f\left(\frac{k}{m}\right), \quad (4.1.1)$$

unde

$$(4.1.2)$$

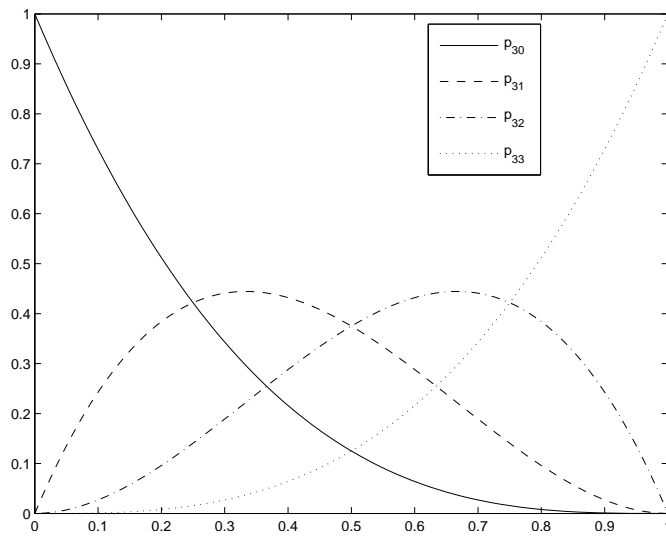
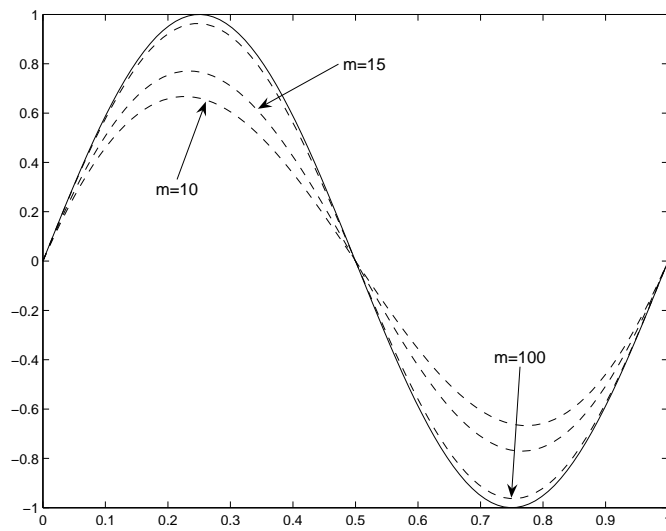
$$p_{m,k}(x) = \binom{m}{k} x^k (1-x)^{m-k} \quad (4.1.3)$$

se numește operatorul lui Bernstein, iar polinomul $B_m f$ se numește polinomul lui Bernstein.

În figura 4.1 apar polinoamele de bază Bernstein pentru $k = 3$, iar în figura 4.2 polinoamele Bernstein pentru funcția $f(x) = \sin 2\pi x$ și $k = 10, 15, 100$.

Karl Theodor Wilhelm Weierstrass (1813-1897), matematician german, considerat a fi unul dintre părinții analizei moderne. A ¹avut contribuții importante în teoria funcțiilor de variabile reale, funcții eliptice, calcul variațional, studiul formelor biliniare și pătratice.



Figura 4.1: Polinoamele de bază Bernstein pentru $k = 3$ Figura 4.2: Polinoamele Bernstein pentru funcția $f(x) = \sin 2\pi x$ (linie continuă) și $k = 10, 15, 100$

Polinomul lui Bernstein poate fi obținut și pe cale probabilistică astfel: dacă f este mărginită pe $[0, 1]$, se consideră variabila aleatoare X cu distribuția

$$X : \begin{pmatrix} f\left(\frac{k}{m}\right) \\ p_{m,k}(x) \end{pmatrix},$$

a cărei valoare medie $M(x) = \sum_{k=0}^m p_{m,k}(x) f\left(\frac{k}{m}\right)$ este chiar $(B_m f)(x)$.

Proprietăți

Teorema 4.1.2 1. B_m este un operator liniar.

2. B_m este un operator pozitiv, adică

$$f(x) \geq 0 \Rightarrow (B_m f)(x) \geq 0, \forall x \in [0, 1].$$

3. B_m reproduce polinoamele până la gradul 1 inclusiv, adică

$$(B_m e_k)(x) = e_k(x), \quad k = 0, 1$$

și în plus

$$(B_m e_2)(x) = e_2(x) + \frac{x(1-x)}{m}.$$

4. Dacă $m \leq f \leq M$ pe $[0, 1]$, atunci $m \leq B_m f \leq M$ pe $[0, 1]$.

5. Dacă $f \in C[0, 1]$, atunci $B_m f \Rightarrow f$ pe $[0, 1]$ când $m \rightarrow \infty$.

Demonstrație. 1. Rezultă imediat din definiție.

2. Imediată ținând cont că $\forall x \in [0, 1] p_m(x) \geq 0$.

3. Se folosește binomul lui Newton

$$(u+v)^m = \sum_{k=0}^m \binom{m}{k} u^k v^{m-k}. \quad (4.1.4)$$

Diferențiem succesiv în raport cu u egalitatea de mai sus

$$u(u+v)^{m-1} = \sum_{k=0}^m \frac{k}{m} u^k v^{m-k} \binom{m}{k} \quad (4.1.5)$$

$$\left(1 - \frac{1}{m}\right) u^2 (u+v)^{m-2} = \sum_{k=0}^m \left[\frac{k^2}{m^2} - \frac{k}{m^2} \right] \binom{m}{k} u^k v^{m-k}. \quad (4.1.6)$$

Punând în ultimele trei relații $u = x$ și $v = 1 - x$ vom obține

$$\begin{aligned} \sum_{k=0}^m p_{mk}(x) &= 1, & \sum_{k=0}^m \frac{k}{m} p_{mk}(x) &= x, \\ \sum_{k=0}^m \left(\frac{k}{m}\right)^2 p_{mk}(x) &= x^2 + \frac{x(1-x)}{m}. \end{aligned}$$

4. Din $f - m \geq 0$, prin pozitivitate, rezultă $B_m(f - m) \geq 0$. Analog și cealaltă inegalitate.
5. $f \in C[0, 1] \Rightarrow f$ uniform continuă pe $[0, 1]$, adică $\forall \varepsilon > 0 \exists \delta = \delta_\varepsilon$ astfel încât $\forall x, x' \in [0, 1]$ cu $|x - x'| < \delta$ să avem $|f(x) - f(x')| < \frac{\varepsilon}{2}$. Folosind a treia proprietate putem scrie

$$\begin{aligned} f(x) - (B_m f)(x) &= \sum_{k=0}^m p_{mk}(x) \left[f(x) - f\left(\frac{k}{m}\right) \right] \text{ și} \\ |f(x) - (B_m f)(x)| &\leq \sum_{k=0}^m p_{mk}(x) \left| f(x) - f\left(\frac{k}{m}\right) \right| = \\ &= \sum_{k \in I_m} p_{mk}(x) \left| f(x) - f\left(\frac{k}{m}\right) \right| + \sum_{k \in J_m} p_{mk}(x) \left| f(x) - f\left(\frac{k}{m}\right) \right| \end{aligned}$$

unde

$$I_m = \left\{ k : \left| \frac{k}{m} - x \right| < \delta \right\} \text{ iar } J_m = \left\{ k : \left| \frac{k}{m} - x \right| \geq \delta \right\}.$$

Dacă $M = \max_{x \in [0, 1]} |f(x)|$ obținem în continuare

$$\begin{aligned} |f(x) - (B_m f)(x)| &\leq \frac{\varepsilon}{2} \sum_{k=0}^m p_{mk}(x) + 2M \sum_{k \in J_m} p_{mk}(x) \\ &\leq \frac{\varepsilon}{2} + 2M \sum_{k \in J_m} p_{mk}(x) \end{aligned}$$

Deoarece

$$\left| \frac{k}{m} - x \right| \geq \delta \Rightarrow \frac{\left(\frac{k}{m} - x\right)^2}{\delta^2}$$

obținem succesiv, folosind proprietatea 3

$$\begin{aligned} \sum_{k \in J_m} p_{mk}(x) &\leq \frac{1}{\delta^2} \sum_{k \in J_m} \left(\frac{k}{m} - x \right)^2 p_{mk}(x) \leq \frac{1}{\delta^2} \sum_{k=0}^m \left(\frac{k}{m} - x \right)^2 p_{mk}(x) \\ &= \frac{1}{\delta^2} \left(\underbrace{\sum_{k=0}^m \frac{k^2}{m^2} p_{mk}(x)}_{x^2 + \frac{x(1-x)}{m}} - 2x \underbrace{\sum_{k=0}^m \frac{k}{m} p_{mk}(x)}_x + x^2 \underbrace{\sum_{k=0}^m p_{mk}(x)}_1 \right) \\ &= \frac{1}{\delta^2} \frac{x(1-x)}{m} \leq \frac{1}{4m\delta^2}, \quad x \in [0, 1]. \end{aligned}$$

Deci, $|f(x) - (B_m f)(x)| \leq \frac{\varepsilon}{2} + \frac{M}{2m\delta^2}$, adică

$$|f(x) - (B_m f)(x)| \leq \varepsilon, \quad \text{dacă } m > \frac{M}{\varepsilon\delta^2}, \quad x \in [0, 1].$$

Această proprietate este o demonstrație constructivă a teoremei lui Weierstrass.

□

Definiția 4.1.3 *Formula*

$$f = B_m f + R_m f \quad (4.1.7)$$

se numește formula de aproximare a lui Bernstein iar $R_m f$ termenul rest.

Teorema 4.1.4 *Dacă $f \in C^2[0, 1]$, atunci*

$$(R_m f)(x) = -\frac{x(1-x)}{2m} f''(\xi), \quad \xi \in [0, 1] \quad (4.1.8)$$

și

$$|(R_m f)(x)| \leq \frac{1}{8m} \|f''\|_\infty. \quad (4.1.9)$$

Demonstrație. Conform teoremei 4.1.2, proprietățile 1 și 2, $B_m f$ reproduce polinoamele de grad 1, adică $B_m f = f$, pentru orice $f \in \mathbb{P}_1$. Aplicând teorema lui Peano se obține

$$(R_m f)(x) = \int_0^1 K(x; t) f''(t) dt,$$

unde

$$K(x; t) = (x - t)_+ - \sum_{k=0}^m p_{m,k}(x) \left(\frac{k}{m} - t \right)_+.$$

Cum $K(x; t) \leq 0$ pentru $x, t \in [0, 1]$ corolarul la teorema lui Peano ne conduce la

$$(R_m f)(x) = -\frac{x(1-x)}{2m} f''(\xi), \quad \xi \in [0, 1].$$

Inegalitatea (4.1.9) rezultă din faptul că $x(1-x) \leq \frac{1}{4}$ pe $[0, 1]$. \square

Importanța teoremei lui Weierstrass a făcut ca ea să fie mereu în atenția matematicienilor. S-au dat numeroase demonstrații constructive și generalizări. Una dintre ele este teorema lui Stone și Weierstrass.

Fie D un compact, $C(D)$ spațiul funcțiilor definite pe compactul D înzestrat cu norma $\|f\| = \max_D |f(x)|$. $C(D)$ este o algebră.

Enunțăm fără demonstrație:

Teorema 4.1.5 (Stone-Weierstrass, 1948) *Dacă A este o subalgebră a lui $C(D)$ care separă punctele lui D , atunci închiderea \bar{A} a lui A este fie $C(D)$, fie mulțimea funcțiilor continue care se anulează într-un punct al lui D .*

Observația 4.1.6. 1. Dacă A conține funcțiile constante, atunci $\bar{A} = C(D)$.

2. O familie S de funcții definite pe E separă punctele lui E dacă $\forall u, v \in E, u \neq v, \exists h \in S$ astfel încât $h(u) \neq h(v)$. \diamond

Pentru detalii asupra demonstrațiilor teoremei lui Weierstrass și asupra teoremei Stone-Weierstrass a se vedea [39].

4.2. B-spline

4.2.1. Noțiuni și rezultate de bază

Fie $m, n \in \mathbb{N}$ și punctele t_0, t_1, \dots, t_m , cu $t_i \leq t_{i+1}$. Fie $[a, b] \subset \mathbb{R}$ astfel încât $t_k \leq a$ și $t_{n-k} \geq b$. Pentru $j = \bar{1}, m - \bar{i}$ punem

$$\omega_{i,j}(x) = \begin{cases} \frac{x - t_i}{t_{i+j} - t_i}, & \text{dacă } t_i < t_{i+j} \\ 0, & \text{altfel.} \end{cases} \quad (4.2.1)$$

Definiția 4.2.1 Pentru $i = \overline{0, m - k - 1}$ funcțiile B-spline de rang i și grad k se definesc prin recurență astfel

$$B_{i,0}(x) = \begin{cases} 1, & \text{dacă } x \in [t_i, t_{i+1}], \\ 0, & \text{în caz contrar.} \end{cases} \quad (4.2.2)$$

$$B_{i,k}(x) = \omega_{i,k}(x)B_{i,k-1}(x) + (1 - \omega_{i+1,k}(x))B_{i+1,k-1}(x), \quad \text{pentru } k \geq 1. \quad (4.2.3)$$

Propoziția 4.2.2 Au loc următoarele proprietăți:

- (a) $B_{i,k}$ este un polinom de grad k pe porțiuni;
- (b) $B_{i,k}(x) = 0$ pentru $x \notin [t_i, t_{i+k+1}]$;
- (c) $B_{i,k}(x) > 0$ pentru $x \in (t_i, t_{i+k+1})$; $B_{i,k}(t_i) = 0$ înafara cazului când $t_i = t_{i+1} = \dots = t_{i+k} < t_{i+k+1}$ (nod de ordin $k + 1$) când avem $B_{i,k}(t_i) = 1$.
- (d)

$$\sum_{i=0}^{m-k-1} B_{i,k}(x) = 1, \quad \forall x \in [a, b].$$

- (e) Fie $x \in (t_i, t_{i+k+1})$; $B_{i,k}(x) = 1 \iff t_{i+1} = \dots = t_{i+k} = x$;

- (f) $B_{i,k}$ este continuă (și indefinit derivabilă) la dreapta $\forall x \in \mathbb{R}$ (precizăm că $B_{i,k}(x) = 0, \forall x \in [t_0, t_m]$).

Demonstrație. (a), (b), (c), (d) și (f) sunt evidente pentru $k = 0$; (a), (b), (c) și (f) pentru $B_{i,k}$ se deduc prin recurență după k . Să demonstrăm (d). Fie $x \in [a, b]$; există un j , $k \leq j \leq m - k - 1$ astfel încât $x \in [t_j, t_{j+1})$. Dacă $x = t_j$ și $B_{j,k}(x) = 1$ proprietatea este evidentă (conform (c)). În celelalte cazuri avem

$$\begin{aligned} \sum_{i=0}^{m-k-1} B_{i,k}(x) &= \sum_{i=j-k}^j B_{i,k}(x) \\ &= \sum_{i=j-k}^j [\omega_{i,k}(x)B_{i,k-1}(x) + (1 - \omega_{i+1,k}(x))B_{i+1,k-1}(x)] \\ &= \omega_{j-k,k}(x)B_{j-k,k-1}(x) + \sum_{i=j+1-k}^j B_{i,k-1}(x) \\ &\quad + (1 - \omega_{j+1,k}(x))B_{j+1,k-1}(x). \end{aligned}$$

Dar $B_{j-k,k-1}(x) = 0$ și $B_{j+1,k-1}(x) = 0$, întrucât $x \in [t_j, t_{j+1})$ (conform lui (b)) și

$$\sum_{i=j+1-k}^j B_{i,k-1}(x) = 1$$

(din ipoteza inducției).

(e) este o consecință imediată a lui (d) și (c). \square

Observația 4.2.3. 1. Deoarece B-splinele sunt nenegative ((b) și (c)) și suma lor este 1, spunem că ele formează o *partiție a unității*.

2. Fiecare $B_{i,k}$ are un suport unic.

3. În cazul când $t_{m-k} = \dots = t_m = b$, pentru ca (d) și (e) să fie adevărate pe întreg $[a, b]$ punem $B_{m-k-1,k}(b) = 1$.

4. B-splinele se pot defini și pentru un șir infinit de noduri t_i , în definiția fiecărui $B_{i,k}$ intervenind un număr finit de noduri.

5. Dacă t_i este un nod de multiplicitate $\geq k + 2(t_i = t_{i+k+1})$ are loc $B_{i,k} \equiv 0$; reciproca este de asemenea adevărată. \diamond

Propoziția 4.2.4 (Identitatea lui Marsden) Fie $n = m - k$ și

$$\Psi_{i,k}(t) = \begin{cases} (t_{i+1} - t) \dots (t_{i+k} - t), & \text{pentru } k \geq 1; \\ 1, & \text{pentru } k = 0. \end{cases} \quad (4.2.4)$$

Are loc relația

$$(x - t)^k = \sum_{i=0}^{n-1} \Psi_{i,k}(t) B_{i,k}(x). \quad (4.2.5)$$

Demonstrație. Se face prin inducție după k . Cazul $k = 0$ este trivial (propoziția 4.2.2 (d)). Presupunem că

$$\sum_{i=0}^{n-1} \Psi_{i,k}(t) B_{i,k}(x) = (x - t)^{k-1}$$

este adevărată. Avem $B_{0,k-1}(x) = 0$ pentru $x \in [a, b]$.

$$\begin{aligned}
\sum_{i=0}^{n-1} \Psi_{i,k}(t) B_{i,k}(x) &= \\
&= \sum_{i=0}^{n-1} \Psi_{i,k}(t) [\omega_{i,k}(x) B_{i,k-1}(x) + (1 - \omega_{i+1,k}(x)) B_{i+1,k-1}(x)] \\
&= \Psi_{0,k}(t) \omega_{0,k}(x) B_{0,k-1}(x) \\
&\quad + \sum_{i=1}^{n-1} B_{i,k-1}(x) [\Psi_{i,k}(t) \omega_{i,k}(x) + \Psi_{i-1,k}(t) (1 - \omega_{i,k}(x))].
\end{aligned}$$

Dacă $t_i = t_{i+k}$, atunci $B_{i,k-1} \equiv 0$ și dacă $t_i < t_{i+k}$ avem

$$\begin{aligned}
\Psi_{i,k}(t) \omega_{i,k}(x) + \Psi_{i-1,k}(t) (1 - \omega_{i,k}(x)) &= \\
&= \Psi_{i,k-1}(t) [(t_{i+k} - t) \omega_{i,k}(x) + (t_i - t) (1 - \omega_{i,k}(x))] \\
&= \Psi_{i,k-1}(t) (\omega_{i,k}(x) (t_{i+k} - t_i) + t_i - t) = \Psi_{i,k-1}(t) (x - t).
\end{aligned}$$

Deci

$$\sum_{i=0}^{n-1} \Psi_{i,k}(t) B_{i,k}(x) = (x - t) \sum_{i=1}^{n-1} \Psi_{i,k-1}(t) B_{i,k-1}(x) = (x - t)^k,$$

ultima egalitate având loc pe baza ipotezei inducției. \square

Observația 4.2.5. Fie $t = (t_0, \dots, t_m)$ și $\mathbb{P}_{k,t}([a, b]) = \mathbb{P}_{k,t}$ spațiul funcțiilor polinomiale pe porțiuni pe $[a, b]$ de grad $\leq k$ și cu racord de clasă C^{k-p_j} în t_j dacă t_j este un nod de multiplicitate p_j . Prin convenție, un racord de clasă C^{k-p_j} , cu $k - p_j < 0$ nu impune nici un fel de condiții asupra lui t_j .

Dacă toate nodurile sunt de multiplicitate $\leq k+1$, atunci funcțiile $B_{i,k}$, $i = \overline{0, n-1}$ formează o bază a lui $\mathbb{P}_{k,t}$ și $\dim \mathbb{P}_{k,t} = n = m - k$. Dacă anumite noduri au multiplicitate $\geq k+2$, atunci $\mathbb{P}_{k,t} = \langle B_{i,k} | i = \overline{0, n-1} \rangle$, dar $\{B_{i,k}\}$ nu formează o bază. \diamond

4.2.2. Algoritmul de evaluare a unui B-spline

Fie $x \geq t_k$ și $S(x) = \sum_{i=0}^{n-1} a_i B_{i,k}(x)$. Aplicând direct definiția inductivă a B-splineelor se obține:

Propoziția 4.2.6 *Are loc*

$$S(x) = \sum a_i^{(0)}(x) B_{i,k}(x) = \sum a_i^{(1)}(x) B_{i,k-1}(x) = \dots = \sum a_i^{(k)}(x) B_{i,0}(x)$$

cu

$$\begin{aligned}
a_i^{(0)}(x) &= a_i \\
a_i^{(r+1)}(x) &= \omega_{i,k-r}(x) a_i^{(r)}(x) + (1 - \omega_{i,k-r}(x)) a_{i-1}^{(r)}(x).
\end{aligned}$$

Observația 4.2.7. 1. Algoritmul este valabil și pentru $x < t_k$, considerând că $B_{i,k} \equiv 0$ pentru $i \leq 0$.

2. Dacă evaluăm $S(x)$ pentru $x \in [t_j, t_{j+1})$, avem $S(x) = a_j^k(x)$, deoarece $B_{j,0}(x)$ este funcția caracteristică a intervalului $[t_j, t_{j+1})$ și pentru a calcula $a_j^{(k)}(x)$ este suficient să evaluăm $a_i^{(r)}$ numai pentru $j - k + r - 1 < i \leq j$, celelalte $B_{i,k-r}(x)$ fiind nule pentru $x \in [t_j, t_{j+1})$. Avem așadar

$$S(x) = \sum_{i=j-k}^j a_i^{(0)}(x) B_{i,k}(x) = \sum_{i=j-k+1}^j a_i^{(1)}(x) B_{i,k-1}(x) = \cdots = a_j^{(k)}(x). \quad \diamond$$

Algoritmul se prezintă, așadar, sub formă triunghiulară, fiecare element obținându-se prin combinația convexă a două elemente din linia anterioară:

$$\begin{array}{cccccccc}
 a_{j-k} & & a_{j-k+1} & & \cdots & & a_{j-1} & & a_j \\
 & a_{j-k+1}^{(1)} & & \cdots & & \cdots & & a_j^{(1)} & \\
 & & \ddots & & & & & \ddots & \\
 & & & a_{j-1}^{(k-1)} & & a_{j-1}^{(k-1)} & & & \\
 & & & & a_j^{(k)} & & & &
 \end{array}$$

(Algoritmul Cox-DeBoor)

Nodurile t formează o diviziune

$$\Delta : t_0 \leq t_1 \leq \cdots \leq t_k \leq a \leq \cdots \leq b \leq t_n \leq t_{n+1} \leq t_{n+k}.$$

Vom presupune că multiplicitatea oricărui nod $\leq k + 1$. O alegere frecventă este

$$t_0 = t_1 = \cdots = t_k = a < t_{k+1} \leq \cdots \leq t_{n-1} < b = t_n = \cdots = t_{n+k}.$$

În cazul particular $t_0 = \cdots = t_k = a < b = t_{k+1} = \cdots = t_{2k+1}$, funcțiile B-spline devin polinoame fundamentale Bernstein

$$p_{i,k}(x) = \binom{k}{i} x^i (1-x)^{k-i}.$$

4.2.3. Aplicații în grafica pe calculator

Curbe B-spline și Bézier

Fie $P_0, \dots, P_{n-1} \in \mathbb{R}^s$, $s \in \mathbb{N}^*$.

Definiția 4.2.8 (a) Se numește curbă B-spline asociată poligonului de control (P_0, \dots, P_{n-1}) curba parametrică (γ) definită prin

$$(\gamma) \quad S(t) = \sum_{i=0}^{n-1} P_i B_{i,k}(t), \quad a \leq t \leq b.$$

(b) Se numește curbă Bézier ² asociată poligonului de control (P_0, \dots, P_{n-1}) curba parametrică

$$B(t) = \sum_{i=0}^k P_i p_{i,k}(t), \quad p_{i,k}(x) = \binom{k}{i} x^i (1-x)^{k-i}.$$

Propoziția 4.2.9 Curbele B-spline au următoarele proprietăți:

1. γ nu trece în general prin punctele P_i ; dacă $t_0 = \dots = t_k$ și $t_n = \dots = t_{n+k} = b$ avem $S(a) = P_0$ și $S(b) = P_{n-1}$; în acest caz γ este tangentă în P_0 și P_{n-1} la laturile poligonului de control.
2. γ este în învelitoarea convexă a punctelor P_0, \dots, P_{n-1} . Mai exact dacă $t_i \leq t \leq t_{i+1}$, $S(t)$ este în învelitoarea convexă a punctelor P_{i-k}, \dots, P_i . Curba Bézier este situată în învelitoarea convexă a punctelor P_0, \dots, P_{n-1} .
3. În cazul când nodurile t_i ($k+1 \leq i \leq n-1$) sunt simple, $\gamma \in C^{k-1}$ și este formată din n arce parametrice polinomiale, de grad $\leq k$.
4. γ este invariantă la deplasările lui \mathbb{R}^s : dacă h este o deplasare a lui \mathbb{R}^s , atunci punctele $h(P_i)$ sunt punctele de control ale curbei $h(S(t))$.
5. $S(t)$ este regularizantă: dacă P este poligonul de control al lui γ și dacă H este un hiperplan, avem $\text{card}(H \cap \gamma) \leq \text{card}(H \cap P)$.

Observația 4.2.10. Curbele B-spline au un caracter local în sensul că:

Pi re Bezier (1910-1999) Fiu  i nepot de ingineri, a studiat ingineria mecanic  la Universitatea din Paris, unde  n 1977 a ob inut doctoratul  n Matematic  . Timp de 42 de ani (1933-1975) a lucrat la Uzinele Renault. A introdus curbele care  i poart  numele. Este considerat fondator al domeniului CAD-CAM, realiz nd primul sistem de acest tip (UNISURF -  ncep nd cu 1960). Laureat al mai multor premii prestigioase, printre care „Steven Anson Coons” al ACM (Association for Computing Machinery).



- (a) Dacă X este un punct al curbei corespunzând unei valori t_0 a parametrului, atunci poziția lui X nu depinde decât de cel mult $k + 1$ puncte P_i , deoarece dacă $t_j \leq t_0 < t_{j+1}$ avem

$$X = S(t_0) = \sum_{i=0}^{n-1} P_i B_{i,k}(t_0) = \sum_{i=j-k}^k P_i B_{i,k}(t_0).$$

- (b) Fiecare P_j nu influențează curba $S(t)$ decât pentru valorile lui t pentru care $B_{j,k}(t) \neq 0$, adică $t_j \leq t < t_{j+k+1}$ (în particular dacă modificăm P_j numai o porțiune a curbei se modifică).
- (c) Caracterul local nu are loc în cazul curbelor Bézier: modificarea unui punct de control atrage modificarea întregii curbe. \diamond

Algoritmi pentru curbe B-spline

Algoritmul 4.1 de evaluare a curbelor B-spline este o simplă traducere a algoritmului de evaluare pentru funcții B-spline.

Algoritmul 4.1 Algoritmul Cox-deBoor pentru evaluarea unei curbe B-spline

Intrare: Poligonul de control P_i , $i = \overline{0, n-1}$, punctul t în care se face evaluarea și diviziunea t_i

Ieșire: Valoarea $S(t)$

- 1: Fie $S(t) = \sum_{i=0}^{n-1} P_i B_{i,k}(t)$ și să presupunem că $t_j \leq t < t_{j+1}$.
 - 2: **for** $i := j - k$ **to** j **do**
 - 3: $P_i^0(t) := P_i$;
 - 4: **end for**
 - 5: **for** $r := 0$ **to** $k - 1$ **do**
 - 6: **for** $i := j - k + r + 1$ **to** j **do**
 - 7: $P_i^{r+1} := \omega_{i,k-r}(t)P_i^r + (1 - \omega_{i,k-r}(t))P_{i-1}^r$
 - 8: $= \frac{(t - t_i)P_i^r(t) + (t_{i+k-r} - t)P_{i-1}^r}{t_{i+k-r} - t_i}$;
 - 9: **end for**
 - 10: **end for**
 - 11: $S_j(t) := P_k^k(t)$;
-

În cazul particular al unei curbe Bézier se obține un algoritm mult mai simplu, care nu depinde de j (vezi algoritmul 4.2). El este datorat lui Paul de Faget de Casteljau.

Dăm și două posibilități de alegere a punctelor diviziunii.

Algoritmul 4.2 Algoritmul de Casteljaou pentru evaluarea unei curbe Bezier**Intrare:** Poligonul de control $P_i, i = \overline{0, k}$ și punctul t în care se face evaluarea**Ieșire:** Valoarea $B(t)$

```

1: Fie  $B(t) = \sum_{i=0}^k P_i p_{i,k}(t)$ .
2: for  $i := 0$  to  $k$  do
3:    $P_i^0(t) := P_i$ ;
4: end for
5: for  $r := 0$  to  $k - 1$  do
6:   for  $i := r + 1$  to  $k$  do
7:      $P_i^{r+1} := (1 - t)P_{i-1}^r(t) + tP_i^r(t)$ ;
8:   end for
9: end for
10:  $B(t) := P_k^k(t)$ ;

```

Varianta 1: $t_j = j$;**Varianta 2:**

$$t_j = \begin{cases} 0, & 0 \leq j \leq k; \\ j - k, & k + 1 \leq j \leq n; \\ n - k, & j > n. \end{cases}$$

4.2.4. Exemple

Să desenăm curba quadratică corespunzând nodurilor $t_0 = t_1 = t_2 = 1, t_3 = 2, t_4 = 3, t_5 = 4$ și $t_6 = t_7 = t_8 = 5$ și al cărei poligon de control este dat în figura 4.3. Graficul curbei apare de asemenea în figura 4.3.

Dacă se alege o repartiție uniformă a nodurilor și dorim o curbă de clasă C^1 , putem alege $t_i = i, i \in \mathbb{Z}$. Avem $B_i(t) = B_{i+r}(t+r), \forall r \in \mathbb{Z}$ (figura 4.4), adică elementele bazei B-spline se deduc una din alta print-o translație întreagă. Să considerăm $S(t) = \sum_{i=-\infty}^{\infty} P_i B_{i,r}(t)$ și să punem $P_{6+i} = P_i$; avem $S(t-6) = S(t)$, iar curba este definită pe un interval arbitrar de lungime 6. Se obține în acest mod o curbă închisă (figura 4.5).

În figura 4.6 se dau două curbe Bézier de grad 3 pentru două forme diferite ale poligonului de control. Urmează acum un exemplu mai detaliat de aplicare a algoritmului Cox-deBoor. Fie nodurile $a = 0 = t_0 = t_1 = t_2 = t_3, t_4 = 1, t_5 = 2, b = 3 = t_6 = t_7 = t_8 = t_9$. Gradul curbei va fi $k = 3$. Poligonul de control are 6 puncte și este dat în figura 4.7. Avem $S(t) = \sum_{i=0}^5 P_i B_{i,3}(t)$; să calculăm $S(\tau)$ pentru

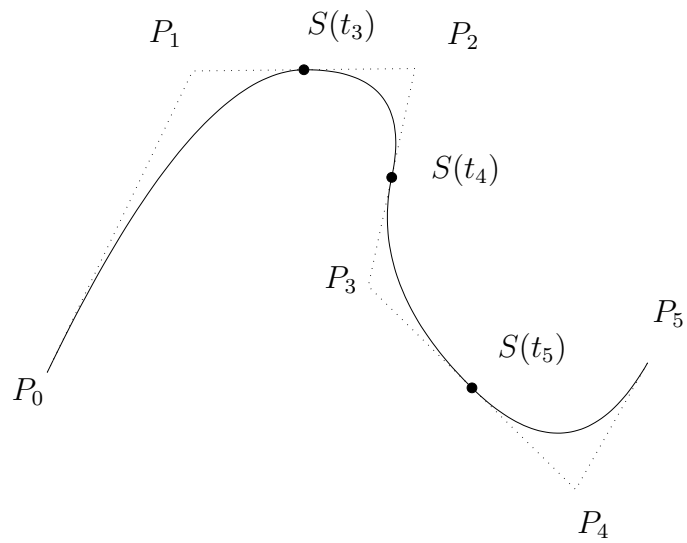
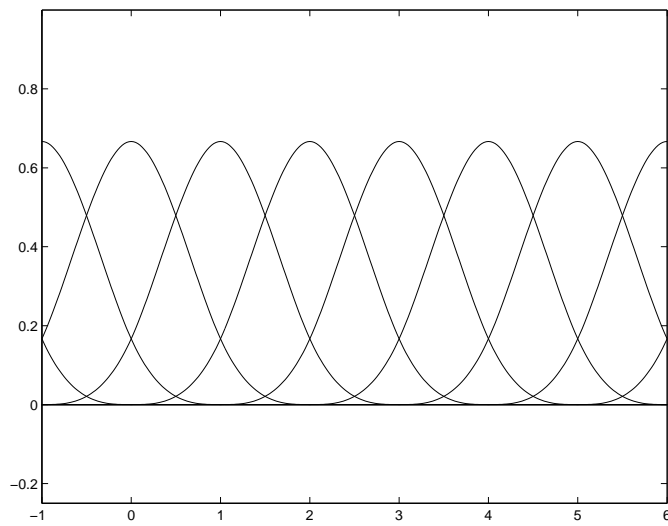


Figura 4.3: Curba B-spline cuadratică și poligonul ei de control

Figura 4.4: Bază pentru noduri din \mathbb{Z}

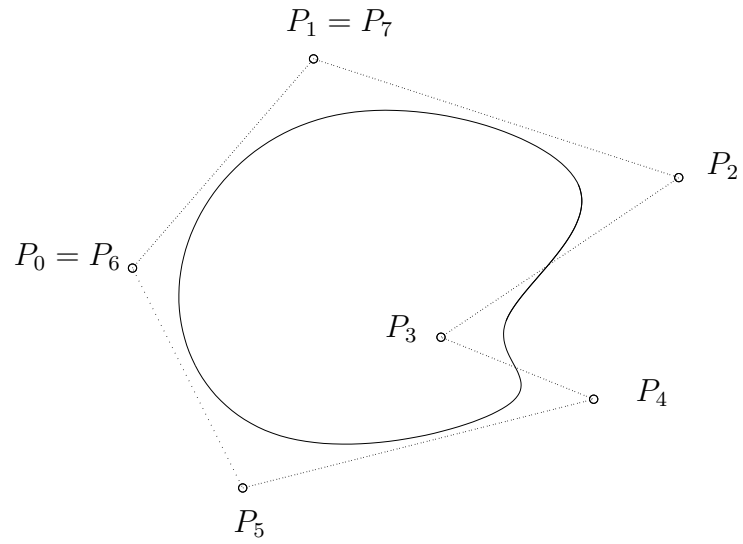


Figura 4.5: Curbă B-spline închisă, periodică

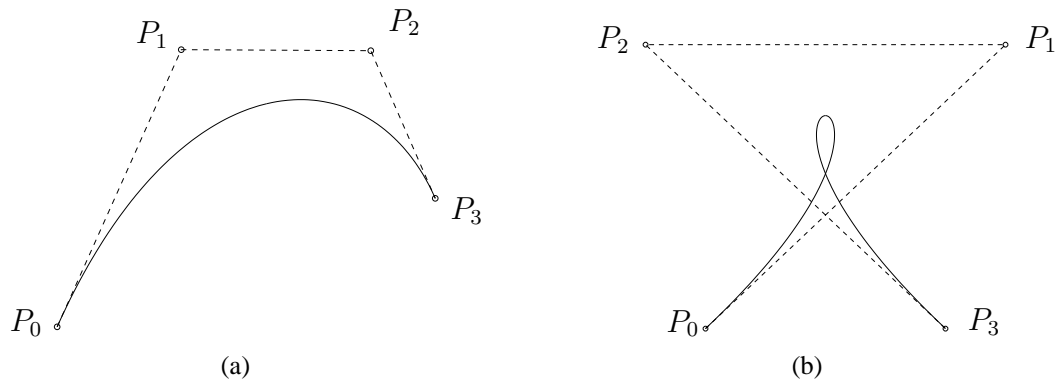


Figura 4.6: Două curbe Bézier de grad 3

$\tau = \frac{3}{2}$. Deoarece $t_4 < \tau < t_5$, calculele se organizează tabelar sub forma

$$\begin{array}{ccccccc} P_1 & & P_2 & & P_3 & & P_4 \\ & P_2^1 & & P_3^1 & & P_4^1 & \\ & & P_3^2 & & P_4^2 & & \\ & & & P_4^3 & & & \end{array}$$

unde $P_4^3 = S\left(\frac{3}{2}\right)$, iar

$$\begin{array}{lll} P_2^1 = \frac{3}{4}P_2 + \frac{1}{4}P_1 & P_3^1 = \frac{1}{2}P_1 + \frac{1}{2}P_2 & P_4^1 = \frac{1}{4}P_4 + \frac{3}{4}P_3 \\ P_3^2 = \frac{3}{4}P_3^1 + \frac{1}{4}P_2^1 & P_4^2 = \frac{1}{4}P_4^1 + \frac{3}{4}P_3^1 & \\ P_4^3 = \frac{1}{2}P_4^2 + \frac{1}{2}P_3^2 & & \end{array}$$

Graficul curbei apare în figura 4.7. Se observă că $S(t)$ este tangentă în punctul $P_j^k(t)$ la segmentul $P_k^{k-1}P_j^{k-1}$. Să ilustrăm acum caracterul local al curbelor B-spline. În figura

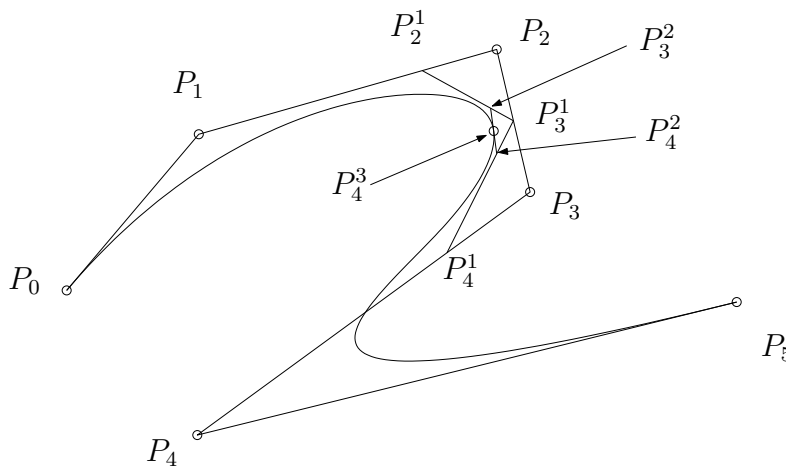


Figura 4.7: Un exemplu de aplicare a algoritmului Cox-deBoor

4.8 este reprezentată o curbă B-spline de gradul 3 corespunzătoare unui poligon de control cu 9 puncte (linie continuă). Se modifică coordonatele punctului P_3 , iar poligonul de control modificat și B-spline-ul corespunzător apar cu linie întreruptă.

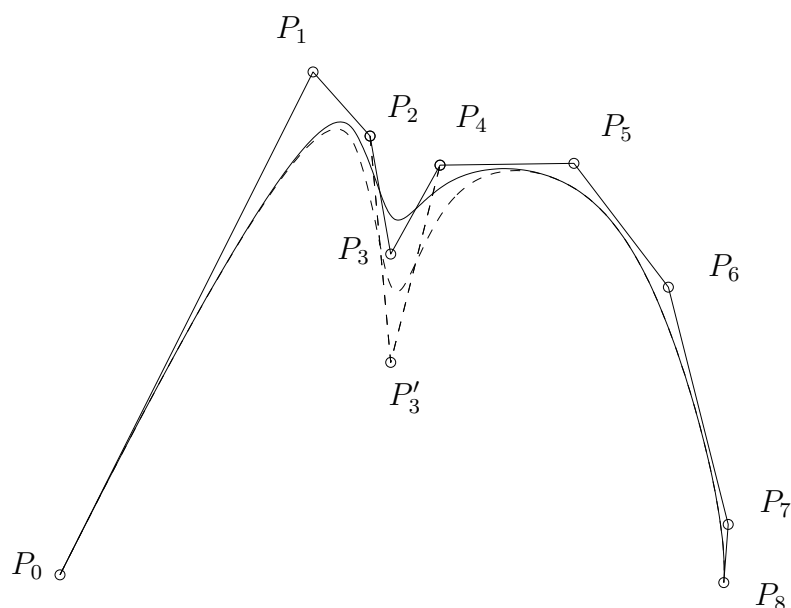


Figura 4.8: Caracterul local al curbelor B-spline. Se ilustrează efectul modificării coordonatelor unui punct (P_3)

4.3. Funcții spline cu variație diminuată

Funcțiile spline cu variație diminuată au fost introduse în 1964 de Isaac J. Schoenberg³ ca o generalizare a polinoamelor Bernstein.

Fie diviziunea

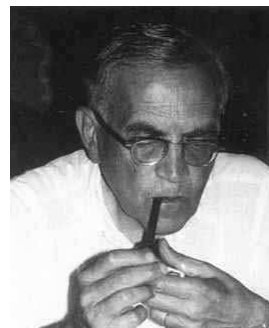
$$\Delta : t_0 \leq t_1 \leq \dots \leq t_k \leq a \leq \dots \leq b \leq t_n \leq t_{n+1} \leq t_{n+k}.$$

Vom presupune că multiplicitatea oricărui nod nu depășește $k + 1$.

Definiția 4.3.1 Fie $f : [t_0, t_{n+k}] \mapsto \mathbb{R}$ și diviziunea $\Delta = (t_i)_{i=\overline{0, n+k}}$ ca mai sus. Punem

$$\xi_i = \frac{t_{i+1} + \dots + t_i + k}{k}, \quad i = \overline{0, n-1} \quad (4.3.1)$$

Isaac J. Schoenberg (1913-1990) - matematician născut la Galați. A studiat la universitățile din Iași, Berlin și Göttingen. În 1926 își susține doctoratul la Universitatea din Iași. Din 1930 a activat în Statele Unite (1941-1966 University of Pennsylvania, 1966-1973 University of Wisconsin). Contribuții în domeniul Teoriei aproximării. Rezultatele sale în domeniul funcțiilor spline l-au făcut celebru (de fapt el a introdus termenul de funcție spline). A fost și un om sensibil, de aleasă cultură.



și definim operatorul S_Δ

$$(S_\Delta f)(x) = \sum_{i=0}^{n-1} f(\xi_i) B_{i,k}(x). \quad (4.3.2)$$

Acest operator se numește operator spline cu variație diminuată sau operatorul lui Schoenberg, iar $S_\Delta f$ se numește funcție spline cu variație diminuată.

$S_\Delta f$ fiind o combinație liniară de funcții B-spline este un spline de gradul k . Pentru punctele ξ avem

$$a = \xi_1 < \dots < \xi_{n-1} = b.$$

Propoziția 4.3.2 S_Δ este un operator liniar și pozitiv și $\forall f \in \mathbb{P}_1$ avem $S_\Delta f = f$.

Demonstrație. Liniaritatea și pozitivitatea rezultă din definiție folosind proprietățile B-splinelor. Trebuie să arătăm că S_Δ reproduce pe 1 și pe x . $S_\Delta 1 = 1$ rezultă din partiția unității. Din identitatea lui Marsden (4.2.4) egalând coeficienții lui t^{k-1} obținem

$$\sum_{i=0}^{n-1} \xi_i B_{i,k}(x) = x,$$

de unde $S_\Delta x = x$. \square

Propoziția 4.3.3 $S_\Delta f$ are următoarea proprietate

$$V(S_\Delta f - \ell) \leq V(f - \ell), \text{ pe } [a, b] \forall \ell \in \mathbb{P}_1$$

unde $V(g)$ este numărul variațiilor de semn ale lui g .

Numele de spline cu variație diminuată vine tocmai de la această proprietate.

Teorema 4.3.4 Are loc inegalitatea

$$\|g - S_\Delta g\|_\infty \leq \frac{k^2}{2} \|\Delta\|^2 \|g''\|_\infty. \quad (4.3.3)$$

Demonstrație. Fie $x_0 \in (a, b)$; trebuie estimată cantitatea $|g(x_0) - S_\Delta f(x_0)|$. Presupunem că $x_0 \in [t_i, t_{i+1})$, ceea ce implică $\xi_{i-k} \leq x_0 < \xi_i$ și atrage $B_{j,k}(x_0) = 0$ pentru $j \notin [i-k, i]$. Fie $P(x) = g(x_0) + g'(x_0)(x - x_0)$ ecuația tangentei la g în x_0 . Avem $S_\Delta P = P$ și deci

$$S_\Delta(g - P) = S_\Delta g - S_\Delta P = S_\Delta g - P = (S_\Delta g - g) + g - P$$

și pentru $x = x_0$

$$g(x_0) - S_{\Delta}g(x_0) = -S_{\Delta}(g - P)(x_0) = - \sum_{j=i-k}^i (g - P)(\xi_j) B_{j,k}(x_0)$$

de unde

$$|g(x_0) - S_{\Delta}g(x_0)| \leq \sup\{|(g - P)(x)| : x \in (\xi_{i-k}, \xi_i)\}.$$

Dar din formula lui Taylor $(g - P)(x) = \frac{(x-x_0)^2}{2} g''(\xi)$ cu $\xi \in (x_0, x)$ și $|x - x_0| \leq k\|\Delta\|$ pentru $\xi \in (\xi_{i-k}, \xi_i)$, de unde

$$|g(x_0) - S_{\Delta}g(x_0)| \leq \frac{k^2}{2} \|\Delta\|^2 \|g''\|_{\infty}.$$

□

Corolarul 4.3.5

$$\|x^2 - S_{\Delta}e_2\|_{\infty} \leq k^2 \|\Delta\|^2.$$

Definiția 4.3.6 Fie $f : [a, b] \rightarrow \mathbb{R}$. Formula

$$F = S_{\Delta}f + R_{\Delta}f \quad (4.3.4)$$

se numește formula de aproximare spline cu variație diminuată, iar $R_{\Delta}f$ este termenul rest.

Teorema 4.3.7 Dacă $f \in C^2[a, b]$ atunci

$$(R_{\Delta}f)(x) = \int_a^b \varphi(x; t) f''(t) dt \quad (4.3.5)$$

unde

$$\varphi(x; t) = (x - t)_+ - \sum_{i=0}^{n-1} B_{i,k}(x) (\xi_i - t)_+$$

respectiv

$$(R_{\Delta}f)(x) = -\frac{1}{2} [S_{\Delta}e_2(x) - x^2] f''(\xi).$$

Demonstrație. Se aplică teorema lui Peano și se ține cont că $\varphi(x; t)$ are semn constant pentru $t, x \in [a, b]$. □

Observația 4.3.8. În cazul particular când $t_0 = \dots = t_k = a < b = t_{k+1} \dots = t_{2k+1}$ se obține operatorul Bernstein ($S_{\Delta}f = B_k f$). ◇

4.4. Operatori liniari și pozitivi

Definiția 4.4.1 Operatorul $L : X \rightarrow Y$ se numește liniar și pozitiv dacă pentru orice $f, g \in X$ și orice scalari α, β

$$\begin{aligned} L(\alpha f + \beta g) &= \alpha Lf + \beta Lg \\ f \geq 0 &\Rightarrow Lf \geq 0. \end{aligned}$$

Observația 4.4.2. X și Y sunt spații liniare ordonate. De obicei $X = \{f : E_1 \subset \mathbb{R} \rightarrow \mathbb{R}\}$, $Y = \{g = Lf : E_2 \subset \mathbb{R} \rightarrow \mathbb{R}\}$. \diamond

Teorema 4.4.3 Dacă $L : X \rightarrow Y$ este un operator liniar și pozitiv și $f, g \in X$, atunci

(i) $f \leq g \Rightarrow Lf \leq Lg$ (monotonie);

(ii) $|Lf| \leq L|f|$.

Observația 4.4.4. Pentru a pune în evidență faptul că un operator liniar L se aplică unei funcții f , ca funcție de variabila independentă t , iar valoarea funcției obținute se consideră pe punctul x vom scrie $L(f(t); x)$ în loc de $(Lf)(x)$. \diamond

Teorema care urmează, numită prima teoremă a lui Korovkin sau teorema Bohman-Popoviciu ⁴-Korovkin, este un criteriu de convergență uniformă pentru șirurile construite cu ajutorul operatorilor liniari și pozitivi.

Teorema 4.4.5 Fie (L_m) , $L_m : C[a, b] \rightarrow C[a, b]$, un șir de operatori liniari și pozitivi. Dacă

$$\begin{aligned} (L_m e_0)(x) &= 1 + u_m(x); \\ (L_m e_1)(x) &= x + v_m(x); \\ (L_m e_2)(x) &= x^2 + w_m(x) \end{aligned} \tag{4.4.1}$$

Tiberiu Popoviciu (1906-1975) mare matematician român, membru al Academiei române, fondatorul școlii de Analiză numerică și Teoria aproximării din Cluj-Napoca. Doctorat în Franța, la École Normale Supérieure (1933). Contribuții importante și de pionierat în domeniul Analizei matematice, Teoriei funcțiilor, Teoriei aproximării și Analizei numerice. A contribuit de asemenea la începuturile și dezvoltarea informaticii românești și clujene.



și $\lim_{m \rightarrow \infty} u_m(x) = 0$, $\lim_{m \rightarrow \infty} v_m(x) = 0$ și $\lim_{m \rightarrow \infty} w_m(x) = 0$, uniform pe $[a, b]$, atunci pentru orice $f \in C[a, b]$

$$\lim_{m \rightarrow \infty} (L_m f)(x) = f(x)$$

uniform pe $[a, b]$.

