

MARILENA POPA

ROMULUS MILITARU

METODE NUMERICE

Note de curs

1. REZOLVAREA NUMERICĂ A SISTEMELOR DE ECUAȚII LINIARE

Introducere. Rezolvarea sistemelor algebrice liniare și operațiile de calcul matriceal (evaluarea determinantilor, inversarea matriceală, calculul valorilor și vectorilor proprii) sunt incluse în domeniul algebrei liniare – implicată în diverse probleme științifice, de exemplu:

– problemele care depind de un număr finit de grade de libertate, reprezentate prin ecuații diferențiale ordinare sau cu derivate parțiale sunt transformate, cu ajutorul diferențelor finite, în sisteme de ecuații liniare;

– problemele neliniare sunt frecvent soluționate (aproximate) prin procese de liniarizare;

– programarea liniară implică rezolvarea unor sisteme de ecuații algebrice liniare;

– foarte multe probleme ingineresti din domeniul rețelelor electrice, analiza structurilor, proiectarea clădirilor, vapoarelor, avioanelor, transportul lichidelor și gazelor prin conducte etc. necesită, pentru soluționare, rezolvarea unor sisteme liniare.

Formularea problemei. Fie $A \in \mathbb{R}^{n \times n}$ (o matrice reală cu n linii și n coloane), $b \in \mathbb{R}^n$ (un vector – matrice coloană – cu n componente reale) și vectorul necunoscut $x \in \mathbb{R}^n$. Atunci un sistem de ecuații liniare se scrie sub una din formele:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ \dots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{cases}$$

sau matriceal: $Ax = b$ sub formă compactă și

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \quad (1)$$

sub formă dezvoltată.

O ultimă variantă de scriere a unui sistem liniar este:

$$\sum_{j=1}^n a_{ij}x_j = b_i, i = 1, 2, \dots, n$$

Observații.

1. Un sistem liniar (1) este **consistent** dacă are cel puțin o soluție și **inconsistent** dacă nu are nici o soluție. Orice sistem liniar considerat în continuare poate avea soluție unică (compatibil determinat), nici o soluție (incompatibil) sau o infinitate de soluții (compatibil nedeterminat) – nu există alte posibilități. Bineînțeles că interesează primul caz.

2. Sistemele (1) se pot clasifica și după vectorul b , în:

a) **sisteme omogene** – dacă $b = 0$. Orice sistem omogen de forma $Ax = 0$ este un sistem consistent (deoarece are soluția $x = 0$ – neinteresantă). Un sistem omogen are o soluție nebanală dacă și numai dacă $\det A = 0$ (adică A este singulară). În această situație soluția depinde de cel puțin un parametru.

b) **sisteme neomogene** – dacă $b \neq 0$. Sistemul (1) este compatibil determinat pentru $(\forall) b \neq 0$ dacă și numai dacă sistemul omogen $Ax=0$ nu are decât soluția banală (adică A este nesingulară).

3. În sistemele obținute din aplicațiile fizice, numerele ce constituie matricea A și vectorul b sunt afectate de erori inerente (provenite din măsurători) sau erori de rotunjire (1/7 nu poate fi reprezentat exact pe nici un calculator electronic – care are lungime fixă a cuvântului – deoarece reprezentarea lui în binar are o infinitate de biți). Dacă erorile mici în cadrul coeficienților lui A și b , sau în procesul de calcul au un efect redus asupra vectorului, soluție, un astfel de sistem este **bine condiționat**, iar dacă efectul este considerabil, un astfel de sistem este **slab condiționat**.

Metodele numerice de rezolvare a sistemului liniar (1) sunt de două tipuri: directe și iterative.

Rezolvarea directă a sistemelor de ecuații liniare

Metodele directe constau în reducerea sistemului (1), într-un număr finit de etape, la un sistem **echivalent**, (cu aceeași soluție) direct rezolvabil – prin substituție directă sau inversă.

Aceste metode furnizează soluția exactă x a sistemului (1) în cazurile (ideale) în care erorile de rotunjire sunt absente și necesită, în

acest scop, efectuarea unui număr de operații aritmetice elementare de ordinul n^3 .

Din acest motiv, metodele directe se utilizează pentru rezolvarea sistemelor „uzuale“, de dimensiune $n \leq 100$.

În continuare vom prezenta câteva metode directe de rezolvare a sistemelor de ecuații liniare.

1.1. Metoda Gauss. Tehnici de pivotare

Procedura de reducere a sistemului (1) la un sistem echivalent are la bază o schemă de eliminare succesivă a necunoscutelor, introdusă de Gauss în 1823, prin care se efectuează triangularizarea superioară a matricei sistemului prin transformări elementare.

Transformările elementare din metoda Gauss sunt de tipul:

- schimbarea a două linii între ele;

- adunarea unei linii, înmulțite cu un număr, la altă linie

și se efectuează asupra matricei extinse $(A | b)$ astfel încât în „ $n-1$ “ etape în locul matricei A să avem o matrice superior triunghiulară.

Observație. Transformările elementare de tipul celor de mai sus sunt permise numai asupra liniilor deoarece nu influențează obținerea soluției (se referă la ecuațiile sistemului care prin schimbare sau adunare – multiplicare cu un număr – nu influențează soluția).

Încercăm prezentarea algoritmică (pe pași) a metodei Gauss:

- la primul pas se folosește prima ecuație a sistemului la eliminarea necunoscutei x_1 din celelalte $n-1$ ecuații obținându-se un nou sistem $A^{(2)}x = b^{(2)}$ (unde $A^{(1)} = A$ și $b^{(1)} = b$);

- la al doilea pas se folosește a doua ecuație transformată (deci din sistemul anterior $A^{(2)}x = b^{(2)}$) pentru eliminarea necunoscutei x_2 din ultimele $n-2$ ecuații, obținându-se sistemul $A^{(3)}x = b^{(3)}$ ș.a.m.d.

$$\begin{cases} x_n = \frac{b_n^{(n)}}{a_{nn}^{(n)}} \\ b_i^{(n)} - \sum_{j=i+1}^n a_{ij}^{(n)} x_j \\ x_i = \frac{\quad}{a_{ii}^{(n)}}, \quad i = n-1, n-2, \dots, 1 \end{cases} \quad (4)$$

Observații.

1. Formulele (2) pot fi completate astfel încât sistemul (3) obținut să aibă forma:

$$\begin{cases} x_1 + a_{12}^{(n)} x_2 + \dots + a_{1n}^{(n)} x_n = b_1^{(n)} \\ x_2 + \dots + a_{2n}^{(n)} x_n = b_2^{(n)} \\ \dots \\ x_n = b_n^{(n)} \end{cases},$$

ceea ce înseamnă că în relațiile (4) dispar numitorii, sau formulele (2) pot fi modificate astfel încât sistemul (3) să aibă forma matriceală:

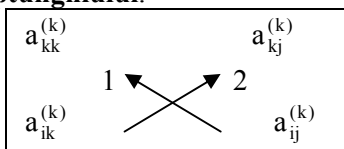
$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1^{(n)} \\ b_2^{(n)} \\ \vdots \\ b_n^{(n)} \end{pmatrix}$$

ceea ce înseamnă că relațiile (4) devin:

$$x_i = b_i^{(n)}, \quad 1 \leq i \leq n$$

Această modificare a metodei Gauss este cunoscută sub numele de metoda Gauss-Jordan.

2. Ultimele din formulele (2), pentru $a_{ij}^{(k+1)}$, $k+1 \leq i, j \leq n$ și $b_i^{(k+1)}$, $k+1 \leq i \leq n$ sunt cunoscute și sub numele de **regula dreptunghiului**.



$$a_{ij}^{(k+1)} = \frac{a_{ij}^{(k)} \cdot a_{kk}^{(k)} - a_{ik}^{(k)} \cdot a_{kj}^{(k)}}{a_{kk}^{(k)}}$$

1.2. Metoda factorizării LR

Definiția 1. Fie $A \in \mathbb{R}^{n \times n}$. O relație de forma:

$$A = L \cdot R \quad (5)$$

unde $L \in \mathbb{R}^{n \times n}$ este o matrice inferior triunghiulară iar R este o matrice superior triunghiulară, se numește factorizare LR a lui A .

Precizare: o matrice $T \in \mathbb{R}^{n \times n}$ se numește superior (inferior) triunghiulară dacă $t_{ij} = 0$ pentru $i > j$ ($t_{ij} = 0$ pentru $i < j$) unde $T = (t_{ij})$, $1 \leq i, j \leq n$.

În contextul definiției 1, sistemul liniar (1) devine:

$$LRx = b$$

care conduce la rezolvarea directă a două sisteme liniare cu matrice triunghiulare și anume:

$$Ly = b \quad (6)$$

$$Rx = y \quad (7)$$

Observații:

1. Sistemul (6) are soluția „intermediară“ y ale cărei componente se obțin prin substituție directă (înainte) datorită formei matricei L – inferior triunghiulară.

2. Sistemul (7) are soluția „finală“ x ale cărei componente se obțin prin substituție inversă (înapoi) datorită formei matricei R – superior triunghiulară.

Pentru simplificarea referirilor ulterioare, notăm prin $A_k = (a_{ij})$, $1 \leq i, j \leq k$ submatricele **lider principale** de ordinul k ale lui A , $k = 1, 2, \dots, n$ și introducem condiția:

i) **submatricele A_k sunt nesingulare** (evident $A_1 = a_{11}$, iar $A_n = A$)

Observație. Condiția i) altfel exprimată este: **toți minorii diagonali principali ai matricei A sunt nenuli.**

$$a_{11} \neq 0, \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \neq 0, \text{ ș.a.m.d. } \det A \neq 0$$

(de fapt toți **primii** minori diagonali de fiecare ordin sunt nenuli)

Este valabilă următoarea:

Teoremă. Dacă matricea A satisface condiția i), atunci există o factorizare LR a lui A , în care matricea L este nesingulară.