Demonstrație. Fie $M = \max_{a \leq x \leq b} |f(x)|$ și t un punct oarecare, dar fix, din intervalul $[a, b]$. Funcția f fiind uniform continuă pe $[a, b]$, $\forall \varepsilon > 0$, $\exists \delta > 0$, astfel încât pentru $x \in [a, b]$ și $|t - x| < \delta$ să avem $|f(t) - f(x)| < \frac{\varepsilon}{2}$. Fie $J = \{x \in [a, b] : |t - x| \geq \delta\}$. Pentru $x \in J$ avem

$$|f(t) - f(x)| \leq 2M \leq \frac{2M}{\delta^2} (t - x)^2$$

prin urmare

$$|f(t) - f(x)| \leq \frac{\varepsilon}{2} + \frac{2M}{\delta^2} (t - x)^2, \quad x, t \in [a, b]. \quad (4.4.2)$$

Cu notația din observația 4.4.4 ($L_m(f(t); x)$ în loc de $(L_m f)(x)$) pentru $t = x$ avem $L_m(f(x); x) = f(x) \cdot L_m(1; x) = f(x)(L_m e_0)(x)$. Din enunț (egalitățile (4.4.1)) se observă că putem scrie

$$\begin{aligned} (L_m f)(x) - f(x) &= L_m(f(t); x) - f(x)L_m(1; x) + f(x)u_m(x) \\ &= L_m(f(t) - f(x); x) + f(x)u_m(x). \end{aligned}$$

astfel că avem

$$|(L_m f)(x) - f(x)| \leq |L_m(f(t) - f(x); x)| + |f(x)||u_m(x)|. \quad (4.4.3)$$

Folosind liniaritatea, teorema 4.4.3 și inegalitatea (4.4.2) deducem că

$$\begin{aligned} |L_m(f(t) - f(x); x)| &\leq L_m(|f(t) - f(x)|; x) < L_m\left(\frac{\varepsilon}{2} + \frac{2M}{\delta^2}(t - x)^2; x\right) \\ &= \frac{\varepsilon}{2}L_m(1; x) + \frac{2M}{\delta^2}L_m((t - x)^2; x). \end{aligned}$$

Ultima inegalitate și (4.4.3) ne dau

$$\begin{aligned} |(L_m f)(x) - f(x)| &< \frac{\varepsilon}{2}L_m(1; x) + \frac{2M}{\delta^2}L_m((t - x)^2; x) + |f(x)||u_m(x)| \\ &< \frac{\varepsilon}{2} + \left(|f(x)| + \frac{\varepsilon}{2}\right)|u_m(x)| + \frac{2M}{\delta^2}L_m((t - x)^2; x). \end{aligned} \quad (4.4.4)$$

Din enunț mai rezultă că

$$L_m((t - x)^2; x) = x^2 u_m(x) - 2x v_m(x) + w_m(x). \quad (4.4.5)$$

Deoarece $u_m \rightrightarrows 0$, $v_m \rightrightarrows 0$, $w_m \rightrightarrows 0$,

$$\left(|f(x) + \frac{\varepsilon}{2}\right) |u_m(x)| + \frac{2M}{\delta^2} L_m((t-x)^2; x) < \frac{\varepsilon}{2}, \text{ pentru } m > N_\varepsilon$$

din (4.4.4) se obține

$$|(L_m f)(x) - f(x)| < \frac{\varepsilon}{2}.$$

□

Corolarul 4.4.6 Dacă (L_m) , $L_m : C[a, b] \rightarrow C[a, b]$, $m \in \mathbb{N}^*$ este un șir de operatori liniari și pozitivi și

$$\lim_{m \rightarrow \infty} (L_m)(1; x) = 1, \quad \lim_{m \rightarrow \infty} (L_m)((t-x)^2; x) = 0 \quad (4.4.6)$$

pentru orice $f \in C[a, b]$ șirul $(L_m f)$ converge uniform la f pe $[a, b]$.

Concluzia rezultă (4.4.4) și (4.4.5).

Observația 4.4.7. Funcțiile e_0, e_1, e_2 din teorema 4.4.5 se numesc funcții de probă. Numărul de funcții de probă nu poate fi modificat. ◇

Un rezultat asemănător cu cel din teorema 4.4.5 are loc și pentru funcțiile 2π -periodice cu funcțiile de probă $1, \cos x, \sin x$, iar corolarul corespunzător se aplică pentru 1 și $\sin^2 \frac{t-x}{2}$. Fie $C_{2\pi}$ spațiul funcțiilor continue 2π -periodice.

Teorema 4.4.8 (teorema a doua a lui Korovkin) Fie (L_m) un șir de operatori liniari pozitivi $L_m : C_{2\pi} \rightarrow C_{2\pi}$. Dacă $\forall x \in \mathbb{R}$

$$\begin{aligned} L_m(1; x) &= 1 + u_m(x); \\ L_m(\sin t; x) &= \sin x + v_m(x); \\ L_m(\cos t; x) &= \cos x + w_m(x) \end{aligned}$$

și

$$\lim_{m \rightarrow \infty} u_m(x) = \lim_{m \rightarrow \infty} v_m(x) = \lim_{m \rightarrow \infty} w_m(x) = 0$$

uniform pe \mathbb{R}_+ , atunci $\forall f \in C_{2\pi}$ $L_m f \rightrightarrows f$ pe \mathbb{R} .

Demonstrație. Este analogă cu a teoremei 4.4.5, analoaga inegalității (4.4.2) fiind

$$|f(t) - f(x)| < \frac{\varepsilon}{2} + 2M \sin^{-2} \frac{\delta}{2} \sin^2 \frac{t-x}{2}.$$

□

Exemple

Exemplul 4.4.9 (Operatorul lui Bernstein). Fie (B_m) șirul operatorilor Bernstein. Avem

$$B_m(1; x) = 1, \quad B_m(t, x) = x, \quad B_m(t^2; x) = x^2 + \frac{x(1-x)}{m}, \quad x \in [0, 1].$$

Deci $u_m(x) = 0$, $v_m(x) = 0$, $w_m(x) = \frac{x(1-x)}{m}$. Deoarece $\lim_{m \rightarrow \infty} \frac{x(1-x)}{m} = 0$ uniform pe $[a, b]$, $B_m f \rightrightarrows f$ pe $[0, 1]$, rezultat echivalent cu teorema 4.1.2. \diamond

Exemplul 4.4.10 (Operatorul lui Schoenberg). Operatorul spline cu variație dimiunată S_Δ este liniar și pozitiv. De asemenea

$$\begin{aligned} S_\Delta(1; x) &= 1; \\ S_\Delta(t; x) &= x; \\ S_\Delta(t^2; x) &= x^2 + E(x). \end{aligned}$$

Deoarece $E(x) \rightrightarrows 0$ când $\|\Delta\| \rightarrow 0$, $S_\Delta f \rightrightarrows f$. \diamond

Exemplul 4.4.11 (Operatorul Hermite-Fejér). Pornind de la operatorul de interpolare Hermite cu noduri duble rădăcini ale polinomului Cebâșev de speța I, T_{m+1}

$$x_k = \cos \frac{2k+1}{2(m+1)}\pi, \quad k = \overline{0, m},$$

$$(H_{2m+1}f)(x) = \sum_{k=0}^m h_{k0}(x)f(x_k) + \sum_{k=0}^m h_{k1}(x)f(x_k)$$

și omițând a doua sumă, Fejér⁵ a obținut operatorul

$$(F_{2m+1}f)(x) = \sum_{k=0}^m h_k(x)f(x_k),$$

unde

$$h_k(x) = h_{k0}(x) = (1 - x_k x) \left[\frac{T_{m+1}(x)}{(m+1)(x - x_k)} \right]^2.$$



Leopold Fejér (1880-1959) Matematician maghiar, lider al generației sale. A avut contribuții importante în domeniul aproximării și interpolării în real și complex și teoriei seriilor Fourier. Rezultatul său privind convergența seriilor Fourier a fost obținut când era încă student. A funcționat ca profesor și la universitatea din Cluj.

Propoziția 4.4.12 *Operatorul lui Fejér are următoarele proprietăți:*

1. $(F_{2m+1}f)(x_k) = f(x_k)$, $(F_{2m+1}f)'(x_k) = 0$ $\sum_{k=0}^m h_k(x) = 1$.
2. F_{2m+1} este un operator liniar și pozitiv.
3. Dacă $f \in C[-1, 1]$, atunci $F_{2m+1}f \Rightarrow f$ pe $[-1, 1]$, când $m \rightarrow \infty$.

Demonstrație. 1. Rezultă din proprietățile polinoamelor de interpolare Hermite.

2. Liniaritatea rezultă din definiție, iar pozitivitatea din

$$1 - x_k x \geq 1 - |x_k| > 0, \quad x \in [-1, 1].$$

3. Din proprietatea 1 rezultă că $F_{2m+1}(1; x) = 1$, pentru $x \in [-1, 1]$.

$$\begin{aligned} F_{2m+1}((t-x)^2; x) &= \sum_{k=0}^m (1 - x_k x) \left[\frac{T_{m+1}(x)}{(m+1)(x-x_k)} \right]^2 (x_k - x)^2 \\ &= \frac{1}{(m+1)^2} T_{m+1}^2(x) \sum_{k=0}^m (1 - x_k x). \end{aligned}$$

Datorită simetriei nodurilor x_k față de origine, $\sum_{k=0}^m x_k = 0$ și se obține

$$F_{2m+1}((t-x)^2; x) = \frac{1}{m+1} T_{m+1}^2(x) \leq \frac{1}{m+1};$$

adică

$$\lim_{m \rightarrow \infty} F_{2m+1}((t-x)^2; x) = 0, \quad \text{uniform pe } [-1, 1].$$

Se aplică apoi corolarul 4.4.6 la teorema 4.4.5.

□

Fiind un operator polinomial, operatorul lui Fejér poate da o demonstrație constructivă a teoremei lui Weierstrass. ◇

4.5. Cea mai bună aproximare uniformă

Fie $f \in C[a, b]$ și presupunem că un polinom de grad cel mult n minimizează norma $\|f - p_n\|_\infty$. Un astfel de polinom se numește *polinomul de cea mai bună aproximare uniformă* a lui f .

În teoria aproximării uniforme rolul central este jucat de *punctele de alternanță Cebâșev* pentru funcție $R(x) = f(x) - p_n(x)$. Punctele de alternanță de grad m sunt nodurile unei grile (diviziuni)

$$a \leq x_1 \leq \dots \leq x_m \leq b$$

având următoarele proprietăți:

$$(1) |R(x_i)| = \max_{a \leq x \leq b} |R(x)|, \quad i = \overline{1, m}.$$

$$(2) R(x_i)R(x_{i+1}) < 0, \quad i = \overline{1, m-1}.$$

Vom nota mulțimea tuturor diviziunilor de acest tip pe $[a, b]$ prin $\mathcal{A}(m, a, b, R)$.

Teorema 4.5.1 (Cebîșev) Pentru ca un polinom p_n de grad cel mult n să fie cea mai bună aproximantă uniformă a lui $f \in C[a, b]$, este necesar și suficient ca $\mathcal{A}(n+2, a, b, R)$ să fie nevidă.

Demonstrația suficienței. Presupunem că există un polinom $q_n \in \mathbb{P}_n$ astfel încât

$$\|f - q_n\|_\infty \leq \|f - p_n\|_\infty.$$

Atunci în punctele de alternanță Cebîșev

$$|f(x_i) - q_n(x_i)| < |f(x_i) - p_n(x_i)|$$

ceea ce implică faptul că funcția $g(x) \equiv (f(x) - p_n(x)) - (f(x) - q_n(x))$ are în punctele x_i același semn ca $R(x) = f(x) - p_n(x)$. Deoarece semnele lui $R(x_i)$ alternează, există un zero în interiorul fiecărui subinterval $[x_i, x_{i+1}]$. Deci, g are $n+1$ zerori pe $[a, b]$. Aceasta nu se poate întâmpla dacă g nu este identic nul. \square

Exemplul 4.5.2. Vom determina polinomul de cea mai bună aproximare uniformă de grad întâi pentru funcția $f(x) = \sqrt{x}$ pe $[a, b] \subset \mathbb{R}_+$.

Polinomul căutat are forma $P_1^* = c_0 + c_1x$. Eroarea de aproximare este $e_1(x) = c_0 + c_1x - \sqrt{x}$. Derivata ei, $e_1'(x) = c_1 - \frac{1}{2\sqrt{x}}$ se anulează în $x_1 = \frac{1}{4c_1^2}$. Conform teoremei lui Cebîșev abaterea maximă se realizează în trei puncte din $[a, b]$ și obținem sistemul neliniar

$$\begin{cases} c_0 + c_1a - \sqrt{a} = E_1 \\ c_0 + \frac{1}{4c_1} - \frac{2}{c_1} = -E_1 \\ c_0 + c_1b - \sqrt{b} = E_1 \end{cases},$$

cu soluțiile

$$c_0 = \frac{1}{2} \left(\sqrt{a} - \frac{a}{\sqrt{a} + \sqrt{b}} + \frac{\sqrt{a} + \sqrt{b}}{4} \right),$$

$$c_1 = \frac{1}{\sqrt{a} + \sqrt{b}},$$

$$E_1 = c_0 + c_1a - \sqrt{a}. \quad \diamond$$

CAPITOLUL 5

Aproximarea funcțiilor liniare

Cuprins

5.1. Introducere	149
5.2. Derivare numerică	154
5.3. Integrare numerică	156
5.3.1. Formula trapezului și formula lui Simpson	157
5.3.2. Formule Newton-Cotes cu ponderi și formule de tip Gauss	161
5.3.3. Proprietăți ale cuadraturilor gaussiene	164
5.4. Cuadraturi adaptive	170
5.5. Cuadraturi iterate. Metoda lui Romberg	171
5.6. Cuadraturi adaptive II	174

5.1. Introducere

Fie X un spațiu liniar, L_1, \dots, L_m funcționale liniare reale, liniar independente, definite pe X și $L : X \rightarrow \mathbb{R}$ o funcțională liniară reală astfel încât L, L_1, \dots, L_m să fie liniar independente.

Definiția 5.1.1 O formulă de aproximare a funcției L în raport cu funcțiile L_1, \dots, L_m este o formulă de forma

$$L(f) = \sum_{i=1}^m A_i L_i(f) + R(f), \quad f \in X. \quad (5.1.1)$$

Parametrii reali A_i se numesc coeficienții formulei, iar $R(f)$ termenul rest.

Pentru o formulă de aproximare de forma (5.1.1), dându-se funcțiile L_i , se pune problema determinării coeficienților A_i și a studiului termenului rest corespunzător valorilor obținute pentru coeficienți.

Observația 5.1.2. Forma funcțiilor L_i depinde de informațiile deținute asupra lui f (ele exprimând de fapt aceste informații), dar și de natura problemei de aproximare, adică de forma lui L . \diamond

Exemplul 5.1.3. (E1) Dacă $X = \{f \mid f : [a, b] \rightarrow \mathbb{R}\}$, $L_i(f) = f(x_i)$, $i = \overline{0, m}$, $x_i \in [a, b]$ și $L(f) = f(\alpha)$, $\alpha \in [a, b]$, formula de interpolare Lagrange

$$f(\alpha) = \sum_{i=0}^m \ell_i(\alpha) f(x_i) + (Rf)\alpha$$

este o formulă de tip (5.1.1), cu coeficienții $A_i = \ell_i(\alpha)$, iar una din reprezentările posibile pentru rest este

$$(Rf)(\alpha) = \frac{u(\alpha)}{(m+1)!} f^{(m+1)}(\xi), \quad \xi \in [a, b]$$

dacă există $f^{(m+1)}$ pe $[a, b]$. \diamond

Exemplul 5.1.4. Dacă X și L_i sunt ca în exemplul 5.1.3 și există $f^{(k)}(\alpha)$, $\alpha \in [a, b]$, $k \in \mathbb{N}^*$, iar $L(f) = f^{(k)}(\alpha)$ se obține o formulă de aproximare a valorii derivatei de ordinul k a lui f în punctul α

$$f^{(k)}(\alpha) = \sum_{i=0}^m A_i f(x_i) + R(f),$$

numită și formulă de derivare numerică. \diamond

Exemplul 5.1.5. Dacă X este un spațiu de funcții definite pe $[a, b]$, integrabile pe $[a, b]$ și pentru care există $f^{(j)}(x_k)$, $k = \overline{0, m}$, $j \in I_k$, cu $x \in [a, b]$ și I_k mulțimi de indici date

$$L_{kj}(f) = f^{(j)}(x_k), \quad k = \overline{0, m}, \quad j \in I_k,$$

iar

$$L(f) = \int_a^b f(x)dx,$$

se obține formula

$$\int_a^b f(x)dx = \sum_{k=0}^m \sum_{j \in I_k} A_{kj} f^{(j)}(x_k) + R(f),$$

numită *formulă de integrare numerică*. ◇

Definiția 5.1.6 Dacă $\mathbb{P}_r \subset X$, numărul $r \in \mathbb{N}$ cu proprietatea că $\text{Ker}(R) = \mathbb{P}_r$ se numește grad de exactitate al formulei de aproximare (5.1.1).

Observația 5.1.7. Deoarece R este o funcțională liniară proprietatea $\text{Ker}(R) = \mathbb{P}_r$ este echivalentă cu $R(e_k) = 0$, $k = \overline{0, r}$ și $R(e_{r+1}) \neq 0$, unde $e_k(x) = x^k$. ◇

Putem acum să formulăm *problema generală de aproximare*: dându-se o funcțională liniară L pe X , m funcționale liniare L_1, L_2, \dots, L_m pe X și valorile lor ("datele") $\ell_i = L_i f$, $i = \overline{1, m}$ aplicate unei anumite funcții f și un subspațiu liniar $\Phi \subset X$ cu $\dim \Phi = m$, dorim să găsim o formulă de aproximare de tipul

$$Lf \approx \sum_{i=1}^m a_i L_i f \tag{5.1.2}$$

care să fie exactă (adică să aibă loc egalitate), pentru orice $f \in \Phi$.

Este natural (deoarece dorim să interpolăm) să facem următoarea

Ipoteză: „problema de interpolare“

Să se găsească $\varphi \in \Phi$ astfel încât

$$L_i \varphi = s_i, \quad i = \overline{1, m} \tag{5.1.3}$$

are o soluție unică $\varphi(\cdot) = \varphi(s, \cdot)$, pentru $s = [s_1, \dots, s_m]^T$, arbitrar.

Putem exprima ipoteza noastră mai explicit în termenii unei baze date $\varphi_1, \varphi_2, \dots, \varphi_m$ a lui Φ și a matricei Gram ¹ asociată

$$G = [L_i \varphi_j] = \begin{pmatrix} L_1 \varphi_1 & L_1 \varphi_2 & \dots & L_1 \varphi_m \\ L_2 \varphi_1 & L_2 \varphi_2 & \dots & L_2 \varphi_m \\ \dots & \dots & \dots & \dots \\ L_m \varphi_1 & L_m \varphi_2 & \dots & L_m \varphi_m \end{pmatrix} \in \mathbb{R}^{m \times m}. \quad (5.1.4)$$

Cerem ca

$$\det G \neq 0. \quad (5.1.5)$$

Este ușor de văzut că această condiție este independentă de alegerea particulară a bazei. Pentru a arăta că solvabilitatea unică a lui (5.1.3) și condiția (5.1.5) sunt echivalente, exprimăm φ din (5.1.3) ca o combinație liniară a funcțiilor de bază

$$\varphi = \sum_{j=1}^{nm} c_j \varphi_j \quad (5.1.6)$$

și observăm că condițiile de interpolare

$$L_i \left(\sum_{j=1}^m c_j \varphi_j \right) = s_i, \quad i = \overline{1, m}$$

pot fi scrise (ținând cont de liniaritatea lui L_i) sub forma

$$\sum_{j=1}^m c_j L_i \varphi_j = s_i, \quad i = \overline{1, m},$$

adică

$$Gc = s, \quad c = [c_1, c_2, \dots, c_m]^T, \quad s = [s_1, s_2, \dots, s_m]^T. \quad (5.1.7)$$

Aceasta are o soluție unică pentru s arbitrar dacă și numai dacă are loc (5.1.5). Avem două abordări pentru rezolvarea acestei probleme.

Jórgen Pedersen Gram (1850-1916), matematician danez, a studiat la Universitatea din Copenhaga. După absolvire a intrat la o companie de asigurări ca asistent calculator și apoi a promovat ¹treptat până a ajuns director. A avut contribuții în domeniul dezvoltării în serie a funcțiilor și aproximării Cebîșev și în sensul celor mai mici pătrate. „Determinantul Gram” a fost introdus de el în legătură cu studiile sale asupra liniar independenței.



Metoda interpolării. Rezolvăm problema generală de aproximare prin interpolare

$$Lf \approx L\varphi(\ell; \cdot), \quad \ell = [\ell_1, \ell_2, \dots, \ell_m]^T, \quad \ell_i = L_i f \quad (5.1.8)$$

Cu alte cuvinte aplicăm L nu lui f , ci soluției $\varphi(\ell; \cdot)$ a problemei de aproximare (5.1.3) în care $s = \ell$. Ipoteza noastră ne garantează că $\varphi(\ell; \cdot)$ este unic determinat. În particular, dacă $f \in \Phi$, atunci (5.1.8) are loc cu egalitate, deoarece $\varphi(\ell; \cdot) = f(\cdot)$, în mod trivial. Astfel, aproximanta noastră (5.1.8) satisface condițiile de exactitate cerute pentru (5.1.2). Rămâne doar să arătăm că (5.1.8) produce o aproximare de forma (5.1.2). Pentru aceasta să observăm că interpolantul în (5.1.8) este

$$\varphi(\ell; \cdot) = \sum_{j=1}^m c_j \varphi_j(\cdot)$$

unde vectorul $c = [c_1, c_2, \dots, c_m]^T$ satisface (5.1.7) cu $s = \ell$

$$Gc = \ell, \quad \ell = [L_1 f, L_2 f, \dots, L_m f]^T.$$

Scriind

$$\lambda_j = L\varphi_j, \quad j = \overline{1, m}, \quad \lambda = [\lambda_1, \lambda_2, \dots, \lambda_m]^T, \quad (5.1.9)$$

avem din liniaritatea lui L

$$L\varphi(\ell; \cdot) = \sum_{j=1}^m c_j L\varphi_j = \lambda^T c = \lambda^T G^{-1} \ell = [(G^T)^{-1} \lambda]^T \ell,$$

adică

$$L\varphi(\ell; \cdot) = \sum_{i=1}^m a_i L_i f, \quad a = [a_1, a_2, \dots, a_m]^T = (G^T)^{-1} \lambda. \quad (5.1.10)$$

Metoda coeficienților nedeterminați. Aici determinăm coeficienții din (5.1.3) astfel încât egalitatea să aibă loc $\forall f \in \Phi$, care, conform liniarității lui L și L_i este echivalentă cu egalitatea pentru $f = \varphi_1, f = \varphi_2, \dots, f = \varphi_m$, adică

$$\left(\sum_{j=1}^m a_j L_j \right) \varphi_i = L\varphi_i, \quad i = \overline{1, m},$$

sau conform (5.1.8)

$$\sum_{j=1}^m a_j L_j \varphi_i = \lambda_i, \quad i = \overline{1, m}.$$

Evident, matricea sistemului este G^T , deci

$$a = [a_1, a_2, \dots, a_m]^T = (G^T)^{-1}\lambda,$$

în concordanță cu (5.1.10). Astfel, metoda interpolării și cea a coeficienților nedeterminați sunt matematic echivalente – ele conduc la exact aceeași aproximare.

S-ar părea că, cel puțin în cazul polinoamelor (adică $\Phi = \mathbb{P}_d$), prima metodă este mai puternică, deoarece poate conduce la o expresie a erorii de interpolare (aplicând funcționala formulei de interpolare $f = a_n f + r_n f$). Dar și în cazul metodei coeficienților nedeterminați, din condiția de exactitate, se poate exprima restul cu ajutorul teoremei lui Peano (teorema 3.3.2).

5.2. Derivare numerică

Pentru simplitate vom considera doar derivata de ordinul I. Se pot aplica tehnici analoge și pentru alte derivate. Vom rezolva problema prin interpolare: în loc să derivăm $f \in C^{m+1}[a, b]$, vom deriva polinomul său de interpolare:

$$f(x) = (L_m f)(x) + (R_m f)(x). \quad (5.2.1)$$

Scriem polinomul de interpolare în forma Newton

$$\begin{aligned} (L_m f)(x) = (N_m f)(x) = & f_0 + (x - x_0)f[x_0, x_1] + \dots + \\ & + (x - x_0) \dots (x - x_{m-1})f[x_0, x_1, \dots, x_m] \end{aligned} \quad (5.2.2)$$

și restul sub forma

$$(R_m f)(x) = (x - x_0) \dots (x - x_m) \frac{f^{(m+1)}(\xi(x))}{(m+1)!}. \quad (5.2.3)$$

Derivând (5.2.2) în raport cu x și punând $x = x_0$ obținem

$$\begin{aligned} (L_m f)(x_0) = & f[x_0, x_1] + (x_0 - x_1)f[x_0, x_1, x_2] + \dots + \\ & + (x_0 - x_1)(x_0 - x_2) \dots (x_0 - x_{m-1})f[x_0, x_1, \dots, x_m]. \end{aligned} \quad (5.2.4)$$

Presupunând că f este continuu derivabilă pe un interval convenabil se obține pentru rest

$$(R_m f)'(x_0) = (x_0 - x_1) \dots (x_0 - x_m) \frac{f^{(m+1)}(\xi(x_0))}{(m+1)!}. \quad (5.2.5)$$

Deci, derivând (5.2.4) obținem

$$f'(x_0) = (L_m f)'(x_0) + \underbrace{(R_m f)'(x_0)}_{e_m}. \quad (5.2.6)$$

Dacă $H = \max_i |x_0 - x_i|$ eroarea are forma $e_m = O(H^m)$, când $H \rightarrow 0$.

Putem obține formule de aproximare de grad arbitrar, dar ele sunt de utilitate practică limitată.

Observația 5.2.1. Derivarea numerică este o operație critică și de aceea este bine să fie evitată pe cât posibil, deoarece chiar dacă aproximanta este bună, nu rezultă că derivata aproximantei este o aproximație bună a derivatei (vezi figura 5.1). Aceasta rezultă și din exemplul 5.2.2 \diamond

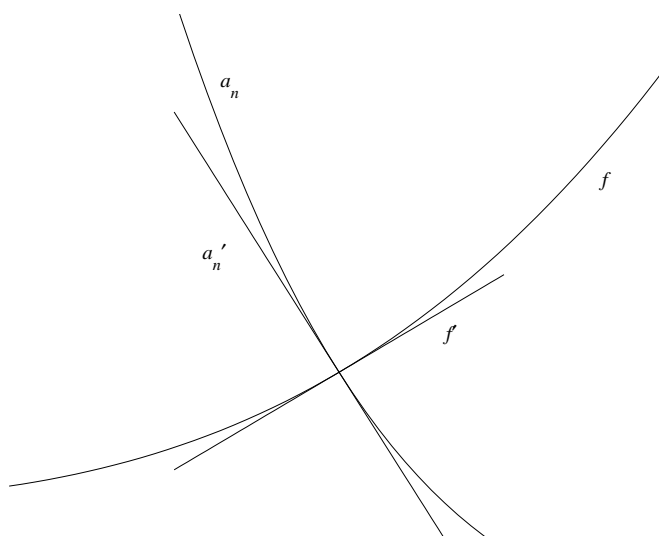


Figura 5.1: Neajunsurile derivării numerice

Exemplul 5.2.2. Fie funcția

$$f(x) = g(x) + \frac{1}{n} \sin n^2(x - a), \quad g \in C^1[a, b].$$

Se constată că $d(f, g) \rightarrow 0$ ($n \rightarrow \infty$), dar $d(f', g') = n \not\rightarrow 0$. \diamond

Formulele de derivare numerică sunt utile pentru deducerea unor metode numerice, în special pentru ecuații diferențiale ordinare și ecuații cu derivate parțiale.

Se pot folosi și alte procedee de aproximare: Taylor, Hermite, spline, metoda celor mai mici pătrate.

5.3. Integrare numerică

Problema este de a calcula integrala definită a unei funcții date pe un interval mărginit $[a, b]$. Dacă f are o comportare bună, aceasta este o problemă de rutină, pentru care metodele cele mai simple de integrare cum ar fi regula trapezelor sau regula repetată a lui Simpson sunt satisfăcătoare, prima având anumite avantaje asupra celei de-a doua în cazul când f este periodică, cu perioada $b - a$. Complicațiile apar atunci când f are singularități (dar rămâne integrabilă), sau când intervalul de integrare este nemărginit (care este o altă manifestare a comportării singulare). Descompunând integrala pe subintervale, dacă este necesar, în mai multe integrale, se poate presupune că singularitatea, dacă locul ei este cunoscut, este la unul sau la ambele capete ale intervalului $[a, b]$. Astfel de integrale improprii pot fi tratate ca și cuadraturi cu ponderi, adică să încorporăm singularitatea într-o pondere, care devine un factor al integrandului, lăsând celălalt factor să aibă o comportare bună. Cel mai important exemplu este formula lui Gauss relativă la o astfel de pondere. În fine, este posibil să se accelereze convergența unei scheme de cuadratură prin recombinații convenabile. Un astfel de exemplu este metoda lui Romberg.

Fie $f : [a, b] \rightarrow \mathbb{R}$ integrabilă pe $[a, b]$, $F_k(f)$, $k = \overline{0, m}$ informații despre f (de regulă funcționale liniare) și $w : [a, b] \rightarrow \mathbb{R}_+$ o funcție pondere integrabilă pe $[a, b]$.

Definiția 5.3.1 *O formulă de forma*

$$\int_a^b w(x)f(x)dx = Q(f) + R(f), \quad (5.3.1)$$

unde

$$Q(f) = \sum_{j=0}^m A_j F_j(f),$$

se numește formulă de integrare numerică a funcției f sau formulă de cuadratură. Parametrii A_j , $j = \overline{0, m}$ se numesc coeficienții formulei, iar $R(f)$ termenul rest al ei. Q se numește funcțională de cuadratură.

Definiția 5.3.2 *Numărul natural $d = d(Q)$ cu proprietatea că $\forall f \in \mathbb{P}d$, $R(f) = 0$ și $\exists g \in \mathbb{P}d + 1$ astfel încât $R(g) \neq 0$ se numește grad de exactitate al formulei de cuadratură.*

Deoarece R este liniar, rezultă că o formulă de cuadratură are gradul de exactitate d dacă și numai dacă $R(e_j) = 0$, $j = \overline{0, d}$ și $R(e_{d+1}) \neq 0$.

Dacă gradul de exactitate al unei formule de cuadratură este cunoscut, restul se poate determina cu ajutorul teoremei lui Peano.

5.3.1. Formula trapezului și formula lui Simpson

Aceste formule au fost denumite de Gautschi în [20] „caii de bătaie” ai integrării numerice. Ele își fac bine munca când intervalul de integrare este mărginit și integrandul este neproblematic. Formula trapezelor este surprinzător de eficientă chiar și pentru intervale infinite.

Ambele reguli se obțin aplicând cele mai simple tipuri de interpolare subintervalului diviziunii

$$a = x_0 < x_1 < x_2 < \cdots < x_{n-1} < x_n = b, \quad x_k = a + kh, \quad h = \frac{b-a}{n}. \quad (5.3.2)$$

În cazul regulii trapezului se interpolează liniar pe fiecare subinterval $[x_k, x_{k+1}]$ și se obține

$$\int_{x_k}^{x_{k+1}} f(x)dx = \int_{x_k}^{x_{k+1}} (L_1 f)(x)dx + \int_{x_k}^{x_{k+1}} (R_1 f)(x)dx, \quad f \in C^1[a, b], \quad (5.3.3)$$

cu

$$(L_1 f)(x) = f_k + (x - x_k)f[x_k, x_{k+1}].$$

Integrând avem

$$\int_{x_k}^{x_{k+1}} f(x)dx = \frac{h}{2}(f_k + f_{k+1}) + R_1(f),$$

unde

$$R_1(f) = \int_{x_k}^{x_{k+1}} K_1(t)f''(t)dt,$$

iar

$$\begin{aligned} K_1(t) &= \frac{(x_{k+1} - t)^2}{2} - \frac{h}{2}[(x_k - t)_+ + (x_{k+1} - t)_+] = \\ &= \frac{(x_k - t)^2}{2} - \frac{h(x_{k+1} - t)}{2} = \\ &= \frac{1}{2}(x_{k+1} - t)(x_k - t) \leq 0. \end{aligned}$$

Deci

$$R_1(f) = -\frac{h^3}{12}f''(\xi_k), \quad \xi_k \in (x_k, x_{k+1})$$

și

$$\int_{x_k}^{x_{k+1}} f(x)dx = \frac{h}{2}(f_k + f_{k+1}) - \frac{1}{12}h^3f''(\xi_k). \quad (5.3.4)$$

Această formulă se numește *regula (elementară a) trapezului*. Însușind pentru toate subintervalele se obține *regula trapezelor* sau *formula compusă a trapezului* sau *formula repetată a trapezului*.

$$\int_a^b f(x)dx = h \left(\frac{1}{2}f_0 + f_1 + \cdots + f_{n-1} + \frac{1}{2}f_n \right) - \frac{1}{12}h^3 \sum_{k=0}^{n-1} f''(\xi_k).$$

Deoarece f'' este continuă pe $[a, b]$, restul se poate scrie sub forma

$$R_{1,n}(f) = -\frac{(b-a)h^2}{12} f''(\xi) = -\frac{(b-a)^3}{12n^2} f''(\xi). \quad (5.3.5)$$

Deoarece f'' este mărginită în modul pe $[a, b]$ avem

$$R_{1,n}(f) = O(h^2),$$

când $h \rightarrow 0$ și deci regula trapezelor converge când $h \rightarrow 0$ (sau echivalent, $n \rightarrow \infty$), atunci când $f \in C^2[a, b]$.

Dacă în locul interpolării liniare se utilizează interpolarea pătratică se obține *regula lui Simpson repetată*. Varianta ei elementară, numită *regula lui Simpson*² sau *formula lui Simpson* este

$$\int_{x_k}^{x_{k+2}} f(x)dx = \frac{h}{3}(f_k + 4f_{k+1} + f_{k+2}) - \frac{1}{90}h^5 f^{(4)}(\xi_k), \quad x_k \leq \xi_k \leq x_{k+2}, \quad (5.3.6)$$

unde s-a presupus că $f \in C^4[a, b]$.

Să demonstrăm formula pentru restul formulei lui Simpson. Deoarece gradul de exactitate este 3, conform teoremei lui Peano avem

$$R_2(f) = \int_{x_k}^{x_{k+2}} K_2(t) f^{(4)}(t) dt.$$

unde

$$K_2(t) = \frac{1}{3!} \left\{ \frac{(x_{k+2} - t)^4}{4} - \frac{h}{3} [(x_k - t)_+^3 + 4(x_{k+1} - t)_+^3 + (x_{k+2} - t)_+^3] \right\},$$



Thomas Simpson (1710-1761), matematician englez autodidact, autor al mai multor texte, populare la vremea respectivă. Simpson și-a publicat formula în 1743, dar ea a fost cunoscută deja, printre alții, de Cavalieri (1639), Gregory (1668) și Cotes (1722).

adică

$$K_2(t) = \frac{1}{6} \begin{cases} \frac{(x_{k+2}-t)^4}{4} - \frac{h}{3} [4(x_{k+1}-t)^3 + (x_{k+2}-t)^3], & t \in [x_k, x_{k+1}] \\ \frac{(x_{k+2}-t)^4}{4} - \frac{h}{3}(x_{k+2}-t)^3, & t \in [x_{k+1}, x_{k+2}] \end{cases}$$

Se arată că pentru $t \in [x_k, x_{k+2}]$, $K_2(t) \leq 0$ și deci putem aplica corolarul la teorema lui Peano.

$$R_2(f) = \frac{1}{4!} f^{(4)}(\xi_k) R_2(e_4),$$

$$\begin{aligned} R_2(e_4) &= \frac{x_{k+2}^5 - x_k^5}{5} - \frac{h}{3} \left[x_k^4 + 4 \left(\frac{x_{k+2} + x_k}{2} \right)^4 + x_{k+2}^4 \right] \\ &= 3h \left[\frac{x_{k+2}^4 + x_{k+2}^3 x_k + x_{k+2}^2 x_k^2 + x_{k+2} x_k^3 + x_k^4}{5} \right. \\ &\quad \left. - \frac{4x_k^4 + x_k^4 + 4x_k^3 x_{k+2} + 6x_k^2 x_{k+2}^2 + 4x_k x_{k+2}^3 + x_{k+2}^4 + 4x_{k+2}^4}{24} \right] \\ &= \frac{h}{40} (-x_k^4 + 4x_k^3 x_{k+2} + 6x_k^2 x_{k+2}^2 + 4x_k x_{k+2}^3 - x_{k+2}^4) \\ &= -\frac{h}{40} (x_{k+2} - x_k)^4. \end{aligned}$$

Deci,

$$R_2(f) = -\frac{h^5}{90} f^{(4)}(\xi_k).$$

Pentru regula repetată a lui Simpson obținem

$$\int_a^b f(x) dx = \frac{h}{3} (f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 + \cdots + 4f_{n-1} + f_n) + R_{2,n}(f) \quad (5.3.7)$$

cu

$$R_{2,n}(f) = -\frac{1}{180} (b-a) h^4 f^{(4)}(\xi) = -\frac{(b-a)^5}{2880n^4} f^{(4)}(\xi), \quad \xi \in (a, b). \quad (5.3.8)$$

Se observă că $R_{2,n}(f) = O(h^4)$, de unde rezultă convergența când $n \rightarrow \infty$. Se observă și creșterea ordinului cu 1, ceea ce face ca regula repetată a lui Simpson să fie foarte populară și larg utilizată.

Regula trapezelor lucrează foarte bine pentru polinoame trigonometrice. Presupunem fără a restrânge generalitatea că $[a, b] = [0, 2\pi]$ și fie

$$\mathbb{T}_m[0, 2\pi] = \{t(x) : t(x) = a_0 + a_1 \cos x + a_2 \cos 2x + \cdots + a_m \cos mx +$$

$$+b_1 \sin x + b_2 \sin 2x + \cdots + b_m \sin mx\}$$

Atunci

$$R_{n,1}(f) = 0, \forall f \in \mathbb{T}_{n-1}[0, 2\pi] \quad (5.3.9)$$

Aceasta se verifică luând

$$f(x) = e_\nu(x) = e^{i\nu x} = \cos \nu x + i \sin \nu x, \quad \nu = 0, 1, 2, \dots$$

$$\begin{aligned} R_{n,1}(e_\nu) &= \int_0^{2\pi} e_\nu(x) dx - \frac{2\pi}{n} \left[\frac{1}{2} e_\nu(0) + \sum_{k=1}^{n-1} e_\nu\left(\frac{2k\pi}{n}\right) + \frac{1}{2} e_\nu(2\pi) \right] = \\ &= \int_0^{2\pi} e^{i\nu x} dx - \frac{2\pi}{n} \sum_{k=0}^{n-1} e^{\frac{2\pi i k \nu}{n}}. \end{aligned}$$

Când $\nu = 0$, aceasta este evident 0 și în caz contrar, deoarece

$$\int_0^{2\pi} e^{i\nu x} dx = (i\nu)^{-1} e^{i\nu x} \Big|_0^{2\pi} = 0,$$

$$R_{n,1}(e_\nu) = \begin{cases} -2\pi & \text{dacă } \nu = 0 \pmod{n}, \nu > 0 \\ -\frac{2\pi}{n} \frac{1 - e^{i\nu n \cdot 2\pi/n}}{1 - e^{i\nu \cdot 2\pi/n}} = 0 & \text{dacă } \nu \neq 0 \pmod{n} \end{cases} \quad (5.3.10)$$

În particular, $R_{n,1}(e_\nu) = 0$ pentru $\nu = 0, 1, \dots, n-1$, ceea ce demonstrează (5.3.9). Luând partea reală și imaginară în (5.3.10) se obține

$$R_{n,1}(\cos \nu \cdot) = \begin{cases} -2\pi, & \nu = 0 \pmod{n}, \nu \neq 0 \\ 0, & \text{în caz contrar} \end{cases}$$

$$R_{n,1}(\sin \nu \cdot) = 0.$$

De aceea, dacă f este 2π -periodică și are o dezvoltare Fourier uniform convergentă

$$f(x) = \sum_{\nu=0}^{\infty} [a_\nu(f) \cos \nu x + b_\nu(f) \sin \nu x], \quad (5.3.11)$$

unde $a_\nu(f), b_\nu(f)$ sunt coeficienții Fourier ai lui f , atunci

$$\begin{aligned} R_{n,1}(f) &= \sum_{\nu=0}^{\infty} [a_\nu(f) R_{n,1}(\cos \nu \cdot) + b_\nu(f) R_{n,1}(\sin \nu \cdot)] = \\ &= -2\pi \sum_{l=1}^{\infty} a_{ln}(f) \end{aligned} \quad (5.3.12)$$

Din teoria seriilor Fourier se știe că coeficienții Fourier ai lui f tind către zero cu atât mai repede cu cât f este mai netedă. Mai exact, dacă $f \in C^r[\mathbb{R}]$, atunci $a_\nu(f) = O(\nu^{-r})$ când $\nu \rightarrow \infty$ (și similar pentru $b_\nu(f)$). Deoarece conform (5.3.12)

$$R_{n,1}(f) \simeq -2\pi a_n(f),$$

rezultă că

$$R_{n,1}(f) = O(n^{-r}) \text{ când } n \rightarrow \infty \quad (5.3.13)$$

($f \in C^r[\mathbb{R}]$, 2π periodică), care dacă $r > 2$ este mai bună decât $R_{n,1}(f) = O(n^{-2})$, valabilă pentru funcții f neperiodice. În particular, dacă $r = \infty$, atunci regula trapezului converge mai repede decât orice putere a lui n^{-1} . Trebuie observat totuși că f trebuie să fie netedă pe întreg \mathbb{R} . Pornind de la o funcție $f \in C^r[0, 2\pi]$ și prelungind-o pe întreaga axă reală prin periodicitate, în general nu se obține o funcție $f \in C^r[\mathbb{R}]$.

O altă împrejurare în care regula trapezelor excelează este pentru funcții definite pe \mathbb{R} și care au proprietatea următoare pentru un anumit $r \geq 1$

$$f \in C^{2r+1}[\mathbb{R}], \quad \int_{\mathbb{R}} |f^{(2r+1)}(x)| dx < \infty, \\ \lim_{x \rightarrow -\infty} f^{(2\rho-1)}(x) = \lim_{x \rightarrow +\infty} f^{(2\rho-1)}(x) = 0, \quad \rho = 1, 2, \dots, r. \quad (5.3.14)$$

În acest caz se poate arăta că

$$\int_{\mathbb{R}} f(x) dx = h \sum_{k=-\infty}^{\infty} f(kh) + R(f; h) \quad (5.3.15)$$

are o eroare $R(f, h)$ ce satisface $R(f; h) = O(h^{2r+1})$, $h \rightarrow \infty$. Deci, dacă (5.3.14) are loc pentru orice $r \in \mathbb{N}$, atunci eroarea tinde la zero mai repede decât orice putere a lui h .

5.3.2. Formule Newton-Cotes cu ponderi și formule de tip Gauss

O formulă de cuadratură cu ponderi este o formulă de tipul

$$\int_a^b f(t)w(t)dt = \sum_{k=1}^n w_k f(t_k) + R_n(f) \quad (5.3.16)$$

unde w este nenegativă, integrabilă pe (a, b) .

Intervalul (a, b) poate fi mărginit sau nemărginit. Dacă este nemărginit trebuie să ne asigurăm că integrala din (5.3.16) este bine definită, cel puțin în cazul când f este polinom. Realizăm aceasta cerând ca toate momentele funcției pondere

$$\mu_s = \int_a^b t^s w(t) dt, \quad s = 0, 1, 2, \dots \quad (5.3.17)$$

să existe și să fie finite.

Spunem că (5.3.16) este de *tip interpolator*, dacă are gradul de exactitate $d = n - 1$. Formulele de tip interpolator sunt chiar formulele obținute prin interpolare, adică pentru care

$$\sum_{k=1}^n w_k f(t_k) = \int_a^b L_{n-1}(f; t_1, \dots, t_n, t) w(t) dt \quad (5.3.18)$$

sau echivalent

$$w_k = \int_a^b \ell_k(t) w(t) dt, \quad k = 1, 2, \dots, n, \quad (5.3.19)$$

unde

$$\ell_k(t) = \prod_{\substack{l=1 \\ l \neq k}}^n \frac{t - t_l}{t_k - t_l} \quad (5.3.20)$$

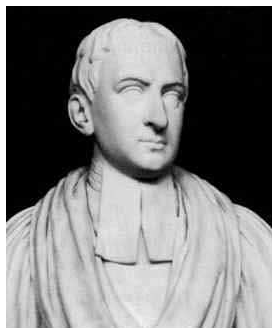
sunt polinoamele fundamentale Lagrange asociate nodurilor t_1, t_2, \dots, t_n . Faptul că (5.3.16) are gradul de exactitate $d = n - 1$ este evident, deoarece pentru orice $L_{n-1}(f; \cdot) \equiv f(\cdot)$ în (5.3.18). Reciproc, dacă (5.3.16) are gradul de exactitate $d = n - 1$, atunci luând $f(t) = \ell_r(t)$ în (5.3.17) ne dă

$$\int_a^b \ell_r(t) w(t) dt = \sum_{k=1}^n w_k \ell_r(t_k) = w_r, \quad r = 1, 2, \dots, n,$$

adică (5.3.19).

Observăm că dacă se dau n noduri distincte t_1, \dots, t_n este posibil întotdeauna să construim o formulă de tip (5.3.16) care este exactă pentru orice polinom de grad $\leq n - 1$. În cazul $w(t) \equiv 1$ pe $[-1, 1]$ și t_k sunt echidistante pe $[-1, 1]$ problema a fost intuită de Newton în 1687 și rezolvată în detaliu de Cotes ³ în jurul anului 1712. Prin extensie vom numi formula (5.3.16) cu t_k prescrise și w_k date de (5.3.19) *formulă de tip Newton-Cotes*.

Chestiunea care se pune în mod natural este dacă nu am putea face aceasta mai bine, adică dacă nu am putea obține gradul de exactitate $d > n - 1$ printr-o alegere judicioasă



Roger Cotes (1682-1716), fiul precoce al unui pastor de țară englez, a fost însărcinat cu pregătirea celei de-a doua ediții a lucrării lui Newton *Principia*. El a dezvoltat ideea lui Newton referitoare la integrarea numerică și a publicat coeficienții pentru formulele de integrare numerică cu n -puncte, pentru $n < 11$ (numere Cotes). La moartea sa prematură la vârsta de 33 de ani, Newton a spus despre el: „Dacă ar fi trăit, am fi putut ști ceva.”

a nodurilor t_k (ponderile w_k fiind date în mod necesar de (5.3.19)). Răspunsul este surprinzător de simplu și de direct. Pentru a-l formula, considerăm polinomul nodurilor

$$u_n(t) = \prod_{k=1}^n (t - t_k). \quad (5.3.21)$$

Teorema 5.3.3 *Dându-se un întreg k , $0 \leq k \leq n$, formula de cuadratură (5.3.16) are gradul de exactitate $d = n - 1 + k$ dacă și numai dacă sunt satisfăcute următoarele condiții:*

- (a) formula (5.3.16) este de tip interpolator ;
 (b) polinomul nodurilor u_n din (5.3.21) satisface

$$\int_a^b u_n(t)p(t)w(t) dt = 0, \quad \forall p \in \mathbb{P}_{k-1}.$$

Condiția din (b) impune k restricții asupra nodurilor t_1, t_2, \dots, t_n din (5.3.16). (Dacă $k = 0$, nu avem nici o restricție, deoarece, așa cum știm, putem atinge gradul de exactitate $d = n - 1$). Într-adevăr u_n trebuie să fie ortogonale pe \mathbb{P}_{k-1} relativ la funcția pondere w . Deoarece $w(t) \geq 0$, avem în mod necesar $k \leq n$. Altfel, u_n trebuie să fie ortogonal pe \mathbb{P}_n , în particular pe el însuși, ceea ce este imposibil. Astfel $k = n$ este optimal, obținându-se o formulă de cuadratură cu gradul maxim de exactitate $d_{max} = 2n - 1$. Condiția (b) impune ortogonalitatea lui u_n pe toate polinoamele de grad mai mic, adică $u_n(\cdot) = \pi_n(\cdot, w)$ este polinomul ortogonal în raport cu ponderea w . Această formulă optimală se numește *formulă de cuadratură de tip Gauss* asociată cu funcția pondere w . Deci nodurile ei sunt zerourile lui $\pi_n(\cdot, w)$, iar ponderile (coeficienții) w_k sunt dați de (5.3.19) adică

$$\begin{aligned} \pi_n(t_k; w) &= 0 \\ w_k &= \int_a^b \frac{\pi_n(t, w)}{(t - t_k)\pi_n'(t_k, w)} w(t) dt, \quad k = 1, 2, \dots, n \end{aligned} \quad (5.3.22)$$

Formula a fost dezvoltată de Gauss în 1814 în cazul special $w(t) \equiv 1$ pe $[-1, 1]$ și extinsă la funcții pondere mai generale de către Christoffel ⁴ în 1877. De aceea se mai numește *formulă de cuadratură Gauss-Christoffel*.