O factorizare LR a lui A poate fi calculată direct printr-o așa-numită **procedură compactă**, în care cele n^2 egalități scalare, din (5), se rezolvă succesiv în raport cu elementele necunoscute l_{ik} , $i \geq k$ și r_{kj} , $k \leq j$ ale matricilor L și respectiv R . Unicitatea procedurii este asigurată precizând apriori elementele diagonale ale matricei L sau elementele diagonale ale matricei R . Din acest punct de vedere, în practică se utilizează două tipuri de factorizări:

a) **Factorizarea Doolittle** – impune L cu diagonala unitate, adică $l_{kk} = 1$, $k = 1, 2, \dots, n$.

b) **Factorizarea Crout** – impune R cu diagonala unitate, adică $r_{kk} = 1$, $k = 1, 2, \dots, n$.

Ne vom ocupa de detalierea acestei factorizări, astfel:

$$\text{Pasul 1. Fie } L = \begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{pmatrix} \text{ și } R = \begin{pmatrix} 1 & r_{12} & \dots & r_{1n} \\ & 1 & \dots & r_{2n} \\ & & \dots & \dots \\ & & & 1 \end{pmatrix}$$

Explicitând egalitatea $A = L \cdot R$, sub forma:

$$a_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} r_{kj}, \quad i, j = 1, 2, \dots, n$$

elementele matricilor L și R se obțin astfel:

$$\left\{ \begin{array}{l} l_{i1} = a_{i1}, \quad 1 \leq i \leq n \rightarrow \text{prima coloană din } L \\ r_{1j} = \frac{a_{1j}}{l_{11}}, \quad 2 \leq j \leq n \rightarrow \text{prima linie din } R \\ l_{ik} = a_{ik} - \sum_{h=1}^{k-1} l_{ih} \cdot r_{hk}, \quad 2 \leq k \leq i \leq n \\ r_{kj} = \left(a_{kj} - \sum_{h=1}^{k-1} l_{kh} \cdot r_{hj} \right) / l_{kk}, \quad 2 \leq k < j \leq n \end{array} \right. \quad (8)$$

Justificarea formulelor (8):

- a_{ik} se obține înmulțind linia i din matricea L cu coloana k din matricea R ;

- dacă $i > k$, atunci $(l_{i1} \ l_{i2} \ \dots \ l_{i,k-1} \ l_{ik} \ \dots \ l_{ii} \ 0 \ \dots \ 0) \cdot$

$$\begin{pmatrix} r_{1k} \\ r_{2k} \\ \vdots \\ r_{k-1,k} \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \rightarrow a_{ik}$$

$$\Rightarrow a_{ik} = \sum_{h=1}^{k-1} l_{ih} r_{hk} + l_{ik} ;$$

- dacă $i = k$, atunci $(l_{i1} \ l_{i2} \ \dots \ l_{i,k-1} \ l_{ii} \ 0 \ \dots \ 0) \cdot$

$$\begin{pmatrix} r_{1k} \\ r_{2k} \\ \vdots \\ r_{k-1,k} \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \rightarrow a_{ik}$$

$$\Rightarrow a_{ik} = \sum_{h=1}^{k-1} l_{ih} r_{hk} + l_{ii} ; \text{ deci am regăsit formulele (8) pentru } i \geq k;$$

$$- \text{dac\`a } k < j, \text{ atunci } (1_{k1} \ 1_{k2} \ \dots \ 1_{k,k-1} \ 1_{kk} \ 0 \ \dots \ 0) \cdot \begin{pmatrix} r_{1j} \\ r_{2j} \\ \vdots \\ r_{k-1,j} \\ r_{kj} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\Rightarrow a_{kj} = \sum_{h=1}^{k-1} l_{kh} r_{hj} + l_{kk} r_{kj}, \text{ deci am reg\`asit formulele (8) pentru } k < j.$$

Pasul 2. Se rezolv\`a sistemul (6):

$$\begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

cu formulele de **substitu\`ie direct\`a** (înapoi)

$$\begin{cases} y_1 = \frac{b_1}{l_{11}} \\ y_k = \left(b_k - \sum_{j=1}^{k-1} l_{kj} y_j \right) / l_{kk}, k = 2, \dots, n \end{cases} \quad (9)$$

Pasul 3. Se rezolv\`a sistemul (7):

$$\begin{pmatrix} 1 & r_{12} & \dots & r_{1n} \\ & 1 & \dots & r_{2n} \\ & & \dots & \dots \\ & & & 1 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

cu formulele de **substitu\`ie invers\`a** (înapoi)

$$\begin{cases} x_n = y_n \\ x_k = y_k - \sum_{j=k+1}^n r_{kj} x_j, k = n-1, n-2, \dots, 1 \end{cases} \quad (10)$$

1.4. Rezolvarea iterativă a sistemelor de ecuații liniare. Metodele Jacobi și Seidel Gauss

Introducere

Metodele iterative constau în construcția unui șir $x^{(k)}$, (evident $x^{(k)} \in \mathbb{R}^n$) $k = 0, 1, \dots$, **convergent** către soluția **exactă** x a sistemului:

$$Ax = b, \text{ unde } A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$$

Oprirea **procesului iterativ**, adică trunchierea șirului $x^{(k)}$, are loc la un indice s , determinat pe parcursul calculului în funcție de precizia impusă, astfel încât termenul curent $x^{(s)}$ să constituie o **aproximație satisfăcătoare** a soluției căutate x .

Metoda Seidel-Gauss (metoda iterațiilor succesive) constă în construcția șirului $x^{(k+1)}$, $k = 0, 1, \dots$ astfel: $x^{(0)} \in \mathbb{R}^n$ arbitrar, iar

$$\begin{cases} x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right), i = \overline{2, n-1}, \text{ unde} \\ x_1^{(k+1)} = \frac{1}{a_{11}} \left(b_1 - \sum_{j=2}^n a_{1j} x_j^{(k)} \right) \\ x_n^{(k+1)} = \frac{1}{a_{nn}} \left(b_n - \sum_{j=1}^{n-1} a_{nj} x_j^{(k+1)} \right) \end{cases} \quad (27)$$

Se observă că iterația „nouă“ folosește iterații „vechi“, numai dacă nu există iterații „noi“ calculate.

Observații: Condiția **suficientă** pentru ca șirul $x^{(k+1)}$ să convergă la soluția x este ca A să fie **diagonal dominantă** pe linii

$$(26) \text{ sau pe coloane } |a_{jj}| > \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|, j = \overline{1, n}.$$

2. Pot exista sisteme pentru care condiția amintită nu este îndeplinită, dar algoritmul Seidel-Gauss să converge. Utilizând norme vectoriale adecvate se poate obține o condiție necesară și suficientă de convergență, calculele depășind cadrul prezentei lucrări. Pentru detalii vezi [19] sau [28].

3. Practic, iterațiile se opresc atunci când pentru un ε , impus de utilizator, avem îndeplinite condițiile:

$$\left| x_i^{(k+1)} - x_i^{(k)} \right| < \varepsilon, \quad i = \overline{1, n} \quad (28)$$

(sau trecând la maxim $d(x^{(k+1)}, x^{(k)}) < \varepsilon$). Am putea determina din teorema Banach, indicele la care ne putem opri, în funcție de q .

În acest caz, vom defini ca soluție vectorul $x^{(k+1)}$.

Metoda Jacobi (metoda iterațiilor simultane) constă în construcția șirului $x^{(k+1)}$, $k = 0, 1, \dots$ astfel

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right), \quad i = \overline{1, n}; \quad x^{(0)} \in \mathbb{R}^n \text{ arbitrar} \quad (29)$$

Observații:

1. Sunt valabile observațiile precedente.
2. Se observă că iterația „nouă“ folosește toate componentele iterației „vechi“ ceea ce face ca formula să fie mai simplă, dar metoda Jacobi este mai lent convergentă decât metoda Seidel-Gauss în aceleași condiții inițiale ($x^{(0)}, \varepsilon$).

2. ECUAȚII ȘI SISTEME DE ECUAȚII NELINIARE

2.1. Metoda bisecției

Fiind dată funcția $f : [a, b] \rightarrow \mathbb{R}$ și ecuația $f(x) = 0$ metoda bisecției se bazează pe proprietatea funcției f de a fi continuă pe $[a, b]$ și $f(a) \cdot f(b) \leq 0$. Atunci f are pe $[a, b]$ un număr impar de zerouri.

Dacă $f\left(\frac{a+b}{2}\right) = 0$, $x_1 = \frac{a+b}{2}$ este soluție. Dacă nu, se notează $a_1 = a$ și $b_1 = \frac{a+b}{2}$ și $f(a) \cdot f\left(\frac{a+b}{2}\right) < 0$ sau $a_1 = \frac{a+b}{2}$ și $b_1 = b$ în cazul $f\left(\frac{a+b}{2}\right) \cdot f(b) < 0$. Continuând, se obține succesiunea de intervale: $[a_1, b_1], [a_2, b_2] \dots [a_n, b_n]$ cu proprietatea că:

$$f(a_n) \cdot f(b_n) < 0, \quad n = 1, 2, 3, \dots$$

$$\text{cu } b_n - a_n = \frac{1}{2^n} (b - a).$$

Obținem astfel $(a_n)_{n \in \mathbb{N}}$ șir crescător și $(b_n)_{n \in \mathbb{N}}$ șir descrescător cu:

$$x_1 = \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n$$

În ipoteza că au fost separate mai multe zerouri, algoritmul expus poate fi aplicat pentru determinarea fiecărei soluții separate.

Metoda înjumătățirii intervalului este cea mai simplă (dar și cea mai slab convergentă) metodă pentru determinarea unei soluții a ecuației $f(x) = 0$.

Algoritmul poate fi aplicat în două variante:

- Au fost separate una sau mai multe soluții;
- Nu a fost separată nici o soluție. În această situație se notează cu $a_1 = -\infty$, $b_1 = 0$ dacă $f(-\infty) \cdot f(0) < 0$ sau $a_1 = 0$, $b_1 = +\infty$

dacă $f(0)f(\infty) < 0$, în practică simbolul ∞ fiind înlocuit cu un număr foarte mare.