Elvin Bruno Christoffel (1829-1900), a activat pentru perioade scurte la Berlin și Zürich și cea mai mare parte a vieții la Strasbourg. Este cunoscut pentru lucrările sale de geometrie, în particular de analiza tensorilor, care a jucat un rol important în teoria relativității a lui Einstein.



Demonstrația teoremei 5.3.3. Vom demonstra întâi necesitatea lui (a) și (b). Deoarece gradul de exactitate este $d = n - 1 + k \geq n - 1$, condiția (a) este trivială. Condiția (b) rezultă de asemenea imediat, deoarece pentru orice $p \in \mathbb{P}_{k-1}$, $u_n p \in \mathbb{P}_{n-1+k}$. Deci

$$\int_a^b u_n(t)p(t)w(t)dt = \sum_{k=1}^n w_k u_k(t_k)p(t_k),$$

care se anulează, căci $u_n(t_k) = 0$ pentru $k = 1, 2, \dots, n$.

Pentru a demonstra suficiența lui (a) și (b) trebuie să arătăm că pentru orice $p \in \mathbb{P}_{n-1+k}$ avem $R_n(p) = 0$ în (5.3.16). Dându-se orice astfel de p , îl împărțim cu u_n , astfel încât

$$p = qu_n + r, \quad q \in \mathbb{P}_{k-1}, \quad r \in \mathbb{P}_{n-1}$$

unde q este câtul și r restul. Rezultă că

$$\int_a^b p(t)w(t)dt = \int_a^b q(t)u_n(t)w(t)dt + \int_a^b r(t)w(t)dt.$$

Prima integrală din dreapta este 0, conform lui b , deoarece $q \in \mathbb{P}_{k-1}$, în timp ce a doua, conform lui (a), deoarece $r \in \mathbb{P}_{n-1}$ este egală cu

$$\sum_{k=1}^n w_k r(t_k) = \sum_{k=1}^n w_k [p(t_k) - q(t_k)u_n(t_k)] = \sum_{k=1}^n w_k p(t_k)$$

ceea ce încheie demonstrația. \square

Cazul $k = n$ va fi discutat în secțiunea 5.3.3. Vom menționa două cazuri importante când $k < n$, care sunt de interes practic. Primul este formula de cuadratură Gauss-Radau în care o extremitate de interval, de exemplu a , este finită și servește ca nod, să zicem $t_1 = a$. Gradul maxim de exactitate care se poate obține este $d = 2n - 2$ și corespunde lui $k = n - 1$ în teorema (5.3.3). Partea (b) a teoremei ne spune că nodurile rămase t_2, \dots, t_n trebuie să fie rădăcinile polinomului $\pi_{n-1}(\cdot, w_a)$, unde $w_a(t) = (t - a)w(t)$. La fel, în formulele Gauss-Lobatto, ambele capete sunt finite și servesc ca noduri, să zicem $t_1 = a$, $t_n = b$, iar nodurile rămase t_2, \dots, t_{n-1} sunt zerourile lui $\pi_{n-2}(\cdot; w_{a,b})$, $w_{a,b}(t) = (t - a)(b - t)w(t)$, obținându-se astfel gradul de exactitate $d = 2n - 3$.

5.3.3. Proprietăți ale cuadraturilor gaussiene

Regula de cuadratură a lui Gauss, dată de (5.3.16) și (5.3.22), pe lângă faptul că este optimă (adică are grad maxim de exactitate) are și unele proprietăți interesante.

(i) Toate nodurile sunt reale, distincte și situate în intervalul deschis (a, b) . Aceasta este o proprietate cunoscută satisfăcută de zerourile polinoamelor ortogonale.

(ii) Toți coeficienții (ponderile) w_k sunt pozitivi. Demonstrația se bazează pe o observație ingenioasă a lui Stieltjes

$$0 < \int_a^b \ell_j^2(t)w(t)dt = \sum_{k=1}^n w_k \ell_j^2(t_k) = w_j, \quad j = 1, 2, \dots, n,$$

prima egalitate rezultând din faptul că gradul de exactitate este $d = 2n - 1$.

(iii) Dacă $[a, b]$ este mărginit, atunci formula lui Gauss converge pentru orice funcție continuă. Adică $R_n(f) \rightarrow 0$, când $n \rightarrow \infty$, pentru orice $f \in C[a, b]$. Aceasta este o consecință a teoremei de aproximare a lui Weierstrass, din care rezultă că dacă $\widehat{p}_{2n-1}(f; \cdot)$ este polinomul de cea mai bună aproximare a lui f pe $[a, b]$ în sensul normei uniforme atunci

$$\lim_{n \rightarrow \infty} \|f(\cdot) - \widehat{p}_{2n-1}(f; \cdot)\|_{\infty} = 0.$$

Deoarece $R_n(\widehat{p}_{2n-1}) = 0$ (căci $d = 2n - 1$), avem succesiv

$$\begin{aligned} |R_n(f)| &= |R_n(f - \widehat{p}_{2n-1})| = \\ &= \left| \int_a^b [f(t) - \widehat{p}_{2n-1}(f; t)]w(t)dt - \sum_{k=1}^n w_k [f(t_k) - \widehat{p}_{2n-1}(f; t_k)] \right| \leq \\ &\leq \int_a^b |f(t) - \widehat{p}_{2n-1}(f; t)|w(t)dt + \sum_{k=1}^n w_k |f(t_k) - \widehat{p}_{2n-1}(f; t_k)| \leq \\ &\leq \|f(\cdot) - \widehat{p}_{2n-1}(f; \cdot)\|_{\infty} \left[\int_a^b w(t)dt + \sum_{k=1}^n w_k \right]. \end{aligned}$$

Aici pozitivitatea ponderilor w_k a intervenit crucial. Observând că

$$\sum_{k=1}^n w_k = \int_a^b w(t)dt = \mu_0,$$

concluzionăm că

$$|R_n(f)| \leq 2\mu_0 \|f - \widehat{p}_{2n-1}\|_{\infty} \rightarrow 0, \quad \text{când } n \rightarrow \infty.$$

Proprietatea care urmează este baza unui algoritm eficient de obținere a unor formule de cuadratură gaussiană.

(iv) Fie $\alpha_k = \alpha_k(w)$ și $\beta_k = \beta_k(w)$ coeficienții din formula de recurență pentru polinoamele ortogonale

$$\begin{aligned} \pi_{k+1}(t) &= (t - \alpha_k)\pi_k(t) - \beta_k\pi_{k-1}(t), & k = 0, 1, 2, \dots \\ \pi_0(t) &= 1, & \pi_{-1}(t) = 0, \end{aligned} \tag{5.3.23}$$

unde

$$\begin{aligned}\alpha_k &= \frac{(t\pi_k, \pi_k)}{(\pi_k, \pi_k)} \\ \beta_k &= \frac{(\pi_k, \pi_k)}{(\pi_{k-1}, \pi_{k-1})},\end{aligned}\tag{5.3.24}$$

cu β_0 definit (ca de obicei) prin

$$\beta_0 = \int_a^b w(t) dt (= \mu_0).$$

Matricea Jacobi de ordinul n pentru funcția pondere w este o matrice simetrică tridiagonală definită prin

$$J_n(w) = \begin{bmatrix} \alpha_0 & \sqrt{\beta_1} & & & 0 \\ \sqrt{\beta_1} & \alpha_1 & \sqrt{\beta_2} & & \\ & \sqrt{\beta_2} & & \ddots & \\ & & \ddots & \ddots & \sqrt{\beta_{n-1}} \\ 0 & & & \sqrt{\beta_{n-1}} & \alpha_{n-1} \end{bmatrix}.$$

Teorema 5.3.4 *Nodurile t_k ale unei formule de tip Gauss sunt valorile proprii ale lui J_n*

$$J_n v_k = t_k v_k, \quad v_k^T v_k = 1, \quad k = 1, 2, \dots, n,\tag{5.3.25}$$

iar ponderile w_k sunt exprimabile cu ajutorul componentelor v_k , ale vectorilor proprii normalizați corespunzători prin

$$w_k = \beta_0 v_{k,1}^2, \quad k = 1, 2, \dots, n\tag{5.3.26}$$

Astfel, pentru a obține o formulă de cuadratură Gaussiană trebuie rezolvată o problemă de vectori și valori proprii pentru o matrice tridiagonală simetrică. Pentru această problemă există metode foarte eficiente. Astfel, abordarea bazată pe valori și vectori proprii este mai eficientă decât cea clasică. În plus, abordarea clasică se bazează pe două probleme prost condiționate: rezolvarea ecuațiilor polinomiale (coeficienții sunt obținuți prin aplicarea de relații de recurență, deci sunt deja perturbați) și rezolvarea unui sistem de ecuații având matricea Vandermonde.

Demonstrația teoremei 5.3.4. Fie $\tilde{\pi}_k(\cdot) = \tilde{\pi}_k(\cdot, w)$ polinomul ortogonal normalizat, deci $\pi_k = \sqrt{(\pi_k, \pi_k)_{d\lambda}} \tilde{\pi}_k$. Inserând aceasta în (5.3.23), împărțind cu $\sqrt{(\pi_k, \pi_k)_{d\lambda}}$ și utilizând (5.3.24), se obține

$$\tilde{\pi}_{k+1}(t) = (t - \alpha_k) \frac{\tilde{\pi}_k}{\sqrt{\beta_{k+1}}} - \beta_k \frac{\tilde{\pi}_{k-1}}{\sqrt{\beta_{k+1}\beta_k}},$$

sau înmulțind cu $\sqrt{\beta_{k+1}}$ și reordonând,

$$t\tilde{\pi}_k(t) = \alpha_k \tilde{\pi}_k(t) + \sqrt{\beta_k} \tilde{\pi}_{k-1}(t) + \sqrt{\beta_{k+1}} \tilde{\pi}_{k+1}(t), \quad k = 0, 1, \dots, n-1. \quad (5.3.27)$$

Cu ajutorul matricei lui Jacobi J_n , putem scrie această relație sub forma vectorială

$$t\tilde{\pi}(t) = J_n \tilde{\pi}(t) + \sqrt{\beta_n} \tilde{\pi}_n(t) e_n, \quad (5.3.28)$$

unde $\tilde{\pi}(t) = [\tilde{\pi}_0(t), \tilde{\pi}_1(t), \dots, \tilde{\pi}_{n-1}(t)]^T$ și $e_n(t) = [0, 0, \dots, 0, 1]^T$ sunt vectori din \mathbb{R}^n . Deoarece t_k sunt zerouri ale lui $\tilde{\pi}_n$ rezultă din (5.3.28) că

$$t_k \tilde{\pi}(t_k) = J_n \tilde{\pi}(t_k), \quad k = 1, 2, \dots, n. \quad (5.3.29)$$

Aceasta demonstrează prima relație a teoremei 5.3.4, deoarece $\tilde{\pi}$ este un vector nenul cu prima componentă

$$\tilde{\pi}_0 = \beta_0^{-1/2}. \quad (5.3.30)$$

Pentru a demonstra a doua relație, observăm din (5.3.29) că vectorul propriu normalizat v_k este

$$v_k = \frac{1}{[\tilde{\pi}(t_k)^T \tilde{\pi}(t_k)]} \tilde{\pi}(t_k) = \left(\sum_{\mu=1}^n \tilde{\pi}_{\mu-1}^2(t_k) \right)^{-1/2} \tilde{\pi}(t_k).$$

Comparând prima componentă din ambi membri și ridicând la pătrat pe baza formulei (5.3.30)

$$\frac{1}{\sum_{\mu=1}^n \tilde{\pi}_{\mu-1}^2(t_k)} = \beta_0 v_{k,1}^2, \quad k = 1, 2, \dots, n. \quad (5.3.31)$$

Pe de altă parte, luând $f(t) = \tilde{\pi}_{\mu-1}(t)$ în formula de cuadratură de tip Gauss (5.3.16), utilizând ortogonalitatea și din nou (5.3.30) se obține că

$$\beta_0^{1/2} \delta_{\mu-1,0} = \sum_{k=0}^n w_k \tilde{\pi}_{\mu-1}(t_k)$$

sau în formă matricială

$$Pw = \beta_0^{1/2} e_1, \quad (5.3.32)$$

unde $\delta_{\mu-1,0}$ este simbolul lui Kronecker, $P \in \mathbb{R}^{n \times n}$ este matricea vectorilor proprii, $w \in \mathbb{R}^n$ este vectorul coeficienților gaussieni și $e_1 = [1, 0, \dots, 0]^T \in \mathbb{R}^n$. Deoarece coloanele lui P sunt ortogonale, avem

$$P^T P = D, \quad D = \text{diag}(d_1, d_2, \dots, d_n), \quad d_k = \sum_{\mu=1}^n \tilde{\pi}_{\mu-1}^2(t_k).$$

Înmulțim acum (5.3.32) la stânga cu P^T și obținem

$$Dw = \beta_0^{1/2} P^T e_1 = \beta_0^{1/2} * \beta_0^{-1/2} e = e, \quad e = [1, 1, \dots, 1]^T.$$

Deci, $w = D^{-1}e$, adică,

$$w_k = \frac{1}{\sum_{\mu=1}^n \tilde{\pi}_{\mu-1}^2(t_k)}, \quad k = 1, 2, \dots, n.$$

Comparând cu (5.3.31) se obține rezultatul dorit. \square

Pentru detalii privind aspectele algoritmice legate de polinoame ortogonale și cuadraturi gaussiene a se vedea [21].

(v) Markov ⁵ a observat că formula de cuadratură a lui Gauss poate fi obținută cu ajutorul formulei de interpolare a lui Hermite cu noduri duble:

$$f(x) = (H_{2n-1}f)(x) + u_n^2(x)f[x, x_1, x_1, \dots, x_n, x_n],$$

$$\begin{aligned} \int_a^b w(x)f(x)dx &= \int_a^b w(x)(H_{2n-1}f)(x)dx + \\ &+ \int_a^b w(x)u_n^2(x)f[x, x_1, x_1, \dots, x_n, x_n]dx. \end{aligned}$$

Dar gradul de exactitate $2n - 1$ implică

$$\int_a^b w(x)(H_{2n-1}f)(x)dx = \sum_{i=1}^n w_i(H_{2n-1}f)(x_i) = \sum_{i=1}^n w_i f(x_i),$$

$$\int_a^b w(x)f(x)dx = \sum_{i=1}^n w_i f(x_i) + \int_a^b w(x)u_n^2(x)f[x, x_1, x_1, \dots, x_n, x_n]dx,$$



Andrei Andreievici Markov (1856-1922), matematician rus, activ în Sankt Petersburg. A avut contribuții importante în teoria probabilităților, teoria numerelor și teoria constructivă a aproximării.

deci

$$R_n(f) = \int_a^b w(x)u_n^2(x)f[x, x_1, x_1, \dots, x_n, x_n]dx.$$

Cum $w(x)u^2(x) \geq 0$, aplicând teorema de medie pentru integrale și teorema de medie pentru diferențe divizate avem

$$\begin{aligned} R_n(f) &= f[\eta, x_1, x_1, \dots, x_n, x_n] \int_a^b w(x)u^2(x)dx = \\ &= \frac{f^{(2n)}(\xi)}{(2n)!} \int_a^b w(x)[\pi_n(x, w)]^2 dx, \quad \xi \in [a, b]. \end{aligned}$$

Vom încheia secțiunea cu o tabelă cu funcțiile pondere clasice, polinoamele lor ortogonale corespunzătoare și coeficienții din formula de recurență α_k, β_k (vezi tabela 5.1).

polinoamele	notația	ponderea	intervalul	α_k	β_k
Legendre	$P_n(l_n)$	1	[-1,1]	0	2 ($k=0$) $(4-k^2)^{-1}$ ($k>0$)
Cebășev #1	T_n	$(1-t^2)^{-\frac{1}{2}}$	[-1,1]	0	π ($k=0$) $\frac{1}{2}\pi$ ($k=1$) $\frac{1}{4}$ ($k>0$)
Cebășev #2	$u_n(Q_n)$	$(1-t^2)^{\frac{1}{2}}$	[-1,1]	0	$\frac{1}{2}\pi$ ($k=0$) $\frac{1}{4}$ ($k>0$)
Jacobi	$P_n^{(\alpha, \beta)}$	$(1-t)^\alpha(1-t)^\beta$ $\alpha > -1, \beta > -1$	[-1,1]	vezi observația 5.3.5	vezi observația 5.3.5
Laguerre	$L_n^{(\alpha)}$	$t^\alpha e^{-t}$ $\alpha > -1$	[0, ∞)	$2k + \alpha + 1$	$\Gamma(1 + \alpha)$ ($k=0$) $k(k + \alpha)$ ($k>0$)
Hermite	H_n	e^{-t^2}	\mathbb{R}	0	$\sqrt{\pi}$ ($k=0$)

Tabela 5.1: Polinoame ortogonale

Observația 5.3.5. Pentru polinoamele Jacobi avem

$$\alpha_k = \frac{\beta^2 - \alpha^2}{(2k + \alpha + \beta)(2k + \alpha + \beta + 2)}$$

și

$$\begin{aligned} \beta_0 &= 2^{\alpha+\beta+1} B(\alpha + 1, \beta + 1), \\ \beta_k &= \frac{4k(k + \alpha)(k + \alpha + \beta)}{(2k + \alpha + \beta - 1)(2k + \alpha + \beta)^2(2k + \alpha + \beta + 1)}, \quad k > 0. \quad \diamond \end{aligned}$$

5.4. Cuadraturi adaptive

La metodele de integrare numerică erorile nu depind numai de dimensiunea intervalului utilizat, ci și de valoarea derivatelor de un anumit ordin ale funcției care urmează a fi integrată. Aceasta implică faptul că metodele nu vor lucra bine pentru funcții cu derivatele de un anumit ordin mari – în special funcții care au fluctuații mari pe unele subintervale sau pe tot intervalul. Este rezonabil să utilizăm subintervale mici acolo unde derivatele sunt mari și subintervale mari acolo unde derivatele sunt mici. O metodă care face aceasta într-o manieră sistematică se numește cuadratură adaptivă.

Abordarea generală într-o cuadratură adaptivă este de a utiliza două metode diferite pe fiecare subinterval, de a compara rezultatul și de a subdiviza intervalul dacă diferențele sunt mari. Există situația nefericită în care se utilizează două metode proaste, rezultatele sunt proaste, dar diferența dintre ele este mică. Un mod de a evita o astfel de situație este de a ne asigura că o metodă supraestimează rezultatul, iar alta îl subestimează. Vom da un exemplu de structură generală de cuadratură adaptiv-recursivă (algoritmul 5.1). Să presupunem că

$metint(a, b : real; f : funcție, n : integer) : real$

este o funcție care aproximează $\int_a^b f(x)dx$ folosind o cuadratură repetată cu n subintervale. Pentru m se alege o valoare mică (4 sau 5). Structura algoritmului: DIVIDE AND

Algoritmul 5.1 Cuadratură adaptivă

Intrare: f - funcția de integrat, a, b - limitele de integrare, ε - toleranța, $metint$ - o cuadratură repetată

Ieșire: valoarea integralei

function $adapt(f, a, b, \varepsilon, metint)$

if $|metint(a, b, f, 2m) - metint(a, b, f, m)| < \varepsilon$ **then**

$adapt := metint(a, b, f, 2m);$

else

$adapt := adapt(f, a, (a + b)/2, \varepsilon, metint) + adapt(f, (a + b)/2, b, \varepsilon, metint);$

end if

CONQUER.

Spre deosebire de alte metode, la care se decide cât de mult se muncește pentru a asigura precizia dorită, la o cuadratură adaptivă se calculează doar atât cât este necesar. Aceasta înseamnă că eroarea absolută ε trebuie aleasă astfel încât să nu se intre într-un ciclu infinit pentru a atinge o precizie imposibil de atins. Numărul de pași depinde de natura funcției de integrat. Posibilități de îmbunătățire: $metint(a, b, f, 2m)$ este apelat de două ori, precizia poate fi scalată cu raportul dintre dimensiunea intervalului curent și dimensiunea întregului interval. Pentru detalii suplimentare recomandăm [19].

5.5. Cuadraturi iterate. Metoda lui Romberg

Un dezavantaj al cuadraturilor adaptive este acela că calculează repetat valorile funcției în noduri, iar atunci când este rulat un astfel de program apare un consum suplimentar de timp de calcul datorită recursivității sau gestiunii stivei într-o implementare iterativă. Cuadraturile iterative înlătură aceste inconveniente. Ele aplică la primul pas o cuadratură repetată și apoi subdivid intervalele în părți egale folosind la fiecare pas aproximantele calculate anterior. Vom exemplifica această tehnică printr-o metodă care pornește de la formula repetată a trapezului și îmbunătățește convergența utilizând extrapolarea Richardson.

Primul pas al procesului presupune aplicarea formulei repetate a trapezului cu $n_1 = 1$, $n_2 = 2, \dots, n_p = 2^{p-1}$, unde $p \in \mathbb{N}^*$. Valoarea pasului h_k corespunzătoare lui n_k va fi

$$h_k = \frac{b-a}{n_k} = \frac{b-a}{2^{k-1}}.$$

Cu aceste notații regula trapezului devine

$$\int_a^b f(x)dx = \frac{h_k}{2} \left[f(a) + f(b) + 2 \sum_{i=1}^{2^{n-1}-1} f(a + ih_k) \right] - \frac{b-a}{12} h_k^2 f''(\mu_k) \quad (5.5.1)$$

$\mu_k \in (a, b)$.

Notăm cu $R_{k,1}$ rezultatul aproximării conform (5.5.1).

$$R_{1,1} = \frac{h_1}{2} [f(a) + f(b)] = \frac{b-a}{2} [f(a) + f(b)] \quad (5.5.2)$$

$$\begin{aligned} R_{2,1} &= \frac{h_2}{2} [f(a) + f(b) + 2f(a + h_2)] = \\ &= \frac{b-a}{4} \left[f(a) + f(b) + 2f\left(a + \frac{b-a}{2}\right) \right] = \\ &= \frac{1}{2} \left[R_{1,1} + h_1 f\left(a + \frac{1}{2}h_1\right) \right]. \end{aligned}$$

și în general

$$R_{k,1} = \frac{1}{2} \left[R_{k-1,1} + h_{k-1} \sum_{i=1}^{2^{k-2}} f\left(a + \left(i - \frac{1}{2}\right) h_{k-1}\right) \right], \quad k = \overline{2, n} \quad (5.5.3)$$

Urmează îmbunătățirea prin extrapolare Richardson ⁶ (deși a fost introdusă de Richardson și Gaunt se pare că este datorată lui Arhimede ⁷)

$$I = \int_a^b f(x)dx = R_{k-1,1} - \frac{(b-a)}{12} h_k^2 f''(a) + O(h_k^4).$$

Vom elimina termenul în h_k^2 combinând două ecuații

$$I = R_{k-1,1} - \frac{(b-a)}{12} h_k^2 f''(a) + O(h_k^4),$$

$$I = R_{k,1} - \frac{b-a}{48} h_k^2 f''(a) + O(h_k^4).$$

Obținem

$$I = \frac{4R_{k,1} - R_{k-1,1}}{3} + O(h^4).$$

Definim

$$R_{k,2} = \frac{4R_{k,1} - R_{k-1,1}}{3}. \quad (5.5.4)$$

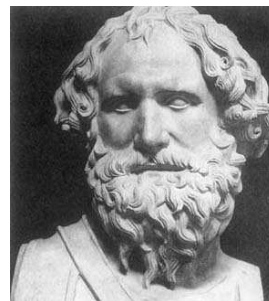
Se aplică extrapolarea Richardson și acestor valori. În general dacă $f \in C^{2n+2}[a, b]$, atunci pentru $k = \overline{1, n}$ putem scrie

$$\int_a^b f(x)dx = \frac{h_k}{2} \left[f(a) + f(b) + 2 \sum_{i=1}^{2^{k-1}-1} f(a + ih_k) \right] + \quad (5.5.5)$$

Lewis Fry Richardson (1881-1953), matematician englez. A avut contribuții la predicția numerică a vremii, propunând rezolvarea ecuațiilor hidro și termodinamice ale meteorologiei cu metode bazate pe diferențe finite. A realizat un studiu de referință asupra turbulenței atmosferice, introducând cantitățile numite astăzi „numerele Richardson”. La 50 de ani și-a luat licența în psihologie și a dezvoltat o teorie științifică a relațiilor internaționale. A fost ales membru al Royal Society în 1926.



Arhimede (287 î.e.n. - 212 î.e.n.), matematician grec din Siracuză, unul dintre cei mai importanți ai întregii antichități. A pus la punct o metodă de integrare care i-a permis să determine arii, volume și suprafețe ale multor corpuri. A dat o aproximare precisă a lui π , metode de aproximare precisă a rădăcinii pătrate și un sistem de reprezentare a numerelor mari. În mecanică, Arhimede a descoperit teoremele fundamentale referitoare la centrul de greutate al figurilor plane și solidelor și principiul care îi poartă numele, referitor la un corp scufundat într-un lichid. Mașinile sale de război au contribuit la apărarea orașului său în timpul asediului romanilor, care a durat trei ani. A murit ucis de un soldat roman la sfârșitul asediului.



$$+ \sum_{i=1}^k K_i h_k^{2i} + O(h_k^{2k+2}),$$

unde K_i nu depinde de h_k .

Formula (5.5.5) se justifică în modul următor. Fie $a_0 = \int_a^b f(x)dx$ și

$$A(h) = \frac{h}{2} \left[f(a) + 2 \sum_{k=1}^{n-1} f(a + kh) + f(b) \right], \quad h = \frac{b-a}{k}.$$

Dacă $f \in C^{2k+1}[a, b]$, $k \in \mathbb{N}^*$ are loc următorul rezultat datorat lui Euler ⁸ și Maclaurin ⁹

$$A(h) = a_0 + a_1 h^2 + a_2 h^4 + \dots + a_k h^{2k} + O(h^{2k+1}), \quad h \rightarrow 0 \quad (5.5.6)$$

unde

$$a_k = \frac{B_{2k}}{(2k)!} [f^{(2k-1)}(b) - f^{(2k-1)}(a)], \quad k = 1, 2, \dots, K.$$

Leonhard Euler (1707-1783), matematician elvețian, a urmat cursurile lui Jakob Bernoulli la Universitatea din Basel, luând și lecții particulare de la Johann Bernoulli. După ce la 20 de ani nu a reușit să obțină o catedră de fizică la Basel, a emigrat la Sankt Petersburg; mai târziu s-a mutat la Berlin și apoi din nou la Sankt Petersburg. Indiscutabil, Euler a fost cel mai prolific matematician al secolului al XVIII-lea, lucrând în aproape toate ramurile calculului diferențial și integral și fiind unul dintre fondatorii calculului variațional. A elaborat lucrări de pionierat în științele aplicate: hidrodinamică, mecanica materialelor deformabile și solidului rigid, optică, astronomie. Nici chiar orbirea sa la vârsta de 59 de ani nu i-a afectat productivitatea fenomenală. Se pare că opera sa nu a fost încă editată în întregime, apărând până acum 71 de volume.



Colin Maclaurin (1698-1768), matematician scoțian. A aplicat calculul infinitezimal la probleme de geometrie. Este cunoscut pentru dezvoltarea în serie în jurul originii dar a avut și contribuții la teoria ecuațiilor.



Cantitățile B_k sunt numerele lui Bernoulli ¹⁰, adică coeficienții dezvoltării

$$\frac{z}{e^z - 1} = \sum_{k=0}^{\infty} \frac{B_k}{k!} z^k, \quad |z| < 2\pi.$$

Formula (5.5.6) se numește *formula Euler-MacLaurin*.

Eliminând succesiv puterile lui h din (5.5.5) se obține

$$R_{k,j} = \frac{4^{j-1} R_{k,j-1} - R_{k-1,j-1}}{4^{j-1} - 1}, \quad k = \overline{2, n}, \quad j = \overline{2, i}.$$

Calcululele se pot aranja tabelar, astfel:

$$\begin{array}{ccccccc} R_{1,1} & & & & & & \\ R_{2,1} & R_{2,2} & & & & & \\ R_{3,1} & R_{3,2} & R_{3,3} & & & & \\ \vdots & \vdots & \vdots & \ddots & & & \\ R_{n,1} & R_{n,2} & R_{n,3} & \dots & R_{n,n} & & \end{array}$$

Deoarece $(R_{n,1})$ este convergent și $(R_{n,n})$ este convergent, mai rapid decât $(R_{n,1})$. Drept criteriu de oprire se poate folosi $|R_{n-1,n-1} - R_{n,n}| \leq \varepsilon$.

5.6. Cuadraturi adaptive II

Coloana a doua din metoda lui Romberg corespunde aproximării prin metoda lui Simpson. Notăm

$$S_{k,1} = R_{k,2}.$$

Coloana a treia este deci o combinație a două aproximante de tip Simpson:

$$S_{k,2} = S_{k,1} + \frac{S_{k,1} - S_{k-1,1}}{15} = R_{k,2} + \frac{R_{k,2} - R_{k-1,2}}{15}.$$

¹⁰Jacob Bernoulli (1654-1705), fratele mai mare al lui Johann Bernoulli, activ în Basel. A fost unul dintre primii matematicieni care a apreciat importanța introducerii de către Leibniz și Newton a calculului diferențial și integral, pe care l-a îmbogățit cu contribuții originale, în competiție (nu întotdeauna amicală) cu fratele său mai mic. Este unul dintre fondatorii Teoriei probabilităților, prin lucrarea sa *Ars conjectandi* și prin „legea numerelor mari.



Relația

$$S_{k,2} = S_{k,1} + \frac{S_{k,1} - S_{k-1,1}}{15}, \quad (5.6.1)$$

va fi folosită la elaborarea unui algoritm de cuadratură adaptivă. Fie $c = (a + b)/2$. Formula elementară a lui Simpson este

$$S = \frac{h}{6} (f(a) + 4f(c) + f(b)).$$

Pentru două subintervale se obține

$$S_2 = \frac{h}{12} (f(a) + 4f(d) + 2f(c) + 4f(e) + f(b)),$$

unde $d = (a + c)/2$ și $e = (c + b)/2$. Cantitatea Q se obține aplicând (5.6.1) celor două aproximante:

$$Q = S_2 + (S_2 - S)/15.$$

Putem să dam acum un algoritm recursiv pentru aproximarea integralei. Funcția *adquad* evaluează integrandul aplicând regula lui Simpson. Ea apelează recursiv *quadstep* și aplică extrapolarea. Descrierea se dă în algoritmul 5.2.

Algoritm 5.2 Cuadratură adaptivă bazată pe metoda lui Simpson și extrapolare

Intrare: funcția f , intervalul $[a, b]$, eroarea ε

Ieșire: Valoarea aproximativă a integralei

function *adquad*(f, a, b, ε) : *real*

$c := (a + b)/2;$

$fa = f(a); fb := f(b); fc := f(c);$

$Q := \text{quadstep}(f, a, b, \varepsilon, fa, fc, fb);$

return $Q;$

function *quadstep*($f, a, b, \varepsilon, fa, fc, fb$) : *real*

$h := b - a; c := (a + b)/2;$

$fd := f((a + c)/2); fe := f((c + b)/2);$

$Q1 := h/6 * (fa + 4 * fc + fb);$

$Q2 := h/12 * (fa + 4 * fb + 2 * fc + 4 * fe + fb);$

if $|Q1 - Q2| < \varepsilon$ **then**

$Q := Q2 + (Q2 - Q1)/15;$

else

$Qa := \text{quadstep}(f, a, c, \varepsilon, fa, fd, fc);$

$Qb := \text{quadstep}(f, c, b, \varepsilon, fc, fe, fb);$

$Q := Qa + Qb;$

end if

return $Q;$

CAPITOLUL 6

Rezolvarea numerică a ecuațiilor neliniare

Cuprins

6.1. Ecuații neliniare	177
6.2. Iterații, convergență și eficiență	178
6.3. Metoda șirurilor Sturm	181
6.4. Metoda falsei poziții	183
6.5. Metoda secantei	185
6.6. Metoda lui Newton	188
6.7. Metoda aproximațiilor succesive	191
6.8. Metoda lui Newton pentru rădăcini multiple	192
6.9. Ecuații algebrice	193
6.10. Metoda lui Newton în \mathbb{R}^n	194
6.11. Metode quasi-Newton	196
6.11.1. Interpolare liniară	197
6.11.2. Metode de modificare	198

6.1. Ecuații neliniare

Problema discutată în acest capitol se poate scrie generic sub forma

$$f(x) = 0, \tag{6.1.1}$$

dar admite diverse interpretări, depinzând de semnificația lui x și f . Cel mai simplu caz este cel al unei singure ecuații cu o singură necunoscută, caz în care f este o funcție dată de o variabilă reală sau complexă și încercăm să găsim valorile acestei variabile pentru care f se anulează. Astfel de valori se numesc *rădăcini* ale ecuației (6.1.1) sau *zerouri* ale funcției f .

Dacă x din (6.1.1) este un vector, să zicem $x = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d$ și f este de asemenea un vector ale cărui componente sunt funcții de cele d variabile x_1, x_2, \dots, x_d , atunci (6.1.1) reprezintă un sistem de ecuații.

Se spune că sistemul este neliniar dacă cel puțin una dintre componentele lui f depinde neliniar de cel puțin una din variabilele x_1, x_2, \dots, x_d . Dacă toate componentele lui f sunt funcții liniare de x_1, \dots, x_d avem de-a face cu un sistem de ecuații algebrice liniare. Mai general (6.1.1) ar putea reprezenta o ecuație funcțională, dacă x este un element al unui spațiu de funcții și f este un operator (liniar sau neliniar) ce acționează pe acest spațiu. În fiecare din aceste situații zeroul din dreapta lui (6.1.1) poate avea diverse interpretări: numărul zero în primul caz, vectorul nul în al doilea și funcția identic nulă în cel de-al treilea.

Mare parte din acest capitol este consacrată unei ecuații neliniare scalare. Astfel de ecuații apar frecvent în analiza sistemelor în vibrație, unde rădăcinile corespund frecvențelor critice (rezonanță). Cazul special al ecuațiilor algebrice, unde f din (6.1.1) este un polinom, este de importanță considerabilă și merită un tratament special.

6.2. Iterații, convergență și eficiență

Nici chiar cele mai simple ecuații - de exemplu cele algebrice - nu admit soluții care să fie exprimabile prin expresii raționale sau radicali. Din acest motiv este imposibil, în general, să calculăm rădăcinile ecuațiilor neliniare printr-un număr finit de operații aritmetice. De aceea este nevoie de o metodă iterativă, adică de o procedură care generează o secvență infinită de aproximații $\{x_n\}_{n \in \mathbb{N}}$ astfel încât

$$\lim_{n \rightarrow \infty} x_n = \alpha, \quad (6.2.1)$$

unde α este o rădăcină a ecuației. În cazul unui sistem x_k și α sunt vectori de dimensiune adecvată, iar convergența trebuie înțeleasă în sensul convergenței pe componente.

Deși convergența unui proces iterativ este de dorit, pentru a putea fi practică, este necesar ceva mai mult decât convergența. Ceea ce se dorește este o convergență rapidă. Conceptul de bază pentru măsurarea vitezei de convergență este ordinul de convergență.

Definiția 6.2.1 Spunem că x_n converge către α (cel puțin) liniar dacă

$$|x_n - \alpha| \leq e_n \quad (6.2.2)$$

unde $\{e_n\}$ este un șir pozitiv ce satisface

$$\lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n} = c, \quad 0 < c < 1. \quad (6.2.3)$$

Dacă (6.2.2) și (6.2.3) au loc cu egalitate în (6.2.2) atunci c se numește eroare asimptotică.

În această definiție, expresia „cel puțin“ se leagă de faptul că avem doar inegalitate în (6.2.2), ceea ce dorim în practică. De fapt, strict vorbind, marginea e_n converge liniar, însemnând că în cele din urmă (pentru n suficient de mare) fiecare din aceste margini ale erorii este aproximativ o fracție constantă din precedentă.

Definiția 6.2.2 Se spune că x_n converge către α (cel puțin) cu ordinul $p \geq 1$ dacă (6.2.2) are loc cu

$$\lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n^p} = c, \quad c > 0 \quad (6.2.4)$$

Astfel convergența de ordinul 1 coincide cu convergența liniară, în timp ce convergența de ordinul $p > 1$ este mai rapidă. De notat că în acest ultim caz nu se pune nici o restricție asupra constantei c : odată ce e_n este suficient de mic, exponentul p va avea grijă de convergență. Și în acest caz, dacă avem egalitate în (6.2.2), c se numește eroare asimptotică.

Aceleași definiții se aplică și șirurilor vectoriale, cu modulul înlocuit cu orice normă vectorială.

Clasificarea convergenței în raport cu ordinul este destul de rudimentară, deoarece sunt tipuri de convergență la care definițiile (6.2.1) și (6.2.2) nu se aplică. Astfel, un șir $\{e_n\}$ poate converge către zero mai încet decât liniar, de exemplu dacă $c = 1$ în (6.2.3). Acest tip de convergență se numește subliniară. La fel, $c = 0$ în (6.2.3) conduce la convergență superliniară, dacă (6.2.4) nu are loc pentru nici un $p > 1$.

Este instructiv să examinăm comportarea lui e_n , dacă în loc de relația la limită avem egalitate pentru un anumit n , să zicem

$$\frac{e_{n+1}}{e_n^p} = c, \quad n = n_0, n_0 + 1, n_0 + 2, \dots \quad (6.2.5)$$

Pentru n_0 suficient de mare, relația este aproape adevărată. Printr-o simplă inducție se obține că

$$e_{n_0+k} = c^{\frac{p^k-1}{p-1}} e_{n_0}^{p^k}, \quad k = 0, 1, 2, \dots, \quad (6.2.6)$$

care desigur are loc pentru $p > 1$, dar și pentru $p = 1$ când $p \downarrow 1$:

$$e_{n_0+k} = c^k e_{n_0}, \quad k = 0, 1, 2, \dots, \quad (p = 1) \quad (6.2.7)$$

Presupunând că e_{n_0} este suficient de mare astfel încât aproximarea x_{n_0} are un număr de zecimale corecte, scriem $e_{n_0+k} = 10^{-\delta_k} e_{n_0}$. Atunci δ_k , în conformitate cu (6.2.2) reprezintă numărul suplimentar de cifre zecimale corecte din aproximația x_{n_0+k} (în contrast cu x_{n_0}). Logaritmând (6.2.6) și (6.2.7) obținem

$$\delta_k = \begin{cases} k \log \frac{1}{c}, & \text{dacă } p = 1 \\ p^k \left[\frac{1-p^{-k}}{p-1} \log \frac{1}{c} + (1-p^{-k}) \log \frac{1}{e_{n_0}} \right], & \text{dacă } p > 1 \end{cases}$$

deci când $k \rightarrow \infty$

$$\delta_k \sim c_1 k \quad (p = 1), \quad \delta_k \sim c_p p^k \quad (p > 1), \quad (6.2.8)$$

unde $c_1 = \log \frac{1}{c} > 0$, dacă $p = 1$ și

$$c_p = \frac{1}{p-1} \log \frac{1}{c} + \log \frac{1}{e_{n_0}}$$

(presupunem că n_0 este suficient de mare și deci e_{n_0} suficient de mic, pentru a avea $c_p > 0$). Aceasta ne arată că numărul de cifre zecimale corecte crește liniar odată cu k când $p = 1$ și exponențial când $p > 1$. În ultimul caz, $\delta_{k+1}/\delta_k \sim p$, însemnând că (pentru k mare) numărul de cifre zecimale corecte, crește pe iterație, cu un factor p .

Dacă fiecare iterație necesită m unități de lucru (o „unitate de lucru” este efortul necesar pentru a calcula o valoare a funcției sau a unei anumite derivate a sa), atunci *indicele de eficiență* al iterației poate fi definit prin

$$\lim_{k \rightarrow \infty} [\delta_{k+1}/\delta_k]^{1/m} = p^{1/m}.$$

Aceasta ne dă o bază comună de comparare între diversele metode iterative. Metodele liniare au indicele de eficiență 1.

Calcululele practice necesită o regulă de oprire care să termine iterația atunci când s-a obținut (sau se crede că s-a obținut) precizia dorită. Ideal, ne oprim atunci când $\|x_n - \alpha\| < tol$, tol dat. Deoarece α nu este cunoscut se obișnuiește să se înlocuiască $x_n - \alpha$ cu $x_n - x_{n-1}$ și se impune cerința ca

$$\|x_n - x_{n-1}\| \leq tol \quad (6.2.9)$$

unde

$$tol = \|x_n\| \varepsilon_r + \varepsilon_a \quad (6.2.10)$$

cu $\varepsilon_r, \varepsilon_a$ valori date ale erorii. Ca o măsură de siguranță, am putea cere ca (6.2.9) să aibă loc pentru mai multe valori consecutive ale lui n , nu doar pentru una singură. Alegând $\varepsilon_r = 0$ sau $\varepsilon_a = 0$ se obține un test de eroare absolută sau relativă. Este totuși prudent să utilizăm un test mixt, cum ar fi, să zicem $\varepsilon_e = \varepsilon_a = \varepsilon$. Atunci, dacă $\|x_n\|$ este mic sau moderat de mare, se controlează efectiv eroarea absolută, în timp ce pentru $\|x_n\|$ foarte mare se controlează eroarea relativă. Testele de mai sus se pot combina cu $\|f(x)\| \leq \varepsilon$. În algoritmiile din acest capitol vom presupune că avem o funcție *crit_oprire* care implementează testul de oprire.

6.3. Metoda șirurilor Sturm

Există situații în care este de dorit să selectăm o rădăcină particulară dintre mai multe și să avem scheme iterative care converg către ea. Acesta este cazul, de exemplu, pentru polinoame ortogonale, ale căror rădăcini sunt reale și distincte. De asemenea am putea dori să facem o selecție a unei anumite rădăcini: cea mai mare, a doua ca mărime, a treia ș.a.m.d. și să o calculăm fără a mai calcula și altele. Aceasta este posibil dacă combinăm înjumătățirea cu teorema lui Sturm ¹.

Să considerăm ecuația

$$f(x) := \pi_d(x) = 0, \quad (6.3.1)$$

unde π_d este un polinom de grad d , ortonormal în raport cu o anumită măsură. Știm că π_d este polinomul caracteristic al unei matrice simetrice tridiagonale și poate fi calculat recursiv printr-o relație de recurență de forma

$$\begin{aligned} \pi_0(x) &= 1, & \pi_1(x) &= x - \alpha_0 \\ \pi_{k+1}(x) &= (x - \alpha_k)\pi_k(x) - \beta_k\pi_{k-1}(x), & k &= 1, 2, \dots, d-1 \end{aligned} \quad (6.3.2)$$

cu β_k pozitiv. Recurența (6.3.2) este utilă nu numai pentru calculul lui $\pi_d(x)$, dar și pentru că are următoarea proprietate utilă, datorată lui Sturm.

Propoziția 6.3.1 Fie $\sigma(x)$ numărul de schimbări de semn (zerourile nu contează) în secvența de numere

$$\pi_d(x), \pi_{d-1}(x), \dots, \pi_1(x), \pi_0(x). \quad (6.3.3)$$

Atunci, pentru orice două numere a, b cu $a < b$, numărul de zerouri a lui π_d pe intervalul $a < x \leq b$ este egal cu $\sigma(a) - \sigma(b)$.

Deoarece $\pi_k(x) = x^k + \dots$, este evident că $\sigma(-\infty) = d$, $\sigma(+\infty) = 0$, astfel încât numărul de rădăcini reale ale lui π_d este $\sigma(-\infty) - \sigma(+\infty) = d$. Mai mult dacă $\xi_1 > \xi_2 > \dots > \xi_d$ desemnează zerourile lui π_d în ordine descrescătoare, avem comportarea lui σ ca în figura 6.1.



Jacques Charles François Sturm (1803-1855), matematician și fizician elvețian, cunoscut pentru teorema sa asupra șirurilor Sturm, descoperită în 1829 și pentru teoria sa asupra ecuației diferențiale Sturm-Liouville. A avut contribuții semnificative și în domeniul geometriei proiective și diferențiale.

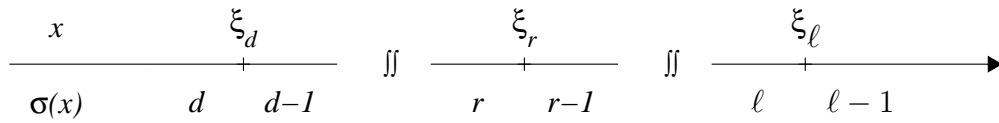


Figura 6.1: Ilustrarea metodei lui Sturm

Este ușor de văzut că:

$$\sigma(x) \leq r - 1 \iff x \geq \xi_r. \quad (6.3.4)$$

Într-adevăr, presupunem că $x \geq \xi_r$. Atunci $\{\# \text{zerouri} \leq x\} \geq d + 1 - r$. Deci conform teoremei lui Sturm, $\sigma(-\infty) - \sigma(x) = d - \sigma(x) = \{\# \text{zerouri} \leq x\} \leq d_1 - r$, adică $\sigma(x) \leq r - 1$. Reciproc, dacă $\sigma(x) \leq r - 1$, atunci, din teorema lui Sturm, $\{\# \text{zerouri} \leq x\} = d - \sigma(x) \geq d + 1 - r$, ceea ce implică $x \geq \xi_r$ (vezi figura 6.1).

Ideea de bază este de a controla procesul de înjumătățire nu ca mai sus, ci mai degrabă verificând inegalitatea (6.3.4) pentru a vedea dacă suntem în stânga sau în dreapta lui ξ_r . Pentru a inițializa procedura, avem nevoie de valorile $a_1 = a$, $b_1 = b$ astfel încât $a < \xi_d$ și $b > \xi_1$. Acestea se obțin trivial ca și capete ale intervalului de ortogonalitate al lui π_d , dacă acesta este finit. Mai general putem aplica teorema lui Gershgorin matricei lui Jacobi J_d asociate polinomului (6.3.2)

$$J_n = \begin{bmatrix} \alpha_0 & \sqrt{\beta_1} & & & 0 \\ \sqrt{\beta_1} & \alpha_1 & \sqrt{\beta_2} & & \\ & \sqrt{\beta_2} & \alpha_2 & \ddots & \\ & & \ddots & \ddots & \sqrt{\beta_{n-1}} \\ 0 & & \sqrt{\beta_{n-1}} & \alpha_{n-1} & \end{bmatrix}$$

ținând cont că zerourile lui π_d sunt valori proprii ale lui J_d .

Teorema lui Gershgorin afirmă că valorile proprii ale matricei $A = [a_{ij}]$ de ordin d sunt localizate în reuniunea discurilor

$$\left\{ z \in \mathbb{C} : |z - a_{ii}| \leq r_i, r_i = \sum_{j \neq i} |a_{ij}| \right\}, \quad i = \overline{1, d}.$$

În acest mod, a poate fi ales ca fiind cel mai mic și b cel mai mare dintre cele d numere $\alpha_0 + \sqrt{\beta_1}$, $\alpha_1 + \sqrt{\beta_1} + \sqrt{\beta_2}$, \dots , $\alpha_{d-2} + \sqrt{\beta_{d-2}} + \sqrt{\beta_{d-1}}$, $\alpha_{d-1} + \sqrt{\beta_{d-1}}$. Metoda șirurilor lui Sturm continuă după cum urmează, pentru orice r cu $1 \leq r \leq d$:

for $n := 1, 2, 3, \dots$ **do**


```

 $x_n := \frac{1}{2}(a_n + b_n);$ 
if  $\sigma(x_n) > r - 1$  then
     $a_{n+1} := x_n; b_{n+1} := b_n;$ 
else
     $a_{n+1} := a_n; b_{n+1} := x_n;$ 
end if
end for

```

Deoarece inițial $\sigma(a) = d > r - 1$, $\sigma(b) = 0 \leq r - 1$, rezultă din construcție că

$$\sigma(a_n) > r - 1, \quad \sigma(b_n) \leq r - 1, \quad n = 1, 2, 3, \dots$$

însemnând că $\xi_r \in [a_n, b_n]$, pentru orice $n = 1, 2, 3, \dots$. Mai mult, deoarece în metoda înjumătățirii, $b_n - a_n = \frac{b - a}{2^{n-1}}$, metoda converge cel puțin liniar către rădăcina ξ_r .

6.4. Metoda falsei poziții

Ca în metoda înjumătățirii, presupunem că avem două numere $a < b$ astfel încât

$$f \in C[a, b], \quad f(a)f(b) < 0 \tag{6.4.1}$$

și generăm un șir descendent de intervale $[a_n, b_n]$, $n = 1, 2, 3, \dots$ cu $a_1 = a$, $b_1 = b$ astfel încât $f(a_n)f(b_n) < 0$. Spre deosebire de metoda înjumătățirii, pentru a determina următorul interval nu luăm mijlocul lui $[a_n, b_n]$, ci soluția $x = x_n$ a ecuației liniare

$$(L_1 f)(x; a_n, b_n) = 0.$$

Aceasta pare să fie o alegere mai flexibilă decât în metoda înjumătățirii deoarece x_n va fi mai apropiat de capătul de care $|f|$ este mai mic.

Procedura decurge după cum urmează:

```

for  $n := 1, 2, \dots$  do
     $x_n := a_n - \frac{a_n - b_n}{f(a_n) - f(b_n)} f(a_n);$ 
    if  $f(a_n)f(x_n) > 0$  then
         $a_{n+1} := x_n; b_{n+1} := b_n;$ 
    else
         $a_{n+1} := a_n; b_{n+1} := x_n;$ 
    end if
end for

```

Iterația se poate termina când $\min(x_n - a_n, b_n - x_n) \leq tol$, unde tol este o valoare dată.

Convergența se analizează mai ușor dacă presupunem că f este convexă sau concavă pe $[a, b]$. Dacă f este convexă, avem

$$f''(x) > 0, \quad x \in [a, b], \quad f(a) < 0, \quad f(b) > 0. \quad (6.4.2)$$

Șirul

$$x_{n+1} = x_n - \frac{x_n - b}{f(x_n) - f(b)} f(x_n), \quad n \in \mathbb{N}^*, \quad x_1 = a \quad (6.4.3)$$

este monoton crescător și mărginit superior de α , deci convergent către o limită x , iar $f(x) = 0$ (vezi figura 6.2).

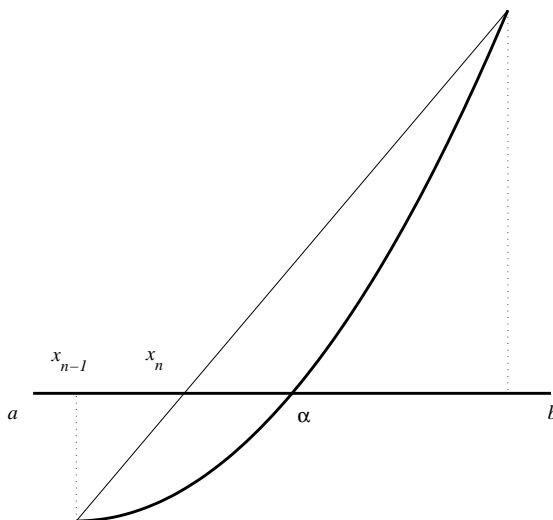


Figura 6.2: Metoda falsei poziții

Viteza de convergență se determină scăzând α din ambii membri ai lui (6.4.3) și utilizând faptul că $f(\alpha) = 0$:

$$x_{n+1} - \alpha = x_n - \alpha - \frac{x_n - b}{f(x_n) - f(b)} [f(x_n) - f(\alpha)].$$

Împărțind cu $x_n - \alpha$ avem

$$\frac{x_{n+1} - \alpha}{x_n - \alpha} = 1 - \frac{x_n - b}{f(x_n) - f(b)} \frac{f(x_n) - f(\alpha)}{x_n - \alpha}.$$

Făcând $n \rightarrow \infty$ și utilizând faptul că $x_n \rightarrow \alpha$, obținem

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{x_n - \alpha} = 1 - (b - \alpha) \frac{f'(\alpha)}{f(b)}. \quad (6.4.4)$$

Deci metoda converge liniar, cu eroarea asimptotică

$$c = 1 - (b - a) \frac{f'(\alpha)}{f(b)}.$$

Datorită ipotezei convexității avem $c \in (0, 1)$. Analog se face demonstrația în cazul când f este concavă. Dacă f nu este nici convexă nici concavă pe $[a, b]$, ci $f \in C^2[a, b]$ și $f''(\alpha) \neq 0$, f'' are semn constant pe o vecinătate a lui α și pentru un n suficient de mare x_n ajunge în acea vecinătate și se poate proceda ca mai sus.

Dezavantaje. (i) Convergența lentă; (ii) Faptul că unul din capete poate rămâne fix. Dacă f este turtită în vecinătatea rădăcinii și a este apropiat de α și b depărtat convergența poate fi foarte lentă.

6.5. Metoda secantei

Este o variantă a metodei falsei poziții, în care nu se mai cere ca f să aibă valori de semne contrare, nici măcar la capetele intervalului inițial.

Se alege două valori arbitrare de pornire x_0, x_1 și se continuă cu

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n), \quad n \in \mathbb{N}^*. \quad (6.5.1)$$

Aceasta preîntâmpină apariția unei false poziții și sugerează o convergență mai rapidă. Din păcate, nu mai are loc convergența „globală“ pe $[a, b]$ ci doar convergența „locală“, adică numai dacă x_0 și x_1 sunt suficient de apropiate de rădăcină.

Vom avea în continuare nevoie de o relație între trei erori consecutive

$$\begin{aligned} x_{n+1} - \alpha &= x_n - \alpha - \frac{f(x_n)}{f[x_{n-1}, x_n]} = (x_n - \alpha) \left(1 - \frac{f(x_n) - f(\alpha)}{(x_n - \alpha)f[x_{n-1}, x_n]} \right) \\ &= (x_n - \alpha) \left(1 - \frac{f[x_n, \alpha]}{f[x_{n-1}, x_n]} \right) = (x_n - \alpha) \frac{f[x_{n-1}, x_n] - f[x_n, \alpha]}{f[x_{n-1}, x_n]} \\ &= (x_n - \alpha)(x_{n-1} - \alpha) \frac{f[x_n, x_{n-1}, \alpha]}{f[x_{n-1}, x_n]}. \end{aligned}$$

Deci

$$(x_{n+1} - \alpha) = (x_n - \alpha)(x_{n-1} - \alpha) \frac{f[x_n, x_{n-1}, \alpha]}{f[x_{n-1}, x_n]}, \quad n \in \mathbb{N}^* \quad (6.5.2)$$

Din (6.5.2) rezultă imediat că dacă α este o rădăcină simplă ($f(\alpha) = 0$, $f'(\alpha) \neq 0$) și $x_n \rightarrow \alpha$ și dacă $f \in C^2$ pe o vecinătate a lui α , convergența este superliniară. Cât este ordinul de convergență?

Înlocuim raportul diferențelor divizate din (6.5.2) cu o constantă, ceea ce este aproape adevărat când n este mare. Punând $c_k = |x_k - \alpha|$, avem

$$e_{n+1} = e_n e_{n-1} C, \quad C > 0$$

Înmulțind ambii membri cu C și punând $E_n = C e_n$ obținem

$$E_{n+1} = E_n E_{n-1}, \quad E_n \rightarrow 0.$$

Logaritmând și punând $y_n = \frac{1}{E_n}$ obținem

$$y_{n+1} = y_n + y_{n-1}, \quad (6.5.3)$$

care este recurența pentru șirul lui Fibonacci. Soluția este

$$y_n = c_1 t_1^n + c_2 t_2^n,$$

c_1, c_2 constante și

$$t_1 = \frac{1}{2}(1 + \sqrt{5}), \quad t_2 = \frac{1}{2}(1 - \sqrt{5}).$$

Deoarece $y_n \rightarrow \infty$, avem $c_1 \neq 0$ și $y_n \sim c_1 t_1^n$, căci $|t_2| < 1$. Revenind la substituție $\frac{1}{E_n} \sim e^{c_1 t_1^n}$, $\frac{1}{e_n} \sim C e^{c_1 t_1^n}$ și deci

$$\frac{e_{n+1}}{e_n} \sim \frac{C^{t_1} e^{c_1 t_1^{n+1}}}{C e^{c_1 t_1^{n+1}}} = C^{t_1-1}, \quad n \rightarrow \infty.$$

Ordinul de convergență este $t_1 = \frac{1 + \sqrt{5}}{2} \approx 1.61803 \dots$ (secțiunea de aur).

Teorema 6.5.1 Fie α un zero simplu al lui f și fie $I_\varepsilon = \{x \in \mathbb{R} : |x - \alpha| < \varepsilon\}$ și presupunem că $f \in C^2[I_\varepsilon]$. Definim pentru ε suficient de mic

$$M(\varepsilon) = \max_{\substack{s \in I_\varepsilon \\ t \in I_\varepsilon}} \left| \frac{f''(s)}{2f'(t)} \right|. \quad (6.5.4)$$

Presupunem că

$$\varepsilon M(\varepsilon) < 1 \quad (6.5.5)$$

Atunci metoda secantei converge către rădăcina unică $\alpha \in I_\varepsilon$ pentru orice valori de pornire $x_0 \neq x_1$ cu $x_0 \in I_\varepsilon$, $x_1 \in I_\varepsilon$.

Observația 6.5.2. Se observă că $\lim_{\varepsilon \rightarrow 0} M(\varepsilon) = \left| \frac{f''(\alpha)}{2f'(\alpha)} \right| < \infty$, deci (6.5.5) poate fi satisfăcută pentru ε suficient de mic. Natura locală a convergenței este cuantificată prin cerința ca $x_0, x_1 \in I_\varepsilon$. \diamond

Demonstrație. Se observă că α este singurul zero al lui f în I_ε . Aceasta rezultă din formula lui Taylor pentru $x = \alpha$:

$$f(x) = f(\alpha) + (x - \alpha)f'(\alpha) + \frac{(x - \alpha)^2}{2}f''(\xi)$$

unde $f(\alpha) = 0$ și $\xi \in (x, \alpha)$ (sau (α, x)). Astfel dacă $x \in I_\varepsilon$, atunci și $\xi \in I_\varepsilon$ și avem

$$f(x) = (x - \alpha)f'(\alpha) \left[1 + \frac{x - \alpha}{2} \frac{f''(\xi)}{f'(\alpha)} \right]$$

Aici, dacă $x \neq \alpha$, toți trei factorii sunt diferiți de 0, căci

$$\left| \frac{x - \alpha}{2} \frac{f''(\xi)}{f'(\alpha)} \right| \leq \varepsilon M(\varepsilon) < 1.$$

Deci f se poate anula pe I_ε numai în $x = \alpha$. Să arătăm că $x_n \in I_\varepsilon$ pentru orice n , în afară de cazul când $f(x_n) = 0$, în care $x_n = \alpha$ și metoda converge într-un număr finit de pași. Vom demonstra aceasta prin inducție: presupunem că $x_{n-1}, x_n \in I_\varepsilon$ și $x_n \neq x_{n-1}$. Acest lucru este adevărat pentru $n = 1$ din ipoteză. Deoarece $f \in C^2[I_\varepsilon]$

$$f[x_{n-1}, x_n] = f'(\xi_1), \quad f[x_{n-1}, x_n, \alpha] = \frac{1}{2}f''(\xi_2), \quad \xi_i \in I_\varepsilon, \quad i = 1, 2,$$

din (6.5.2) rezultă

$$|x_{n+1} - \alpha| \leq \varepsilon^2 \left| \frac{f''(\xi_2)}{2f'(\xi_1)} \right| \leq \varepsilon \varepsilon M(\varepsilon) < \varepsilon,$$

adică $x_{n+1} \in I_\varepsilon$. Mai mult, din relația între trei erori consecutive, (6.5.2), rezultă $x_{n+1} \neq x_n$ în afară de cazul când $f(x_n) = 0$ (și atunci $x_n = \alpha$). Utilizând (6.5.2) avem

$$|x_{n+1} - \alpha| \leq |x_n - \alpha| \varepsilon M(\varepsilon)$$

care aplicată repetat ne dă

$$|x_{n+1} - \alpha| \leq |x_n - \alpha| \varepsilon M(\varepsilon) \leq \dots \leq [\varepsilon M(\varepsilon)]^{n-1} |x_1 - \alpha|.$$

Cum $\varepsilon M(\varepsilon) < 1$, rezultă că metoda este convergentă și $x_n \rightarrow \alpha$ când $n \rightarrow \infty$. \square

Deoarece este nevoie de o singură evaluare a lui f pe pas, indicele de eficiență este $p = \frac{1+\sqrt{5}}{2} \approx 1.61803\dots$. Pseudocodul metodei este dat în algoritmul 6.1.

Algoritmul 6.1 Metoda secantei pentru ecuații neliniare în \mathbb{R}

Intrare: Funcția f , valorile de pornire x_0 și x_1 , numărul maxim de iterații, $Nmax$, informații de toleranță tol

Ieșire: O aproximație a rădăcinii sau un mesaj de eroare

```

1:  $x_c := x_1; \quad x_v = x_0;$ 
2:  $f_c := f(x_1); \quad f_v := f(x_0);$ 
3: for  $k := 1$  to  $Nmax$  do
4:    $x_n := x_c - f_c * (x_c - x_v) / (f_c - f_v);$ 
5:   if  $crit\_oprire(tol)$  then
6:     return  $x_n; \{Succes\}$ 
7:   end if
8:    $x_v := x_c; \quad f_v := f_c; \quad x_c := x_n; \quad f_c = f(x_n);$ 
9: end for
10:  $error("S-a depășit numărul de iterații").$ 

```

6.6. Metoda lui Newton

Poate fi privită ca un caz la limită al metodei secantei, când $x_{n-1} \rightarrow x_n$. Obținem iterația

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (6.6.1)$$

O altă interpretare mult mai fructuoasă este liniarizarea ecuației $f(x) = 0$ în $x = x_n$:

$$f(x) \approx f(x_n) + (x - x_n)f'(x_n) = 0.$$

Privită în acest mod metoda lui Newton se poate generaliza la ecuații neliniare de toate tipurile (sisteme neliniare, ecuații funcționale, caz în care f' trebuie înțeleasă ca derivată Fréchet), iar iterația este

$$x_{n+1} = x_n - [f'(x_n)]^{-1} f(x_n) \quad (6.6.2)$$

Studiul erorii în metoda lui Newton este la fel ca cel al erorii în metoda secantei.

$$\begin{aligned}
x_{n+1} - \alpha &= x_n - \alpha - \frac{f(x_n)}{f'(x_n)} \\
&= (x_n - \alpha) \left[1 - \frac{f(x_n) - f(\alpha)}{(x_n - \alpha)f'(x_n)} \right] \\
&= (x_n - \alpha) \left(1 - \frac{f[x_n, \alpha]}{f[x_n, x_n]} \right) = (x_n - \alpha)^2 \frac{f[x_n, x_n, \alpha]}{f[x_n, x_n]}
\end{aligned} \quad (6.6.3)$$

De aceea, dacă $x_n \rightarrow \alpha$, atunci

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^2} = \frac{f''(\alpha)}{2f'(\alpha)}$$

și ordinul de convergență al metodei lui Newton este 2 dacă $f''(\alpha) \neq 0$. Referitor la convergența locală a metodei lui Newton avem

Teorema 6.6.1 Fie α o rădăcină simplă a ecuației $f(x) = 0$ și $I_\varepsilon = \{x \in \mathbb{R} : |x - \alpha| \leq \varepsilon\}$. Presupunem că $f \in C^2[I_\varepsilon]$. Definim

$$M(\varepsilon) = \max_{\substack{s \in I_\varepsilon \\ t \in I_\varepsilon}} \left| \frac{f''(s)}{2f'(t)} \right| \quad (6.6.4)$$

Dacă ε este suficient de mic astfel încât

$$2\varepsilon M(\varepsilon) < 1, \quad (6.6.5)$$

atunci pentru orice $x_0 \in I_\varepsilon$, metoda lui Newton este bine definită și converge pătratic către singura rădăcină $\alpha \in I_\varepsilon$.

Criteriul de oprire pentru metoda lui Newton

$$|x_n - x_{n-1}| < \varepsilon$$

se bazează pe următoarea propoziție:

Propoziția 6.6.2 Fie (x_n) șirul de aproximante generat prin metoda lui Newton. Dacă α este o rădăcină simplă din $[a, b]$, $f \in C^2[a, b]$ și metoda este convergentă, atunci există un $n_0 \in \mathbb{N}$ astfel încât

$$|x_n - \alpha| \leq |x_n - x_{n-1}|, \quad n > n_0.$$

Demonstrație. Vom arăta întâi că

$$|x_n - \alpha| \leq \frac{1}{m_1} |f(x_n)|, \quad m_1 \leq \inf_{x \in [a, b]} |f'(x)|. \quad (6.6.6)$$

Utilizând teorema lui Lagrange, $f(\alpha) - f(x_n) = f'(\xi)(\alpha - x_n)$, cu $\xi \in (\alpha, x_n)$ (sau (x_n, α)). Din relațiile $f(\alpha) = 0$ și $|f'(x)| \geq m_1$ pentru $x \in (a, b)$ rezultă că $|f(x_n)| \geq m_1 |\alpha - x_n|$, adică chiar (6.6.6).

Pe baza formulei lui Taylor avem

$$f(x_n) = f(x_{n-1}) + (x_n - x_{n-1})f'(x_{n-1}) + \frac{1}{2}(x_n - x_{n-1})^2 f''(\mu), \quad (6.6.7)$$

cu $\mu \in (x_{n-1}, x_n)$ sau $\mu \in (x_n, x_{n-1})$. Ținând cont de modul de obținere a unei aproximații de metoda lui Newton, avem $f(x_{n-1}) + (x_n - x_{n-1})f'(x_{n-1}) = 0$ și din (6.6.7) se obține

$$|f(x_n)| = \frac{1}{2}(x_n - x_{n-1})^2 |f''(\mu)| \leq \frac{1}{2}(x_n - x_{n-1})^2 \|f''\|_\infty,$$

iar pe baza formulei (6.6.6) rezultă că

$$|\alpha - x_n| \leq \frac{\|f''\|_\infty}{2m_1}(x_n - x_{n-1})^2.$$

Cum am presupus că metoda este convergentă, există un n_0 natural cu proprietatea că

$$\frac{\|f''\|_\infty}{2m_1}(x_n - x_{n-1}) < 1, \quad n > n_0$$

și deci

$$|x_n - \alpha| \leq |x_n - x_{n-1}|, \quad n > n_0.$$

□

Interpretarea geometrică a metodei lui Newton apare în figura 6.3, iar descrierea metodei în algoritmul 6.2.

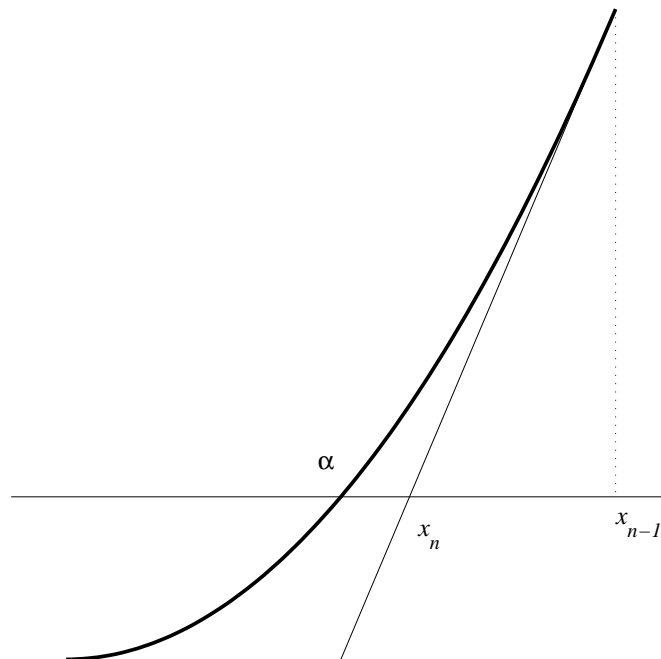


Figura 6.3: Metoda lui Newton

Alegerea valorii de pornire este, în general, o problemă dificilă. În practică, se alege o valoare, iar dacă după un număr maxim fixat de iterații nu s-a obținut precizia dorită, testată prin unul din criteriile uzuale, se încearcă cu altă valoare de pornire. De exemplu, dacă rădăcina este izolată într-un interval $[a, b]$ și $f''(x) \neq 0$, $x \in (a, b)$, un criteriu de alegere este $f(x_0)f''(x_0) > 0$.

Algoritmul 6.2 Metoda lui Newton pentru ecuații neliniare în \mathbb{R}

Intrare: Funcția f , derivata f' , valoarea de pornire x_0 , numărul maxim de iterații, $Nmax$, informații de toleranță tol

Ieșire: O aproximație a rădăcinii sau un mesaj de eroare

```

1: for  $k := 0$  to  $Nmax$  do
2:    $x_{k+1} := x_k - \frac{f(x_k)}{f'(x_k)}$ ;
3:   if  $crit\_oprive(tol)$  then
4:     return  $x_{k+1}; \{ Succes \}$ 
5:   end if
6: end for
7:  $error("S-a depășit numărul de iterații")$ .

```

6.7. Metoda aproximațiilor succesive

Adesea, în aplicații, ecuațiile neliniare apar sub forma unei *probleme de punct fix*: să se determine x astfel încât

$$x = \varphi(x). \quad (6.7.1)$$

Un număr α ce satisface această ecuație se numește punct fix al lui φ . Orice ecuație $f(x) = 0$ se poate scrie (în multe moduri diferite) în forma echivalentă (6.7.1). De exemplu, dacă $f'(x) \neq 0$ în intervalul de interes, putem lua

$$\varphi(x) = x - \frac{f(x)}{f'(x)}. \quad (6.7.2)$$

Dacă x_0 este o aproximație inițială a unui punct fix α a lui (6.7.1) atunci metoda aproximațiilor succesive generează un șir de aproximații

$$x_{n+1} = \varphi(x_n). \quad (6.7.3)$$

Dacă acest șir converge și φ este continuă, atunci șirul converge către un punct fix a lui φ . De notat că (6.7.3) este chiar metoda lui Newton dacă φ este dată de (6.7.2). Astfel metoda lui Newton poate fi privită ca o iterație de tip punct fix, dar nu și metoda secantei. Pentru o iterație de forma (6.7.3), presupunând că $x_n \rightarrow \alpha$ când $n \rightarrow \infty$, ordinul de convergență este ușor de determinat. Să presupunem că în punctul fix α avem

$$\varphi'(\alpha) = \varphi''(\alpha) = \dots = \varphi^{(p-1)}(\alpha) = 0, \quad \varphi^{(p)}(\alpha) \neq 0 \quad (6.7.4)$$

Presupunem că $\varphi \in C^p$ pe o vecinătate a lui α . Avem atunci, conform teoremei lui Taylor

$$\begin{aligned} \varphi(x_n) &= \varphi(\alpha) + (x_n - \alpha)\varphi'(\alpha) + \dots + \frac{(x_n - \alpha)^{p-1}}{(p-1)!}\varphi^{(p-1)}(\alpha) \\ &\quad + \frac{(x_n - \alpha)^p}{p!}\varphi^{(p)}(\xi_n) = \varphi(\alpha) + \frac{(x_n - \alpha)^p}{p!}\varphi^{(p)}(\xi_n), \end{aligned}$$

unde $\xi_n \in (\alpha, x_n)$ (sau (x_n, α)). Deoarece $\varphi(x_n) = x_{n+1}$ și $\varphi(\alpha) = \alpha$ obținem

$$\frac{x_{n+1} - \alpha}{(x_n - \alpha)^p} = \frac{1}{p!} \varphi^{(p)}(\xi_n).$$

Când $x_n \rightarrow \alpha$, deoarece ξ_n este între x_n și α , deducem pe baza continuității că

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^p} = \frac{1}{p!} \varphi^{(p)}(\alpha) \neq 0. \quad (6.7.5)$$

Aceasta ne arată că ordinul de convergență este exact p și eroarea asimptotică este

$$c = \frac{1}{p!} \varphi^{(p)}(\alpha). \quad (6.7.6)$$

Combinând aceasta cu condiția uzuală de convergență locală se obține:

Teorema 6.7.1 Fie α un punct fix al lui φ și $I_\varepsilon = \{x \in \mathbb{R} : |x - \alpha| \leq \varepsilon\}$. Presupunem că $\varphi \in C^p[I_\varepsilon]$ și satisface (6.7.4). Dacă

$$M(\varepsilon) := \max_{t \in I_\varepsilon} |\varphi'(t)| < 1 \quad (6.7.7)$$

atunci iterația (6.7.3) converge către α , $\forall x_0 \in I_\varepsilon$. Ordinul de convergență este p , iar eroarea asimptotică este dată de (6.7.6).

6.8. Metoda lui Newton pentru rădăcini multiple

Dacă α este o rădăcină multiplă de ordinul m , atunci ordinul de convergență a metodei lui Newton este doar 1. Într-adevăr, fie

$$\varphi(x) = x - \frac{f(x)}{f'(x)}.$$

Deoarece

$$\varphi'(x) = \frac{f(x)f''(x)}{[f'(x)]^2}$$

procesul va fi convergent dacă $\varphi'(\alpha) = 1 - 1/m < 1$.

O modalitate de a evita aceasta este să rezolvăm ecuația

$$u(x) := \frac{f(x)}{f'(x)} = 0$$

care are aceleași rădăcini ca și f , dar simple. Metoda lui Newton pentru problema modificată are forma

$$x_{k+1} = x_k - \frac{u(x_k)}{u'(x_k)} = \frac{f(x_k)f'(x_k)}{[f'(x_k)]^2 - f(x_k)f''(x_k)}. \quad (6.8.1)$$

Deoarece α este o rădăcină simplă a lui u , convergența lui (6.8.1) este pătratică. Singurul dezavantaj teoretic al lui (6.8.1) este derivata a doua necesară suplimentar și complexitatea mai mare a calculului lui x_{k+1} din x_k . În practică aceasta este o slăbiciune, deoarece numitorul lui (6.8.1) poate lua valori foarte mici în vecinătatea lui α când $x_k \rightarrow \alpha$.

Convergența pătratică a metodei lui Newton se poate realiza nu numai prin modificarea problemei ci și prin modificarea metodei. În vecinătatea unei soluții multiple de ordinul m , α , avem

$$f(x) = (x - \alpha)^m \varphi(x) \approx (x - \alpha)^m \cdot c, \quad (6.8.2)$$

de unde rezultă

$$\frac{f(x)}{f'(x)} \approx \frac{x - \alpha}{m} \Rightarrow \alpha \approx x - m \frac{f(x)}{f'(x)}.$$

Metoda modificată corespunzătoare

$$x_{k+1} := x_k - m \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, 2, \dots \quad (6.8.3)$$

converge pătratic către rădăcina multiplă de ordinul m când se întrebuițează o valoare corectă a lui m în (6.8.3). Eficiența variantei (6.8.3) a metodei lui Newton depinde de utilizarea unei valori de aproximare bune pentru m , dacă această valoare nu este cunoscută din alte surse.

În ipoteza

$$|x_k - \alpha| < |x_{k-1} - \alpha| \wedge |x_k - \alpha| < |x_{k-2} - \alpha|$$

putem înlocui în (6.8.2) α prin x_k

$$\begin{aligned} f(x_{k-1}) &\approx (x_{k-1} - x_k)^m \cdot c \\ f(x_{k-2}) &\approx (x_{k-2} - x_k)^m \cdot c. \end{aligned}$$

În continuare se obține m :

$$m \approx \frac{\log [f(x_{k-1})/f(x_{k-2})]}{\log [(x_{k-1} - x_k)/(x_{k-2} - x_k)]}.$$

Această valoare poate fi utilizată în (6.8.3).

6.9. Ecuații algebrice

Există multe metode special concepute pentru a rezolva ecuații algebrice. Aici vom descrie numai metoda lui Newton aplicată în acest context, concentrându-ne asupra unui mod eficient de a evalua simultan valoarea polinomului și a primei derivate.

Metoda lui Newton aplicată ecuațiilor algebrice. Considerăm o ecuație algebrică de grad d

$$f(x) = 0, \quad f(x) = x^d + a_{d-1}x^{d-1} + \dots + a_0, \quad (6.9.1)$$

în care coeficientul dominant se presupune (fără a restrânge generalitatea) a fi egal cu 1 și unde putem presupune, fără a restrânge generalitatea că $a_0 \neq 0$. Pentru simplitate vom presupune că toți coeficienții sunt reali. Pentru a aplica metoda lui Newton ecuației (6.9.1) este nevoie de a metodă bună de evaluare a polinomului și derivatei.

Schema lui Horner este bună pentru așa ceva:

```
bd := 1;   cd := 1;
for k = d - 1 downto 1 do
    bk := tbk+1 + ak;
    ck := tck+1 + bk;
end for
b0 := tb1 + a0;
```

Atunci $f(t) = b_0$, $f'(t) = c_1$.

Deci procedăm astfel:

Se aplică metoda lui Newton, calculând simultan $f(x_n)$ și $f'(x_n)$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

Se aplică apoi metoda lui Newton polinomului $\frac{f(x)}{x - \alpha}$. Pentru rădăcini complexe se începe cu x_0 complex și toate calculele se fac în aritmetică complexă. Este posibil să se împartă cu factori pătratici și să se folosească aritmetica reală – se ajunge astfel la metoda lui Bairstow. Folosind metoda aceasta de scădere a gradului erorile pot fi mari. O modalitate de îmbunătățire este de a utiliza rădăcinile astfel calculate ca aproximații inițiale și a aplica metoda lui Newton polinomului original.

6.10. Metoda lui Newton în \mathbb{R}^n

Metoda lui Newton este ușor de generalizat la sisteme neliniare

$$F(x) = 0, \quad (6.10.1)$$

unde $F : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$, iar $x, F(x) \in \mathbb{R}^n$. Sistemul (6.10.1) se scrie pe componente

$$\begin{cases} F_1(x_1, \dots, x_n) = 0 \\ \vdots \\ F_n(x_1, \dots, x_n) = 0 \end{cases}$$

Fie $F'(x^{(k)})$ jacobianul lui F în $x^{(k)}$:

$$J := F'(x^{(k)}) = \begin{bmatrix} \frac{\partial F_1}{\partial x_1}(x^{(k)}) & \cdots & \frac{\partial F_1}{\partial x_n}(x^{(k)}) \\ \vdots & \ddots & \vdots \\ \frac{\partial F_n}{\partial x_1}(x^{(k)}) & \cdots & \frac{\partial F_n}{\partial x_n}(x^{(k)}) \end{bmatrix}. \quad (6.10.2)$$

Cantitatea $1/f'(x)$ se înlocuiește în acest caz cu inversa jacobianului în $x^{(k)}$:

$$x^{(k+1)} = x^{(k)} - [F'(x^{(k)})]^{-1} F(x^{(k)}). \quad (6.10.3)$$

Scriem iterația sub forma

$$x^{(k+1)} = x^{(k)} + w^{(k)}. \quad (6.10.4)$$

Se observă că w_k este soluția sistemului de n ecuații liniare cu n necunoscute

$$F'(x^{(k)})w^{(k)} = -F(x^{(k)}). \quad (6.10.5)$$

Este mai eficient și mai convenabil ca, în loc să inversăm jacobianul la fiecare pas, să rezolvăm sistemul (6.10.5) și să folosim iterația în forma (6.10.4).

Teorema 6.10.1 *Fie α o soluție a ecuației $F(x) = 0$ și presupunem că în bila închisă $B(\delta) \equiv \{x : \|x - \alpha\| \leq \delta\}$, există matricea Jacobi a lui $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, este nesingulară și satisface condiția Lipschitz*

$$\|F'(x) - F'(y)\|_\infty \leq c\|x - y\|_\infty, \quad \forall x, y \in B(\delta), \quad c > 0.$$

Punem $\gamma = c \max\{\|[F'(x)]^{-1}\|_\infty : \|\alpha - x\|_\infty \leq \delta\}$ și $0 < \varepsilon < \min\{\delta, \gamma^{-1}\}$. Atunci pentru orice aproximație inițială $x^{(0)} \in B(\varepsilon) := \{x : \|x - \alpha\|_\infty \leq \varepsilon\}$ metoda lui Newton este convergentă, iar vectorii $e^{(k)} := \alpha - x^{(k)}$ satisfac următoarele inegalități:

$$(a) \|e^{(k+1)}\|_\infty \leq \gamma \|e^{(k)}\|_\infty^2$$

$$(b) \|e^{(k)}\|_\infty \leq \gamma^{-1} (\gamma \|e^{(0)}\|_\infty)^{2^k}.$$

Demonstrație. Dacă F' este continuă pe segmentul ce unește punctele $x, y \in \mathbb{R}^n$, conform teoremei lui Lagrange

$$F(x) - F(y) = J_k(x - y),$$

unde

$$J_k = \begin{bmatrix} \frac{\partial F_1}{\partial x_1}(\xi_1) & \cdots & \frac{\partial F_1}{\partial x_n}(\xi_1) \\ \vdots & \ddots & \vdots \\ \frac{\partial F_n}{\partial x_1}(\xi_n) & \cdots & \frac{\partial F_n}{\partial x_n}(\xi_n) \end{bmatrix} \Rightarrow$$

$$\begin{aligned} e^{(k+1)} &= e^{(k)} - [F'(x^{(k)})]^{-1}(F(\alpha) - F(x^{(k)})) = e^{(k)} - [F'(x^{(k)})]^{-1}J_k e^{(k)} \\ &= [F'(x^{(k)})]^{-1}(F'(x^{(k)}) - J_k)e^{(k)} \end{aligned}$$

și de aici rezultă imediat (a). Din condiția Lipschitz

$$\|F'(x^{(k)}) - J_k\|_\infty \leq c \max_{j=1, n} \|x^{(k)} - \xi^{(j)}\| \leq c \|x^{(k)} - \alpha\|$$

Deci, dacă $\|\alpha - x^{(k)}\|_\infty \leq \varepsilon$, atunci $\|\alpha - x^{(k+1)}\|_\infty \leq (\gamma\varepsilon)\varepsilon \leq \varepsilon$. Deoarece (a) este adevărată pentru orice k , se obține (b) imediat. \square

Metoda lui Newton pentru sisteme de ecuații neliniare este descrisă în algoritmul 6.3.

Algoritmul 6.3 Metoda lui Newton pentru sisteme de ecuații neliniare

Intrare: Funcția F , derivata Fréchet F' , vectorul de pornire $x^{(0)}$, numărul maxim de iterații, $Nmax$, informații de toleranță tol

Ieșire: O aproximație a rădăcinii sau un mesaj de eroare

- 1: **for** $k := 0$ **to** $Nmax$ **do**
 - 2: Calculează matricea jacobian $J = F'(x^{(k)})$;
 - 3: Rezolvă sistemul $Jw = -F(x^{(k)})$;
 - 4: $x^{(k+1)} := x^{(k)} + w$;
 - 5: **if** $crit_oprire(tol)$ **then**
 - 6: **return** $x^{(k+1)}$; {Succes}
 - 7: **end if**
 - 8: **end for**
 - 9: $error$ ("S-a depășit numărul de iterații").
-

6.11. Metode quasi-Newton

O slăbiciune semnificativă a metodei lui Newton pentru rezolvarea sistemelor de ecuații neliniare este necesitatea ca la fiecare pas să calculăm matricea jacobiană și să rezolvăm un sistem $n \times n$ cu această matrice. Pentru a ilustra dimensiunile unei astfel de slăbiciuni, să evaluăm volumul de calcule asociat cu o iterație a metodei lui Newton. Matricea jacobiană asociată unui sistem de n ecuații neliniare scris în forma $F(x) = 0$ necesită evaluarea celor n^2 derivate parțiale ale celor n funcții componente ale lui F . În cele mai multe situații, evaluarea exactă a derivatelor parțiale este neconvenabilă și de multe ori imposibilă. Efortul total pentru o iterație a metodei lui Newton va fi de cel puțin $n^2 + n$ evaluări de funcții scalare (n^2 pentru evaluarea jacobianului și n pentru evaluarea lui F) și $O(n^3)$ operații aritmetice pentru a rezolva sistemul liniar. Acest volum de

calculul este prohibitiv, exceptând valori mici ale lui n și funcții scalare ușor de evaluat. Este firesc ca atenția să fie îndreptată spre reducerea numărului de evaluări și evitarea rezolvării unui sistem liniar la fiecare pas.

La metoda secantei aproximația următoare $x^{(k+1)}$ se obține ca soluție a ecuației liniare

$$\bar{l}_k = f(x^{(k)}) + (x - x^{(k)}) \frac{f(x^{(k)} + h_k) - f(x^{(k)})}{h_k} = 0.$$

Aici funcția \bar{l}_k poate fi interpretată în două moduri:

1. ca aproximare a ecuației tangentei

$$l_k(x) = f(x^{(k)}) + (x - x^{(k)})f'(x^{(k)});$$

2. ca interpolare liniară între punctele $x^{(k)}$ și $x^{(k+1)}$.

Se pot obține diverse generalizări ale metodei secantei la sisteme de ecuații neliniare în funcție de modul în care se interpretează \bar{l}_k . Prima interpretare conduce la metode de tip Newton discretizate, iar a doua la metode bazate pe interpolare.

Metodele de tip Newton discretizate se obțin dacă în metoda lui Newton (6.10.3) $F'(x)$ se înlocuiește cu o aproximare discretă $A(x, h)$. Derivatele parțiale din matricea jacobiană (6.10.2) se vor înlocui prin diferențele divizate

$$A(x, h)e_i := [F(x + h_i e_i) - F(x)]/h_i, \quad i = \overline{1, n}, \quad (6.11.1)$$

unde $e_i \in \mathbb{R}^n$ este al i -lea vector al bazei canonice și $h_i = h_i(x)$ este mărimea pasului de discretizare. O alegere posibilă a pasului este de exemplu

$$h_i := \begin{cases} \varepsilon |x_i|, & \text{dacă } x_i \neq 0; \\ \varepsilon, & \text{altfel,} \end{cases}$$

cu $\varepsilon := \sqrt{\text{eps}}$, unde eps este epsilon-ul mașinii.

6.11.1. Interpolare liniară

La interpolare fiecare dintre planele tangente se înlocuiește cu un (hiper)plan care interpolatează funcțiile componente F_i ale lui F în $n + 1$ puncte date $x^{k,j}$, $j = \overline{0, n}$, într-o vecinătate a lui $x^{(k)}$, adică se determină vectorii $a^{(i)}$ și scalarii α_i , astfel încât pentru

$$L_i(x) = \alpha_i + a^{(i)T} x, \quad i = \overline{1, n} \quad (6.11.2)$$

are loc

$$L_i(x^{k,j}) = F_i(x^{k,j}), \quad i = \overline{1, n}, \quad j = \overline{0, n}.$$

Următoarea aproximație $x^{(k+1)}$ se obține ca punct de intersecție între cele n hiperplane (6.11.2) din \mathbb{R}^{n+1} cu hiperplanul $y = 0$. $x^{(k+1)}$ rezultă ca soluție a sistemului de ecuații liniare

$$L_i(x) = 0, \quad i = \overline{1, n}. \quad (6.11.3)$$

În funcție de alegerea punctelor de interpolare se obțin diferite metode, dintre care cele mai cunoscute sunt metoda lui Brown și metoda lui Brent. Metoda lui Brown combină aproximarea lui F' și rezolvarea sistemului prin eliminare gaussiană. În metoda lui Brent se întrebuițează la rezolvarea sistemului metoda QR. Ambele metode aparțin unei clase, care, la fel ca metoda lui Newton converg pătratic, dar au nevoie doar de $(n^2 + 3n)/2$ evaluări de funcții pe iterație.

Într-un studiu comparativ, Moré și Cosnard [28] au ajuns la concluzia că metoda Brent este adeseori de preferat metodei lui Brown și că pentru sisteme de ecuații neliniare, la care evaluarea lui F necesită un efort mai mic, metoda lui Newton discretizată este cea mai eficientă metodă de rezolvare.

6.11.2. Metode de modificare

Din punct de vedere al efortului de calcul, sunt deosebit de convenabile metodele în care la fiecare pas se întrebuițează o aproximare A_k a lui $F'(x^{(k)})$, care se obține din A_{k-1} printr-o modificare de rang 1, adică prin adăugarea unei matrice de rang 1:

$$A_{k+1} := A_k + u^{(k)} [v^{(k)}]^T, \quad u^{(k)}, v^{(k)} \in \mathbb{R}^n, \quad k = 0, 1, 2, \dots$$

Pe baza formulei Sherman-Morrison (vezi [14])

$$(A + uv^T)^{-1} = A^{-1} - \frac{1}{1 + v^T A^{-1} u} A^{-1} uv^T A^{-1},$$

pentru $B_{k+1} := A_{k+1}^{-1}$ are loc relația de recurență

$$B_{k+1} = B_k - \frac{B_k u^{(k)} [v^{(k)}]^T B_k}{1 + [v^{(k)}]^T B_k u^{(k)}}, \quad k = 0, 1, 2, \dots,$$

atât timp cât $1 + [v^{(k)}]^T B_k u^{(k)} \neq 0$. Astfel, nu mai este necesară rezolvarea unui sistem liniar la fiecare pas; ea se înlocuiește cu înmulțiri matrice-vector, ceea ce corespunde unei reduceri a efortului de calcul de la $O(n^3)$ la $O(n^2)$. Acest avantaj va fi plătit prin aceea că nu vom mai avea o convergență pătratică ca la metoda lui Newton, ci doar una superliniară:

$$\lim_{k \rightarrow \infty} \frac{\|x^{(k+1)} - \alpha\|}{\|x^{(k)} - \alpha\|} = 0. \quad (6.11.4)$$

În metoda lui Broyden alegerea vectorilor $u^{(k)}$ și $v^{(k)}$ are loc după principiul aproximației secantei. În cazul scalar aproximația $a_k \approx f'(x^{(k)})$ se face unic prin

$$a_{k+1}(x^{(k+1)} - x^{(k)}) = f(x^{(k+1)}) - f(x^{(k)}).$$

Pentru $n > 1$, din contră, aproximația

$$A_{k+1}(x^{(k+1)} - x^{(k)}) = F(x^{(k+1)}) - F(x^{(k)}) \quad (6.11.5)$$

(așa numita ecuație quasi-Newton) nu mai este unic determinată; orice altă matrice de forma

$$\bar{A}_{k+1} := A_{k+1} + pq^T$$

cu $p, q \in \mathbb{R}^n$ și $q^T(x^{(k+1)} - x^{(k)}) = 0$ verifică de asemenea ecuația (6.11.5). Pe de altă parte,

$$y_k := F(x^{(k)}) - F(x^{(k-1)}) \text{ și } s_k := x^{(k)} - x^{(k-1)}$$

conțin numai informații despre derivata parțială a lui F în direcția s_k , dar nici o informație în direcții ortogonale pe s_k . Pe această direcție trebuie ca efectul lui A_{k+1} și A_k să coincidă

$$A_{k+1}q = A_kq, \quad \forall q \in \{v : v \neq 0, v^T s_k = 0\}. \quad (6.11.6)$$

Pornind de la prima aproximație $A_0 \approx F'(x^{(0)})$, se generează șirul A_1, A_2, \dots utilizând formulele (6.11.5) și (6.11.6) (Broyden [7], Dennis și Moré [14]).

Pentru șirul $B_0 = A_0^{-1} \approx [F'(x^{(0)})]^{-1}$, B_1, B_2, \dots cu ajutorul formulei Sherman-Morrisson se obține relația de recurență

$$B_{k+1} := B_k + \frac{(s_{k+1} - B_k y_{k+1}) s_{k+1}^T B_k}{s_{k+1}^T B_k y_{k+1}}, \quad k = 0, 1, 2, \dots$$

care necesită doar înmulțiri matrice vector și a cărei complexitate este doar $O(n^2)$. Cu ajutorul matricelor B_k se poate defini metoda lui Broyden prin

$$x^{(k+1)} := x^{(k)} - B_k F(x^{(k)}), \quad k = 0, 1, 2, \dots$$

Această metodă converge superliniar în sensul lui (6.11.4), dacă pașii s_k se apropie asimptotic (când $k \rightarrow \infty$) de vectorii de actualizare (corecție) ai metodei lui Newton. Se poate recunoaște în aceasta semnificația centrală a principiului linearizării locale la rezolvarea ecuațiilor neliniare.

Metoda lui Broyden este descrisă în algoritmul 6.4.

Algoritmul 6.4 Metoda lui Broyden pentru sisteme de ecuații neliniare

Intrare: Funcția F , vectorul de pornire $x^{(0)}$, numărul maxim de iterații, $Nmax$, informații de toleranță tol

Ieșire: O aproximație a rădăcinii sau un mesaj de eroare

```

1:  $B_0 := F'(x^{(0)}); \quad v := F(x); \quad B := B_0^{-1};$ 
2:  $s := -Bv; \quad x := x + s;$ 
3: for  $k := 1$  to  $Nmax$  do
4:    $w := v; \quad v := F(x); \quad y := v - w;$ 
5:    $z := -By; \{z = -B_{k-1}y_k\}$ 
6:    $p := -s^T z; \{p = s_k^T B_{k-1}y_k\}$ 
7:    $C := pI + (s + z)s^T; \{C = s_k^T B_{k-1}^{-1}y_k I + (s_k + B_{k-1}y_k)s_k^T\}$ 
8:    $B := (1/p)CB; \{B = B_k\}$ 
9:    $s := -Bv; \{s = -B_k F(x^{(k)})\}$ 
10:   $x := x + s;$ 
11:  if  $crit\_oprive(tol)$  then
12:    return  $x; \{succes\}$ 
13:  end if
14: end for
15:  $error("S-a depășit numărul maxim de iterații")$ 

```

CAPITOLUL 7

Vectori și valori proprii

Cuprins

7.1. Valori proprii și rădăcini ale polinoamelor	202
7.2. Terminologie și descompunere Schur	203
7.3. Iterația vectorială	206
7.4. Metoda QR – teoria	209
7.5. Metoda QR – practica	213
7.5.1. Metoda QR clasică	213
7.5.2. Deplasare spectrală	218
7.5.3. Metoda QR cu pas dublu	220

În acest capitol ne ocupăm de determinarea valorilor (și vectorilor) proprii ale unei matrice pătratice $A \in \mathbb{R}^{n \times n}$, adică a valorilor $\lambda \in \mathbb{C}$ și vectorilor $x \in \mathbb{C}^n$ pentru care

$$Ax = \lambda x. \quad (7.0.1)$$

Definiția 7.0.1 Numărul $\lambda \in \mathbb{C}$ se numește valoare proprie a matricei $A \in \mathbb{R}^{n \times n}$, când există un vector $x \in \mathbb{C}^n \setminus \{0\}$ numit vector propriu astfel încât $Ax = \lambda x$.

Observația 7.0.2. 1. Cerința $x \neq 0$ este importantă, căci vectorul nul este un vector propriu corespunzător oricărei valori proprii.

2. Chiar dacă A este reală, ea poate avea valori proprii complexe. În acest caz ele apar în perechi conjugate. \diamond

7.1. Valori proprii și rădăcini ale polinoamelor

Orice problemă de calcul al valorilor proprii se poate reduce la calculul zerourilor unui polinom: valorile proprii ale unei matrice $A \in \mathbb{R}^{n \times n}$ sunt rădăcinile *polinomului caracteristic*

$$p_A \lambda = \det(A - \lambda I), \quad \lambda \in \mathbb{C},$$

căci determinatul este nul exact atunci când sistemul $(A - \lambda I)x = 0$ are o soluție nebanală, adică atunci când λ este o valoare proprie.

O metodă de rezolvare a problemelor proprii ar putea fi calculul polinomului caracteristic și apoi determinarea rădăcinilor. Natural, calculul unui determinant *în general* fiind o problemă complexă și instabilă, transformarea matricei ar fi mai potrivită. Reciproc, problema găsirii rădăcinilor unui polinom poate fi formulată ca o problemă de determinare a valorilor proprii. Fie $p \in \mathbb{P}_n$ un polinom cu coeficienți reali, pe care îl putem scrie (cu ajutorul rădăcinilor sale, z_1, \dots, z_n , eventual complexe) sub forma

$$p(x) = a_n x^n + \dots + a_0 = a_n (x - z_1) \dots (x - z_n), \quad a_n \in \mathbb{R}, \quad a_n \neq 0.$$

Pe spațiul vectorial \mathbb{P}_{n-1} „înmulțirea modulo p ”

$$\mathbb{P}_{n-1} \ni q \rightarrow r \quad xq(x) = \alpha p(x) + r(x), \quad r \in \mathbb{P}_n \quad (7.1.1)$$

este o *transformare liniară* și deoarece

$$x^n = \frac{1}{a_n} p(x) - \sum_{j=0}^{n-1} \frac{a_j}{a_n} x^j, \quad x \in \mathbb{R},$$

vom reprezenta p relativ la baza $1, x, \dots, x^{n-1}$ prin așa-numita matrice *companion* a lui Frobenius (de dimensiune $n \times n$)

$$M = \begin{bmatrix} 0 & & & -\frac{a_0}{a_n} \\ 1 & 0 & & -\frac{a_1}{a_n} \\ & \ddots & \ddots & \vdots \\ & & 1 & 0 \\ & & & 1 & -\frac{a_{n-2}}{a_n} \\ & & & & 1 & -\frac{a_{n-1}}{a_n} \end{bmatrix}, \quad (7.1.2)$$

Fie $v_j = (v_{jk} : k = \overline{1, n}) \in \mathbb{C}^n$, $j = \overline{1, n}$ alese astfel ca

$$\ell_j(x) = \frac{p(x)}{x - z_j} = a_n \prod_{k \neq j} (x - z_k) = \sum_{k=1}^n v_{jk} x^{k-1}, \quad j = \overline{1, n},$$

atunci

$$\sum_{k=1}^n (Mv_j - z_j v_j)_k x^{k-1} = x \ell_j(x) - z_j \ell_j(x) = (x - z_j) \ell_j(x) = p(x) \approx 0,$$

și cu aceasta $Mv_j = z_j v_j$, $j = \overline{1, n}$.

Valorile proprii ale lui M sunt deci rădăcini ale lui p .

Matricea Frobenius dată de (7.1.2) este doar o modalitate dintre multe altele prin care se poate reprezenta „înmulțirea” din (7.1.1); orice altă bază a lui \mathbb{P}_{n-1} furnizează o matrice M ale cărei valori proprii să fie rădăcini ale lui p . Singurul mijloc auxiliar de manipulare a polinoamelor de care avem nevoie este realizarea unei „împărțiri cu rest”.