2.2. Metoda secantei

Considerăm din nou ecuația

$$f(x) = 0, x \in [a, b] \quad (1)$$

cu f continuă și derivabilă de două ori pe $[a, b]$ și în plus $f(a) \cdot f(b) < 0$, adică a fost separată o soluție a ecuației în acest interval, iar $f''(x)$ păstrează semn constant pe $[a, b]$.

Vom prezenta în continuare metoda secantei pentru aproximarea acestei soluții, notată cu \bar{x} .

Pentru a obține o primă aproximație x_1 pentru \bar{x} vom scrie ecuația dreptei care trece prin punctele $(a, f(a))$ și $(b, f(b))$:

$$\frac{x - a}{b - a} = \frac{y - f(a)}{f(b) - f(a)} \quad (2)$$

Aproximația x_1 va fi dată de abscisa punctului în care dreapta taie axa Ox , adică pentru $y = 0$, relația (2) se scrie sub forma:

$$x_1 = a - \frac{f(a)(b - a)}{f(b) - f(a)} \quad (3)$$

Vom construi din nou o dreaptă similară cu (2) folosind un punct: $(x_1, f(x_1))$ și un al doilea punct în funcție de intervalul în care se află soluția, (a, x_1) sau (x_1, b) .

Dacă $f(a) \cdot f(x_1) < 0$ atunci al doilea punct al dreptei va fi $(a, f(a))$ și suntem conduși la formula generală:

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - a)}{f(x_n) - f(a)} \quad (4)$$

șir descrescător (din construcție) și mărginit, deci convergent.

Dacă $f(x_1) \cdot f(b) < 0$, al doilea punct va fi $(b, f(b))$ și formula iterativă va fi:

$$x_{n+1} = x_n - \frac{f(x_n)(b - x_n)}{f(b) - f(x_n)} \quad (5)$$

obținând prin construcție un șir crescător și mărginit.

În ambele cazuri șirul dat de (4) sau (5) va avea ca limită soluția căutată, adică:

$$\bar{x} = \lim_{n \rightarrow \infty} x_n$$

Trecând la limită în oricare din relația de recurență (4) sau (5) se obține:

$$f(\bar{x}) = 0$$

2.3. Metoda Newton (tangentei)

2.3.2. Metoda Newton pentru sisteme neliniare

Prin metoda lui Newton se construiește un șir de aproximații $x^{(n)}$ ale soluției sistemului (6), cunoscând că acest sistem are o soluție \bar{x} și având o aproximație inițială $x^{(0)}$ a acestei soluții.

În definirea acestui șir de aproximații intervine derivata $F'(x)$ a operatorului F .

Din relația (7) rezultă că numărătorul tinde la 0 mai repede decât h . Pentru $\|h\|$ suficient de mic, putem face aproximarea

$$F(x+h) - F(x) - F'(x)h \cong 0$$

Să aplicăm această relație pentru o aproximație $x^{(n)}$ a soluției ecuației (1) și pentru $h = \bar{x} - x^{(n)}$. Obținem:

$$F(\bar{x}) - F(x^{(n)}) - F'(x^{(n)})(\bar{x} - x^{(n)}) \cong 0, \text{ și cum:}$$

$$F(\bar{x}) = 0 - \text{căci } \bar{x} \text{ este soluția ecuației (6), obținem:}$$

$$F(x^{(n)}) + F'(x^{(n)})(\bar{x} - x^{(n)}) \cong 0$$

Datorită aproximării, rezolvând această ecuație nu vom obține soluția exactă (\bar{x}), ci o valoare $x^{(n+1)}$, care o definim ca aproximație de ordinul $n+1$.

$$F(x^{(n)}) + F'(x^{(n)})(x^{(n+1)} - x^{(n)}) = 0$$

Aproximația $x^{(n+1)}$, presupunând că există $[F'(x^{(n)})]^{-1}$, este de forma:

$$x^{(n+1)} = x^{(n)} - [F'(x^{(n)})]^{-1} F(x^{(n)}) \quad (8)$$

Procesul iterativ definit de (8) poartă numele de metoda lui Newton.

Observație. În particular, metoda lui Newton, se poate aplica pentru rezolvarea ecuațiilor de forma:

$$f(x) = 0$$

unde $f : \mathbb{R} \rightarrow \mathbb{R}$ este o funcție derivabilă, cu derivata nenulă într-un anumit interval care conține o soluție \bar{x} a acestei ecuații. Șirul (8) devine, în acest caz:

$$x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{f'(x^{(n)})}$$

și metoda admite următoarea interpretare geometrică: „ $x^{(n+1)}$ este abscisa punctului de intersecție cu axa Ox a tangentei la graficul funcției $f(x)$ în punctul $(x^{(n)}, f(x^{(n)}))$ “. Din acest motiv, metoda lui Newton se mai numește și metoda tangentei.