7.2. Terminologie și descompunere Schur

Asa cum ne arată exemplul

$$A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad p_A(\lambda) = \lambda^2 + 1 = (\lambda + i)(\lambda - i),$$

o matrice reală poate avea valori proprii complexe. De aceea (cel puțin în teorie) este avantajos să ne ocupăm de matrice complexe $A \in \mathbb{C}^{n \times n}$.

Definiția 7.2.1 Două matrice $A, B \in \mathbb{C}^{n \times n}$ se numesc similare dacă există o matrice $T \in \mathbb{C}^{n \times n}$, astfel încât

$$A = TBT^{-1}.$$

Lema 7.2.2 Dacă $A, B \in \mathbb{C}^{n \times n}$ sunt similare, ele au aceleași valori proprii.

Demonstrație. Fie $\lambda \in \mathbb{C}$ o valoare proprie a lui $A = TBT^{-1}$ și $x \in \mathbb{C}^n$ vectorul propriu corespunzător. Atunci avem

$$B(T^{-1}x) = T^{-1}ATT^{-1}x = T^{-1}Ax = \lambda T^{-1}x$$

și deci, λ este, de asemenea, valoare proprie a lui B . \square

Din algebra liniară se cunoaște următorul rezultat important.

Teorema 7.2.3 (Forma normală Jordan) Orice matrice $A \in \mathbb{C}^{n \times n}$ este similară cu o matrice

$$J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_k \end{bmatrix}, \quad J_\ell = \begin{bmatrix} \lambda_\ell & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_\ell \end{bmatrix} \in \mathbb{C}^{n_\ell \times n_\ell}, \quad \sum_{\ell=1}^k n_\ell = n,$$

numită forma normală Jordan a lui A .

Definiția 7.2.4 O matrice se numește diagonalizabilă, dacă toate blocurile sale Jordan J_ℓ au dimensiunea 1, adică $n_\ell = 1$, $\ell = \overline{1, n}$. O matrice se numește nederogatorie dacă pentru fiecare valoare proprie λ_ℓ există exact un bloc Jordan în care ea apare pe diagonală.

Observația 7.2.5. Dacă o matrice $A \in \mathbb{R}^{n \times n}$ are n valori proprii simple, atunci ea este diagonalizabilă și de asemenea nederogatorie și potrivită pentru a fi tratată numeric. \diamond

Teorema 7.2.6 (Descompunere Schur) Pentru orice matrice $A \in \mathbb{C}^{n \times n}$ există o matrice unitară $U \in \mathbb{C}^{n \times n}$ și o matrice triunghiulară superior

$$R = \begin{bmatrix} \lambda_1 & * & \dots & * \\ & \ddots & \ddots & \vdots \\ & & \ddots & * \\ & & & \lambda_n \end{bmatrix} \in \mathbb{C}^{n \times n},$$

astfel încât $A = URU^*$.

Observația 7.2.7. 1. Elementele diagonale ale lui R sunt, natural, valorile proprii ale lui A . Deoarece A și R sunt similare, ele au aceleași valori proprii.

2. Între A și R are loc o formă mai puternică de similaritate: ele sunt *unitar-similare*. \diamond

Demonstrația teoremei 7.2.6. Demonstrația se face prin inducție. Cazul $n = 1$ este trivial. Presupunem teorema adevărată pentru $n \in \mathbb{N}$ și fie $A \in \mathbb{C}^{(n+1) \times (n+1)}$. Fie $\lambda \in \mathbb{C}$ o valoare proprie a lui A și $x \in \mathbb{C}^{n+1}$, $\|x\|_2 = 1$, vectorul propriu corespunzător. Luăm vectorul $u_1 = x$ și alegem u_2, \dots, u_{n+1} astfel încât u_1, \dots, u_{n+1} să formeze o bază ortonormală a lui \mathbb{C}^{n+1} , sau echivalent, matricea $U = [u_1, \dots, u_{n+1}]$ să fie unitară. Așadar,

$$U^*AUe_1 = U^*Au_1 = U^*Ax = \lambda U^*x = \lambda e_1,$$

adică

$$U^*AU = \begin{bmatrix} \lambda & * \\ 0 & B \end{bmatrix}, \quad B \in \mathbb{C}^{n \times n}.$$

Conform ipotezei inducției există o matrice unitară $V \in \mathbb{C}^{n \times n}$, astfel încât $B = VSV^*$, unde $S \in \mathbb{C}^{n \times n}$ este o matrice triunghiulară superior. De aceea

$$A = U \begin{bmatrix} \lambda_1 & * \\ 0 & VSV^* \end{bmatrix} U^* = \underbrace{U \begin{bmatrix} 1 & 0 \\ 0 & V \end{bmatrix}}_{=:U} \underbrace{\begin{bmatrix} \lambda_1 & * \\ 0 & S \end{bmatrix}}_{=:R} \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & V^* \end{bmatrix}}_{=:U^*} U^*$$

și demonstrația este completă. \square

Să dăm acum două consecințe nemijlocite ale descompunerii Schur.

Corolarul 7.2.8 *Oricărei matrice hermitiene $A \in \mathbb{C}^{n \times n}$ îi corespunde o matrice ortogonală $U \in \mathbb{C}^{n \times n}$ astfel încât*

$$A = U \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} U^*, \quad \lambda_j \in \mathbb{R}, \quad j = \overline{1, n}.$$

Demonstrație. Matricea R din teorema 7.2.6 verifică $R = U^*AU$. Deoarece

$$R^* = (U^*AU) = U^*A^*U = U^*AU = R,$$

R trebuie să fie diagonală, cu elemente reale pe diagonală (fiind hermitiană). \square

Altfel spus, corolarul 7.2.8 ne asigură că orice matrice hermitiană este unitar diagonalizabilă și posedă o bază formată din vectori proprii ortonormali. În plus, toate valorile proprii ale unei matrice hermitiene sunt reale. Este interesant, nu numai din punct de vedere teoretic, ce matrice sunt unitar diagonalizabile.

Teorema 7.2.9 *O matrice $A \in \mathbb{C}^{n \times n}$ este unitar diagonalizabilă, adică există o matrice unitară $U \in \mathbb{C}^{n \times n}$ astfel încât*

$$A = U \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} U^*, \quad \lambda_j \in \mathbb{R}, \quad j = 1, \dots, n. \quad (7.2.1)$$

dacă și numai dacă A este normală, adică

$$AA^* = A^*A. \quad (7.2.2)$$

Demonstrație. Punem $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Conform lui (7.2.1), A are forma $A = U\Lambda U^*$, deci

$$AA^* = U\Lambda U^*U\Lambda^*U^* = U|\Lambda|^2U^* \text{ și } A^*A = U^*\Lambda^*U^*U\Lambda U^* = U|\Lambda|^2U^*,$$

adică (7.2.2). Pentru reciprocă, folosim descompunerea Schur a lui A sub forma $R = U^*AU$. Atunci

$$|\lambda_1|^2 = (R^*R)_{11} = (RR^*)_{11} = |\lambda_1|^2 + \sum_{k=2}^n |r_{1k}|^2,$$

de unde rezultă $r_{12} = \dots = r_{1n} = 0$. Prin inducție, se vede că pentru $j = \overline{2, n}$

$$(R^*R)_{jj} = |\lambda_j|^2 + \sum_{k=1}^{j-1} |r_{kj}|^2 = (RR^*)_{jj} = |\lambda_j|^2 + \sum_{k=j+1}^n |r_{jk}|^2,$$

din care cauză R trebuie să fie diagonală. \square

Pentru matricea reală *descompunerea Schur reală* este un pic mai complicată.

Teorema 7.2.10 Pentru orice matrice $A \in \mathbb{R}^{n \times n}$ există o matrice ortogonală $U \in \mathbb{R}^{n \times n}$ astfel încât

$$A = U \begin{bmatrix} R_1 & * & \dots & * \\ & \ddots & \ddots & \vdots \\ & & \ddots & * \\ & & & R_k \end{bmatrix} U^*, \quad (7.2.3)$$

în care fie $R_j \in \mathbb{R}^{1 \times 1}$, fie $R_j \in \mathbb{R}^{2 \times 2}$, cu două valori proprii complexe conjugate, $j = \overline{1, k}$.

Descompunerea Schur reală transformă A într-o matrice Hessenberg superioară

$$U^T A U = \begin{bmatrix} * & \dots & \dots & * \\ * & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ & & * & * \end{bmatrix}.$$

Demonstrație. Dacă A are doar valori proprii reale, atunci se procedează ca la descompunerea Schur complexă. Altfel, fie $\lambda = \alpha + i\beta$, $\beta \neq 0$, o valoare proprie complexă a lui A și $x + iy$ vectorul propriu corespunzător. Atunci

$$A(x + iy) = \lambda(x + iy) = (\alpha + i\beta)(x + iy) = (\alpha x - \beta y) + i(\beta x + \alpha y)$$

sau matricial

$$A \underbrace{\begin{bmatrix} x & y \end{bmatrix}}_{\in \mathbb{R}^{n \times 2}} = \begin{bmatrix} x & y \end{bmatrix} \underbrace{\begin{bmatrix} \alpha & \beta \\ -\beta & \alpha \end{bmatrix}}_{:=R}.$$

Deoarece $R = \alpha^2 + \beta^2 > 0$, căci $\beta \neq 0$, $\text{span}\{x, y\}$ este un subspațiu bidimensional A -invariant al lui \mathbb{R}^n . Atunci alegem u_1, u_2 astfel încât să formeze o bază a acestui spațiu, completată cu u_3, \dots, u_n la o bază ortonormală a lui \mathbb{R}^n și raționând analog cu cazul complex obținem

$$U^T A U = \begin{bmatrix} R & * \\ 0 & B \end{bmatrix}$$

și inducția se derulează ca la descompunerea Schur complexă. \square

7.3. Iterația vectorială

Iterația vectorială (numită și metoda puterii) este cea mai simplă metodă atunci când dorim o valoare proprie și vectorul propriu corespunzător.

Pornind de la un vector $y^{(0)} \in \mathbb{C}^n$ se construiește șirul y^k , $k \in \mathbb{N}$ prin intermediul iterației

$$\begin{aligned} z^{(k)} &= Ay^{(k-1)}, \\ y^{(k)} &= \frac{z^{(k)}}{z_{j^*}^{(k)} \|z^{(k)}\|_\infty}, \quad j = \min \left\{ 1 \leq j \leq n : |z_j^{(k)}| \geq \left(1 - \frac{1}{k}\right) \|z^{(k)}\|_\infty \right\} \end{aligned} \quad (7.3.1)$$

și se afirmă că, în condiții determinate, acest șir converge către vectorul propriu dominant.

Propoziția 7.3.1 Fie $A \in \mathbb{C}^{n \times n}$ o matrice diagonalizabilă ale cărei valori proprii $\lambda_1, \dots, \lambda_n$ verifică condiția

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|.$$

Atunci șirul $y^{(k)}$, $k \in \mathbb{N}$, converge către un multiplu al vectorului normat x_1 corespunzător valorii proprii λ_1 , pentru aproape orice $y^{(0)}$.

Demonstrație. Fie x_1, \dots, x_n vectorii proprii ortonormali ai lui A corespunzători valorilor proprii $\lambda_1, \dots, \lambda_n$ – existența lor rezultă din diagonalizabilitatea lui A . Scriem

$$y^{(0)} = \sum_{j=1}^n \alpha_j x_j, \quad \alpha_j \in \mathbb{C}, \quad j = \overline{1, n},$$

și afirmăm că

$$A^k y^{(0)} = \sum_{j=1}^n \alpha_j A^k x_j = \sum_{j=1}^n \alpha_j \lambda_j^k x_j = \lambda_1^k \sum_{j=1}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1}\right)^k x_j.$$

De aici rezultă, deoarece $|\lambda_1| > |\lambda_j|$, $j = \overline{2, n}$ că

$$\lim_{k \rightarrow \infty} \lambda_1^{-k} A^k y^{(0)} = \alpha_1 x_1 + \lim_{k \rightarrow \infty} \sum_{j=2}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1}\right)^k x_j = \alpha_1 x_1,$$

precum și

$$\lim_{k \rightarrow \infty} |\lambda_1|^{-k} \|A^k y^{(0)}\|_\infty = \left\| \sum_{j=1}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1}\right)^k \alpha_j x_j \right\| = |\alpha_1| \|x_1\|_\infty.$$

Dacă $\alpha_1 = 0$ și deci $y^{(0)}$ aparține hiperplanului

$$x_1^\perp = \{x \in \mathbb{C}^n, x^* x_1 = 0\},$$

atunci ambele limite sunt nule și nu se poate spune nimic despre convergența șirului y^k , $k \in \mathbb{N}$; acest hiperplan este o mulțime de măsură nulă, așa că în continuare vom presupune că $\alpha_1 \neq 0$.

Datorită lui (7.3.1), $y^k = \gamma_k A^k y^{(0)}$, $\gamma_k \in \mathbb{C}$ și pe lângă aceasta $\|y^k\|_\infty = 1$, așadar

$$\lim_{k \rightarrow \infty} |\lambda_1|^k |\gamma_k| = \lim_{k \rightarrow \infty} \frac{1}{|\lambda_1|^{-1} \|A^{(k)} y^{(0)}\|} = \frac{1}{|\alpha_1| \|x_1\|_\infty}.$$

Astfel

$$y^{(k)} = \gamma_k A^k y^{(0)} = \underbrace{\frac{\gamma_k \lambda_1^k}{|\gamma_k \lambda_1^k|}}_{=: e^{-2\pi i \theta_k}} \underbrace{\frac{\alpha_1 x_1}{|\alpha_1| \|x_1\|_\infty}}_{=: \alpha x_1} + O\left(\frac{|\lambda_2|^k}{|\lambda_1|^k}\right), \quad k \in \mathbb{N}, \quad (7.3.2)$$

unde $\theta_k \in [0, 1]$. Aici intervine normarea mai ciudată din (7.3.1): fie j indicele minim pentru care $|(\alpha x_1)_j| = \|\alpha x_1\|_\infty$; atunci conform lui (7.3.2) pentru un k suficient de mare de asemenea în (7.3.1) $j^* = j$. Deci are loc

$$\lim_{k \rightarrow \infty} y_j^{(k)} = 1 \Rightarrow \lim_{k \rightarrow \infty} e^{2\pi i \theta_k} = \lim_{k \rightarrow \infty} \frac{y_j^{(k)}}{(\alpha x_1)_j} = \frac{1}{(\alpha x_1)_j}.$$

Înlocuind aceasta în (7.3.2) obținem convergența șirului $y^{(k)}$, $k \in \mathbb{N}$. \square

Se poate aplica de asemenea iterația vectorială pentru a găsi toate valorile proprii și toți vectorii proprii, în măsura în care valorile proprii ale lui A sunt diferite în modul. Pentru aceasta se determină cea mai mare în modul valoare proprie λ_1 a lui A și vectorul propriu corespunzător x_1 și se continuă cu

$$A^{(1)} = A - \lambda_1 x_1 x_1^T.$$

Matricea diagonalizabilă $A^{(1)}$ are aceiași vectori proprii ortonormali ca și A , doar că x_1 este vectorul propriu corespunzător valorii proprii 0 și nu mai joacă nici un rol în iterație, atât timp cât nu se pornește chiar cu un multiplu al lui x_1 . Aplicând încă odată iterația vectorială lui $A^{(1)}$ se obține a doua valoare proprie ca mărime în modul λ_2 și vectorul propriu corespunzător; iterația

$$A^{(j)} = A^{(j-1)} - \lambda_j x_j x_j^T, \quad j = \overline{1, n}, \quad A^{(0)} = A$$

calculează succesiv toate valorile proprii și toți vectorii proprii ai lui A , presupunând că valorile proprii sunt diferite în modul.

Observația 7.3.2 (Dezavantajele iterației vectoriale). 1. Metoda funcționează în general doar când există un vector propriu dominant, adică când există pentru

o valoare proprie dominantă exact un vector propriu. Dacă se consideră, de exemplu, matricea

$$A = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix},$$

atunci ea transformă vectorul $[x_1 \ x_2]^T$ în vectorul $[x_2 \ x_1]^T$ și convergența are loc exact atunci când iterația pornește cu unul din vectorii proprii.

2. Metoda dă rezultate numai pentru vectori de pornire „potrivii”. Sună nemaipomenit că toți vectorii de pornire cu excepția celor dintr-un hiperplan sunt potriviți, dar nu este chiar așa de simplu. Dacă valoarea proprie cea mai mare în modul a unei matrice reale este complexă atunci se poate itera la nesfârșit cu valori de pornire reale, căci oricum nu se va găsi vectorul propriu.
3. Ar trebui făcute tot timpul calcule în complex, ceea ce ar ridica serios complexitatea (la adunare un număr dublu de calcule, la înmulțire de șase ori mai multe).
4. Viteza de convergență depinde de raportul

$$\frac{|\lambda_2|}{|\lambda_1|} < 1$$

care poate fi oricât de aproape de 1. Dacă dominanța vectorului propriu dominant nu este bine reliefată, atunci convergența poate fi extrem de lentă. \diamond

Ținând cont de toate acestea, conchidem că iterația vectorială nu este o metodă prea bună pentru problemele de valori proprii.

7.4. Metoda QR – teoria

Metoda practică pentru tratarea problemelor de valori proprii este în zilele noastre metoda QR descoperită de Francis [18] și Kublanovskaya [27], o extensie unitară a metodei LR a lui Rutishauser [32]. Vom începe cu cazul complex.

Metoda este una iterativă extrem de simplă: se pornește cu $A^{(0)} = A$ și se calculează iterativ cu ajutorul descompunerii QR,

$$A^{(k)} = Q_k R_k, \quad A^{(k+1)} = R_k Q_k, \quad k \in \mathbb{N}_0. \quad (7.4.1)$$

Cu puțin noroc, sau așa cum spun matematicienii, în anumite ipoteze, acest șir va converge către o matrice ale cărei elemente diagonale sunt valorile proprii ale lui A .

Lema 7.4.1 *Matricele $A^{(k)}$ construite prin (7.4.1), $k \in \mathbb{N}$, sunt similar ortogonale cu A și au aceleași valori proprii ca A .*

Demonstrație. Are loc

$$A^{(k+1)} = Q_k^* Q_k R_k Q_k = Q_k^* A^{(k)} Q_k = \cdots = \underbrace{Q_k^* \cdots Q_0^*}_{=: U_k^*} A \underbrace{Q_0 \cdots Q_k}_{=: U_k}.$$

□

Pentru a arăta convergența, vom interpreta iterația QR ca o generalizare a iterației vectoriale (7.3.1) (fără normarea ciudată) la spații vectoriale. Pentru aceasta vom scrie baza ortonormală $u_1, \dots, u_m \in \mathbb{C}^n$ a unui subspațiu m -dimensional $U \subset \mathbb{C}^n$, $m \leq n$, ca vectori coloane ai unei matrice unitare $U \in \mathbb{R}^{n \times m}$ și iterăm subspațiul vectorial (respectiv matricele) peste descompunerea QR

$$U_{k+1} R_k = A U_k, \quad k \in \mathbb{N}_0, \quad U_0 \in \mathbb{C}^n. \quad (7.4.2)$$

De aici rezultă imediat că

$$U_{k+1}(R_k \cdots R_0) = A U_k(R_{k-1} \cdots R_0) = A^2 U_{k-1}(R_{k-2} \cdots R_0) = \cdots = A^{k+1} U_0. \quad (7.4.3)$$

Dacă definim acum, pentru $m = n$, $A^{(k)} = U_k^* A U_k$, atunci conform lui (7.4.2) au loc relațiile

$$\begin{aligned} A^{(k)} &= U_k^* A U_k = U_k^* U_{k+1} R_k \\ A^{(k+1)} &= U_{k+1}^* A U_{k+1} = U_{k+1}^* A U_k U_k^* U_{k+1} \end{aligned}$$

și punând $Q_k := U_k^* U_{k+1}$, obținem regula de iterare (7.4.1). Ca matrice de pornire putem alege $U_0 = I$.

Pentru a obține rezultate despre metoda QR mai avem nevoie de un tip de matrice.

Definiția 7.4.2 O matrice de fază $\Theta \in \mathbb{C}^{n \times n}$ este o matrice diagonală de forma

$$\Theta = \begin{bmatrix} e^{-i\theta_1} & & \\ & \ddots & \\ & & e^{-i\theta_n} \end{bmatrix}, \quad \theta_j \in [0, 2\pi), \quad j = \overline{1, n}.$$

Propoziția 7.4.3 Presupunem că matricea $A \in \mathbb{C}^{n \times n}$ are valori proprii distincte în modul, $|\lambda_1| > |\lambda_2| > \cdots > |\lambda_n| > 0$. Dacă matricea X^{-1} din forma normală Jordan $A = X \Lambda X^{-1}$ a lui A are o descompunere LU

$$X^{-1} = ST, \quad S = \begin{bmatrix} 1 & & & \\ * & 1 & & \\ \vdots & \ddots & \ddots & \\ * & \cdots & * & 1 \end{bmatrix}, \quad T = \begin{bmatrix} * & \cdots & * \\ & \ddots & \vdots \\ & & * \end{bmatrix},$$

atunci există matricele de fază Θ_k , $k \in \mathbb{N}_0$, astfel încât șirul de matrice $(\Theta_k U_k, k \in \mathbb{N})$ să convergă.

Observația 7.4.4 (La propoziția 7.4.3). 1. Convergența șirului $(\Theta_k U_k)$, înseamnă mai ales, că dacă bazele ortonormale aparținătoare converg către o bază ortonormală a lui \mathbb{C}^n , avem de asemenea convergența spațiilor vectoriale.

2. Existența descompunerii LU a lui X^{-1} este fără nici o limitare: deoarece X^{-1} este inversabilă în mod trivial, există întotdeauna o permutare P , astfel încât $X^{-1}P^T = (PX)^{-1} = LU$ și PX este într-adevăr o matrice inversabilă. Aceasta înseamnă că matricea $\widehat{A} = P^TAP$, în care A se transformă prin permutarea liniilor și coloanelor are aceleși valori proprii ca și A , verifică ipoteza propoziției 7.4.3.

3. Demonstrația propoziției 7.4.3 este o modificare a demonstrației din [40, pag. 54–56] pentru convergența metodei LR, care își are originea în lucrarea lui Wilkinson¹ [46]. Ce este de fapt metoda LR? Se procedează ca la metoda QR, dar lui A i se aplică eliminarea gaussiană, $A^{(k)} = L_k R_k$ și apoi se construiește $A^{(k+1)} = R_k L_k$. De asemenea, această metodă converge în anumite condiții către o matrice triunghiulară superior. \diamond

Înainte de a demonstra propoziția 7.4.3 să vedem întâi de ce convergența șirului (U_k) atrage convergența metodei QR. Și anume, dacă avem $\|U_{k+1} - U_k\| \leq \varepsilon$ sau echivalent

$$U_{k+1} = U_k + E, \quad \|E\|_2 \leq \varepsilon,$$

atunci

$$Q_k = U_{k+1}^* U_k = (U_k + E)^* U_k = I + E^* U_k = I + F, \quad \|F\|_2 \leq \|E\|_2 \underbrace{\|U_k\|_2}_{=1} \leq \varepsilon,$$

și cu aceasta

$$A^{(k+1)} = R_k Q_k = R_k (I + F) = R_k + G, \quad \|G\|_2 \leq \varepsilon \|R_k\|_2,$$

deci șirul $A^{(k)}$, $k \in \mathbb{N}$, converge, de asemenea, către o matrice triunghiulară superior, numai dacă normele lui R_k , $k \in \mathbb{N}$ sunt uniform mărginite. Chiar așa se întâmplă, căci

$$\|R_k\|_2 = \|Q_k^* A^{(k)}\|_2 = \|A^{(k)}\|_2 = \|Q_{k-1}^* \dots Q_0^* A Q_0 \dots Q_{k-1}\| = \|A\|_2.$$

Mai avem nevoie de un rezultat ajutător despre „unicitatea” descompunerii QR.

James Hardy Wilkinson (1919-1986), matematician englez. Contribuții importante în domeniul Analizei numerice, Algebrei liniare numerice și Informaticii. Membru al Royal Society, laureat al premiului Turing al ACM. Pe lângă numeroasele sale lucrări în domeniul Analizei numerice, a lucrat și la dezvoltarea de biblioteci de rutine numerice. Grupul NAG (Numerical Algorithms Group) și-a început activitatea în 1970 și multe dintre rutinele de algebră liniară numerică s-au datorat lui Wilkinson.



Lema 7.4.5 Fie $U, V \in \mathbb{C}^{n \times n}$ matrice unitare și fie $R, S \in \mathbb{C}^{n \times n}$ matrice triunghiulare superior inversabile. Atunci are loc $UR = VS$ dacă și numai dacă există o matrice de fază

$$\Theta = \begin{bmatrix} e^{-i\theta_1} & & \\ & \ddots & \\ & & e^{-i\theta_n} \end{bmatrix}, \quad \theta_j \in [0, 2\pi), j = \overline{1, n},$$

astfel încât $V = V\Theta^*$, $R = \Theta S$.

Demonstrație. Deoarece $UR = V\Theta^* \Theta S = VS$, ” \Leftarrow ” este trivială. Pentru necesitate, din $UR = VS$ rezultă că $V^*U = SR^{-1}$ trebuie să fie o matrice triunghiulară superior astfel încât $(V^*U)^* = U^*V = RS^{-1}$. Așadar $\Theta = V^*U$ este o matrice diagonală unitară și are loc $U = VV^*U = V\Theta$. \square

Demonstrația propoziției 7.4.3. Fie $A = X\Lambda X^{-1}$ forma normală Jordan a lui A , unde $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Pentru $U_0 = I$ și $k \in \mathbb{N}_0$

$$U_k \left(\prod_{j=k-1}^0 R_j \right) = (X^{-1}\Lambda X)^k = X\Lambda^k X^{-1} = X\Lambda^k S T = X \underbrace{(\Lambda^k S \Lambda^{-k})}_{=: L_k} \Lambda^k T,$$

în care L_k este o matrice triunghiulară inferior cu elementele

$$(L_k)_{jm} = \left(\frac{\lambda_j}{\lambda_m} \right)^k, \quad 1 \leq m \leq j \leq n \quad (7.4.4)$$

astfel încât pentru $k \in \mathbb{N}$

$$|L_k - I| \leq \left(\max_{1 \leq m < j \leq n} |s_{jm}| \right) \left(\max_{1 \leq m < j \leq n} \left| \frac{\lambda_j}{\lambda_m} \right| \right)^k \begin{bmatrix} 0 & & & \\ 1 & \ddots & & \\ \vdots & \ddots & \ddots & \\ 1 & \dots & 1 & 0 \end{bmatrix}, \quad (k \in \mathbb{N}). \quad (7.4.5)$$

Fie $\widehat{U}_k \widehat{R}_k = X L_k$ descompunerea QR a lui $X L_k$, care datorită lui (7.4.5) și lemei 7.4.5 converge până la o matrice de fază către descompunerea QR $X = UR$ a lui X . Dacă aplicăm acum lema 7.4.5 identității

$$U_k \left(\prod_{j=k-1}^0 R_j \right) \widehat{Q}_k \widehat{R}_k \Lambda^k T,$$

atunci există matricele de fază Θ_k , astfel încât

$$U_k = \widehat{Q}_k \Theta_k^* \text{ și } \left(\prod_{j=k-1}^0 R_j \right) = \Theta_k \widehat{R}_k \Lambda^k T,$$

deci există matricele de fază $\widehat{\Theta}_k$, astfel încât $U_k \widehat{\Theta}_k \rightarrow U$, când $k \rightarrow \infty$. \square

Merită să aruncăm o scurtă privire asupra „termenului de eroare” din (7.4.4), ale cărui elemente subdiagonale verifică relația

$$|L_k|_{jm} \leq \left(\frac{|\lambda_j|}{|\lambda_m|} \right) |s_{jm}|, \quad 1 \leq m < j \leq n.$$

Așadar are loc

Convergența unui element subdiagonal către 0 este cu atât mai rapidă cu cât elementul este mai îndepărtat de diagonală.

7.5. Metoda QR – practica

7.5.1. Metoda QR clasică

Am văzut că metoda QR generează un șir de matrice $A^{(k)}$ care în condiții determinate trebuie să convergă către o matrice triunghiulară superior care are pe diagonală valorile proprii. Putem întrebuința această metodă pentru matrice reale.

Exemplul 7.5.1. Fie

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 2 & 1 \end{bmatrix}.$$

Această matrice are valorile proprii

$$\lambda_1 \approx 4.56155, \quad \lambda_2 = -1, \quad \lambda_3 \approx 0.43845.$$

Iterând prin metoda QR se obțin pentru elementele subdiagonale rezultatele din tabela 7.1 (o implementare MATLAB brută). Se poate vedea că, după k iterații, elementele

#iterații	a_{21}	a_{31}	a_{32}
10	6.64251e-007	-2.26011e-009	0.00339953
20	1.70342e-013	-1.52207e-019	8.9354e-007
30	4.36711e-020	-1.02443e-029	2.34578e-010
40	1.11961e-026	-6.89489e-040	6.15829e-014

Tabela 7.1: Rezultate pentru exemplul 7.5.1

$a_{m\ell}^{(k)}$, $\ell < m$, se apropie de 0 la fel ca $|\lambda_\ell/\lambda_k|$.

◇

Exemplul 7.5.2. Matricea

$$\begin{bmatrix} 1 & 5 & 7 \\ 3 & 0 & 6 \\ 4 & 3 & 1 \end{bmatrix}$$

are valorile proprii

$$\lambda_1 \approx 9.7407, \quad \lambda_2 \approx -3.8703 + 0.6480i, \quad \lambda_3 \approx -3.8703 - 0.6480i.$$

În acest caz nu ne putem aștepta ca metoda QR să convergă către o matrice triunghiulară superioară, căci atunci toate valorile proprii ale lui A ar fi reale. De fapt, după 100 de iterații se obține matricea

$$A^{(100)} \approx \begin{bmatrix} 9.7407 & -4.3355 & 0.94726 \\ 8.552e - 039 & -4.2645 & 0.7236 \\ 3.3746e - 039 & -0.79491 & -3.4762 \end{bmatrix},$$

care ne furnizează corect valoarea proprie reală. Pe lângă aceasta, matricea 2×2 din colțul din dreapta jos ne furnizează valorile proprii complexe $-3.8703 \pm 0.6480i$. \diamond

Al doilea exemplu recomandă următoarea strategie: dacă elementele situate sub diagonală nu vor să dispară, ar fi indicat să privim mai atent matricea 2×2 corespunzătoare.

Definiția 7.5.3 Dacă $A \in \mathbb{R}^{n \times n}$ are descompunerea QR $A = QR$, atunci transformarea RQ a lui A se definește prin $A_* = RQ$.

Ce probleme apar la realizarea practică a metodei QR? Deoarece complexitatea descompunerii QR este $O(n^3)$, nu este prea înțelept să utilizăm o metodă care se bazează pe un astfel de pas iterativ. Pentru a evita problema, vom converti matricele inițiale în matrice similare a căror descompunere QR poate fi calculată mai repede. Astfel de matrice sunt matricele Hessenberg superioare a căror descompunere QR se poate calcula cu n iterații Givens, deci un total de $O(n^2)$ operații: de vreme ce numai elementele $h_{j-1,j}$, $j = \overline{2, n}$ trebuie eliminate, vom determina unghiurile ϕ_2, \dots, ϕ_n , astfel ca

$$G(n-1, n; \phi_n) \dots G(1, 2; \phi_2)H = R$$

și are loc

$$H_* = RG^T(1, 2; \phi_2) \dots G^T(n-1, n; \phi_n). \quad (7.5.1)$$

Această este ideea algoritmului 7.1. După [35], motto-ul trebuie să fie „odată Hessenberg, întotdeauna Hessenberg”.

Lema 7.5.4 Dacă $H \in \mathbb{R}^{n \times n}$ este o matrice Hessenberg superioară, atunci H_* este de asemenea matrice Hessenberg superioară.

Demonstrație. Rezultă direct din reprezentarea (7.5.1). Înmulțirea la dreapta cu o matrice Givens $G^T(j, j+1, \phi_{j+1})$, $j = \overline{1, n-1}$ înseamnă o combinație a coloanelor j și $j+1$ și crează valori diferite de zero doar în prima subdiagonală – R este triunghiulară superior. \square

Algoritmul 7.1 Transformarea RQ a unei matrice Hessenberg H, adică $H_* = RQ$ unde $H = QR$ este descompunerea QR a lui H

for $k := 1$ **to** $n - 1$ **do**

$[c_k, s_k] := \mathbf{givens}(H_{kk}, H_{k+1,k});$

$H_{k:k+1,k:n} := \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix}^T H_{k:k+1,k:n};$

end for

for $k := 1$ **to** $n - 1$ **do**

$H_{1:k+1,k:k+1} := H_{1:k+1,k:k+1} \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix};$

end for

Să vedem cum aducem matricea inițială la forma Hessenberg. În acest scop vom utiliza pentru variație transformările Householder. Să presupunem că am găsit deja o matrice Q_k , astfel încât primele k coloane ale matricei transformate să aibă deja forma Hessenberg, adică

$$Q_k A Q_k^T = \left[\begin{array}{ccc|ccc} * & \dots & * & * & * & \dots & * \\ * & \dots & * & * & * & \dots & * \\ & & \ddots & \vdots & \vdots & \ddots & \vdots \\ & & & * & * & \dots & * \\ \hline & & & a_1^{(k)} & * & \dots & * \\ & & & \vdots & \vdots & \ddots & \vdots \\ & & & a_{n-k-1}^{(k)} & * & \dots & * \end{array} \right].$$

Apoi determinăm $\hat{y} \in \mathbb{R}^{n-k-1}$ și $\alpha \in \mathbb{R}$ (care rezultă automat), astfel ca

$$H(\hat{y}) \begin{bmatrix} a_1^{(k)} \\ \vdots \\ \vdots \\ a_{n-k-1}^{(k)} \end{bmatrix} = \begin{bmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow U_{k+1} := \begin{bmatrix} I_{k+1} & \\ & H(\hat{y}) \end{bmatrix}$$

și obținem că

$$\underbrace{U_{k+1}Q_k}_{=:Q_{k+1}} A \underbrace{Q_k U_{k+1}}_{=:Q_{k+1}^T} = \left[\begin{array}{ccc|c|ccc} * & \dots & * & * & * & \dots & * \\ * & \dots & * & * & * & \dots & * \\ & & \ddots & \vdots & \vdots & \ddots & \vdots \\ & & & * & * & \dots & * \\ \hline & & & \alpha & * & \dots & * \\ & & & 0 & * & \dots & * \\ & & & \vdots & \ddots & \vdots & \\ & & & 0 & * & \dots & * \end{array} \right] U_{k+1};$$

matricea unitate I_{k+1} din stânga sus în matricea U_{k+1} are grijă să avem în primele $k + 1$ coloane o structură Hessenberg. Algoritmul 7.2 dă metoda de trecere a unei matrice la forma Hessenberg. În concluzie, metoda noastră QR va fi o metodă în două faze:

1. Transformă pe A în forma Hessenberg utilizând o transformare ortogonală:

$$H^{(0)} = Q A Q^T, \quad Q^T Q = Q Q^T = I.$$

2. Execută iterațiile QR

$$H^{(k+1)} = H_*^{(k)}, \quad k \in \mathbb{N}_0,$$

în speranța că elementele de pe prima subdiagonală converg toate către zero.

Algoritmul 7.2 Reducere la forma Hessenberg superioară

Intrare: Matricea $A \in \mathbb{R}^{n \times n}$

Ieșire: H forma Hessenberg a lui A și dacă se dorește Q astfel încât $H = Q A Q^T$

for $i := 1$ **to** $n - 2$ **do**

$u_i := \text{House}(A_{i+1:n,i});$

$P_i := I - 2u_i u_i^T; \{Q_i = \text{diag}(I_i, P_i)\}$

$A_{i+1:n,i:n} := P_i A_{i+1:n,i:n};$

$A_{1:n,i+1:n} := A_{1:n,i+1:n} P_i;$

end for

if se dorește Q **then**

$Q := I;$

for $i := 1$ **to** $n - 2$ **do**

$Q_{i+1:n,i+1:n} := P_i Q_{i+1:n,i+1:n};$

end for

end if

ε	#iterații	λ_1	λ_2	λ_3
10^{-3}	11	4.56155	-0.999834	0.438281
10^{-4}	14	4.56155	-1.00001	0.438461
10^{-5}	17	4.56155	-0.999999	0.438446
10^{-10}	31	4.56155	-1	0.438447

Tabela 7.2: Rezultatele pentru exemplul 7.5.5

Deoarece elementele subdiagonale converg cel mai încet, putem folosi drept criteriu de oprire maximul modului. Aceasta ne conduce la *metoda QR simplă*, vezi algoritmul 7.3. Desigur, la apariția unor valori proprii complexe această metodă iterează la nesfârșit.

Algoritmul 7.3 Metoda QR simplă

Intrare: Matricea A , toleranța tol

Ieșire: Vectorul valorilor proprii λ și numărul de iterații it

$H := \text{Hessenberg}(A)$; {Forma Hessenberg a lui A }

$it := 0$;

while $\|\text{diag}(H, -1)\|_\infty > tol$ **do**

$H := H_*$; {Transformarea RQ a lui H }

$it := it + 1$;

end while

$\lambda := \text{diag}(H)$;

Exemplul 7.5.5. Aplicăm noua metodă matricei din exemplul 7.5.1. Pentru diverse toleranțe ε date obținem rezultatele din tabela 7.2. Se observă că la fiecare trei iterații se câștigă o zecimală la elementele vecine cu diagonală. \diamond

Putem încerca să accelerăm metoda descompunând problema în subprobleme. Dacă avem o matrice Hessenberg de forma

$$H = \left[\begin{array}{cccc|cccc} * & \dots & \dots & * & & & & \\ * & \ddots & & \vdots & & & & * \\ & \ddots & \ddots & \vdots & & & & \\ & & * & * & & & & \\ \hline & & & & * & \dots & \dots & * \\ & & & & * & \ddots & & \vdots \\ & & & & & \ddots & \ddots & \vdots \\ & & & & & & * & * \end{array} \right] = \begin{bmatrix} H_1 & * \\ & H_2 \end{bmatrix},$$

atunci problema de valori proprii referitoare la H se poate descompune într-o problemă de valori proprii referitoare la H_1 și una referitoare la H_2 .

Conform lui [24], un element subdiagonal $h_{j+1,j}$ este considerat „suficient de mic” dacă

$$|h_{j+1,j}| \leq \text{eps} (|h_{jj}| + |h_{j+1,j+1}|). \quad (7.5.2)$$

Aici vom proceda mai simplu și vom descompune o matrice dacă cel mai mic element în modul situat pe prima subdiagonală devine mai mic decât o toleranță dată. Procedeu este după cum urmează: funcția de determinare a valorilor proprii determină cu ajutorul unei iterații QR o descompunere în matricele H_1 și H_2 și se apelează recursiv pe sine pentru fiecare din aceste matrice.

Dacă una dintre matrice este 1×1 , valoarea proprie este automat determinată, iar dacă este 2×2 , atunci polinomul său caracteristic este

$$\begin{aligned} p_A(x) &= \det(A - xI) = x^2 - \text{trace}(A)x + \det(A) \\ &= x^2 \underbrace{(a_{11} + a_{22})}_{=:b} x + \underbrace{(a_{11}a_{22} - a_{12}a_{21})}_{=:c}. \end{aligned}$$

Dacă discriminantul $b^2 - 4c$ este pozitiv, A are două valori proprii reale și distincte

$$x_1 = \frac{1}{2} \left(-b - \text{sgn}(b)\sqrt{b^2 - 4c} \right) \text{ și } x_2 = \frac{c}{x_1},$$

altfel valorile proprii sunt complexe și anume

$$\frac{1}{2} \left(-b \pm i\sqrt{4c - b^2} \right);$$

astfel avem sub control și cazul valorilor proprii complexe. Presupunem că funcția **Eigen2x2** returnează valorile proprii ale unei matrice 2×2 . Ideea este implementată în algoritmul 7.4. Iterațiile QR efective sunt date de algoritmul 7.5.

Exemplul 7.5.6. Să considerăm din nou matricea din exemplul 7.5.2 căreia îi aplicăm algoritmul 7.4. Rezultatele apar în tabela 7.3. \diamond

7.5.2. Deplasare spectrală

Utilizarea matricelor Hessenberg ne permite, într-adevăr, să executăm fiecare pas de iterație într-un timp mai scurt. Vom încerca acum să micșorăm numărul de iterații, așadar să creștem viteza de convergență deoarece

Rata de convergență a elementelor subdiagonale $h_{j+1,j}$ are ordinul de mărime

$$\left(\frac{\lambda_{j+1}}{\lambda_j} \right)^k, \quad j = \overline{1, n-1}.$$

Algoritmul 7.4 QRSplit1a – metoda QR cu partiționare și tratarea cazurilor 2×2 **Intrare:** Matricea $A \in \mathbb{R}^{n \times n}$ și toleranța tol **Ieșire:** Valorile proprii λ și numărul de iterații it

```

if  $n = 1$  then
     $it := 0; \lambda := A;$ 
    return;
else
    if  $n = 2$  then
         $it := 0;$ 
         $\lambda := \text{Eigen2x2}(A);$ 
        return;
    else
         $H := \text{Hessenberg}(A); \{\text{Forma Hessenberg}\}$ 
         $[H_1, H_2, it] := \text{QRIter}(H, tol);$ 
         $[\lambda_1, it_1] := \text{QRSplit1a}(H_1, tol); \{\text{apeluri recursive}\}$ 
         $[\lambda_2, it_2] := \text{QRSplit1a}(H_2, tol);$ 
         $it := it + it_1 + it_2;$ 
         $\lambda := [\lambda_1, \lambda_2];$ 
    end if
end if

```

Algoritmul 7.5 Iterație QR pe o matrice Hessenberg; utilizat de algoritmul 7.4 – apel $[H_1, H_2, it] = \text{QRIter}(H, t)$ **Intrare:** Matricea Hessenberg H , toleranța tol **Ieșire:** Matricele H_1, H_2 reprezentând descompunerea lui H după elementul subdiagonal minim în modul și It numărul de iterații

```

 $it := 0;$ 
Determină minimul  $m$  al modululelor elementelor de pe prima subdiagonală a lui  $H$  și poziția sa  $j$ ;
while  $m > tol$  do
     $it := it + 1;$ 
     $H := H_*; \{\text{Transformarea RQ a matricei } H\}$ 
    Determină minimul  $m$  al modululelor elementelor de pe prima subdiagonală a lui  $H$  și poziția sa  $j$ ;
end while
 $H_1 := H_{1:j,1:j};$ 
 $H_2 := H_{j+1:n,j+1:n};$ 

```

ε	#iterații	λ_1	λ_2	λ_3
10^{-3}	12	9.7406	$-3.8703 + 0.6479i$	$-3.8703 - 0.6479i$
10^{-4}	14	9.7407	$-3.8703 + 0.6479i$	$-3.8703 - 0.6479i$
10^{-5}	17	9.7407	$-3.8703 + 0.6480i$	$-3.8703 - 0.6480i$
10^{-5}	19	9.7407	$-3.8703 + 0.6480i$	$-3.8703 - 0.6480i$
10^{-5}	22	9.7407	$-3.8703 + 0.6480i$	$-3.8703 - 0.6480i$

Tabela 7.3: Rezultate pentru exemplul 7.5.6

Cuvântul de ordine este aici *deplasare (translație) spectrală*. Se observă că, pentru $\mu \in \mathbb{R}$, matricea $A - \mu I$ are valorile proprii $\lambda_1 - \mu, \dots, \lambda_n - \mu$. Pentru o matrice oarecare inversabilă, B , matricea $B(A - \mu I)B^{-1} + \mu I$ are valorile proprii $\lambda_1, \dots, \lambda_n - \mu$ - se poate așadar deplasa spectrul matricei înainte și înapoi printr-o transformare de similaritate. Se ordonează valorile proprii μ_1, \dots, μ_n astfel ca

$$|\mu_1 - \mu| > |\mu_2 - \mu| > \dots > |\mu_n - \mu|, \quad \{\mu_1, \dots, \mu_n\} = \{\lambda_1, \dots, \lambda_n\}$$

și dacă μ este apropiat de μ_n , atunci dacă metoda QR începe cu $H^0 = A - \mu I$, elementul subdiagonal $h_{n-1,n}^{(n)}$ converge foarte repede către zero. Mai bine este dacă deplasarea spectrală se realizează la fiecare pas individual. Pe lângă aceasta, putem lua ca aproximație pentru μ (în mod euristic) valoarea $h_{nn}^{(k)}$. Se obține astfel următoarea schemă iterativă

$$H^{(k+1)} = (H^{(k)} - \mu_k I)_* + \mu_k I, \quad \mu_k := h_{nn}^{(k)}, \quad k \in \mathbb{N}_0,$$

cu matricea de pornire $H^0 = Q A Q^T$. Algoritmul 7.6 dă o variantă a metodei care tratează și valorile proprii complexe. El utilizează algoritmul 7.7. În ultimul algoritm, $(H - H_{n,n} I_n)_*$ din linia 6 înseamnă transformarea RQ a matricei $H - H_{n,n} I_n$.

Observația 7.5.7. Dacă valoarea de deplasare μ este suficient de apropiată de o valoare proprie λ , atunci matricea se descompune într-un singur pas iterativ. \diamond

7.5.3. Metoda QR cu pas dublu

Se poate arăta că metoda QR cu deplasare spectrală converge *pătratic*, eroarea fiind pentru un $\rho < 1$ doar

$$O(\rho^{2k}) \text{ în loc de } O(\rho^k).$$

Oricât de frumoasă ar fi această idee ea funcționează bine doar pentru valori proprii reale, în cazul valorilor proprii complexe fiind problematică. Cu toate acestea, putem exploata faptul că valorile proprii apar în perechi. Aceasta ne conduce la ideea de „metode cu pas dublu”:

Algoritmul 7.6 Metoda QR cu deplasare spectrală, partiționare și tratarea valorilor proprii complexe

Intrare: Matricea $A \in \mathbb{R}^{n \times n}$ și toleranța tol

Ieșire: Valorile proprii λ ale lui A și numărul de iterații It

```

It := 0;
if  $n = 1$  then
     $\lambda := A$ ;
    return
else if  $n = 2$  then
     $\lambda := \text{Eigen2x2}(A)$ ;
    return
else
     $H := \text{Hessenberg}(A)$ ; {aducere la forma Hessenberg}
     $[H_1, H_2, It] := \text{QRIter2}(H, tol)$ 
     $[\lambda_1, It_1] := \text{QRSplit2}(H_1, tol)$  {apel recursiv}
     $[\lambda_2, It_2] := \text{QRSplit2}(H_2, tol)$  {apel recursiv}
     $It := It + It_1 + It_2$ ;
     $\lambda = [\lambda_1, \lambda_2]$ ;
end if

```

Algoritmul 7.7 Iterație QR și partiționare

```

It := 0;
Determină minimul  $m$  al modululelor elementelor de pe prima subdiagonală a lui  $H$ 
și poziția sa  $j$ ;
while  $m > tol$  do
     $It := It + 1$ ;
     $H := (H - H_{n,n}I_n)_* + H_{n,n}I_n$ ;
    Determină minimul  $m$  al modululelor elementelor de pe prima subdiagonală a lui
     $H$  și poziția sa  $j$ ;
end while
 $H_1 := H_{1:j,1:j}$ ;
 $H_2 := H_{j+1:n,j+1:n}$ ;

```

în loc să deplasăm spectrul cu o valoare proprie, aproximată în mod euristic cu $h_{n,n}^{(k)}$, se execută două deplasări într-un pas, și anume cu valorile proprii ale lui

$$B = \begin{bmatrix} h_{n-1,n-1}^{(k)} & h_{n-1,n}^{(k)} \\ h_{n-1,n}^{(k)} & h_{n,n}^{(k)} \end{bmatrix}.$$

Există două posibilități: fie ambele valori μ și μ' ale lui B sunt reale și atunci se procedează ca mai sus, fie sunt complexe conjugate și avem valorile proprii μ și $\bar{\mu}$. Așa cum vom vedea, al doilea caz se poate de asemenea trata în aritmetica reală. Fie $Q_k, Q'_k \in \mathbb{C}^{n \times n}$ și $R_k, R'_k \in \mathbb{C}^{n \times n}$ matricele descompunerii QR complexe

$$\begin{aligned} Q_k R_k &= H^{(k)} - \mu I, \\ Q'_k R'_k &= R_k Q_k + (\mu - \bar{\mu}) I. \end{aligned}$$

Atunci are loc

$$\begin{aligned} H^{(k+1)} &:= R'_k Q'_k + \mu I = (Q'_k)^* (R_k Q_k + (\mu - \bar{\mu}) I) Q'_k + \mu I \\ &= Q_k^* R_k Q_k Q'_k + \mu I = (Q'_k)^* Q_k^* (H^{(k)} - \mu I) Q_k Q'_k + \mu I \\ &= \underbrace{(Q_k Q'_k)^*}_{=U^*} H^{(k)} \underbrace{Q_k Q'_k}_{=U}. \end{aligned}$$

Utilizând matricea $S = R'_k R_k$ avem

$$\begin{aligned} US &= Q_k Q'_k R'_k R_k = Q_k (R_k Q_k + (\mu - \bar{\mu}) I) R_k \\ &= Q_k R_k Q_k R_k + (\mu - \bar{\mu}) Q_k R_k = (H^{(k)} - \mu I)^2 + (\mu - \bar{\mu}) (H^{(k)} - \mu I) \\ &= (H^{(k)})^2 - 2\mu H^{(k)} + \mu^2 I + (\mu - \bar{\mu}) H^{(k)} - (\mu^2 - \mu \bar{\mu}) I \\ &= (H^{(k)})^2 - (\mu + \bar{\mu}) H^{(k)} + \mu \bar{\mu} I =: X \end{aligned} \quad (7.5.3)$$

Dacă $\mu = \alpha + i\beta$, atunci $\mu + \bar{\mu} = 2\alpha$ și $\mu \bar{\mu} = |\mu|^2 = \alpha^2 + \beta^2$, așadar matricea X din membrul drept al lui (7.5.3), deci are o descompunere QR $X = QR$ reală și conform lemei 7.4.5 există o matrice de fază $\Theta \in \mathbb{C}^{n \times n}$ astfel încât $U = \Theta Q$. Dacă vom itera în real mai departe, vom obține metoda QR cu pas dublu

$$\begin{aligned} Q_k R_k &= (H^{(k)})^2 - (h_{n-1,n-1}^{(k)} + h_{n,n}^{(k)}) H^{(k)} \\ &\quad + \left((h_{n-1,n-1}^{(k)} h_{n,n}^{(k)} - h_{n-1,n}^{(k)} h_{n,n-1}^{(k)}) \right) I, \\ H^{(k+1)} &= Q_k^T H^{(k)} Q_k. \end{aligned} \quad (7.5.4)$$

Observația 7.5.8 (Metoda QR cu pas dublu). 1. Matricea X din (7.5.3) nu mai este o matrice Hessenberg, ea având o diagonală suplimentară. Cu toate acestea se poate calcula descompunerea QR a lui X destul de simplu, cu doar $2n - 3$ rotații Jacobi, în loc de $n - 1$, cât este necesar pentru o matrice Hessenberg.