Observație. Metoda Newton se poate aplica și în cazul sistemelor neliniare:

$$\begin{cases} f_1(x_1, x_2, \dots, x_m) = 0 \\ f_2(x_1, x_2, \dots, x_m) = 0 \\ \dots\dots\dots \\ f_m(x_1, x_2, \dots, x_m) = 0 \end{cases}$$

unde $F : \mathbb{R}^m \rightarrow \mathbb{R}^m$, cu $F = (f_1, f_2, \dots, f_m)^t$. În exemplul dat am constatat că dacă f_i au derivate parțiale de ordinul I continue, atunci

există derivata lui F și $F'(x) = \left(\frac{\partial f_i(x)}{\partial x_j} \right)$, matricea iacobiană.

3. POLINOM CARACTERISTIC.

VECTORI ȘI VALORI PROPRII

Preliminarii (noțiuni generale).

Definiția 1. Dacă matricea $A \in \mathbb{R}^{n \times n}$, atunci polinomul definit prin:

$$p_A(\lambda) = \det(\lambda I - A) \quad (1)$$

se numește **polinomul caracteristic** al matricei A .

Precizări:

1. Polinomul p_A este un polinom unitar, de grad n , cu coeficienți reali și ecuația caracteristică de forma:

$$p_A(\lambda) = 0 \Leftrightarrow \lambda^n + c_1 \lambda^{n-1} + c_2 \lambda^{n-2} + \dots + c_{n-1} \lambda + c_n = 0 \quad (1')$$

are cel mult n rădăcini distincte (reale sau complexe).

2. Dacă λ este o rădăcină a lui p_A , atunci, din $\det(\lambda I - A) = 0$, înseamnă că sistemul liniar omogen definit prin

$$(A - \lambda I) \cdot x = \mathbf{0} \quad (\text{vectorul nul din } \mathbb{R}^n) \quad (2)$$

are și soluții nebanale.

În continuare vom studia rădăcinile lui p_A și soluțiile nebanale ale sistemului (2).

Definiția 2. Dacă p_A este polinomul caracteristic al matricei A , rădăcinile lui p_A se numesc **valori proprii** (sau valori caracteristice) ale matricei A . Dacă λ este o valoare proprie a lui A și $x \neq 0$ verifică sistemul (2) atunci x se numește **vector propriu** (sau vector caracteristic) al lui A corespunzător valorii proprii λ .

Observație. Egalitatea (2) conduce la:

$$Ax = \lambda x \quad (3)$$

Definiția 3. Mulțimea $\sigma(A) = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$ se numește **spectrul** lui A , iar numărul $\rho(A) = \max_{1 \leq i \leq n} |\lambda_i|$ se numește **raza spectrală** a lui A .

Consecință. Orice matrice $A \in \mathbb{R}^{n \times n}$ are cel puțin un vector propriu x , mai precis, fiecărei valori proprii $\lambda_i \in \sigma(A)$ îi corespunde cel puțin un vector propriu $x_i \in \text{Ker}(\lambda_i I_n - A)$.

Metodele numerice pentru calculul valorilor și vectorilor proprii se împart în două clase astfel:

I. – Metode pentru calcularea polinomului caracteristic, în care determinarea coeficienților lui $p_A(\lambda)$ se face în mod direct. Deci o cale posibilă pentru a calcula valorile proprii este să obținem întâi coeficienții polinomului caracteristic:

$$p_A(\lambda) = \lambda^n + c_1\lambda^{n-1} + \dots + c_n$$

și apoi să rezolvăm ecuația caracteristică $p_A(\lambda) = 0$, folosind în acest scop orice metodă aproximativă de rezolvare a ecuațiilor algebrice neliniare.

Totuși, nu trebuie să procedăm astfel decât în cazurile când matricea A este de ordin foarte mic și soluțiile ecuației caracteristice sunt bine separate.

Motivul este că valorile proprii sunt foarte sensibile la perturbațiile coeficienților c_1, c_2, \dots, c_n , datorate erorilor de rotunjire inerente.

II. – Metode pentru determinarea valorilor și vectorilor proprii ai matricei A , fără a dezvolta în prealabil polinomul caracteristic. Aceste metode sunt în esență proceduri iterative de aducere a matricei A la anumite forme simple („canonice“) prin transformări de asemănare sau prin transformări de asemănare ortogonală (forma diagonală, forma tridiagonală, forma superior triunghiulară, forma superior Hessenberg).

3.1. Calculul polinomului caracteristic cu ajutorul minorilor diagonali

Dacă $A \in \mathbb{R}^{n \times n}$ este de forma $A = (a_{ij})$, $1 \leq i, j \leq n$, atunci polinomul ei caracteristic se poate scrie:

$$p_A(\lambda) = \det(\lambda\delta_{ij} - a_{ij})$$

Elementele fiecărei coloane a determinantului de mai sus sunt sume algebrice de câte două elemente. Deci $p_A(\lambda)$ se descompune într-o sumă de 2^n determinanți. Fie Δ unul dintre acești determinanți și să notăm cu j_1, j_2, \dots, j_m coloanele lui Δ care conțin numai elemente ale matricei A . Toate celelalte coloane conțin λ la intersecția cu diagonala principală și în rest 0.