2. Datorită complexității sale ridicate, înmulțirea $Q_k^T H^{(k)} Q_k$ nu mai este o metodă efectivă de iterație; acest lucru se poate remedia, a se vedea de exemplu [18] sau [36, pag. 272–278].
3. Natural, $H^{(k+1)}$ poate fi adusă la forma Hessenberg.
4. Metoda cu dublu pas se aplică numai atunci când matricea A are valori proprii complexe; din contră, la matrice simetrice nu este avantajoasă. \diamond

Metoda QR cu pas dublu, partiționare și tratarea matricelor 2×2 este dată în algoritmul 7.8. El apelează algoritmul 7.9.

Algoritmul 7.8 Metoda QR cu dublu pas, partiționare și tratarea matricelor 2×2

Intrare: Matricea $A \in \mathbb{R}^{n \times n}$ și toleranța tol

Ieșire: Valorile proprii λ ale lui A și numărul de iterații It

```

It := 0;
if  $n = 1$  then
     $\lambda := A$ ;
    return
else if  $n = 2$  then
     $\lambda := \text{Eigen2x2}(A)$ ;
    return
else
     $H := \text{Hessenberg}(A)$ ; {aducere la forma Hessenberg}
     $[H_1, H_2, It] := \text{QRDouble}(H, tol)$ 
     $[\lambda_1, It_1] := \text{QRSplit2}(H_1, tol)$  {apel recursiv}
     $[\lambda_2, It_2] := \text{QRSplit2}(H_2, tol)$  {apel recursiv}
     $It := It + It_1 + It_2$ ;
     $\lambda = [\lambda_1, \lambda_2]$ ;
end if

```

Exemplul 7.5.9. Aplicăm algoritmi 7.6 și 7.8 matricelor din exemplele 7.5.1 și 7.5.2. Se obțin rezultatele din tabela 7.4. Buna comportare a metodei pasului dublu se justifică prin aceea că obține două valori proprii dintr-o dată. \diamond

Algoritmul 7.9 Iterație QR cu dublu pas și transformare Hessenberg

$It := 0;$
 Determină minimul m al modulelor elementelor de pe prima subdiagonală a lui H și poziția sa j ;
while $m > tol$ **do**
 $It := It + 1;$
 $X := H^2 - (H_{n-1,n-1} + H_{n,n})H + (H_{n-1,n-1}H_{n,n} - H_{n,n-1}H_{n-1,n})I_n;$
 Determină descompunerea $X = QR$ a lui X ;
 $H := \text{Hessenberg}(Q^T H Q);$
 Determină minimul m al modulelor elementelor de pe prima subdiagonală a lui H și poziția sa j ;
end while
 $H_1 := H_{1:j,1:j};$
 $H_2 := H_{j+1:n,j+1:n};$

ε	#iterații în real		#iterații în complex	
	alg. 7.6	alg. 7.8	alg. 7.6	alg. 7.8
1e-010	1	1	9	4
1e-020	9	2	17	5
1e-030	26	3	45	5

Tabela 7.4: Comparațiile din exemplul 7.5.9

CAPITOLUL 8

Rezolvarea numerică a ecuațiilor diferențiale ordinare

Cuprins

8.1. Ecuații diferențiale	226
8.2. Metode numerice	227
8.3. Descrierea locală a metodelor cu un pas	228
8.4. Exemple de metode cu un pas	229
8.4.1. Metoda lui Euler	229
8.4.2. Metoda dezvoltării Taylor	231
8.4.3. Metode de tip Euler îmbunătățite	232
8.5. Metode Runge-Kutta	233
8.6. Descrierea globală a metodelor cu un pas	238
8.6.1. Stabilitatea	240
8.6.2. Convergență	243
8.6.3. Asimptotica erorii globale	244
8.7. Monitorizarea erorilor și controlul pasului	247
8.7.1. Estimarea erorii globale	247
8.7.2. Estimarea erorii de trunchiere	249
8.7.3. Controlul pasului	252

8.1. Ecuații diferențiale

Considerăm problema cu valori inițiale sau problema Cauchy ¹: să se determine o funcție cu valori vectoriale $y \in C^1[a, b]$, $y : [a, b] \rightarrow \mathbb{R}^d$, astfel încât

$$(PC) \quad \begin{cases} \frac{dy}{dx} = f(x, y), & x \in [a, b] \\ y(a) = y_0 \end{cases} \quad (8.1.1)$$

Vom evidenția două clase importante de astfel de probleme:

(i) pentru $d = 1$ avem o singură ecuație diferențială scalară de ordinul I

$$\begin{cases} y' = f(x, y) \\ y(a) = y_0 \end{cases}$$

(ii) pentru $d > 1$ avem un sistem de ecuații diferențiale ordinare de ordinul I

$$\begin{cases} \frac{dy^i}{dx} = f^i(x, y^1, y^2, \dots, y^d), & i = \overline{1, d} \\ y^i(a) = y_0^i, & i = \overline{1, d} \end{cases}$$

Reamintim următoarea teoremă clasică referitoare la existență și unicitate.

Teorema 8.1.1 *Presupunem că $f(x, y)$ este continuă în prima variabilă pentru $x \in [a, b]$ și în raport cu cea de-a doua variabilă satisface o condiție Lipschitz uniformă*

$$\|f(x, y) - f(x, y^*)\| \leq L\|y - y^*\|, \quad y, y^* \in \mathbb{R}^d, \quad (8.1.2)$$

unde $\|\cdot\|$ este o anumită normă vectorială. Atunci problema Cauchy (PC) are o soluție unică $y(x)$, $a \leq x \leq b$, $\forall y_0 \in \mathbb{R}^d$. Mai mult, $y(x)$ depinde continuu de a și y_0 .

Augustin Louis Cauchy (1789-1857), matematician francez, considerat părintele analizei moderne. A fundamentat solid analiza pe baza conceptului riguros de limită. Este de asemenea creatorul analizei complexe, în care „formula lui Cauchy” ocupă un loc central. Numele său este legat și de contribuții de pionierat în domeniul ecuațiilor diferențiale și cu derivate parțiale, în particular legate de problema existenței și unicității. La fel ca în cazul multor mari matematicieni din secolele al optzecelea și al nouăzecelea, lucrările sale au tratat probleme din geometrie, algebră, teoria numerelor, mecanică, dar și fizică teoretică.



Condiția Lipschitz (8.1.2) are singur loc dacă toate funcțiile $\frac{\partial f^i}{\partial y^j}(x, y)$, $i, j = \overline{1, d}$ sunt continue în raport cu variabilele y și sunt mărginite pe $[a, b] \times \mathbb{R}^d$. Aceasta este situația în cazul sistemelor de ecuații diferențiale ordinare liniare, unde

$$f^i(x, y) = \sum_{j=1}^d a_{ij}(x)y^j + b_i(x), \quad i = \overline{1, d}$$

și $a_{ij}(x)$, $b_i(x)$ sunt funcții continue pe $[a, b]$.

De multe ori condiția Lipschitz (8.1.2) are loc într-o vecinătate a lui x_0 astfel încât $y(x)$ să rămână într-un compact D .

8.2. Metode numerice

Se face distincție între *metode de aproximare analitice* și *metode discrete*. În cadrul primei categorii se încearcă să se găsească aproximații $y_a(x) \approx y(x)$ ale soluției exacte, valabile pentru orice $x \in [a, b]$. Acestea de obicei au forma unei dezvoltări într-o serie trunchiată, fie după puterile lui x , fie în polinoame Cebâșev, fie într-un alt sistem de funcții de bază. În cazul metodelor discrete, se încearcă să se găsească aproximații $u_n \in \mathbb{R}^d$ ale lui $y(x_n)$ pe o grilă de puncte $x_n \in [a, b]$. Abscisele x_n pot fi predeterminate (de exemplu puncte echidistante pe $[a, b]$), sau mai convenabil sunt generate dinamic ca parte a procesului de integrare.

Dacă se dorește, se pot obține din aceste aproximante discrete $\{u_n\}$ aproximante $y_n(x)$ definite pe întreg intervalul $[a, b]$, fie prin interpolare, sau mai natural, printr-un mecanism de continuare conținut în metoda de aproximare însăși. Ne vom ocupa numai de metode discrete cu un pas, adică metode în care u_{n+1} este determinat cunoscând numai x_n, u_n și pasul h pentru a trece de la x_n la $x_{n+1} = x_n + h$. Într-o metodă cu k pași ($k > 1$) este necesară cunoașterea a $k - 1$ puncte adiționale (x_{n-j}, u_{n-j}) , $j = 1, 2, \dots, k - 1$ pentru a obține o nouă componentă a soluției.

Când se descrie o metodă cu un pas este suficient să arătăm cum se trece de la un punct generic (x, y) , $x \in [a, b]$, $y \in \mathbb{R}^d$ la punctul următor $(x+h, y_{next})$. Ne vom referi la aceasta ca fiind *descrierea locală* a unei metode cu un pas. Aceasta include de asemenea o discuție a preciziei locale, adică cât de apropiat este y_{next} de soluție în $x+h$. O metodă cu un pas pentru rezolvarea problemei Cauchy (8.1.1) generează efectiv o funcție grilă $\{u_n\}$, $u_n \in \mathbb{R}^d$, pe o grilă $a = x_0 < x_1 < x_2 < \dots < x_{N-1} < x_N = b$ ce acoperă intervalul $[a, b]$, prin care se intenționează ca u_n să aproximeze soluția exactă $y(x)$ în $x = x_n$. Punctul (x_{n+1}, u_{n+1}) se obține din punctul (x_n, u_n) aplicând o metodă cu un pas, având un pas $h_n = x_{n+1} - x_n$ adecvat ales. Ne vom referi la aceasta ca *descrierea globală* a unei metode cu un pas. Chestiunile de interes aici sunt comportarea erorii globale $u_n - y(x_n)$, în particular stabilitatea, convergența și alegerea lui h_n pentru a trece de la un punct al grilei x_n , la următorul, $x_{n+1} = x_n + h_n$.

8.3. Descrierea locală a metodelor cu un pas

Dându-se un punct generic $x \in [a, b]$, $y \in \mathbb{R}^d$, definim un pas al metodei cu un pas prin

$$y_{next} = y + h\Phi(x, y; h), \quad h > 0. \quad (8.3.1)$$

Funcția $\Phi : [a, b] \times \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$ poate fi gândită ca un increment aproximativ pe unitatea de pas sau ca o aproximare a diferenței divizate și ea definește metoda. Împreună cu (8.3.1) considerăm soluția $u(t)$ a ecuației diferențiale (8.1.1) ce trece prin punctul (x, y) , adică problema locală cu valoarea inițială

$$\begin{cases} \frac{du}{dt} = f(t, u) \\ u(t) = y \quad t \in [t, t + h] \end{cases} \quad (8.3.2)$$

Vom numi $u(t)$ *soluție de referință*. Se intenționează ca vectorul y_{next} din (8.3.1) să aproximeze $u(x + h)$. Cu cât succes se realizează aceasta se măsoară prin eroarea locală de trunchiere, definită după cum urmează.

Definiția 8.3.1 Eroarea de trunchiere a metodei Φ în punctul (x, y) este definită prin

$$T(x, y; h) = \frac{1}{h}[y_{next} - u(x + h)]. \quad (8.3.3)$$

Eroarea de trunchiere este o funcție cu valori vectoriale de $d + 2$ variabile. Utilizând (8.3.1) și (8.3.2) o putem scrie sub forma

$$T(x, y; h) = \Phi(x, y; h) - \frac{1}{h}[u(x + h) - u(x)], \quad (8.3.4)$$

ceea ce arată că T este diferența între incrementul aproximativ și cel exact pe unitatea de pas.

Definiția 8.3.2 Metoda Φ se numește consistentă dacă

$$T(x, y; h) \rightarrow 0 \text{ când } h \rightarrow 0, \quad (8.3.5)$$

uniform pentru $(x, y) \in [a, b] \times \mathbb{R}^d$.

Conform lui (8.3.4) și (8.3.3) avem consistență dacă și numai dacă

$$\Phi(x, y; 0) = f(x, y), \quad x \in [a, b], \quad y \in \mathbb{R}^d. \quad (8.3.6)$$

O descriere mai fină a preciziei locale este furnizată de definiția următoare, bazată pe conceptul de eroare de trunchiere.

Definiția 8.3.3 Spunem că metoda Φ are ordinul p dacă pentru o anumită normă vectorială $\|\cdot\|$

$$\|T(x, y; h)\| \leq Ch^p, \quad (8.3.7)$$

uniform pe $[a, b] \times \mathbb{R}^d$, cu constanta C independentă de x, y și h .

Această proprietate se mai poate exprima sub forma

$$T(x, y; h) = O(h^p), \quad h \rightarrow 0. \quad (8.3.8)$$

De notat că $p > 0$ implică consistența. De obicei $p \in \mathbb{N}^*$. El se numește *ordin exact*, dacă (8.3.7) nu are loc pentru nici un p mai mare.

Definiția 8.3.4 O funcție $\tau : [a, b] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ care satisface $\tau(x, y) \neq 0$ și

$$T(x, y; h) = \tau(x, y)h^p + O(h^{p+1}), \quad h \rightarrow 0 \quad (8.3.9)$$

se numește funcție de eroare principală.

Funcția de eroare principală determină termenul principal (dominant) al erorii de trunchiere. Numărul p din (8.3.9) este ordinul exact al metodei deoarece $\tau \neq 0$.

Toate definițiile precedente sunt formulate în ideea că $h > 0$ este un număr mic. Cu cât p este mai mare, cu atât metoda este mai precisă.

8.4. Exemple de metode cu un pas

Unele dintre metodele cele mai vechi sunt motivate prin considerații geometrice simple asupra pantei definite de membrul drept al ecuației diferențiale. În această categorie intră metoda lui Euler și metoda lui Euler modificată. Alte metode mai precise și mai sofisticate se bazează pe dezvoltarea Taylor.

8.4.1. Metoda lui Euler

Euler a propus metoda sa în 1768, la începutul istoriei calculului diferențial și integral. Ea constă pur și simplu în a urma panta în punctul generic (x, y) pe un interval de lungime h

$$y_{next} = y + hf(x, y). \quad (8.4.1)$$

(vezi figura 8.1).

Astfel $\Phi(x, y; h) = f(x, y)$ nu depinde de h și conform lui (8.3.6) metoda este consistentă. Pentru eroarea de trunchiere avem conform lui (8.3.3)

$$T(x, y; h) = f(x, y) - \frac{1}{h}[u(x+h) - u(x)], \quad (8.4.2)$$

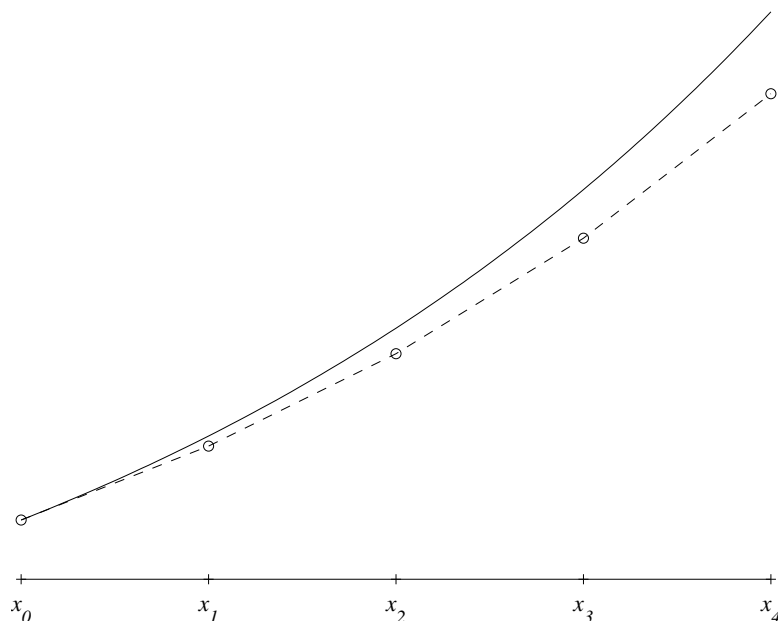


Figura 8.1: Metoda lui Euler – soluția exactă (linie continuă) și soluția aproximativă (linie punctată)

unde $u(t)$ este soluția de referință definită de (8.3.2). Deoarece $u' = f(x, u(x)) = f(x, y)$, putem scrie, utilizând formula lui Taylor

$$\begin{aligned} T(x, y; h) &= u'(x) - \frac{1}{h}[u(x+h) - u(x)] = \\ &= u'(x) - \frac{1}{h}[u(x) + hu'(x) + \frac{1}{2}h^2u''(\xi) - u(x)] = \\ &= -\frac{1}{2}hu''(\xi), \quad \xi \in (x, x+h), \end{aligned} \quad (8.4.3)$$

presupunând că $u \in C^r[x, x+h]$. Aceasta este adevărată dacă $f \in C^1([a, b] \times \mathbb{R}^d)$. Diferențiind acum total (8.3.2) în raport cu t și făcând $t = \xi$, suntem conduși la

$$T(x, y; h) = -\frac{1}{2}h[f_x + f_y f](\xi, u(\xi)), \quad (8.4.4)$$

unde f_x este derivata parțială a lui f în raport cu x și f_y este jacobianul lui f în raport cu variabila y . Dacă în spiritul teoremei 8.1.1, presupunem că f și toate derivatele sale parțiale de ordinul I sunt uniform mărginite în $[a, b] \times \mathbb{R}^d$, există o constantă C , independentă de x, y și h astfel încât

$$\|T(x, y; h)\| \leq Ch. \quad (8.4.5)$$

Astfel, metoda lui Euler are ordinul $p = 1$. Dacă facem aceeași presupunere și despre derivatele parțiale de ordinul doi ale lui f avem

$$u''(\xi) = u''(x) + O(h)$$

și de aceea din (8.4.3) rezultă

$$T(x, y; h) = -\frac{1}{2}h[f_x + f_y f](x, y) + O(h^2), \quad h \rightarrow 0, \quad (8.4.6)$$

arătând că funcția eroare principală este dată de

$$\tau(x, y) = -\frac{1}{2}[f_x + f_y f](x, y). \quad (8.4.7)$$

Exceptând situația când $f_x + f_y f \equiv 0$, ordinul exact al metodei lui Euler este $p = 1$.

8.4.2. Metoda dezvoltării Taylor

Am văzut că metoda lui Euler se bazează pe trunchierea dezvoltării Taylor a soluției de referință după cel de-al doilea termen. Este o idee naturală, propusă încă de Euler, de a utiliza mai mulți termeni din dezvoltarea Taylor. Aceasta necesită calculul succesiv al „derivatele totale“ ale lui f ,

$$\begin{aligned} f^{[0]}(x, y) &= f(x, y) \\ f^{[k+1]}(x, y) &= f_x^{[k]}(x, y) + f_y^{[k]}(x, y)f(x, y), \quad k = 0, 1, 2, \dots \end{aligned} \quad (8.4.8)$$

care determină derivatele succesive ale soluției de referință $u(t)$ a lui (8.3.2) în virtutea relației

$$u^{(k+1)}(t) = f^{[k]}(t, u(t)), \quad k = 0, 1, 2, \dots \quad (8.4.9)$$

Acestea, pentru $t = x$ devin

$$u^{(k+1)}(x) = f^{[k]}(x, y), \quad k = 0, 1, 2, \dots \quad (8.4.10)$$

și sunt utilizate pentru a scrie dezvoltarea Taylor conform cu

$$y_{next} = y + h \left[f^{[0]}(x, y) + \frac{1}{2}h f^{[1]}(x, y) + \dots + \frac{1}{p!}h^{p-1} f^{[p-1]}(x, y) \right], \quad (8.4.11)$$

adică

$$\Phi(x, y; h) = f^{[0]}(x, y) + \frac{1}{2}h f^{[1]}(x, y) + \dots + \frac{1}{p!}h^{p-1} f^{[p-1]}(x, y). \quad (8.4.12)$$

Pentru eroarea de trunchiere, utilizând (8.4.10) și (8.4.12) și presupunând că $f \in C^p([a, b] \times \mathbb{R}^d)$ se obține din teorema lui Taylor

$$\begin{aligned} T(x, y; h) &= \Phi(x, y; h) - \frac{1}{h}[u(x+h) - u(x)] = \\ &= \Phi(x, y; h) - \sum_{k=0}^{p-1} u^{(k+1)}(x) \frac{h^k}{(k+1)!} - u^{(p+1)}(\xi) \frac{h^p}{(p+1)!} = \\ &= -u^{(p+1)}(\xi) \frac{h^p}{(p+1)!}, \quad \xi \in (x, x+h), \end{aligned}$$

așa că

$$\|T(x, y; h)\| \leq \frac{C_p}{(p+1)!} h^p,$$

unde C_p este o margine a derivatei totale de ordin p a lui f . Astfel metoda are ordinul exact p (în afară de cazul când $f^{[p]}(x, y) \equiv 0$) și funcția eroare principală este

$$\tau(x, y) = -\frac{1}{(p+1)!} f^{[p]}(x, y). \quad (8.4.13)$$

Necesitatea calculului multor derivate parțiale în (8.4.8) a fost un factor descurajant în trecut, când se făcea cu mâna. Dar în zilele noastre această sarcină poate cădea în seama calculatorului, astfel încât metoda a devenit din nou o opțiune viabilă.

8.4.3. Metode de tip Euler îmbunătățite

Există prea multă inerție în metoda lui Euler: nu se urmează aceeași pantă pe întreg intervalul de lungime h , deoarece de-a lungul acestui segment de dreaptă panta definită de ecuația diferențială se schimbă. Aceasta sugerează mai multe alternative. De exemplu am putea să reevaluăm panta la mijlocul segmentului – să luăm pulsul ecuației diferențiale – și apoi să urmărim panta actualizată pe întreg intervalul (vezi figura 8.2). Formula este

$$y_{next} = y + hf \left(x + \frac{1}{2}h, y + \frac{1}{2}hf(x, y) \right) \quad (8.4.14)$$

sau

$$\Phi(x, y; h) = f \left(x + \frac{1}{2}h, y + \frac{1}{2}hf(x, y) \right) \quad (8.4.15)$$

Observați „imbricarea” necesară aici. Pentru programarea acestei metode este indicat să se scrie

$$\begin{aligned} K_1(x, y) &= f(x, y) \\ K_2(x, y; h) &= f \left(x + \frac{1}{2}h, y + \frac{1}{2}hK_1 \right) \\ y_{next} &= y + hK_2 \end{aligned} \quad (8.4.16)$$

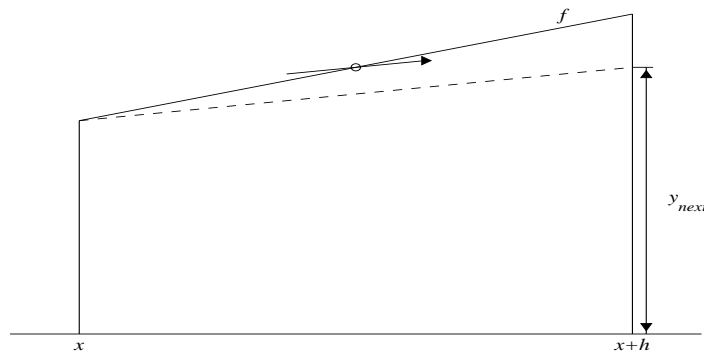


Figura 8.2: Metoda lui Euler modificată

Cu alte cuvinte, încercăm să luăm două pante de test K_1 și K_2 , una în punctul inițial și a doua în apropiere și apoi s-o alegem pe ultima ca pantă finală. Metoda se numește *metoda lui Euler modificată*.

Putem la fel de bine să luăm o a doua pantă de încercare $(x + h, y + hf(x, y))$, dar atunci, deoarece trebuie să așteptăm prea mult înainte de a reevalua panta, luăm ca pantă finală media celor două pante

$$\begin{aligned} K_1(x, y) &= f(x, y) \\ K_2(x, y; h) &= f(x + h, y + hK_1) \\ y_{next} &= y + \frac{1}{2}h(K_1 + K_2). \end{aligned} \quad (8.4.17)$$

Această metodă se numește *metoda lui Heun*. Efectul ambelor modificări este creșterea ordinului cu 1, așa cum se va vedea în continuare.

8.5. Metode Runge-Kutta

Se caută Φ de forma:

$$\begin{aligned} \Phi(x, y; h) &= \sum_{s=1}^r \alpha_s K_s \\ K_1(x, y) &= f(x, y) \\ K_s(x, y) &= f\left(x + \mu_s h, y + h \sum_{j=1}^{s-1} \lambda_{sj} K_j\right), \quad s = 2, 3, \dots, r \end{aligned} \quad (8.5.1)$$

Este natural să impunem în (8.5.1) condițiile

$$\mu_s = \sum_{j=1}^{s-1} \lambda_{sj}, \quad s = 2, 3, \dots, r, \quad \sum_{s=1}^r \alpha_s = 1, \quad (8.5.2)$$

unde primul set de condiții este echivalent cu

$$K_s(x, y; h) = u'(x + \mu_s h) + O(h^2), \quad s \geq 2,$$

iar a doua este condiția de consistență (8.3.6) (adică $\Phi(x, y; h) = f(x, y)$).

Vom numi metoda (8.5.1) *metodă Runge-Kutta explicită în r stadii* deoarece necesită r evaluări ale funcției f din membrul drept al ecuației diferențiale. Condițiile (8.5.2) conduc la un sistem neliniar. Fie $p^*(r)$ ordinul maxim pentru o metodă Runge-Kutta explicită în r stadii. Kutta ² a arătat în 1901 că

$$p^*(r) = r, \quad r = \overline{1, 4}.$$

Se pot considera *metode Runge-Kutta implicite* cu r stadii

$$\begin{aligned} \Phi(x, y; h) &= \sum_{s=1}^r \alpha_s K_s(x, y; h), \\ K_s &= f \left(x + \mu_s h, y + \sum_{j=1}^r \lambda_{sj} K_j \right), \quad s = \overline{1, r}, \end{aligned} \quad (8.5.3)$$

în care ultimele r ecuații formează un sistem de ecuații (în general neliniar) cu necunoscutele K_1, K_2, \dots, K_r . Deoarece fiecare dintre necunoscute este un vector din \mathbb{R}^d , înainte de construirea incrementului aproximativ Φ trebuie să rezolvăm un sistem de rd ecuații cu rd necunoscute. *Metodele Runge-Kutta semiimplicite*, la care limitele de sumare merg de la $j = 1$ la $j = s$ necesită un efort mai mic. Se ajunge la un sistem de r ecuații, fiecare având d necunoscute, componentele lui K_s . Volumul considerabil de calcule necesar în metodele implicite și semiimplicite se justifică numai în împrejurări



Wilhelm Martin Kutta (1867-1944), matematician german, cu preocupări în domeniul matematicilor aplicate. Cunoscut pentru lucrările sale în domeniul rezolvării numerice a ecuațiilor diferențiale ordinare, a avut și contribuții al aplicarea transformărilor conforme la probleme de hidro și aerodinamică (formula Kutta-Jukovski).

speciale, de exemplu la rezolvarea problemelor stiff. Motivul este acela că metodele implicite pot avea ordin mai mare și proprietăți de stabilitate mai bune.

Parametrii se aleg astfel ca ordinul să fie cât mai mare posibil.

Exemplul 8.5.1. Fie

$$\Phi(x, y; h) = \alpha_1 K_1 + \alpha_2 K_2 \quad \diamond$$

unde

$$\begin{aligned} K_1(x, y) &= f(x, y) \\ K_2(x, y; h) &= f(x + \mu_2 h, y + \lambda_{21} h K_1) \\ \lambda_{21} &= \mu_2 \end{aligned}$$

Avem deci 3 parametri α_1, α_2, μ . Un mod sistematic de a determina ordinul maxim p este de a dezvolta atât $\Phi(x, y; h)$ cât și $h^{-1}[u(x+h) - u(x)]$ după puterile lui h și să impunem coincidența a cât mai multor termeni posibili, fără a impune restricții asupra lui f . Pentru a dezvolta Φ avem nevoie de dezvoltarea Taylor a unei funcții vectoriale de mai multe variabile

$$\begin{aligned} f(x + \Delta x, y + \Delta y) &= f + f_x \Delta x + f_y \Delta y + \\ &+ \frac{1}{2} [f_{xx} (\Delta x)^2 + 2f_{xy} \Delta x \Delta y + (\Delta y)^T f_{yy} (\Delta y)] + \dots, \end{aligned} \quad (8.5.4)$$

unde f_y este jacobianul lui f , iar $f_{yy} = [f_{yy}^i]$ este vectorul matricelor hessiene ale lui f . În formula de mai sus toate funcțiile și derivatele parțiale se evaluează în (x, y) . Punând $\Delta x = \mu h$, $\Delta y = \mu h f$ obținem

$$\begin{aligned} K_2(x, y; h) &= f + \mu h (f_x + f_y f) \\ &+ \frac{1}{2} \mu^2 h^2 (f_{xx} + 2f_{xy} f + f^T f_{yy} f) + O(h^3) \end{aligned} \quad (8.5.5)$$

$$\frac{1}{h} [u(x+h) - u(x)] = u'(x) + \frac{1}{2} h u''(x) + \frac{1}{6} h^2 u'''(x) + O(h^3) \quad (8.5.6)$$

unde

$$\begin{aligned} u'(x) &= f \\ u''(x) &= f^{[1]} = f_x + f_y f \\ u'''(x) &= f^{[2]} = f_x^{[1]} + f_y^{[1]} f = f_{xx} + f_x f_y f + f_y f_x + (f_{xy} + (f_y f)_y) f = \\ &= f_{xx} + 2f_{xy} f + f^T f_{yy} f + f_y (f_x + f_y) f \end{aligned}$$

unde în ultima ecuație s-a utilizat

$$(f_y f)_y f = f^T f_{yy} f + f_y^2 f$$

Avem

$$T(x, y; h) = \alpha_1 K_1 + \alpha_2 K_2 - \frac{1}{h} [u(x+h) - u(x)]$$

în care înlocuim (8.5.5) și (8.5.6). Găsim

$$T(x, y; h) = (\alpha_1 + \alpha_2 - 1)f + \left(\alpha_2\mu - \frac{1}{2}\right)h(f_x + f_y f) + \\ + \frac{1}{2}h^2 \left[\left(\alpha_2\mu^2 - \frac{1}{3}\right)(f_{xx} + 2f_{xy}f + f^T f_{yy}f) - \frac{1}{3}f_y(f_x + f_y f) \right] + O(h^3) \quad (8.5.7)$$

Nu putem impune asupra coeficientul lui h^2 condiția ca el să fie zero decât dacă impunem restricții severe asupra lui f . Astfel ordinul maxim este 2 și el se obține pentru

$$\begin{cases} \alpha_1 + \alpha_2 = 1 \\ \alpha_2\mu = \frac{1}{2} \end{cases}$$

Soluția

$$\alpha_1 = 1 - \alpha_2 \\ \mu = \frac{1}{2\alpha_2}$$

depinde de un parametru, $\alpha_2 \neq 0$, arbitrar.

Pentru $\alpha_2 = 1$ avem metoda lui Euler modificată, iar pentru $\alpha_2 = \frac{1}{2}$ metoda lui Heun.

Vom menționa formula Runge-Kutta clasică de ordin $p = 4$.

$$\begin{aligned} \Phi(x, y; h) &= \frac{1}{6}(K_1 + 2K_2 + 2K_3 + K_4) \\ K_1(x, y; h) &= f(x, y) \\ K_2(x, y; h) &= f\left(x + \frac{1}{2}h, y + \frac{1}{2}hK_1\right) \\ K_3(x, y; h) &= f\left(x + \frac{1}{2}h, y + \frac{1}{2}hK_2\right) \\ K_4(x, y; h) &= f(x + h, y + hK_3) \end{aligned} \quad (8.5.8)$$

Dacă f nu depinde de y , atunci (8.5.8) coincide cu formula lui Simpson. Runge ³ a avut ideea de a generaliza formula lui Simpson la ecuații diferențiale ordinare. El a reușit



Carl David Tolmé Runge (1856-1927) matematician german, membru al școlii matematice de la Göttingen și unul dintre pionierii matematicii numerice. Este cunoscut pentru metodele Runge-Kutta din domeniul rezolvării numerice a ecuațiilor diferențiale ordinare, ale căror idei de bază i se datorează. A avut contribuții notabile și în domeniul aproximărilor în planul complex.

doar parțial, formula sa având $r = 4$ și $p = 3$. Metoda (8.5.8) a fost descoperită de Kutta în 1901 printr-o căutare sistematică.

În cazul când f nu depinde de y , atunci (8.5.8) se reduce la formula lui Simpson. Metoda Runge-Kutta clasică de ordinul 4 pentru o grilă de $N + 1$ puncte echidistante este dată de algoritmul 8.1.

Algoritmul 8.1 Metoda Runge-Kutta de ordinul 4

Intrare: capetele a, b ale intervalului; întregul N ; valoarea inițială α .

Ieșire: $N + 1$ abscise t și aproximantele w ale lui valorilor lui y în t .

$$h := (b - a)/N;$$

$$t_0 := a;$$

$$w_0 := \alpha;$$

for $i := 0$ **to** $N - 1$ **do**

$$K_1 := hf(t_i, w_i);$$

$$K_2 := hf(t_i + h/2, w_i + K_1/2);$$

$$K_3 := hf(t_i + h/2, w_i + K_2/2);$$

$$K_4 := hf(t_i + h, w_i + K_3);$$

$$w_{i+1} := w_i + (K_1 + 2 * K_2 + 2 * K_3 + K_4);$$

$$t_{i+1} := t_i + h;$$

end for

Exemplul 8.5.2. Utilizând metoda Runge-Kutta de ordinul 4 pentru a aproxima soluția problemei Cauchy

$$y' = -y + t + 1, \quad t \in [0, 1]$$

$$y(0) = 1,$$

cu $h = 0.1$, $N = 10$ și $t_i = 0.1i$ se obțin rezultatele din tabelul 8.1. ◇

Se obișnuiește să se asocieze unei metode Runge-Kutta cu r stadii (8.5.3) tabloul

$$\begin{array}{c|cccc} \mu_1 & \lambda_{11} & \lambda_{12} & \dots & \lambda_{1r} \\ \mu_2 & \lambda_{21} & \lambda_{22} & \dots & \lambda_{2r} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \mu_r & \lambda_{r1} & \lambda_{r2} & \dots & \lambda_{rr} \\ \hline & \alpha_1 & \alpha_2 & \dots & \alpha_r \end{array} \quad \left(\text{în formă matricială } \frac{\mu}{\alpha^T} \right)$$

numit *tabelă Butcher*. Pentru o metodă explicită $\mu_1 = 0$ și Λ este triunghiulară inferior, cu zerouri pe diagonală. Putem asocia primelor r linii ale tabelului Butcher o formulă de cuadratură $\int_0^{\mu_s} u(t) dt \approx \sum_{j=1}^r \lambda_{sj} u(\mu_j)$, $s = \overline{1, r}$ și ultimei linii formula de cuadratură

t_i	Aproximante	Valori exacte	Eroarea
0.0	1	1	0
0.1	1.00483750000	1.00483741804	8.19640e-008
0.2	1.01873090141	1.01873075308	1.48328e-007
0.3	1.04081842200	1.04081822068	2.01319e-007
0.4	1.07032028892	1.07032004604	2.42882e-007
0.5	1.10653093442	1.10653065971	2.74711e-007
0.6	1.14881193438	1.14881163609	2.98282e-007
0.7	1.19658561867	1.19658530379	3.14880e-007
0.8	1.24932928973	1.24932896412	3.25617e-007
0.9	1.30656999120	1.30656965974	3.31459e-007
1.0	1.36787977441	1.36787944117	3.33241e-007

Tabela 8.1: Rezultate numerice pentru exemplul 8.5.2

$\int_0^1 u(t) dt \approx \sum_{s=1}^r \alpha_s u(\mu_j)$. Dacă gradele de exactitate respective sunt $d_s = q_s - 1$, $1 \leq s \leq r + 1$ ($d_s = \infty$ dacă $\mu_s = 0$ și toți $\lambda_{sj} = 0$) atunci din teorema lui Peano rezultă că în reprezentarea restului apar derivatele de ordinul q_s ale lui u și deci punând $u(t) = y'(x + th)$ se obține

$$\frac{y(x + \mu_s h) - y(x)}{h} - \sum_{j=1}^r \lambda_{sj} y'(x + \mu_j h) = O(h^{q_s}), \quad s = \overline{1, r}$$

și

$$\frac{y(x + h) - y(x)}{h} - \sum_{s=1}^r \alpha_s y'(x + \mu_s h) = O(h^{q_{r+1}}).$$

Pentru metoda Runge-Kutta clasică de ordinul patru (8.5.8) tabela Butcher este:

0	0			
$\frac{1}{2}$	$\frac{1}{2}$	0		
$\frac{1}{2}$	0	$\frac{1}{2}$	0	
1	0	0	1	0
	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{1}{6}$

8.6. Descrierea globală a metodelor cu un pas

Descrierea metodelor se face cel mai bine în termeni de grile și funcții grilă.

O *grilă* pe intervalul $[a, b]$ este o mulțime de puncte $\{x_n\}_{n=0}^N$ astfel încât

$$a = x_0 < x_1 < x_2 < \cdots < x_{N-1} < x_N = b, \quad (8.6.1)$$

cu lungimile grilei definite prin

$$h_n = x_{n+1} - x_n, \quad n = 0, 1, \dots, N-1. \quad (8.6.2)$$

Finețea grilei este măsurată prin

$$|h| = \max_{0 \leq n \leq N-1} h_n. \quad (8.6.3)$$

Vom utiliza litera h pentru a desemna colecția de lungimi $h = \{h_n\}$. Dacă $h_1 = h_2 = \cdots = h_N = (b-a)/N$, grila se numește *uniformă*, iar în caz contrar *neuniformă*. O funcție cu valori vectoriale $v = \{v_n\}$, $v_n \in \mathbb{R}^d$, definită pe grila (8.6.1) se numește *funcție grilă*. Astfel, v_n este valoarea funcției v în punctul x_n al grilei. Orice funcție $v(x)$ definită pe $[a, b]$ induce o funcție grilă prin restricție. Vom desemna mulțimea funcțiilor grilă pe $[a, b]$ prin $\Gamma_h[a, b]$ și pentru fiecare funcție grilă $v = \{v_n\}$ definim norma sa prin

$$\|v\|_\infty = \max_{0 \leq n \leq N} \|v_n\|, \quad v \in \Gamma_h[a, b] \quad (8.6.4)$$

O metodă cu un pas – de fapt orice metodă discretă – este o metodă care produce o funcție grilă $u = \{u_n\}$ astfel încât $u \approx y$, unde $y = \{y_n\}$ este funcția grilă indusă de soluția exactă a problemei Cauchy. Fie metoda

$$\begin{aligned} x_{n+1} &= x_n + h_n \\ u_{n+1} &= u_n + h_n \Phi(x_n, u_n; h_n) \end{aligned} \quad (8.6.5)$$

unde $x_0 = a$, $u_0 = y_0$.

Pentru a clarifica analogia dintre (8.1.1) și (8.6.5) vom introduce operatorii R și R_h care acționează pe $C^1[a, b]$ și respectiv pe $\Gamma_h[a, b]$. Aceștia sunt operatorii reziduali

$$(Rv)(x) := v'(x) - f(x, v(x)), \quad v \in C^1[a, b] \quad (8.6.6)$$

$$(R_h v)_n := \frac{1}{h_n} (v_{n+1} - v_n) - \Phi(x_n, v_n; h_n), \quad n = 0, 1, \dots, N-1, \quad (8.6.7)$$

unde $v = \{v_n\} \in \Gamma_h[a, b]$. (Funcția grilă $\{(R_h v)_n\}$ nu este definită pentru $n = N$, dar putem lua arbitrar $(R_h v)_N = (R_h v)_{N-1}$). Atunci problema Cauchy și analogul său discret (8.6.5) se pot scrie transparent ca

$$Ry = 0 \text{ pe } [a, b], \quad y(a) = y_0 \quad (8.6.8)$$

$$R_h u = 0 \text{ pe } [a, b], \quad u_0 = y_0 \quad (8.6.9)$$

De notat că operatorul rezidual (8.6.7) este strâns înrudit cu eroarea locală de trunchiere (8.3.3) când aplicăm operatorul într-un punct $(x_n, y(x_n))$ pe traiectoria soluției exacte. Atunci, într-adevăr, soluția de referință $u(t)$ coincide cu soluția $y(t)$ și

$$\begin{aligned} (R_h y)_n &= \frac{1}{h_n} [y(x_{n+1}) - y(x_n)] - \Phi(x_n, y(x_n); h_n) = \\ &= -T(x_n, y(x_n); h_n). \end{aligned} \quad (8.6.10)$$

8.6.1. Stabilitatea

Stabilitatea este o proprietate numai a schemei numerice (8.6.5) și nu are nimic de-a face apriori cu puterea de aproximare. Ea caracterizează robustețea schemei în raport cu perturbații mici. Totuși stabilitatea combinată cu consistența conduce la convergența soluției numerice către soluția adevărată. Definim stabilitatea în termeni de operatori reziduali discreți R_h în (8.6.7). Ca de obicei, presupunem că $\Phi(x, y; h)$ este definită pe $[a, b] \times \mathbb{R}^d \times [0, h_0]$, unde $h_0 > 0$ este un număr pozitiv adecvat.

Definiția 8.6.1 *Metoda (8.6.5) se numește stabilă pe $[a, b]$ dacă există o constantă $K > 0$ care nu depinde de h astfel încât pentru o grilă arbitrară h pe $[a, b]$ și pentru două funcții grilă arbitrare $v, w \in \Gamma_h[a, b]$ are loc*

$$\|v - w\|_\infty \leq K (\|v_0 - w_0\|_\infty + \|R_h v - R_h w\|_\infty), \quad v, w \in \Gamma_h[a, b] \quad (8.6.11)$$

pentru orice h cu $|h|$ suficient de mică. În (8.6.11) norma este definită prin (8.6.4).

Vom numi în continuare (8.6.11) *inegalitatea de stabilitate*. Motivația ei este următoarea. Să presupunem că avem două funcții grilă u, w ce satisfac

$$R_h u = 0, \quad u_0 = y_0 \quad (8.6.12)$$

$$R_h w = \varepsilon, \quad w_0 = y_0 + \eta_0, \quad (8.6.13)$$

unde $\varepsilon = \{\varepsilon_n\} \in \Gamma_h[a, b]$ este o funcție grilă cu $\|\varepsilon_n\|$ mic și $\|\eta_0\|$ este de asemenea mic. Putem interpreta $u \in \Gamma_h[a, b]$ ca fiind rezultatul aplicării schemei numerice (8.6.5) cu precizie infinită, în timp ce $w \in \Gamma_h[a, b]$ ar putea fi soluția lui (8.6.5) în aritmetica în virgulă flotantă. Atunci, dacă stabilitatea are loc, avem

$$\|u - w\|_\infty \leq K (\|\eta_0\|_\infty + \|\varepsilon\|_\infty), \quad (8.6.14)$$

adică, schimbarea locală în u este de același ordin de mărime ca și eroarea reziduală locală $\{\varepsilon_n\}$ și eroarea inițială η_0 . Trebuie apreciat, totuși, că prima ecuație în (8.6.13) spune

$$w_{n+1} - w_n - h_n \Phi(x_n, w_n, h_n) = h_n \varepsilon_n,$$

însemnând că erorile de rotunjire trebuie să tindă către zero atunci când $|h| \rightarrow \infty$.

Interesant, pentru stabilitate este necesară doar o condiție Lipschitz asupra lui Φ .

Teorema 8.6.2 Dacă $\Phi(x, y; h)$ satisface o condiție Lipschitz în raport cu variabila y ,

$$\|\Phi(x, y; h) - \Phi(x, y^*; h)\| \leq M\|y - y^*\| \text{ pe } [a, b] \times \mathbb{R}^d \times [0, h_0], \quad (8.6.15)$$

atunci metoda (8.6.5) este stabilă.

Vom pregăti demonstrația prin următoarea leamnă utilă.

Lema 8.6.3 Fie $\{e_n\}$ o secvență de numere, $e_n \in \mathbb{R}$, ce satisface

$$e_{n+1} \leq a_n e_n + b_n, \quad n = 0, 1, \dots, N-1 \quad (8.6.16)$$

unde $a_n > 0$ și $b_n \in \mathbb{R}$. Atunci

$$e_n \leq E_n, \quad E_n = \left(\prod_{k=0}^{n-1} a_k \right) e_0 + \sum_{k=0}^{n-1} \left(\prod_{\ell=k+1}^{n-1} a_\ell \right) b_k, \quad n = 0, 1, \dots, N \quad (8.6.17)$$

Adoptăm aici convenția uzuală că un produs vid are valoarea 1 și o sumă vidă are valoarea 0.

Demonstrația lemei 8.6.3. Se verifică că

$$E_{n+1} = a_n E_n + b_n, \quad n = 0, 1, \dots, N-1, \quad E_0 = e_0.$$

Scăzând aceasta din (8.6.16) se obține

$$e_{n+1} - E_{n+1} \leq a_n (e_n - E_n), \quad n = 0, 1, \dots, N-1.$$

Acum, $e_0 - E_0 = 0$, așa că $e_1 - E_1 \leq 0$, căci $a_0 > 0$. Prin inducție, mai general, $e_n - E_n \leq 0$, deoarece $a_{n-1} > 0$. \square

Demonstrația teoremei 8.6.2. Fie $h = \{h_n\}$ o grilă arbitrară pe $[a, b]$ și $v, w \in \Gamma_h[a, b]$ două funcții grilă cu valori vectoriale arbitrare. Din definiția lui R_h putem scrie

$$v_{n+1} = v_n + h_n \Phi(x_n, v_n; h_n) + h_n (R_h v)_n, \quad n = 0, 1, \dots, N-1$$

și similar pentru w_{n+1} . Scăzându-le obținem

$$\begin{aligned} v_{n+1} - w_{n+1} &= v_n - w_n + h_n [\Phi(x_n, v_n; h_n) - \Phi(x_n, w_n; h_n)] + \\ &+ h_n [(R_h v)_n - (R_h w)_n], \quad n = 0, 1, \dots, N-1. \end{aligned} \quad (8.6.18)$$

Definim acum

$$e_n = \|v_n - w_n\|, \quad d_n = \|(R_h v)_n - (R_h w)_n\|, \quad \delta = \|d_n\|_\infty. \quad (8.6.19)$$

Utilizând inegalitatea triunghiului în (8.6.18) și condiția Lipschitz (8.6.19) pentru Φ obținem

$$e_{n+1} \leq (1 + h_n M)e_n + h_n \delta, \quad n = 0, 1, \dots, N-1 \quad (8.6.20)$$

Aceasta este inegalitatea (8.6.16) cu $a_n = 1 + h_n M$, $b_n = h_n \delta$. Deoarece $k = 0, 1, \dots, n-1$, $n \leq N$ avem

$$\begin{aligned} \prod_{\ell=k+1}^{n-1} a_\ell &\leq \prod_{\ell=0}^{n-1} a_\ell = \prod_{\ell=0}^{N-1} (1 + h_\ell M) \leq \prod_{\ell=0}^{N-1} e^{h_\ell M} \\ &= e^{(h_0+h_1+\dots+h_{N-1})M} = e^{(b-a)M}, \end{aligned}$$

unde în a doua inegalitate a fost utilizată inegalitatea clasică $1 + x \leq e^x$, din lema 8.6.3 se obține că

$$\begin{aligned} e_n &\leq e^{(b-a)M} e_0 + e^{(b-a)M} \sum_{k=0}^{n-1} h_k \delta \leq \\ &\leq e^{(b-a)M} (e_0 + (b-a)\delta), \quad n = 0, 1, \dots, N-1. \end{aligned}$$

De aceea

$$\|e\|_\infty = \|v - w\|_\infty \leq e^{(b-a)M} (\|v_0 - w_0\| + (b-a)\|R_h v - R_h w\|_\infty),$$

care este (8.6.11) cu $K = e^{(b-a)M} \max\{1, b-a\}$. \square

Am demonstrat de fapt stabilitatea pentru orice $|h| \leq h_0$, nu numai pentru h suficient de mic.