Dezvoltând succesiv Δ după coloanele care conțin λ și scoțând -1 factor comun pe celelalte coloane, se obține:

$$\Delta = (-1)^m \cdot \lambda^{n-m} \cdot M_{j_1 j_2 \dots j_m}$$

unde $M_{j_1 j_2 \dots j_m}$ este minorul diagonal al matricei A format cu elementele care se află la intersecția liniilor și coloanelor de indici j_1, j_2, \dots, j_m .

Dacă în dezvoltarea lui $p_A(\lambda)$ se grupează termenii după puterile lui λ , se obține:

$$p_A(\lambda) = \lambda^n - \sigma_1 \lambda^{n-1} + \sigma_2 \lambda^{n-2} - \dots + (-1)^n \sigma_n,$$

unde:

$$\sigma_1 = \sum_{i=1}^n a_{ii} \quad (\sigma_1 \text{ este suma minorilor diagonali de ordinul unu});$$

$$\sigma_2 = \sum_{1 \leq i < j \leq n} \begin{vmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{vmatrix} \quad (\sigma_2 \text{ este suma minorilor diagonali de}$$

ordinul doi); ș.a.m.d.

$$\sigma_n = \det A \quad (\sigma_n \text{ reprezintă singurul minor diagonal de ordinul } n);$$

3.2. Metoda Leverrier

Dacă $A \in R^{n \times n}$, atunci polinomul său caracteristic este:

$$p_A(\lambda) = \lambda^n - \sigma_1 \lambda^{n-1} + \sigma_2 \lambda^{n-2} - \dots + (-1)^n \sigma_n$$

Pentru a-i determina coeficienții procedăm astfel:

1) Determinăm $s_k = \text{Tr}(A^k)$, $1 \leq k \leq n$.

2) $\sigma_1 = s_1$

$$\sigma_k = (s_1 \sigma_{k-1} - s_2 \sigma_{k-2} + \dots + (-1)^{k+1} s_k) / k, \quad 2 \leq k \leq n.$$

Algoritmul lui Leverrier se programează fără dificultate și prezintă un grad mare de generalitate (nu există cazuri particulare).

Pentru n mare, calculele sunt dificile (trebuie calculate puterile matricei date) și se mărește timpul de răspuns.

3.3. Metoda Krylov

Este o metodă pentru calculul polinomului caracteristic bazată pe teorema lui Cayley-Hamilton. Fie $A \in \mathbb{R}^{n \times n}$ și $p_A(\lambda) = \lambda^n + c_1\lambda^{n-1} + \dots + c_{n-1}\lambda + c_n$, polinomul caracteristic al matricei A . Trebuie să determinăm coeficienții c_1, c_2, \dots, c_n . Conform teoremei amintite, avem:

$$p_A(A) = A^n + c_1A^{n-1} + \dots + c_{n-1}A + c_nI_n = O_n \quad (7)$$

Fie $y^{(0)} \in \mathbb{R}^n$, oarecare, nenul. Din relația (7) rezultă:

$$A^n y^{(0)} + c_1 A^{n-1} y^{(0)} + \dots + c_{n-1} A y^{(0)} + c_n y^{(0)} = \mathbf{0} \quad (8)$$

Dacă se notează

$$A^k y^{(0)} = y^{(k)}, \quad k = 1, 2, \dots, n \quad (9)$$

atunci relația (8) se mai scrie:

$$c_1 y^{(n-1)} + c_2 y^{(n-2)} + \dots + c_{n-1} y^{(1)} + c_n y^{(0)} = -y^{(n)} \quad (10)$$

Scrisă pe componente, relația (10) reprezintă un sistem de n ecuații liniare cu n necunoscute c_1, c_2, \dots, c_n . Dacă determinantul sistemului (10) este nenul, atunci soluția unică reprezintă coeficienții polinomului caracteristic. Această soluție se poate obține prin una din metodele de rezolvare a sistemelor de ecuații liniare. Dacă determinantul sistemului (10) este nul, atunci trebuie reluate calculele cu un alt vector inițial $y^{(0)}$.

Observație. Relația (9) pentru obținerea vectorilor $y^{(1)}, y^{(2)}, \dots, y^{(n)}$ conduce la calcule simplificate dacă este scrisă sub forma echivalentă:

$$y^{(k)} = A y^{(k-1)}, \quad k = 1, 2, \dots, n \quad (11)$$

Dacă polinomul caracteristic are rădăcini distincte, atunci vectorii $y^{(0)}, y^{(1)}, \dots, y^{(n-1)}$ obținuți mai sus permit determinarea simplă și a vectorilor proprii. Fie $\lambda_1, \lambda_2, \dots, \lambda_n$ valorile proprii distincte și $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ vectorii proprii corespunzători care formează o bază.

Fie:

$$y^{(0)} = b_1 x^{(1)} + b_2 x^{(2)} + \dots + b_n x^{(n)} \quad (12)$$

expresia vectorului inițial în această bază.