Toate metodele cu un pas utilizate în practică satisfac o condiție Lipschitz dacă f satisface o astfel de condiție și constanta M pentru Φ poate fi aproximată cu ajutorul constantei L pentru f . Este evident pentru metoda lui Euler și nu este dificil de demonstrat pentru celelalte. Este util de observat că Φ nu este nevoie să fie continuă în x ; continuitatea pe porțiuni fiind suficientă atât timp cât (8.6.15) are loc pentru orice $x \in [a, b]$, luând limitele laterale în punctele de discontinuitate.

Vom folosi următoarea aplicație a lemei 8.6.3, relativă la o funcție grilă $v \in \Gamma_h[a, b]$ ce satisface

$$v_{n+1} = v_n + h_n(A_n v_n + b_n), \quad n = 0, 1, \dots, N-1, \quad (8.6.21)$$

unde $A_n \in \mathbb{R}^{d \times d}$, $b_n \in \mathbb{R}^d$, și h_n este o grila arbitrară pe $[a, b]$.

Lema 8.6.4 *Presupunem că în (8.6.21)*

$$\|A_n\| \leq M, \quad \|b_n\| \leq \delta, \quad n = 0, 1, \dots, N-1, \quad (8.6.22)$$

unde constantele M , δ nu depind de h . Atunci, există o constantă $K > 0$, independentă de h , dar depinzând de $\|v_0\|$, astfel încât

$$\|v\|_\infty \leq K. \quad (8.6.23)$$

Demonstrație. Lemma rezultă observând că

$$\|v_{n+1}\| \leq (1 + h_n M) \|v_n\| + h_n \delta, \quad n = 0, 1, \dots, N - 1,$$

care este chiar inegalitatea (8.6.19) din demonstrația teoremei 8.6.2, deci

$$\|v_n\| \leq e^{(b-a)M} \{ \|v_0\| + (b-a)\delta \}. \quad (8.6.24)$$

□

8.6.2. Convergență

Stabilitatea este un concept puternic. Ea implică aproape imediat convergența și este un instrument de deducere a estimării erorii globale. Vom începe prin a defini precis ce înseamnă convergența.

Definiția 8.6.5 Fie $a = x_0 < x_1 < x_2 < \dots < x_N = b$ o grilă pe $[a, b]$ cu lungimea grilei $|h| = \max_{1 \leq n \leq N} (x_n - x_{n-1})$. Fie $u = \{u_n\}$ o funcție grilă definită aplicând metoda (8.6.5) pe $[a, b]$ și $y = \{y_n\}$ funcția grilă indusă de soluția exactă a problemei Cauchy. Metoda (8.6.5) se numește convergentă pe $[a, b]$ dacă are loc

$$\|u - y\|_\infty \rightarrow 0, \quad \text{când } |h| \rightarrow 0 \quad (8.6.25)$$

Teorema 8.6.6 Dacă metoda (8.6.5) este consistentă și stabilă pe $[a, b]$, atunci ea converge. Mai mult, dacă Φ are ordinul p , atunci

$$\|u - y\|_\infty = O(|h|^p) \quad \text{când } |h| \rightarrow 0. \quad (8.6.26)$$

Demonstrație. Din inegalitatea de stabilitate (8.6.11) aplicată funcțiilor grilă $v = h$ și $w = y$ din definiția 8.6.5 avem pentru $|h|$ suficient de mic

$$\|u - y\|_\infty \leq K(\|u_0 - y(x_0)\| + \|R_h u - R_h y\|_\infty) = K \|R_h y\| \quad (8.6.27)$$

deoarece $u_0 = y(x_0)$ și $R_h u = 0$ conform (8.6.5). Dar din (8.6.10)

$$\|R_h y\|_\infty = \|T(\cdot, y; h)\|_\infty \quad (8.6.28)$$

unde T este eroare de trunchiere a metodei Φ . Din definiția consistenței

$$\|T(\cdot, y; h)\|_\infty \rightarrow 0, \quad \text{când } |h| \rightarrow 0,$$

ceea ce demonstrează prima parte a teoremei. Partea a doua rezultă imediat din (8.6.27) și (8.6.28), deoarece ordinul p , prin definiția înseamnă că

$$\|T(\cdot, y; h)\|_\infty = O(|h|^p), \quad \text{când } |h| \rightarrow 0. \quad (8.6.29)$$

□

8.6.3. Asimptotica erorii globale

Deoarece funcția eroare principală descrie contribuția termenului principal al erorii locale de trunchiere este de interes să identificăm termenul dominant în eroarea globală $u_n - y(x_n)$. Pentru a simplifica lucrurile vom presupune că avem o grilă de lungime constantă h , deși nu este dificil să lucrăm cu o grilă de lungime variabilă de forma $h_n = \vartheta(x_n)h$, unde $\vartheta(x)$ este o funcție continuă pe porțiuni și $0 < \vartheta(x) < \theta$ pentru $a \leq x \leq b$. Astfel, considerăm că metoda cu un pas având forma

$$\begin{aligned} x_{n+1} &= x_n + h \\ u_{n+1} &= u_n + h\Phi(x_n, u_n; h); \quad n = 0, 1, \dots, N-1 \\ x_0 &= a, \quad u_0 = y_0, \end{aligned} \quad (8.6.30)$$

definește o funcție grilă $u = \{u_n\}$ pe o grilă uniformă pe $[a, b]$. Suntem interesați în comportarea asimptotică a lui $u_n - y(x_n)$ când $h \rightarrow 0$, unde $y(x)$ este soluția exactă a problemei Cauchy

$$\begin{cases} \frac{dy}{dx} = f(x, y) & x \in [a, b] \\ y(a) = y_0 \end{cases} \quad (8.6.31)$$

Teorema 8.6.7 *Presupunem că*

- (1) $\Phi(x, y, h) \in C^2([a, b] \times \mathbb{R}^d \times [0, h_0])$;
- (2) Φ este o metodă de ordin $p \geq 1$ ce admite o funcție de eroare principală $\tau(x, y) \in C([a, b] \times \mathbb{R}^d)$;
- (3) $e(x)$ este soluția problemei Cauchy

$$\begin{cases} \frac{de}{dx} = f_y(x, y(x))e + \tau(x, y(x)), & a \leq x \leq b \\ e(a) = 0 \end{cases} \quad (8.6.32)$$

Atunci, pentru $n = \overline{0, N}$,

$$u_n - y(x_n) = e(x_n)h^p + O(h^{p+1}), \quad \text{când } h \rightarrow 0. \quad (8.6.33)$$

Înainte de a demonstra teorema, facem următoarele observații:

1. Semnificația precisă a lui (8.6.33) este

$$\|u - y - h^p e\|_\infty = O(h^{p+1}),$$

unde u, y, e sunt funcțiile grilă $u = \{u_n\}$, $y = \{y(x_n)\}$ și $e = \{e(x_n)\}$.

2. Datorită consistenței $\Phi(x, y; 0) = f(x, y)$, condiția (1) din enunț implică $f \in C^2([a, b] \times \mathbb{R}^d)$, ceea ce este mai mult decât suficient pentru a garanta existența și unicitatea soluției $e(x)$ a lui (8.6.32) pe întreg intervalul $[a, b]$.
3. Faptul că anumite componente, dar nu toate ale lui $\tau(x, y)$ ar putea să fie identic nule nu implică faptul că $e(x)$ se anulează de asemenea, deoarece (8.6.32) este un sistem cuplat de ecuații diferențiale.

Demonstrația teoremei 8.6.7. Vom începe cu un calcul ajutător, estimarea lui

$$\Phi(x_n, u_n; h) - \Phi(x_n, y(x_n); h). \quad (8.6.34)$$

Conform teoremei lui Taylor (pentru funcții de mai multe variabile), aplicată celei de-a i -a componente a lui (8.6.34), avem

$$\begin{aligned} \Phi^i(x_n, u_n; h) - \Phi^i(x_n, y(x_n); h) &= \sum_{j=1}^d \Phi^i y^j(x_n, y(x_n); h) [u_n^j - y^j(x_n)] \\ &+ \frac{1}{2} \sum_{j,k=1}^d \Phi^i y^j y^k(x_n, \bar{u}_n; h) [u_n^j - y^j(x_n)] [u_n^k - y^k(x_n)], \end{aligned} \quad (8.6.35)$$

unde \bar{u}_n este pe segmentul ce unește u_n și $y(x_n)$. Utilizând teorema lui Taylor încă o dată, în variabila h , putem scrie

$$\Phi_{y^j}^i(x_n, y(x_n); h) = \Phi_{y^j}^i(x_n, y(x_n); 0) + h \Phi_{y^j h}^i(x_n, y(x_n); \bar{h}),$$

unde $0 < \bar{h} < h$. Deoarece, conform consistenței, $\Phi(x, y; 0) \equiv f(x, y)$ pe $[a, b] \times \mathbb{R}^d$, avem

$$\Phi_{y^j}^i(x, y; 0) = f_{y^j}^i(x, y), \quad x \in [a, b], \quad y \in \mathbb{R}^d,$$

și condiția (1) ne permite să scriem

$$\Phi_{y^j}^i(x_n, y(x_n); h) = f_{y^j}^i(x_n, y(x_n)) + O(h), \quad h \rightarrow 0. \quad (8.6.36)$$

Observând acum că $u_n - y(x_n) = O(h^p)$, în virtutea teoremei 8.6.6 și utilizând (8.6.36) și (8.6.35), obținem conform ipotezei (1),

$$\begin{aligned} \Phi^i(x_n, u_n; h) - \Phi^i(x_n, y(x_n); h) &= \sum_{j=1}^d f_{y^j}^i(x_n, y(x_n)) [u_n^j - y^j(x_n)] + \\ &O(h^{p+1}) + O(h^{2p}). \end{aligned}$$

Dar $O(h^{2p})$ este de ordinul $O(h^{p+1})$, căci $p \geq 1$. Astfel, în notație vectorială,

$$\Phi(x_n, u_n; h) - \Phi(x_n, y(x_n); h) = f_y(x_n, y(x_n)) [u_n - y(x_n)] + O(h^{p+1}). \quad (8.6.37)$$

Acum, pentru a evidenția termenul dominant în eroarea globală, definim funcția grilă $r = \{r_n\}$ prin

$$r = h^{-p}(u - y). \quad (8.6.38)$$

Atunci

$$\begin{aligned} \frac{1}{h}(r_{n+1} - r_n) &= \frac{1}{h} [h^{-p}(u_{n+1} - y(x_{n+1})) - h^{-p}(u_n - y(x_n))] = \\ &= h^p \left[\frac{1}{h}(u_{n+1} - u_n) - \frac{1}{h}(y(x_{n+1}) - y(x_n)) \right] = \\ &= h^{-p} \{ \Phi(x_n, u_n; h) - [\Phi(x_n, y(x_n); h) - T(x_n, y(x_n); h)] \}, \end{aligned}$$

unde am utilizat (8.6.30) și relația (8.6.10) pentru eroarea de trunchiere T . De aceea, exprimând T cu ajutorul funcției eroare principală τ , obținem

$$\begin{aligned} \frac{1}{h}(r_{n+1} - r_n) &= h^{-p} [\Phi(x_n, u_n; h) - \Phi(x_n, y(x_n); h) + \tau(x_n, y(x_n)) h^p \\ &\quad + O(h^{p+1})] \end{aligned}$$

Pentru primii doi termeni din paranteză utilizăm (8.6.37) și definiția lui r (8.6.38) pentru a obține

$$\begin{aligned} \frac{1}{h}(r_{n+1} - r_n) &= f_y(x_n, y(x_n)) r_n + \tau(x_n, y(x_n)) + O(h), \quad n = \overline{0, N-1} \\ r_0 &= 0. \end{aligned} \quad (8.6.39)$$

Acum punând

$$g(x, y) = f_y(x, y(x))y + \tau(x, y(x)) \quad (8.6.40)$$

interpretăm (8.6.39) scriind

$$\left(R_h^{Euler, g} \right)_n = \varepsilon_n, \quad n = \overline{0, N-1}, \quad \varepsilon_n = O(h),$$

unde $R_h^{Euler, g}$ este operatorul rezidual discret (8.6.7) corespunzător metodei lui Euler aplicată lui $e' = g(x, e)$, $e(a) = 0$. Deoarece metoda lui Euler este stabilă pe $[a, b]$ și g fiind liniară în y satisface o condiție Lipschitz uniformă, avem conform inegalității de stabilitate (8.6.11)

$$\|r - e\|_\infty = O(h),$$

și conform lui (8.6.38)

$$\|u - y - h^p e\|_\infty = O(h^{p+1}),$$

așa cum trebuia arătat. \square

8.7. Monitorizarea erorilor și controlul pasului

Vom încerca să realizăm monitorizarea erorilor globale, cel puțin asimptotic, implementând rezultatul din teorema 8.6.7. Aceasta necesită evaluarea matricei jacobiene $f_y(x, y)$ de-a lungul sau în apropierea traiectoriei soluției; acest lucru este natural, deoarece $f_y(x, y)$ guvernează, într-o primă aproximare, efectul perturbațiilor prin ecuația diferențială variațională (8.6.32). În această ecuație tonul este dat de funcția de eroare principală, evaluată de-a lungul traiectoriei, așa că estimarea erorii locale de trunchiere (mai exact a funcției de eroare principală) este de asemenea necesară în această abordare. Pentru simplitate vom presupune că grila are o lungime constantă.

8.7.1. Estimarea erorii globale

Ideea estimării este de a integra „ecuația variațională” (8.6.32) împreună cu ecuația principală (8.6.31). Deoarece avem nevoie de $e(x)$ în (8.6.31) numai în limita unei erori de $O(h)$ (orice termen de eroare $O(h)$ din $e(x_n)$, înmulțit cu h^p , fiind absorbit de termenul $O(h^{p-1})$), putem utiliza în acest scop metoda lui Euler, care va furniza aproximația dorită $v_n \approx e(x_n)$.

Teorema 8.7.1 *Presupunem că*

- (1) $\Phi(x, y; h) \in C^2([a, b] \times \mathbb{R}^d \times [0, h_0])$;
- (2) Φ este o metodă de ordin $p \geq 1$ ce admite o funcție de eroare principală $\tau(x, y) \in C^1([a, b] \times \mathbb{R}^d)$;
- (3) este disponibilă o estimare $r(x, y; h)$ pentru funcție de eroare principală ce satisface

$$r(x, y; h) = \tau(x, y) + O(h), \quad h \rightarrow 0, \quad (8.7.1)$$

uniform pe $[a, b] \times \mathbb{R}^d$;

- (4) odată cu funcția grilă $u = \{u_n\}$ generăm funcția grilă $v = \{v_n\}$ în modul următor

$$\begin{aligned} x_{n+1} &= x_n + h; \\ u_{n+1} &= u_n + h\Phi(x_n, u_n; h) \\ v_{n+1} &= v_n + h[f_y(x_n, v_n)v_n + r(x_n, u_n; h)] \\ x_0 &= a, \quad u_0 = y_0, \quad v_0 = 0. \end{aligned} \quad (8.7.2)$$

Atunci pentru orice $n = \overline{0, N-1}$,

$$u_n - y(x_n) = v_n h^p + O(h^{p+1}), \quad \text{când } h \rightarrow 0. \quad (8.7.3)$$

Demonstrație. Demonstrația începe cu stabilirea următoarelor estimări

$$f_y(x_n, u_n) = f_y(x_n, y(x_n)) + O(h), \quad (8.7.4)$$

$$r(x_n, u_n; h) = \tau(x_n, y(x_n)) + O(h). \quad (8.7.5)$$

Din ipoteza (1) observăm pe baza consistenței $f(x, y) = \Phi(x, y; 0)$ că $f(x, y) \in C^2([a, b] \times \mathbb{R}^d)$. Ținând cont de teorema 8.6.6, avem $u_n = y(x_n) + O(h^p)$, și deci

$$f_y(x_n, u_n) = f_y(x_n, y_n) + O(h^p),$$

relație ce implică (8.7.4), deoarece $p \geq 1$. În continuare, deoarece $\tau(x, y) \in C^1([a, b] \times \mathbb{R}^d)$, conform ipotezei (2) avem

$$\begin{aligned} \tau(x_n, u_n) &= \tau(x_n, y(x_n)) + \tau_y(x_n, \bar{u}_n)(u_n - y(x_n)) \\ &= \tau(x_n, y(x_n)) + O(h^p) \end{aligned}$$

și aplicând apoi ipoteza (3) obținem

$$r(x_n, u_n; h) = \tau(x_n, u_n) + O(h) = \tau(x_n, y(x_n)) + O(h^p) + O(h),$$

din care rezultă imediat (8.7.5).

Fie (a se compara cu (8.6.40))

$$g(x, y) = f_y(x, y(x))y + \tau(x, y(x)). \quad (8.7.6)$$

Ecuația pentru v_{n+1} în (8.7.2) are forma

$$v_{n+1} = v_n + h(A_n v_n + b_n),$$

unde A_n sunt matrice mărginite și b_n sunt vectori mărginiți. Conform lemei 8.6.6, avem mărginirea lui v_n ,

$$v_n = O(1), \quad h \rightarrow 0. \quad (8.7.7)$$

Înlocuind (8.7.4) și (8.7.5) în ecuația lui v_{n+1} și ținând cont de (8.7.7) obținem

$$\begin{aligned} v_{n+1} &= v_n + h[f_y(x_n, y(x_n))v_n + \tau(x_n, y(x_n)) + O(h)] \\ &= v_n + hg(x_n, v_n) + O(h^2). \end{aligned}$$

Astfel, cu notația utilizată în demonstrația teoremei 8.6.7

$$\left(R_h^{Euler, g} v \right)_n = O(h), \quad v_0 = 0.$$

Deoarece metoda lui Euler este stabilă,

$$v_n - e(x_n) = O(h),$$

unde $e(x)$ este, ca mai înainte, soluția ecuației

$$\begin{aligned} e' &= g(x, e) \\ e(a) &= 0 \end{aligned}$$

Deci, conform lui (8.6.33)

$$u_n - y(x_n) = e(x_n)h^p + O(h^{p+1}).$$

□

8.7.2. Estimarea erorii de trunchiere

Pentru a aplica teorema 8.7.1 avem nevoie de estimări $r(x, y; h)$ ale funcției de eroare principală $\tau(x, y)$ care să aibă precizia $O(h)$. Vom descrie două dintre ele, în ordinea crescătoare a eficienței.

Extrapolare Richardson la zero

Aceasta funcționează pentru orice metodă cu un pas Φ , dar de obicei este considerată prea costisitoare. Dacă Φ are ordinul p , procedura este următoarea

$$\begin{aligned} y_h &= y + h\Phi(x, y; h), \\ y_{h/2} &= y + \frac{1}{2}h\Phi\left(x, y; \frac{1}{2}h\right), \\ y_h^* &= y_{h/2} + \frac{1}{2}h\Phi\left(x + \frac{1}{2}h, y_{h/2}; \frac{1}{2}h\right), \\ r(x, y; h) &= \frac{1}{1 - 2^{-p}} \frac{1}{h^{p+1}} (y_h - y_h^*). \end{aligned} \tag{8.7.8}$$

De notat că y_h^* este rezultatul aplicării lui Φ pentru doi pași consecutivi, fiecare de lungime $h/2$, pe câtă vreme y_h este rezultatul aplicării lui Φ pentru un pas de lungime h .

Să verificăm acum că $r(x, y; h)$ dat grupul de formule (8.7.8) este o estimare acceptabilă. Pentru aceasta trebuie să presupunem că $\tau(x, y) \in C^1([a, b] \times \mathbb{R}^d)$. Conform lui (8.3.4) și (8.3.8), utilizând soluția de referință $u(t)$ ce trece prin (x, y) avem

$$\Phi(x, y; h) = \frac{1}{h}[u(x+h) - u(x)] + \tau(x, y)h^p + O(h^{p+1}). \tag{8.7.9}$$

Mai mult,

$$\begin{aligned} \frac{1}{h}(y_h - y_h^*) &= \frac{1}{h}(y_h - y_{h/2}) + \Phi(x, y; h) - \frac{1}{2}h\Phi\left(x + \frac{1}{2}h, y_{h/2}; \frac{1}{2}h\right) \\ &= \Phi(x, y; h) - \frac{1}{2}\Phi\left(x, y; \frac{1}{2}h\right) - \frac{1}{2}h\Phi\left(x + \frac{1}{2}h, y_{h/2}; \frac{1}{2}h\right). \end{aligned}$$

Aplicând (8.7.9) fiecărui termen din membrul drept găsim

$$\begin{aligned} \frac{1}{h}(y_h - y_h^*) &= \frac{1}{h}[u(x+h) - u(x)] + \tau(x, y)h^p + O(h^{p+1}) \\ &\quad - \frac{1}{2} \frac{1}{h/2} \left[u\left(x + \frac{1}{2}h\right) - u(x) \right] - \frac{1}{2} \tau(x, y) \left(\frac{1}{2}h^p \right) + O(h^{p+1}) \\ &\quad - \frac{1}{2} \frac{1}{h/2} \left[u(x+h) - u\left(x + \frac{1}{2}h\right) \right] - \frac{1}{2} \tau\left(x + \frac{1}{2}h, y + O(h)\right) \left(\frac{1}{2}h^p \right) \\ &\quad + O(h^{p+1}) = \tau(x, y)(1 - 2^{-p})h^p + O(h^{p+1}). \end{aligned}$$

În consecință

$$\frac{1}{1 - 2^{-p}} \frac{1}{h}(y_h - y_h^*) = \tau(x, y)h^p + O(h^{p+1}), \quad (8.7.10)$$

așa cum s-a dorit.

Scăzând (8.7.10) din (8.7.9), rezultă incidental că

$$\Phi^*(x, y; h) := \Phi(x, y; h) - \frac{1}{1 - 2^{-p}} \frac{1}{h}(y_h - y_h^*) \quad (8.7.11)$$

definește o metodă cu un pas de ordin $p + 1$.

Procedura (8.7.8) este costisitoare. Pentru un proces Runge-Kutta de ordinul 4 sunt necesare în total 11 evaluări ale lui f pe pas, aproape de trei ori mai mult decât pentru un pas Runge-Kutta. De aceea, extrapolarea Richardson este utilizată numai după fiecare doi pași ai lui Φ , adică se continuă în conformitate cu formulele

$$\begin{aligned} y_h &= y + h\Phi(x, y; h), & (8.7.12) \\ y_{2h}^* &= y_h + h\Phi(x + h, y_h; h) \\ y_{2h} &= y + 2h\Phi(x, y; 2h). \end{aligned}$$

Atunci (8.7.10) ne dă

$$\frac{1}{2(2^p - 1)} \frac{1}{h^{p+1}}(y_{2h} - y_{2h}^*) = \tau(x, y) + O(h), \quad (8.7.13)$$

așa că expresia din membrul drept este un estimator acceptabil al lui $r(x, y; h)$. Dacă cei doi pași din (8.7.12) conduc la o precizie acceptabilă (a se vedea subsecțiunea 8.7.3), atunci pentru un proces Runge-Kutta de ordinul 4 procedura necesită numai trei evaluări adiționale ale lui f , deoarece y_h și y_{2h}^* trebuie calculat oricum. Vom vedea că există scheme mai eficiente.

Metode scufundate (imbricate)

Ideea de bază a aceste abordări este următoarea: se consideră o metodă Φ de ordinul p și o metodă Φ^* de ordinul $p^* = p + 1$ și se definește

$$r(x, y; h) = \frac{1}{h^p} [\Phi(x, y; h) - \Phi^*(x, y; h)]. \quad (8.7.14)$$

Acesta este un estimator acceptabil, așa cum rezultă scăzând relațiile

$$\begin{aligned} \Phi(x, y; h) - \frac{1}{h} [u(x+h) - u(x)] &= \tau(x, y)h^p + O(h^{p+1}) \\ \Phi^*(x, y; h) - \frac{1}{h} [u(x+h) - u(x)] &= O(h^{p+1}) \end{aligned}$$

și împărțind rezultatul cu h^p .

Cheia problemei este de a face această procedură eficientă. Urmând o idee a lui Fehlberg, putem încerca să facem aceasta incluzând un proces Runge-Kutta de ordinul p în altul de ordin $p + 1$. Mai concret, fie Φ o metodă Runge-Kutta explicită în r stadii

$$\begin{aligned} K_1(x, y) &= f(x, y) \\ K_s(x, y; h) &= f\left(x + \mu_s h; y + h \sum_{j=1}^{s-1} \lambda_{sj} K_j\right), \quad s = 2, 3, \dots, r \\ \Phi(x, y; h) &= \sum_{s=1}^r \alpha_s K_s \end{aligned}$$

Atunci pentru Φ^* alegem un proces similar în r^* -stadii, cu $r^* > r$, astfel încât

$$\mu_s^* = \mu_s, \quad \lambda_{sj}^* = \lambda_{sj}, \quad \text{pentru } s = 2, 3, \dots, r.$$

Estimarea (8.7.14) costă atunci din $r^* - r$ evaluări suplimentare ale lui f . Dacă $r^* = r + 1$ putem încă să mai facem economii de evaluări suplimentare, selectând (dacă este posibil)

$$\mu_r^* = 1, \quad \lambda_{rj}^* = \alpha_j \text{ pentru } j = \overline{1, r^* - 1} \quad (r^* = r + 1) \quad (8.7.15)$$

Atunci, într-adevăr, K_r^* va fi identic cu K_1 pentru pasul următor.

Perechi de astfel de formule Runge-Kutta imbricate $(p, p + 1)$ au fost dezvoltate la sfârșitul anilor '60 de E. Fehlberg[16, 17]. Este un grad considerabil de libertate în alegerea acestor parametri. Alegerile lui Fehlberg au fost ghidate de încercarea de a reduce mărimea coeficienților tuturor derivatelor parțiale care intervin în funcția de eroare principală $\tau(x, y)$ a lui Φ . El a reușit să obțină pentru parametrii p, r, r^* valorile date în tabela 8.2

Pentru procesul de ordinul 3 (și numai pentru acesta) se pot alege parametrii astfel ca să aibă loc și (8.7.15).

p	3	4	5	6	7	8
r	4	5	6	8	11	15
r^*	5	6	8	10	13	17

Tabela 8.2: Formule Runge-Kutta imbricate

8.7.3. Controlul pasului

Orice estimare $r(x, y; h)$ a funcției de eroare principală $\tau(x, y)$ implică o estimare

$$h^p r(x, y; h) = T(x, y; h) + O(h^{p+1}) \quad (8.7.16)$$

a erorii locale de trunchiere, care poate fi utilizată pentru a monitoriza eroarea de trunchiere în timpul procesului de integrare. Totuși, trebuie avut în vedere faptul că eroarea locală de trunchiere este chiar diferită de eroarea globală, eroare pe care vrem de fapt să o controlăm. Pentru a obține o mai bună cunoaștere a relației dintre aceste două erori reamintim teorema următoare, care cuantifică continuitatea soluției problemei Cauchy în raport cu valorile inițiale.

Teorema 8.7.2 *Fie $f(x, y)$ continuă în $x \in [a, b]$ și care satisface o condiție Lipschitz cu constanta L , uniform pe $[a, b] \times \mathbb{R}$, adică*

$$\|f(x, y) - f(x, y^*)\| \leq L \|y - y^*\|.$$

Atunci problema Cauchy

$$\begin{aligned} \frac{dy}{dx} &= f(x, y), & x &\in [a, b], \\ y(c) &= y_c \end{aligned} \quad (8.7.17)$$

are o soluție unică pentru orice $c \in [a, b]$ și orice $y_c \in \mathbb{R}^d$. Fie $y(x, s)$ și $y(x; s^)$ soluțiile lui (8.7.17) ce corespund lui $y_c = s$ și respectiv $y_c = s^*$. Atunci, pentru orice normă vectorială $\|\cdot\|$,*

$$\|y(x; s) - y(x; s^*)\| \leq e^{L|x-c|} \|s - s^*\|. \quad (8.7.18)$$

Rezolvarea numerică a problemei (8.6.31) printr-o metodă cu un pas (nu neapărat constant) înseamnă în realitate că se urmărește o secvență de „piste ale soluției“ (expresia este din [20]) prin care în fiecare punct al grilei se sare de la o pistă la următoarea cu o cantitate egală cu eroarea de trunchiere în x_n (vezi figura 8.3). Aceasta rezultă din definiția erorii de trunchiere, soluția de referință fiind una din pistele soluției. Mai concret, a n -a pistă, $n = \overline{0, N}$, este dată de soluția problemei Cauchy

$$\begin{aligned} \frac{dv_n}{dx} &= f(x, v_n), \quad x \in [x_n, b], \\ v_n(x_n) &= u_n, \end{aligned} \quad (8.7.19)$$

și

$$u_{n+1} = v(x_{n+1}) + h_n T(x_n, u_n; h_n), \quad n = \overline{0, N-1}. \quad (8.7.20)$$

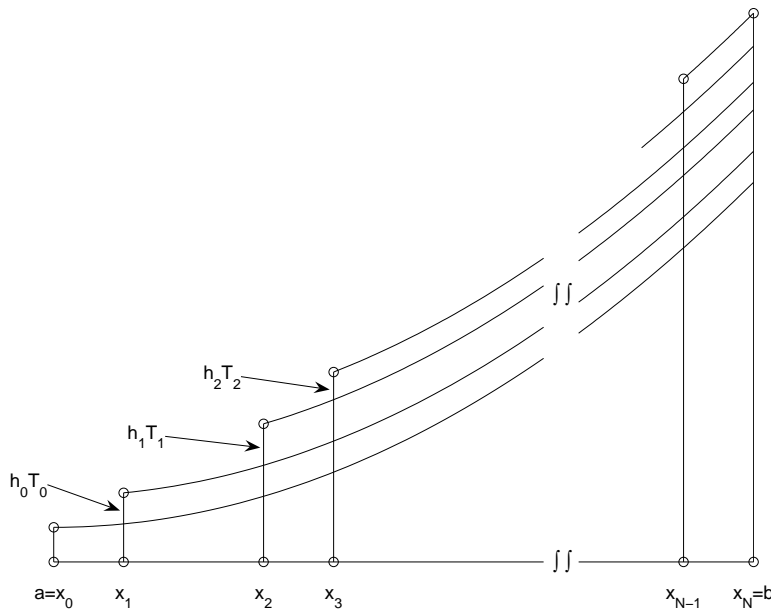


Figura 8.3: Acumularea erorilor într-o metodă cu un pas

Deoarece conform lui (8.7.19) avem $u_{n+1} = v_{n+1}(x_{n+1})$, putem aplica teorema 8.7.2 soluțiilor v_{n+1} și v_n , luând $c = x_{n+1}$, $s = u_{n+1}$, $s^* = u_{n+1} - h_n T(x_n, u_n; h_n)$ (conform lui (8.7.20)) și astfel obținem

$$\|v_{n+1}(x) - v_n(x)\| \leq h_n e^{L|x-x_n|} \|T(x_n, u_n; h_n)\|, \quad n = \overline{0, N-1}. \quad (8.7.21)$$

Acum

$$\sum_{n=0}^{N-1} [v_{n+1}(x) - v_n(x)] = v_N(x) - v_0(x) = v_N(x) - y(x), \quad (8.7.22)$$

și deoarece $v_N(x_N) = u_N$, luând $x = x_N$, obținem din (8.7.21) și (8.7.22) că

$$\begin{aligned} \|u_N - y(x_N)\| &\leq \sum_{n=0}^{N-1} \|v_{n+1}(x_N) - v_n(x_N)\| \\ &\leq \sum_{n=0}^{N-1} h_n e^{L|x_N - x_{n+1}|} \|T(x_n, u_n; h_n)\|. \end{aligned}$$

De aceea, dacă ne asigurăm că

$$\|T(x_n, u_n; h_n)\| \leq \varepsilon_T, \quad n = \overline{0, N-1}, \quad (8.7.23)$$

atunci

$$\|u_N - y(x_N)\| \leq \varepsilon_T \sum_{n=0}^{N-1} (x_{n+1} - x_n) e^{L|x_N - x_{n+1}|}.$$

Interpretând suma din dreapta ca o sumă Riemann pentru o integrală definită, obținem în final aproximarea

$$\|u_N - y(x_N)\| \leq \varepsilon_T \int_a^b e^{L(b-x)} dx = \frac{\varepsilon_T}{L} (e^{L(b-a)} - 1).$$

Astfel, cunoașterea unei estimări pentru L ne va permite să găsim un ε_T

$$\varepsilon_T = \frac{L}{e^{L(b-a)} - 1} \varepsilon, \quad (8.7.24)$$

care să ne garanteze o eroare $\|u_N - y(x_N)\| \leq \varepsilon$. Ceea ce are loc pentru întreaga grilă pe $[a, b]$ are loc, desigur, pentru orice grilă pe subintervalul $[a, x]$, $a \leq x \leq b$. Astfel, în principiu, dându-se precizia dorită ε pentru soluția $y(x)$, putem determina un „nivel de toleranță“ ε_T (din (8.7.24)) și putem asigura precizia dorită păstrând eroarea locală de trunchiere sub limita ε_T (a se compara cu (8.7.23)). De notat că dacă $L \rightarrow 0$, avem $\varepsilon_T \rightarrow \varepsilon/(b-a)$. Această valoare limită pentru ε_T ar fi adecvată pentru o problemă de cuadratură, dar nu pentru o ecuație diferențială veritabilă, unde ε_T , în general, trebuie ales mult mai mic decât eroarea finală ε .

Considerații ca acestea motivează următorul mecanism de control al pasului: fiecare pas de integrare (de la x_n la $x_{n+1} = x_n + h_n$) constă din următoarele părți:

1. Estimăm h_n .
2. Se calculează $u_{n+1} = u_n + h_n \Phi(x_n, u_n; h_n)$ și $r(x_n, u_n; h_n)$.

3. Se testează dacă $h_n^p \|r(x_n, u_n; h_n)\| \leq \varepsilon_T$. Dacă testul este satisfăcut se trece la pasul următor. Dacă nu, repetăm pasul cu un h_n mai mic, până când testul este satisfăcut.

Pentru a estima h_n , presupunem că $n \geq 1$, astfel ca estimatorul din pasul precedent, $r(x_{n-1}, u_{n-1}; h_{n-1})$ (sau cel puțin norma sa) să fie disponibil. Atunci, neglijând termenul $O(h)$,

$$\|\tau(x_{n-1}, u_{n-1})\| \approx \|r(x_{n-1}, u_{n-1}; h_{n-1})\|,$$

și deoarece $\tau(x_n, u_n) \approx \tau(x_{n-1}, u_{n-1})$, în plus.

$$\|\tau(x_n, u_n)\| \approx \|r(x_{n-1}, u_{n-1}; h_{n-1})\|.$$

Ceea ce dorim este

$$\|\tau(x_n, u_n)\| h_n^p \approx \theta \varepsilon_T,$$

unde θ este un factor de siguranță (să zicem $\theta = 0.8$). Eliminând $\tau(x_n, u_n)$ găsim

$$h_n \approx \left\{ \frac{\theta \varepsilon_T}{\|r(x_{n-1}, u_{n-1}; h_{n-1})\|} \right\}^{1/p}.$$

De notat că din pasul precedent avem

$$h_{n-1}^p \|r(x_{n-1}, u_{n-1}; h_{n-1})\| \leq \varepsilon_T,$$

așa că

$$h_n \geq \theta^{1/p} h_{n-1}$$

și tendința este de creștere a pasului.

Dacă $n = 0$, procedăm la fel, alegând o valoare inițială $h_0^{(0)}$ a lui h_0 și calculăm $r(x_0, y_0; h_0^{(0)})$ pentru a obține

$$h_0^{(1)} = \left\{ \frac{\theta \varepsilon_T}{r(x_0, y_0; h_0^{(0)})} \right\}^{1/p}.$$

Procesul se poate repeta odată sau de două ori pentru a obține estimarea finală a lui h_0 și $r(x_0, y_0; h_0^{(0)})$.

Pentru o descriere sintetică a metodelor Runge-Kutta cu pas variabil tabela Butcher se completează cu o linie suplimentară care servește la calculul lui Φ^* (și deci a lui $r(x, y; h)$):

μ_1	λ_{11}	λ_{12}	\dots	λ_{1r}	
μ_2	λ_{21}	λ_{22}	\dots	λ_{2r}	
\vdots	\vdots	\vdots	\dots	\vdots	
μ_r	λ_{r1}	λ_{r2}	\dots	λ_{rr}	
	α_1	α_2	\dots	α_r	
	α_1^*	α_2^*	\dots	α_r^*	α_{r+1}^*

Ca exemplu, în tabela 8.3 dăm tabela Butcher pentru o metodă de ordinul 2-3. Pentru deducerea elementelor tablei a se consulta [41, paginile 451–452].

μ_j	λ_{ij}			
0	0			
$\frac{1}{4}$	$\frac{1}{4}$	0		
$\frac{27}{40}$	$-\frac{189}{800}$	$\frac{729}{800}$	0	
1	$\frac{214}{891}$	$\frac{1}{33}$	$\frac{650}{891}$	0
α_i	$\frac{214}{891}$	$\frac{1}{33}$	$\frac{650}{891}$	0
α_i^*	$\frac{533}{2106}$	0	$\frac{800}{1053}$	$-\frac{1}{78}$

Tabela 8.3: O pereche 2-3

Tabela 8.4 este tabela Butcher pentru metoda Bogacki-Shampine [6]. Ea stă la baza rezolvitorului ode23 din MATLAB.

μ_j	λ_{ij}			
0	0			
$\frac{1}{2}$	$\frac{1}{2}$	0		
$\frac{3}{4}$	0	$\frac{3}{4}$	0	
1	$\frac{2}{9}$	$\frac{3}{9}$	$\frac{4}{9}$	0
α_i	$\frac{2}{9}$	$\frac{3}{9}$	$\frac{4}{9}$	0
α_i^*	$\frac{7}{24}$	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{8}$

Tabela 8.4: Tabela Butcher pentru metoda Bogacki-Shampine

Un alt exemplu important este DORPRI5 sau RK5(4)7FM, o pereche cu ordinele 4-5 și cu 7 stadii (tabela 8.5). Aceasta este o pereche foarte eficientă, ea stând la baza rezolvitorului ode45 din MATLAB, dar și a altor rezolvitori importanți.

Algoritmul 8.2 încearcă să dea sugestii pentru implementarea unei metode Runge-Kutta cu pas variabil când se cunoaște tabela Butcher. $ttol$ este produsul dintre tol și factorul de siguranță (0.8 sau 0.9).

Algoritmul 8.2 Fragment de pseudocod ce ilustrează implementarea unei metode RK cu pas variabil

```

done := false;
loop
   $K_{1,:} := f(x, y)$ ;
  for  $i = 2$  to  $s$  do
     $w := y + hK_{:,1:i-1}\lambda_{i,1:i-1}^T$ ;
     $K_{:,i} := f(x + \mu_i h, w)$ ;
  end for
   $\delta := h \max(|(\alpha^* - \alpha)^T K|)$ ; {estimarea erorii}
   $\beta := (\delta/ttol)^{1/(1+p)}$ ; {raport lung. pas}
  if  $\delta < tol$  then
    {acceptare pas}
     $y := y + h(K\alpha^T)$ ; {actualizare  $y$ }
     $x := x + h$ ;
    if done then
      EXIT {terminare și ieșire}
    end if
     $h := h / \max(\beta, 0.1)$ ; {predicție pas următor}
    if  $x + h > x_{end}$  then
       $x := x_{end} - x$ ; {reducere pas la capăt}
      done := true;
    end if
  else
    {respingere pas}
     $h := h / \min(\beta, 10)$ ; {reducere pas}
    if done then
      done := false;
    end if
  end if
end loop

```

μ_j	λ_{ij}						
0	0						
$\frac{1}{5}$	$\frac{1}{5}$	0					
$\frac{3}{10}$	$\frac{3}{40}$	$\frac{9}{40}$	0				
$\frac{4}{5}$	$\frac{44}{45}$	$-\frac{56}{15}$	$\frac{32}{9}$	0			
$\frac{8}{9}$	$\frac{19372}{6561}$	$-\frac{25360}{2187}$	$\frac{64448}{6561}$	$-\frac{212}{729}$	0		
1	$\frac{9017}{3168}$	$-\frac{355}{33}$	$\frac{46732}{5247}$	$\frac{49}{176}$	$-\frac{5103}{18656}$	0	
1	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	0
α_i	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	0
α_i^*	$\frac{5179}{57600}$	0	$\frac{7571}{16695}$	$\frac{393}{640}$	$-\frac{92097}{339200}$	$\frac{187}{2100}$	$\frac{1}{40}$

Tabela 8.5: Perechea inclusă RK5(4)7FM (DORPRI5)

Aproximări în mai multe variabile

Cuprins

9.1. Aproximarea funcțiilor de mai multe variabile pe un domeniu rectangular	260
9.2. Integrarea numerică a funcțiilor de mai multe variabile	267
9.2.1. Considerații de implementare	273

Problema aproximării funcțiilor de mai multe variabile și cea a aproximării integralelor multiple sunt deosebit de importante atât din punct de vedere teoretic cât și practic, ele intervenind atât în probleme matematice abstracte cât și în modele ale diverselor fenomene sau procese din natură și societate.

Vom nota cu \mathbb{P}_m^n mulțimea polinoamelor în n variabile și de grad global cel mult m . Fie $D \in \mathbb{R}^n$, $f : D \rightarrow \mathbb{R}$. Fie de asemenea $F_1f, F_2f, \dots, F_p f$ informații despre f (de regulă valori ale funcției sau ale unor derivate parțiale ale acestora). Se pune problema ca pe baza informațiilor $F_k, k = \overline{1, p}$, să se determine o funcție F , astfel încât $f \approx F$, într-un sens precizat, în domeniul D .

De exemplu, dacă $P_i \in D, i = \overline{1, s}$, iar $F_i f = f(P_i), i = \overline{1, s}$ se ajunge la o problemă de interpolare polinomială.

Dacă $F_i f = f(P_i), i = \overline{1, s}$, iar $f \in \mathbb{P}_m^n$ se determină astfel încât

$$\sum_{i=1}^s [F(p_i) - f(P_i)]^2 \rightarrow \min$$

se obține o problemă de aproximare discretă în sensul celor mai mici pătrate.

9.1. Aproximarea funcțiilor de mai multe variabile pe un domeniu rectangular

Ne vom ocupa în continuare de aproximarea funcțiilor definite pe un domeniu rectangular și de extinderea unor procedee unidimensionale.

Fie $D = \prod_{k=1}^n [a_k, b_k] \subset \mathbb{R}^n$ și \mathcal{F}_n o mulțime de funcții de n variabile independente definite pe D . Notăm prin $P_i : \mathcal{F}_n \rightarrow \mathcal{G}_i$, $i = \overline{1, n}$ o mulțime de n proiectori ce acționează asupra funcției $f \in \mathcal{F}_n$, fiecare în raport cu variabila x_i , iar prin $R_i = I - P_i$, unde I este operatorul identic, operatorul rest corespunzător. Operatorul $R_i : \mathcal{F}_n \rightarrow \mathcal{H}_i$ este de asemenea projector.

Cu ajutorul operației de *produs tensorial* (compunere) se pot construi operatorii $P_i P_j$, $i, j = \overline{1, n}$, care, în caz că oricare doi comută sunt de asemenea proiectori, deoarece $(P_i P_j)^2 = P_i P_j$. Deoarece suma a doi proiectori $P_i + P_j$ nu este projector ($(P_i + P_j)^2 \neq P_i + P_j$), se lucrează cu *suma booleană* $P_i \oplus P_j = P_i + P_j - P_i P_j$, care este projector.

Datorită proprietății de asociativitate a produsului și a sumei booleene, aceste operații se pot extinde la trei sau mai mulți operatori. Putem defini proiectorii

$$P := P_1 P_2 \dots P_n,$$

$$S := P_1 \oplus P_2 \oplus \dots \oplus P_n = \sum_{i=1}^n \tilde{P}_i,$$

cu

$$\tilde{P}_1 = P_1, \quad \tilde{P}_i = P_i - \left(\sum_{k=1}^{i-1} \tilde{P}_k \right) P_i, \quad i = \overline{2, n},$$

adică produsul și suma booleană a celor n proiectori considerați. De exemplu, pentru trei proiectori avem

$$P_1 P_2 P_3 = P_1 (P_2 P_3),$$

$$P_1 \oplus P_2 \oplus P_3 = P_1 + P_2 + P_3 - P_1 P_2 - P_1 P_3 - P_2 P_3 + P_1 P_2 P_3.$$

Fie \mathcal{P} mulțimea formată din proiectorii P_1, \dots, P_n și toți proiectorii obținuți din aceștia cu ajutorul operațiilor de produs și sumă booleană.

Imaginea Qf a unei funcții $f \in \mathcal{F}_n$ prin orice element $Q \in \mathcal{P}$ poate fi considerată o aproximantă a lui f . De exemplu, dacă P_i este operatorul de interpolare Lagrange $L_{m_i}^{x_i}$ relativ la nodurile $x_{i_0}, \dots, x_{i_{m_i}}$ ce acționează asupra variabilei x_i atunci

$$(L_{m_i}^{x_i} f) = \sum_{k=0}^{m_i} \ell_k(x_i) f(x_1, \dots, x_{i-1}, x_{i_k}, x_{i+1}, \dots, x_{i_{m_i}}),$$

unde

$$\ell_k(x_i) = \frac{(x_i - x_{i_0}) \dots (x_i - x_{i_{k-1}})(x_i - x_{i_{k+1}}) \dots (x_i - x_{i_{m_i}})}{(x_{i_k} - x_{i_0}) \dots (x_{i_k} - x_{i_{k-1}})(x_{i_k} - x_{i_{k+1}}) \dots (x_{i_k} - x_{i_{m_i}})},$$

este o aproximantă a funcției f . Mai mult, putem scrie una din expresiile termenului rest din această aproximare. Astfel, utilizând forma cu diferențe divizate pentru rest, avem

$$(R_{m_i}^{x_i} f) = u(x_i)[x_i, x_{i_0}, \dots, x_{i_{m_i}}, f(x_1, \dots, x_{i-1}, \cdot, x_{i+1}, \dots, x_{i_{m_i}})].$$

Operatorul de diferență divizată acționează asupra funcției f în raport cu variabila x_i , iar $u(x_i) = (x_i - x_{i_0}) \dots (x_i - x_{i_{m_i}})$.

Pentru a compara aproximantele generate de elementele lui \mathcal{P} se definește pe \mathcal{P} o relație de ordine parțială.

Definiția 9.1.1 Pentru $X, Y \in \mathcal{P}$, $X \leq Y \Leftrightarrow XY = X$.

Dacă produsul a oricare doi proiectori este comutativ, adică $P_i P_j = P_j P_i$, $i \neq j$, atunci (\mathcal{P}, \leq) este o latice distributivă. Din relațiile

$$P = P_1 \dots P_n \leq X, \quad X \in \mathcal{P}$$

și

$$Y \leq S = P_1 \oplus \dots \oplus P_n, \quad Y \in \mathcal{P},$$

rezultă că P este element minimal, iar S element maximal al laticei (\mathcal{P}, \leq) .

Definiția 9.1.2 Fie $f \in \mathcal{F}_n$. Aproximanta Pf a funcției f se numește algebric minimală, iar aproximanta Sf se numește algebric maximală.

Pentru exemplificare, fie $n = 2$ și $P_1 = L_r^x$, $P_2 = L_s^y$ operatorii de interpolare Lagrange unidimensionali relativ la nodurile x_0, \dots, x_r , respectiv y_0, \dots, y_s . Dacă $f \in \mathcal{F}_2$, atunci aproximanta algebric minimală este

$$(L_r^x L_s^y f)(x, y) = \sum_{i=0}^r \sum_{j=0}^s \ell_i(x) \tilde{\ell}_j(y) f(x_i, y_j), \tag{9.1.1}$$

iar cea algebric maximală este

$$\begin{aligned} (L_r^x \oplus L_s^y f)(x, y) &= \sum_{i=0}^r \ell_i(x) f(x_i, y) + \sum_{j=0}^s \tilde{\ell}_j(y) f(x, y_j) \\ &\quad - \sum_{i=0}^r \sum_{j=0}^s \ell_i(x) \tilde{\ell}_j(y) f(x_i, y_j), \end{aligned} \tag{9.1.2}$$

unde $\ell_i(x)$ și $\tilde{\ell}_j(y)$ sunt polinoamele fundamentale Lagrange corespunzătoare.

Prin verificare directă pentru produs se obține

$$(L_r^x L_s^y f)(x_i, y_j) = f(x_i, y_j), \quad i = \overline{0, r}, j = \overline{0, s},$$

iar pentru suma booleană

$$\begin{aligned} (L_r^x \oplus L_s^y f)(x_i, y) &= f(x_i, y), & i = \overline{0, r}, y \in [a_2, b_2], \\ (L_r^x \oplus L_s^y f)(x, y_j) &= f(x, y_j), & j = \overline{0, s}, x \in [a_1, b_1]. \end{aligned}$$

Să considerăm grila $\{(x_i, y_j) : i = \overline{0, r}, j = \overline{0, s}\}$. Interpolarea algebric minimală

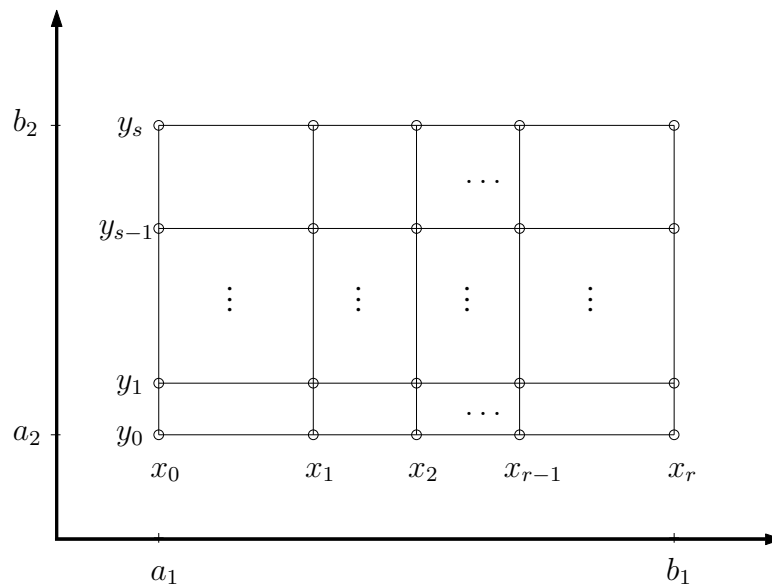


Figura 9.1: Interpretarea geometrică a interpolanților produs tensorial și sumă booleană

este o interpolare discretă (punctuală); ea reproduce valorile funcției numai pe noduri. Interpolarea algebric maximală este o *interpolare transformată* sau *interpolare blending*; funcția interpolatoare reproduce valorile lui f pe segmentele $\{(x_i, y) \in \mathbb{R}^2 : a_2 \leq y \leq b_2\}, i = \overline{0, r}$ și $\{(x, y_j) \in \mathbb{R}^2 : a_1 \leq x \leq b_1\}, j = \overline{0, s}$ (figura 9.1).

Revenind la operatorii rest $R_i = I - P_i, i = \overline{1, n}$ se observă că mulțimea formată din R_i și operatorii obținuți din aceștia prin produs tensorial și sumă booleană, formează și ei în raport cu relația „ \leq ” (dată în definiția 9.1.1) tot o latice (dacă produsul este comutativ). În această latice $R_1 R_2 \dots R_n$ este element minimal, iar $R_1 \oplus R_2 \oplus \dots \oplus R_n$ este element maximal. Avem de asemenea următoarele descompuneri pentru operatorul identic

$$\begin{aligned} I &= P_1 \dots P_n + R_1 \oplus \dots \oplus R_n \\ I &= P_1 \oplus \dots \oplus P_n + R_1 \dots R_n, \end{aligned} \tag{9.1.3}$$

numite *descompunere minimală* și respectiv *maximală*. Dacă $f \in \mathcal{F}_n$, pe baza descompunerilor anterioare avem formulele de aproximare

$$f = P_1 \dots P_n f + R_1 \oplus \dots \oplus R_n f \tag{9.1.4}$$

$$f = P_1 \oplus \dots \oplus P_n f + R_1 \dots R_n f, \tag{9.1.5}$$

numite *formulă de aproximare algebric minimală* și respectiv *formulă de aproximare algebric maximală*.

Pentru exemplificare să considerăm din nou cazul bidimensional și operatorii de interpolare Lagrange L_r^x și L_s^y . Operatorii rest corespunzători vor fi $R_r^x = I - L_r^x$ și $R_s^y = I - L_s^y$. Pentru formula algebric minimală se va obține

$$f = L_r^x L_s^y + R_r^x \oplus R_s^y, \tag{9.1.6}$$

iar pentru cea algebric maximală

$$f = L_r^x \oplus L_s^y + R_r^x R_s^y. \tag{9.1.7}$$

Cunoscând expresiile termenului rest al formulei de interpolare Lagrange și presupunând că există toate derivatele parțiale ale lui f necesare și că acestea satisfac condițiile de continuitate cerute, obținem

$$\begin{aligned} (R_r^x \oplus R_s^y f)(x, y) &= \frac{u_r(x)}{(r+1)!} \frac{\partial^{r+1}}{\partial x^{r+1}} f(\xi_1, y) + \frac{u_s(y)}{(s+1)!} \frac{\partial^{s+1}}{\partial y^{s+1}} f(x, \eta_1) \\ &\quad - \frac{u_r(x)u_s(y)}{(r+1)!(s+1)!} \frac{\partial^{r+s+2}}{\partial x^{r+1} \partial y^{s+1}} f(\xi_2, \eta_2), \end{aligned}$$

cu $\xi_1, \xi_2 \in [\alpha_1, \alpha_2]$, $\eta_1, \eta_2 \in [\beta_1, \beta_2]$,

$$(R_r^x R_s^y f)(x, y) = \frac{u_r(x)u_s(y)}{(r+1)!(s+1)!} \frac{\partial^{r+s+2}}{\partial x^{r+1} \partial y^{s+1}} f(\xi_3, \eta_3),$$

cu $\xi_3 \in [\alpha_1, \alpha_2]$, $\eta_3 \in [\beta_1, \beta_2]$, $\alpha_1 = \min\{x, x_0, \dots, x_r\}$, $\alpha_2 = \max\{x, x_0, \dots, x_r\}$, $\beta_1 = \min\{y, y_0, \dots, y_s\}$, $\beta_2 = \max\{y, y_0, \dots, y_s\}$. De notat că termenul rest se poate exprima și sub formă integrală și cu ajutorul diferențelor divizate.

Exemplul 9.1.3. Dacă luăm $r = s = 1$ și nodurile 0 și 1 după fiecare coordonată, operatorii Lagrange unidimensionali sunt:

$$\begin{aligned} L_1^x(x, y) &= (1-x)f(0, y) + xf(1, y) \\ L_1^y(x, y) &= (1-y)f(x, 0) + yf(x, 1), \end{aligned}$$

iar pentru operatorii produs tensorial și sumă booleană se obțin expresiile

$$(L_1^x L_1^y)(x, y) = (1-x)(1-y)f(0, 0) + (1-x)yf(0, 1) + x(1-y)f(1, 0) + xyf(1, 1)$$

și respectiv

$$(L_1^x \oplus L_1^y)(x, y) = (1-x)f(0, y) + xf(1, y) + (1-y)f(x, 0) + yf(x, 1) - (1-x)(1-y)f(0, 0) - (1-x)yf(0, 1) - x(1-y)f(1, 0) + xyf(1, 1).$$

Prima aproximantă reproduce valorile funcției doar în vârfurile pătratului $[0, 1] \times [0, 1]$, pe când ce-a de-a doua reproduce valorile funcției pe frontiera pătratului. Dacă există derivatele parțiale până la ordinele necesare, expresiile corespunzătoare pentru rest sunt

$$(R_P f)(x, y) = (R_1^x R_1^y)(x, y) = \frac{x(x-1)}{2} \frac{\partial^2}{\partial x^2} f(\xi_1, y) + \frac{y(y-1)}{2} \frac{\partial^2}{\partial y^2} f(x, \eta_1) - \frac{x(x-1)}{2} \frac{y(y-1)}{2} \frac{\partial^4}{\partial x^2 \partial y^2} f(\xi_2, \eta_2).$$

și respectiv

$$(R_S f)(x, y) = (R_1^x \oplus R_1^y)(x, y) = \frac{x(x-1)}{2} \frac{y(y-1)}{2} \frac{\partial^4}{\partial x^2 \partial y^2} f(\xi, \eta),$$

unde $\xi, \eta, \xi_1, \eta_1, \xi_2, \eta_2 \in (0, 1)$. Dacă f este de clasă $C^{2,2}$, au loc delimitările

$$\|R_P f\|_\infty \leq \frac{1}{8} \left[\left\| \frac{\partial^2}{\partial x^2} f \right\|_\infty + \left\| \frac{\partial^2}{\partial y^2} f \right\|_\infty + \frac{1}{8} \left\| \frac{\partial^4}{\partial x^2 \partial y^2} f \right\|_\infty \right],$$

$$\|R_S f\|_\infty \leq \frac{1}{64} \left\| \frac{\partial^4}{\partial x^2 \partial y^2} f \right\|_\infty. \quad \diamond$$

Exemplul 9.1.4. Să considerăm funcția $f : [-2, 2] \times [-2, 2] \rightarrow \mathbb{R}$, $f(x, y) = xe^{-x^2-y^2}$. Graficul ei este dat în figura 9.2. Graficele aproximantelor algebric minimală, respectiv algebric maximală de tip Lagrange și ale resturilor corespunzătoare apar în figura 9.3. S-au considerat cinci noduri echidistante după fiecare coordonată, $x_k, y_k = -2 + k$, $k = \overline{0, 4}$. \diamond

În concluzie, formula de aproximare generată de produsul operatorilor P_1, \dots, P_n este algebric minimală, iar formula generată de suma booleană a acestor operatori este algebric maximală, proprietăți datorate calității operatorilor produs și sumă booleană de a fi element minimal, respectiv maximal în laticea (\mathcal{P}, \leq) .

Se va da în continuare și o altă interpretare a celor două formule de aproximare extremală, având în vedere ordinul de aproximare al operatorilor unidimensionali P_i ,

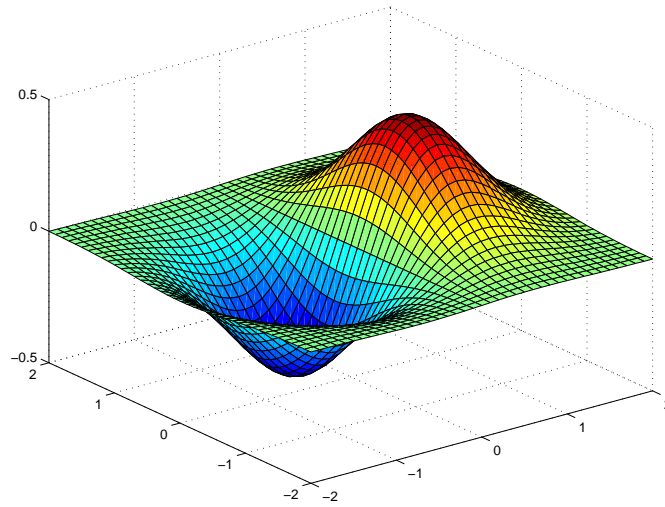
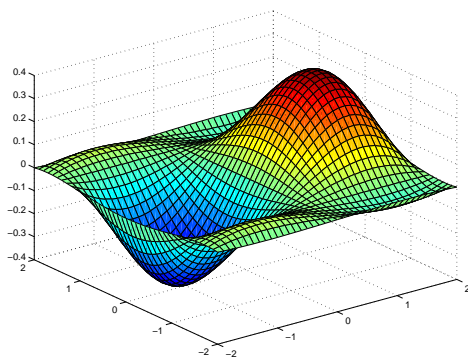
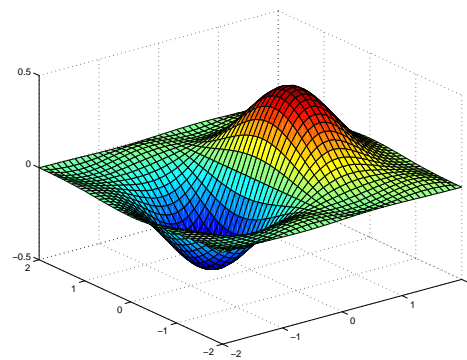


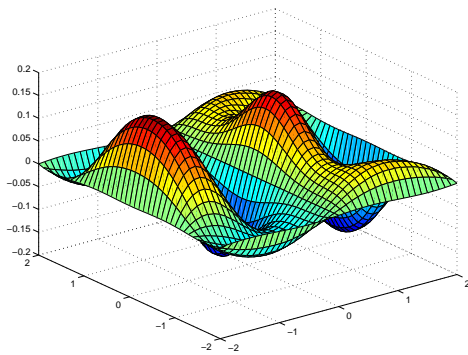
Figura 9.2: Graficul funcției $f : [-2, 2] \times [-2, 2] \rightarrow \mathbb{R}$, $f(x, y) = xe^{-x^2-y^2}$ din exemplul 9.1.4



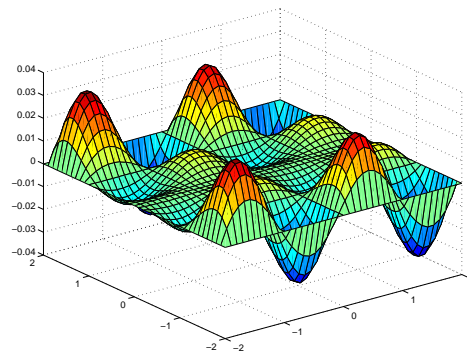
(a) produs tensorial



(b) sumă booleană



(c) rest produs tensorial



(d) rest sumă booleană

Figura 9.3: Aproximarea algebric minimală (figura 9.3(a)), cea maximală (figura 9.3(b)) și resturile corespunzătoare (figura 9.3(c) și respectiv 9.3(d))

$i = \overline{1, n}$ de la care se pornește. Această interpretare va permite extinderea procedeeilor de mai sus și la alți operatori care nu sunt proiectori.

Dacă notăm prin m_i ordinul de aproximare al operatorului P_i , $i = \overline{1, n}$, atunci din (9.1.4) rezultă că ordinul de aproximare al operatorului produs tensorial este

$$\text{ord}(P) = \min\{\text{ord}(P_1), \dots, \text{ord}(P_n)\}, \quad (9.1.8)$$

iar din (9.1.5) rezultă că ordinul de aproximare al operatorului sumă booleană este

$$\text{ord}(S) = \text{ord}(P_1) + \dots + \text{ord}(P_n). \quad (9.1.9)$$

Pentru două variabile formula de aproximare algebric minimală (9.1.6) are ordinul de aproximare $\min(r, s) + 1$, în timp ce ordinul formulei de interpolare algebric maximale (9.1.7) este $r + s + 2$.

Prin urmare, proprietățile de extremalitate ale celor doi operatori sunt caracterizate și prin ordinul de aproximare, operatorul produs având ordinul de aproximare minim, iar operatorul sumă booleană având ordinul de aproximare maxim în mulțimea \mathcal{P} a tuturor operatorilor generați din P_1, \dots, P_n cu ajutorul celor două operații.

Deși operatorul sumă booleană are un ordin mare de aproximare, în expresia sa

$$Sf = (P_1 + \dots + P_n - P_1P_2 - \dots - P_{n-1}P_n + \dots P_1P_2 \dots P_n)f, \quad (9.1.10)$$

vor interveni termeni care sunt funcții de $n - 1, n - 2, \dots, 1$ variabile. O astfel de formulă este aplicabilă doar atunci când se cunosc valorile funcției f sau ale unor derivate parțiale ale ei pe hiperfețele lui D sau pe secțiuni ale domeniului D paralele cu aceste hiperfețe. Acesta este un dezavantaj major al aproximării blending. El poate fi înlăturat aplicând funcției f pe hipersuprafețele respective ale lui D alți operatori liniari de aproximare. Continuând acest prodeu din aproape în aproape se ajunge după cel mult $n - 1$ etape la o aproximare scalară (numerică), adică la o aproximare ce conține numai valori ale lui f pe puncte ale lui D . Evident că restul unei astfel de formule de aproximare, obținută din formula blending conține mai mulți termeni. Pentru ilustrare, să considerăm cazul unei funcții de două variabile, definită pe $[a_1, b_1] \times [a_2, b_2]$. Considerând operatorii P_1^1, P_2^1 , în care indicele superior indică numărul nivelului (etapei) de aproximare, obținem la prima etapă

$$f = (P_1^1 + P_2^1 - P_1^1P_2^1)f + R_1^1R_2^1f, \quad (9.1.11)$$

unde P_1^1f conține pe x_2 ca variabilă liberă, iar P_2^1 pe x_1 . Aplicând pe al doilea nivel lui P_1f operatorul P_2^2 , iar lui P_2^1f operatorul P_1^2 , obținem formula de aproximare scalară:

$$f = (P_1^1P_2^2 + P_1^2P_2^1 - P_1^1P_2^1)f + (P_2^1R_1^2 + P_1^1R_2^2 + R_1^1R_2^1)f. \quad (9.1.12)$$

Restul acestei formule conține trei termeni. Singura proprietate a tuturor operatorilor folosiți este comutativitatea produsului. Avem mai multe posibilități de alegere a operatorilor de pe al doilea nivel. Ei depind în primul rând de informațiile disponibile asupra

lui f . O cale naturală de alegere este aceea ca ordinul de aproximare al formulei originale (9.1.11) să fie păstrat. O astfel de formulă de aproximare se va numi *formulă consistentă*. De exemplu, alegând în (9.1.12) pe post de P_1^1 și P_2^1 operatorii de interpolare Lagrange L_1^x și L_2^y relativi la nodurile a_1, b_1 și a_2, b_2 , iar în locul lui P_1^2 și P_2^2 operatorii de interpolare Hermite H_3^x și H_4^y relativi la nodurile duble a_1, b_1 și respectiv a_2 triplu și b_2 dublu se obține

$$f = (L_1^x H_4^y + H_3^x L_1^y - L_1^x L_1^y)f + (L_1^y R_3^x + L_1^x R_4^y + R_1^x R_1^y)f. \quad (9.1.13)$$

Cum ordinele de aproximare ale operatorilor L_1 , H_3 și H_4 sunt 2, 4 și 5, rezultă că ordinul formulei de aproximare blending

$$f = L_1^x \oplus L_2^y + R_1^x R_2^y f \quad (9.1.14)$$

este 4, iar ordinul de aproximare al lui (9.1.13) este tot 4, căci

$$\begin{aligned} \|(L_1^y R_3^x + L_1^x R_4^y + R_1^x R_1^y)f\| &\leq \|R_3^x f\| + \|R_4^y f\| + \|R_1^x R_2^y f\| = \\ &= c_1(b_1 - a_1)^4 + c_2(b_2 - a_2)^5 + c_3(b_1 - a_1)^2(b_2 - a_2)^2 \\ &= (c_1 + c_2 h + c_3)h^4, \quad \text{unde } h = b_1 - a_1 = b_2 - a_2. \end{aligned}$$

Deci (9.1.13) este consistentă.

Deoarece ordinul de aproximare al formulei inițiale nu poate fi mărit, este preferabil ca operatorii folosiți în următoarele nivele de aproximare să fie astfel aleși încât termenii din expresia restului formulei finale să aibă același ordin de mărime. Formula numerică astfel obținută se va numi *omogenă*.

Formula (9.1.13) nu este omogenă deoarece $\|R_3^x f\| = O(h^4)$, $\|R_1^x R_1^y\| = O(h^4)$, dar $\|R_4^y f\| = O(h^5)$. Pentru a deveni omogenă trebuie ca operatorul H_4^y folosit pe al doilea nivel să fie înlocuit printr-un operator cu ordinul de aproximare 4, de exemplu H_3^y . Astfel, din formula originală (9.1.14), am dedus următoarea formulă numerică omogenă

$$f = (L_1^x H_3^y + H_3^x L_1^y - L_1^x L_1^y)f + (L_1^y R_3^x + L_1^x R_3^y + R_1^x R_1^y)f.$$

9.2. Integrarea numerică a funcțiilor de mai multe variabile

Fie $D \subseteq \mathbb{R}^n$, funcția $f : D \rightarrow \mathbb{R}$, punctele $P_i \in D$, $i = \overline{0, m}$ și w o funcție pondere nenegativă, definită pe D .

Definiția 9.2.1 *Formula*

$$\int \cdots \int_D w(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n = \sum_{i=0}^m A_i f(P_i) + R_m f \quad (9.2.1)$$

se numește formulă de integrare numerică a funcției f sau formulă de cubatură. Parametrii A_i se numesc coeficienții formulei, punctele P_i se numesc nodurile ei, iar R_m termenul rest.

Problema construirii unei formule de cubatură constă în determinarea coeficienților și a nodurilor ei, folosind condiții corespunzătoare. În funcție de condițiile folosite, formulele de cubatură pot fi clasificate în determinate și nedeterminate, iar cele determinate în formule de tip interpolator, de tip Gauss, optimale, etc.

O metodă eficientă de construire a formulelor de cubatură în cazul în care D este un domeniu rectangular, constă în exprimarea parametrilor acestuia cu ajutorul coeficienților, respectiv a nodurilor unei formule de cuadratură unidimensionale. Pentru simplificare ne vom limita la cazul bidimensional. Fie $D = [a, b] \times [c, d]$, Δ_x diviziunea $a = x_0 < x_1 < \dots < x_m = b$, Δ_y diviziunea $c = y_0 < y_1 < \dots < y_n = d$, $w(x, y) = 1, \forall (x, y) \in D$.

În acest caz formula (9.2.1) devine

$$\int_a^b \int_c^d f(x, y) dx dy = \sum_{i=0}^m \sum_{j=0}^n A_{ij} f(x_i, y_j) + R_{m,n}(f). \quad (9.2.2)$$

Se poate obține o astfel de formulă pornind de la formula de interpolare Lagrange (bidimensională)

$$f = L_m^x L_n^y f + R_m^x \oplus R_n^y f.$$

Dacă $f \in C^{m+1, n+1}(D)$, prin integrare termen cu termen se obține

$$\int_a^b \int_c^d f(x, y) dx dy = \sum_{i=0}^m \sum_{j=0}^n A_i B_j f(x_i, y_j) + R_{mn}(f) \quad (9.2.3)$$

unde

$$A_i = \int_a^b \ell_i(x) dx, \quad B_j = \int_c^d \tilde{\ell}_j(y) dy,$$

iar

$$\begin{aligned} R_{m,n}(f) &= \int_a^b \int_c^d (R_m^x \oplus R_n^y) f(x, y) dx dy \\ &= \frac{1}{(m+1)!} \int_a^b \int_c^d u_m(x) f^{(m+1,0)}(\xi_x, y) dx dy \\ &\quad + \frac{1}{(n+1)!} \int_a^b \int_c^d u_n(y) f^{(0,n+1)}(x, \eta_y) dx dy \\ &\quad - \frac{1}{(m+1)!} \frac{1}{(n+1)!} \int_a^b \int_c^d u_m(x) u_n(y) f^{(m+1,n+1)}(\tilde{\xi}_x, \tilde{\eta}_y) dx dy. \end{aligned}$$

Dacă Δ_x și Δ_y sunt diviziuni uniforme ale intervalului $[a, b]$ și respectiv $[c, d]$, atunci formula de cubatură (9.2.3) se numește de tip *Newton-Cotes*.

Cazuri particulare Pentru $m = n = 1$ se obține formula de cubatură a trapezului.

$$\int_a^b \int_c^d f(x, y) dx dy = \frac{(b-a)(d-c)}{4} [f(a, c) + f(a, d) + f(b, c) + f(b, d)] + R_{11}(f),$$

unde

$$R_{11}(f) = -\frac{(b-a)^3(d-c)}{12} f^{(2,0)}(\xi_1, \eta_1) - \frac{(b-a)(d-c)^3}{12} f^{(0,2)}(\xi_2, \eta_2) - \frac{(b-a)^3(d-c)^3}{144} f^{(2,2)}(\xi_3, \eta_3).$$

Pentru $m = n = 2$ se obține formula de cubatură a lui Simpson.

$$\int_a^b \int_c^d f(x, y) dx dy = \frac{(b-a)(d-c)}{36} \left\{ f(a, c) + f(a, d) + f(b, c) + f(b, d) + 4 \left[f\left(\frac{a+b}{2}, c\right) + f\left(\frac{a+b}{2}, d\right) + f\left(a, \frac{c+d}{2}\right) + f\left(b, \frac{c+d}{2}\right) \right] + 16f\left(\frac{a+b}{2}, \frac{c+d}{2}\right) \right\} + R_{22}(f),$$

unde

$$R_{22}(f) = -\frac{(b-a)^5(d-c)}{2880} f^{(4,0)}(\xi_1, \eta_1) - \frac{(b-a)(d-c)^5}{2880} f^{(0,4)}(\xi_2, \eta_2) - \frac{(b-a)^5(d-c)^5}{2880^2} f^{(4,4)}(\xi_3, \eta_3).$$

Partiționând intervalele $[a, b]$ și $[c, d]$ se pot obține formule de cubatură repetate. Vom ilustra pentru formula lui Simpson. Să presupunem că $[a, b]$ este împărțit în m părți egale, iar $[c, d]$ în n părți egale, obținându-se o diviziune cu mn dreptunghiuri. Vom împărți fiecare dreptunghi în patru părți egale ca în figura 9.4 (vârfurile dreptunghiului sunt indicate prin cercuri negre, iar punctele intermediare cu cercuri mai albe).

Fie

$$h = \frac{b-a}{2m}, \quad k = \frac{d-c}{2n}.$$

Punctele diviziunii vor avea coordonatele

$$\begin{aligned} x_i &= x_0 + ih, & x_0 &= a, & i &= \overline{0, 2m} \\ y_j &= y_0 + jk, & y_0 &= b, & j &= \overline{0, 2n}. \end{aligned}$$

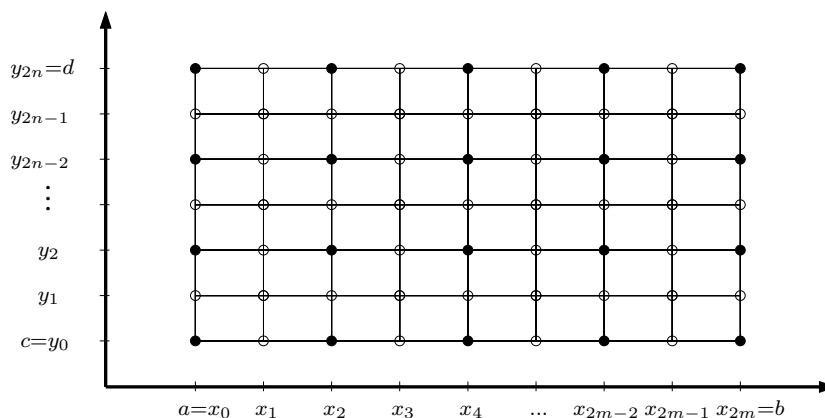


Figura 9.4: Diviziunea pentru formula repetată a lui Simpson

Introducem notația $f(x_i, y_j) = f_{ij}$. Aplicând formula lui Simpson fiecărui dreptunghi al diviziunii avem

$$\int_a^b \int_c^d f(x, y) dx dy = \frac{hk}{9} \sum_{i=0}^m \sum_{j=0}^n [f_{2i,2j} + f_{2i+2,2j} + f_{2i+2,2j+2} + f_{2i,2j+2} + 4(f_{2i+1,2j} + f_{2i+2,2j+1} + f_{2i+1,2j+2} + f_{2i,2j+1}) + 16f_{2i+1,2j+1}] + R_{m,n}(f).$$

Reducând termenii asemenea se obține

$$\int_a^b \int_c^d f(x, y) dx dy = \frac{hk}{9} \sum_{i=0}^m \sum_{j=0}^n \lambda_{i,j} f_{ij} + R_{m,n}(f),$$

unde λ_{ij} sunt dați de matricea

$$\Lambda = \begin{bmatrix} 1 & 4 & 2 & 4 & 2 & \dots & 4 & 2 & 4 & 1 \\ 4 & 16 & 8 & 16 & 8 & \dots & 16 & 8 & 16 & 4 \\ 2 & 8 & 4 & 8 & 4 & \dots & 8 & 4 & 8 & 2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 2 & 8 & 4 & 8 & 4 & \dots & 8 & 4 & 8 & 2 \\ 4 & 16 & 8 & 16 & 8 & \dots & 16 & 8 & 16 & 4 \\ 1 & 4 & 2 & 4 & 2 & \dots & 4 & 2 & 4 & 1 \end{bmatrix}.$$

În [8, 38], pentru $f \in C^{4,4}(D)$ se dă următoarea formă a restului

$$R_{mn}(f) = -\frac{(b-a)(d-c)}{180} [h^4 f^{(4,0)}(\xi_1, \eta_1) + k^4 f^{(0,4)}(\xi_2, \eta_2)],$$

cu $\xi_1, \eta_1, \xi_2, \eta_2 \in D$.

Se pot da formule de cubatură de tip Gauss bidimensionale. De exemplu, dacă $x_i, i = \overline{0, m}$ și $y_j, j = \overline{0, n}$ sunt rădăcinile polinomului Legendre relativ la intervalul $[a, b]$ și respectiv $[c, d]$, se obține formula de cubatură Gauss-Legendre, în care

$$A_i = \frac{[(m+1)!]^4 (b-a)^{2m+3}}{[(2m+2)!]^2 (x_i-a)(b-x_i)[u'(x_i)]^2}, \quad i = \overline{0, m}$$

$$B_j = \frac{[(n+1)!]^4 (b-a)^{2n+3}}{[(2n+2)!]^2 (y_j-c)(d-y_j)[u'(y_j)]^2}, \quad j = \overline{0, n},$$

iar dacă $f \in C^{2m+2, 2n+2}(D)$

$$R_{mn}(f) = (d-c)\lambda_m f^{(2m+2,0)}(\xi_1, \eta_1) + (b-a)\lambda_n f^{(0,2n+2)}(\xi_2, \eta_2) - \lambda_m \lambda_n f^{(2m+2, 2n+2)}(\xi_3, \eta_3),$$

unde

$$\lambda_m = \frac{[(m+1)!]^4 (b-a)^{2m+3}}{[(2m+2)!]^3 (2m+3)}, \quad \lambda_n = \frac{[(n+1)!]^4 (b-a)^{2n+3}}{[(2n+2)!]^3 (2n+3)}.$$

Pentru detalii a se vedea [10].

În cazul $m = n = 0$ se obține o formulă de cubatură cu un singur nod, analoagă formulei dreptunghiului

$$\int_a^b \int_c^d f(x, y) dx dy = (b-a)(d-c) f\left(\frac{a+b}{2}, \frac{c+d}{2}\right) + R_{00}(f),$$

unde

$$R_{00}(f) = \frac{(b-a)^3(d-c)}{24} f^{(2,0)}(\xi_1, \eta_1) + \frac{(b-a)(d-c)^3}{24} f^{(0,2)}(\xi_2, \eta_2) - \frac{(b-a)^3(d-c)^3}{576} f^{(2,2)}(\xi_3, \eta_3).$$

Utilizarea metodelor de aproximare de mai sus nu este limitată la domenii rectangulare. De exemplu, tehnica de la formula de cubatură a lui Simpson se poate modifica pentru a fi aplicabilă la aproximarea unor integrare de forma

$$\int_a^b \int_{c(x)}^{d(x)} f(x, y) dx dy$$

sau

$$\int_c^d \int_{a(y)}^{b(y)} f(x, y) dx dy.$$

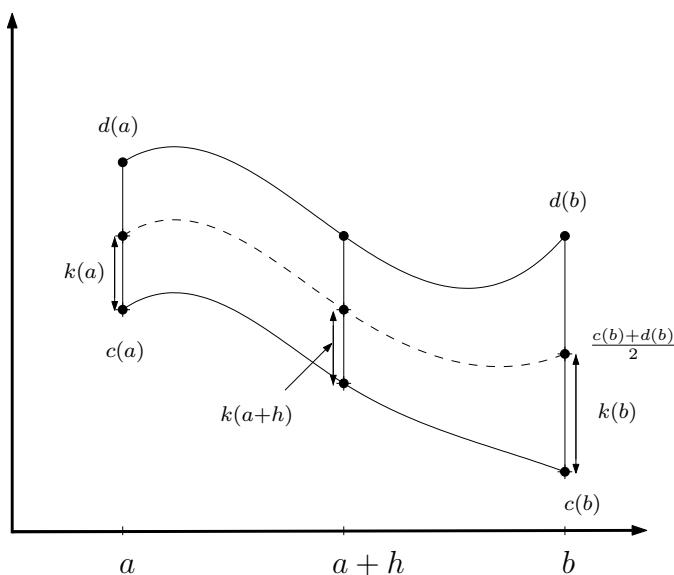


Figura 9.5: Formula lui Simpson pentru un domeniu simplu în raport cu x

Pentru o integrală de primul tip pasul după x va fi $h = \frac{b-a}{2}$, dar cel după y va varia odată cu x (vezi figura 9.5)

$$k(x) = \frac{d(x) + c(x)}{2}.$$

Se obține

$$\begin{aligned} & \int_a^b \int_{c(x)}^{d(x)} f(x, y) \, dx \, dy \approx \\ & \approx \int_a^b \frac{k(x)}{3} [f(x, c(x)) + 4f(x, c(x) + k(x)) + f(x, d(x))] \, dx \\ & \approx \frac{h}{3} \left\{ \frac{k(a)}{3} [f(a, c(a)) + 4f(a, c(a) + k(a)) + f(a, d(a))] \right. \\ & \quad + 4 \frac{k(a+h)}{3} [f(a+h, c(a+h)) + 4f(a+h, c(a+h) + k(a+h)) \\ & \quad + f(a+h, d(a+h))] \\ & \quad \left. + \frac{k(b)}{3} [f(b, c(b)) + 4f(b, c(b) + k(b)) + f(b, d(b))] \right\}. \end{aligned}$$

Dacă domeniul de integrare D este curbiliniu, construim un dreptunghi $R \supset D$, ale cărui laturi să fie paralele cu axele de coordonate (figura 9.6). Se consideră funcția auxiliară

$$f^*(x, y) = \begin{cases} f(x, y), & \text{dacă } (x, y) \in D; \\ 0, & \text{dacă } (x, y) \in R \setminus D;. \end{cases}$$

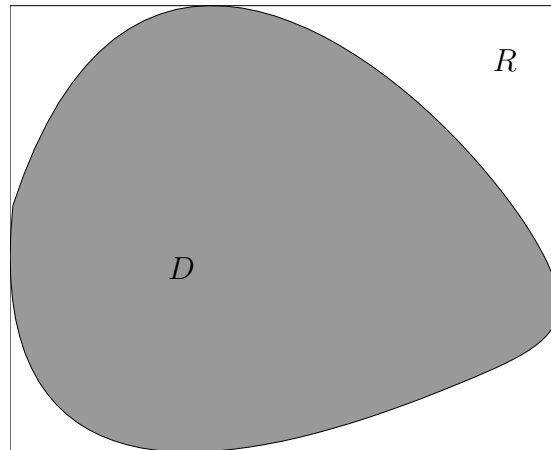


Figura 9.6: Încadrarea unui domeniu curbiliniu într-un dreptunghi

Evident

$$\iint_D f(x, y) \, dx \, dy = \iint_R f^*(x, y) \, dx \, dy$$

Ultima integrală se poate aproxima printr-o tehnică cunoscută.

9.2.1. Considerații de implementare

Fie domeniul de integrare

$$D = [a, b] \times [c, d].$$

Integrala de aproximat se poate scrie sub forma

$$\int_a^b \int_c^d f(x, y) \, dx \, dy = \int_a^b \left(\int_c^d f(x, y) \, dy \right) dx = \int_a^b F(x) \, dx,$$

unde

$$F(x) = \int_c^d f(x, y) \, dy.$$

Să presupunem că *adquad* este o rutină de cuadratură unidimensională adaptivă de forma

function *adquad*(*f* : funcție, *a, b* : real, ε : real) : real

Ideea este de a folosi această rutină pentru a calcula valorile lui F definite mai sus și de a folosi din nou rutina pentru a integra F . Descrierea este dată în algoritmul 9.1.

Algoritmul 9.1 Aproximarea unei integrale duble pe dreptunghi

Intrare: funcția f , intervalele $[a, b]$ și $[c, d]$, rutina $adquad$, eroarea ε

Ieșire: Valoarea aproximativă a integralei

function $dblquad(f, a, b, c, d, adquad, \varepsilon) : real$

$Q := adquad(integint, a, b, \varepsilon, f, c, d, quadf);$

function $integint(x, f, c, d, \varepsilon, adquad) : real$

{integrala interioară}

for each x **do**

$F(x) := adquad(f(x, \cdot), c, d, \varepsilon);$

end for

Bibliografie

- [1] O. Agratini, *Positive Approximation Processes*, Hiperboreea Press, Turda, 2001.
- [2] Octavian Agratini, Ioana Chiorean, Gheorghe Coman, Radu Trîmbițaș, *Analiză numerică și teoria aproximării*, vol. III, Presa Universitară Clujeană, Cluj-Napoca, 2002, coordonatori D. D. Stancu și Gh. Coman.
- [3] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, H. van der Vorst, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*, 2nd ed., SIAM, Philadelphia, PA, 1994, disponibilă prin [www](http://www.netlib.org/templates), <http://www.netlib.org/templates>.
- [4] Å. Björk, *Numerical Methods for Least Squares Problem*, SIAM, Philadelphia, 1996.
- [5] E. Blum, *Numerical Computing: Theory and Practice*, Addison-Wesley, 1972.
- [6] P. Bogacki, L. F. Shampine, *A 3(2) pair of Runge-Kutta formulas*, Appl. Math. Lett. **2** (1989), no. 4, 321–325.
- [7] C. G. Broyden, *A Class of Methods for Solving Nonlinear Simultaneous Equations*, Math. Comp. **19** (1965), 577–593.
- [8] L. Burden, J. D. Faires, *Numerical Analysis*, PWS Kent, Boston, 1986.
- [9] P. G. Ciarlet, *Introduction à l'analyse numérique matricielle et à l'optimisation*, Masson, Paris, Milan, Barcelone, Mexico, 1990.
- [10] Gheorghe Coman, *Analiză numerică*, Editura Libris, Cluj-Napoca, 1995.

- [11] I. Cuculescu, *Analiză numerică*, Editura Tehnică, București, 1967.
- [12] P. J. Davis, P. Rabinowitz, *Numerical Integration*, Blaisdell, Waltham, Massachusetts, 1967.
- [13] James Demmel, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- [14] J. E. Dennis, J. J. Moré, *Quasi-Newton Methods, Motivation and Theory*, SIAM Review **19** (1977), 46–89.
- [15] J. Dormand, *Numerical Methods for Differential Equations. A Computational Approach*, CRC Press, Boca Raton New York, 1996.
- [16] E. Fehlberg, *Klassische Runge-Kutta-Formeln fünfter und siebenter Ordnung mit Schrittweiten-Kontrolle*, Computing **4** (1969), 93–106, Corrigendum: *ibid.* 5, 184.
- [17] E. Fehlberg, *Klassische Runge-Kutta-Formeln vierter und niedriger Ordnung mit Schrittweiten-Kontrolle und ihre Anwendung auf Wärmeleitungsprobleme*, Computing **6** (1970), 61–71, Corrigendum: *ibid.* 5, 184.
- [18] J. G. F. Francis, *The QR transformation: A unitary analogue to the LR transformation*, Computer J. **4** (1961), 256–272, 332–345, parts I and II.
- [19] W. Gander, W. Gautschi, *Adaptive quadrature - revisited*, BIT **40** (2000), 84–101.
- [20] W. Gautschi, *Numerical Analysis, an Introduction*, Birkhäuser, Basel, 1997.
- [21] Walther Gautschi, *Orthogonal polynomials: applications and computation*, Acta Numerica **5** (1996), 45–119.
- [22] D. Goldberg, *What every computer scientist should know about floating-point arithmetic*, Computing Surveys **23** (1991), no. 1, 5–48.
- [23] H. H. Goldstine, J. von Neumann, *Numerical inverting of matrices of high order*, Amer. Math. Soc. Bull. **53** (1947), 1021–1099.
- [24] Gene H. Golub, Charles van Loan, *Matrix Computations*, 3rd ed., John Hopkins University Press, Baltimore and London, 1996.
- [25] Nicholas J. Higham, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.
- [26] E. Isaacson, H. B. Keller, *Analysis of numerical methods*, John Wiley, New York, 1966.

- [27] V. N. Kublanovskaya, *On some algorithms for the solution of the complete eigenvalue problem*, USSR Comp. Math. Phys. **3** (1961), 637–657.
- [28] J. J. Moré, M. Y. Cosnard, *Numerical Solutions of Nonlinear Equations*, ACM Trans. Math. Softw. **5** (1979), 64–85.
- [29] W. H. Press, S. A. Teukolsky, W. T. Vetterling, B. P. Flannery, *Numerical Recipes in C*, Cambridge University Press, Cambridge, New York, Port Chester, Melbourne, Sidney, 1996, disponibilă prin [www, http://www.nr.com/](http://www.nr.com/).
- [30] I. A. Rus, *Ecuatii diferențiale, ecuații integrale și sisteme dinamice*, Transilvania Press, Cluj-Napoca, 1996.
- [31] I. A. Rus, P. Pavel, *Ecuatii diferențiale*, Editura Didactică și Pedagogică, București, 1982, ediția a doua.
- [32] H. Rutishauser, *Solution of the eigenvalue problems with the LR transformation*, Nat. Bur. Stand. App. Math. Ser. **49** (1958), 47–81.
- [33] Y. Saad, *Iterative Methods for Sparse Linear Systems*, PWS Publishing, Boston, 1996, disponibilă via [www](http://www-users.cs.umn.edu/~saad/books.html) la adresa <http://www-users.cs.umn.edu/~saad/books.html>.
- [34] A. Sard, *Linear Approximation*, American Mathematical Society, Providence, RI, 1963.
- [35] Thomas Sauer, *Numerische Mathematik II*, Universität Erlangen-Nurnberg, Erlangen, 2000, Vorlesungskript.
- [36] R. Schwarz, H., *Numerische Mathematik*, B. G. Teubner, Stuttgart, 1988.
- [37] D. D. Stancu, *Analiză numerică – Curs și culegere de probleme*, Lito UBB, Cluj-Napoca, 1977.
- [38] D. D. Stancu, G. Coman, P. Blaga, *Analiză numerică și teoria aproximării*, vol. II, Presa Universitară Clujeană, Cluj-Napoca, 2002, D. D. Stancu, Gh. Coman, (coord.).
- [39] D. D. Stancu, Gh. Coman, O. Agratini, R. Trîmbițaș, *Analiză numerică și Teoria aproximării*, vol. I, Presa Universitară Clujeană, Cluj-Napoca, 2001.
- [40] J. Stoer, R. Burlisch, *Einführung in die Numerische Mathematik*, vol. II, Springer Verlag, Berlin, Heidelberg, 1978.
- [41] J. Stoer, R. Burlisch, *Introduction to Numerical Analysis*, 2nd ed., Springer Verlag, 1992.

-
- [42] A. H. Stroud, *Approximate Calculation of Multiple Integrals*, Prentice Hall Inc., Englewood Cliffs, NJ, 1971.
- [43] Lloyd N. Trefethen, David Bau III, *Numerical Linear Algebra*, SIAM, Philadelphia, 1996.
- [44] C. Überhuber, *Computer-Numerik*, vol. 1, 2, Springer Verlag, Berlin, Heidelberg, New-York, 1995.
- [45] C. Ueberhuber, *Numerical Computation. Methods, Software and Analysis*, vol. I, II, Springer Verlag, Berlin, Heidelberg, New York, 1997.
- [46] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.

- algoritm
 - precis, 18
 - regresiv stabil, 19
 - stabil, 18
- algoritmul
 - Cox-deBoor, 133
 - de Casteljaou, 134
- algoritmul lui Strassen, 50
- analiza progresivă a erorilor, 20
- analiza regresivă a erorilor, 20
- aproximantă, 4
- aproximantă algebric maximală, 259
- aproximantă algebric minimală, 259
- axioma fundamentală a aritmeticii in virgulă flotantă, 8

- condiționare
 - a unui sistem liniar, 28
- condiționare
 - a unei probleme, 11
 - a unui algoritm, 13
 - număr de condiționare, 12
- convergență liniară, 176
- cuadraturi adaptive, 168
- curbă
 - B-spline, 131
 - Bézier, 132

- derivare numerică, 152

- descompunere maximală, 261
- descompunere minimală, 261
- diferență divizată cu noduri multiple, 112
- diferența finită
 - regresivă, 109
- diferența regresivă, 109
- diferențe divizate, 105
- diferențe finite
 - progresive, *Vezi* diferențe progresive
- diferențe progresive, 107

- ecuații normale, 69
- element de cea mai bună aproximare, 64
- element pivot, 34
- eliminare gaussiană, 32
- eroare, 4
- eroare absolută, 4
- eroare asimptotică, 177
- eroare de trunchiere, 226
- eroare relativă, 4

- flops, 35
- formulă de derivare numerică, 148
- formulă de integrare numerică, 154
- formulă de tip Gauss-Christoffel, *Vezi* formulă de cuadratură de tip Gauss

- formulă de aproximare
 - algebric maximală, 261
 - algebric minimală, 261
 - consistentă, 265
 - omogenă, 265
- formulă de cuadratură, *Vezi* formulă de integrare numerică
- formulă de cubatură, 266
 - a lui Simpson, 267
 - a trapezului, 267
 - Newton-Cotes, 266
- formula de aproximare
 - a lui Bernstein, 126
- formula Euler-MacLaurin, 172
- formula lui Newton progresivă, 108
- formula lui Newton regresivă, 110
- formula lui Simpson, 156
- formula lui Simpson repetată, 156
- formula repetată a trapezului, 156
- formula trapezului, *Vezi* regula trapezului
- formule de cuadratură de tip Gauss, 161
- formule Newton-Cotes, 160
- funcție
 - B-spline, 127
 - grilă, 237
- grad de exactitate, 149, 154
- grilă, 237
- identitatea lui Marsden, 129
- IEEE 754 (standardul), 9
- indice de eficiență, 178
- inegalitatea de stabilitate, 238
- integrare numerică, 154
- interpolare blending, *Vezi* interpolare transfinită
- interpolare Hermite, 87
- interpolare Lagrange, 85
 - metoda lui Aitken, 104
 - metoda lui Neville, 103
- interpolare polinomială, 83
- interpolare spline, 112
- interpolare transfinită, 260
- iterație vectorială, 204
- matrice
 - companion, 200
 - complement Schur, 37
 - descompunere LU a unei \sim , 37
 - descompunere LUP a unei \sim , 39
 - descompunere Schur a unei \sim , 202
 - descompunere Schur reală a unei \sim , 204
 - diagonalizabilă, 202
 - factorizare Cholesky a unei \sim , 42
 - forma normală Jordan a unei \sim , 201
 - hermitiană, 22
 - Hessenberg superioară, 204
 - nederogatorie, 202
 - normală, 22
 - ortogonală, 22
 - polinom caracteristic al unei \sim , 200
 - simetrică, 22
 - similare, 201
 - transformarea RQ a unei \sim , 212
 - unitară, 22
 - valoare proprie a unei \sim , 199
 - vector propriu al unei \sim , 199
- metodă cu un pas
 - consistentă, 226
 - convergentă, 241
 - funcția de eroare principală, 227
 - stabilă, 238
- metodă cu un pas
 - ordin, 227
 - ordin exact, 227
- metoda
 - aproximațiilor succesive, 189
 - dezvoltării Taylor, 229
 - eliminării a lui Gauss, *Vezi* eliminare gaussiană

- falsei poziții, 181
- Gauss-Seidel, 57
- lui Broyden, 197
- lui Euler, 227
- lui Euler modificată, 231
- lui Heun, 231
- lui Jacobi, 57
- lui Newton, 186
- lui Romberg, 169
- lui Sturm, 179
- puterii, *Vezi* iterație vectorială QR, 207
 - cu deplasare spectrală, 218
 - cu pas dublu, 220
 - simplă, 215
- rafinării iterative, 50
- relaxării, 58
- secantei, 183
- SOR, 58
- metode quasi-Newton, 194
- metode Runge-Kutta, 231

- normă matricială, 23
- normă Frobenius, 28
- normă matricială
 - indusă, 23
 - subordonată, 23
- notația
 - O , 17
 - Ω , 17
 - Θ , 17
- număr de condiționare, 30
- număr de condiționare al unei matrice, 30

- operator
 - liniar și pozitiv, 140
 - spline cu variație diminuată, see operatorul lui Schoenberg 138
- operatorul
 - Bernstein, 122, 143
 - Hermite-Fejér, 143
 - lui Schoenberg, 138, 143
- ordin de convergență, 177

- partiția unității, 129
- pivotare maximală pe coloana, 34
- pivotare parțială, *Vezi* pivotare maximală pe coloană
- pivotare scalată pe coloană, 34
- polinoame ortogonale, 75
 - Cebîșev de speța a doua, 82
 - Cebîșev de speța I, 79
 - Hermite, 83
 - Jacobi, 83
 - Laguerre, 82
 - Legendre, 77
- polinom
 - de cea mai bună aproximare uniformă, 145
- polinomul
 - Bernstein, 122
- principiul efectelor egale, 5
- problemă de interpolare, 65
- problemă de cea mai buna aproximare, 64
- problemă incorect pusă, 16
- problemă numerică, 2
- problemă prost condiționată, 15
- produs tensorial, 258
- puncte de alternața Cebîșev, 145

- rafinare iterativă, *Vezi* metoda rafinării iterative
- regula trapezelor, *Vezi* formula repetată a trapezului
- regula trapezului, 155
- rezolvarea numerică a ecuațiilor diferențiale
 - metode cu un pas, 226
- rotație Givens, 47

spline

- Not-a-knot, 117

- spline complete, 116

- spline cubice, 115

- spline naturale, 117

- suma booleană, 258

- tabelă Butcher, 235

- tabela Butcher, 253

- Teorema lui Peano, 84

transformare

- Householder, 44

virgulă flotantă

- anulare, 8

- exponent, 6

- normalizarea unui număr, 6

- reprezentare, 6

- semnificant, 6