Observație. Menționăm faptul că vectorul inițial a condus la obținerea coeficienților c_1, c_2, \dots, c_n ai polinomului caracteristic.

În continuare, din (9) sau (11) și definiția vectorilor proprii, obținem:

$$\left\{ \begin{array}{l} y^{(1)} = \sum_{k=1}^n b_k \lambda_k x^{(k)} \\ \dots \\ y^{(n-1)} = \sum_{k=1}^n b_k \lambda_k^{n-1} x^{(k)} \end{array} \right. \quad (13)$$

Să notăm:

$$q_j(\lambda) = \frac{p_A(\lambda)}{\lambda - \lambda_j} = \lambda^{n-1} + q_{1j}\lambda^{n-2} + \dots + q_{n-1,j}, j = 1, 2, \dots, n. \quad (14)$$

În ipoteza făcută (valori proprii distincte), avem.

$$\left\{ \begin{array}{l} q_j(\lambda_k) = 0, k \neq j \\ q_j(\lambda_j) \neq 0 \end{array} \right. \quad (15)$$

și din relațiile de mai sus deducem:

$$y^{(n-1)} + q_{1j}y^{(n-2)} + \dots + q_{n-1,j}y^{(0)} = b_1q_j(\lambda_1)x^{(1)} + b_2q_j(\lambda_2)x^{(2)} + \dots + b_jq_j(\lambda_j)x^{(j)} + \dots + b_nq_j(\lambda_n)x^{(n)} = b_jq_j(\lambda_j)x^{(j)}$$

Dacă $b_j \neq 0$, atunci $b_jq_j(\lambda_j)x^{(j)}$ este de asemenea un vector propriu corespunzător valorii proprii λ_j , deci

$$y^{(n-1)} + q_{1j}y^{(n-2)} + \dots + q_{n-1,j}y^{(0)} \quad (16)$$

este un vector propriu corespunzător valorii proprii λ_j .

3.4. Metoda Fadeev

Este o metodă, care pornind de la relația (17), cu un efort minim, permite obținerea atât a coeficienților polinomului caracteristic al matricei date $A \in \mathbb{R}^{n \times n}$, cât și a inversei acesteia A^{-1} .

Algoritmul metodei Fadeev este caracterizat de relațiile:

$$\left\{ \begin{array}{l} A_1 = A; \quad c_1 = -\text{Tr}(A_1); \quad B_1 = A_1 + c_1 I_n \\ A_2 = AB_1; \quad c_2 = -\frac{1}{2}\text{Tr}(A_2); \quad B_2 = A_2 + c_2 I_n \\ \dots \\ A_n = AB_{n-1}; \quad c_n = -\frac{1}{n}\text{Tr}(A_n); \quad B_n = A_n + c_n I_n \end{array} \right. \quad (18)$$

Se demonstrează ușor următoarele afirmații:

$$\begin{cases} B_n = O_n; \\ A^{-1} = -\frac{1}{c_n} B_{n-1}, \text{ dac\u0103 } c_n \neq 0 \end{cases}$$

3.5. Metoda Danilevski

Fie $A \in R^{n \times n}$. Ne propunem s\u0103 ob\u021binem p_A .

Defini\u021bia 7. Spunem c\u0103 matricea $F \in R^{n \times n}$ are forma normal\u0103 Frobenius, dac\u0103:

$$F = \begin{pmatrix} f_1 & f_2 & \dots & f_{n-1} & f_n \\ 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \end{pmatrix} \quad (19)$$

Proprietatea 2.

$$p_F(\lambda) = \lambda^n - f_1 \lambda^{n-1} - f_2 \lambda^{n-2} - \dots - f_{n-1} \lambda - f_n \quad (20)$$

(deci o matrice cu form\u0103 normal\u0103 Frobenius are \u021n prima linie opu\u0219ii coeficien\u021bilor polinomului s\u0103u caracteristic).

Observa\u021bie. Demonstrarea rela\u021biei (20) se face dezvolt\u0103nd determinantul matricei $(\lambda I - F)$ dup\u0103 prima linie.

Metoda lui Danilevski, pentru aflarea polinomului caracteristic al matricei date, A , const\u0103 \u021n aducerea matricei A , la forma normal\u0103 Frobenius prin procedee de asem\u0103nare. Acest lucru se realizeaz\u0103 \u021n $n-1$ etape, la fiecare etap\u0103, ob\u021bin\u0103nd c\u0103te o linie din matricea F , dat\u0103 de (19), de la ultima linie p\u0103n\u0103 la prima.

Astfel:

Etapa 1. Presupunem $a_{n,n-1} \neq 0$ pentru a realiza ultima linie din matricea F .

Observa\u021bie. Cazul $a_{n,n-1} = 0$ (posibil practic) va fi tratat la cazuri particulare.

Cu presupunerea $a_{n,n-1} \neq 0$, vom prelucra matricea A astfel:

- \u021mp\u0103r\u021bim elementele coloanei $n-1$ prin $a_{n,n-1}$

- apoi, din fiecare coloană j , $1 \leq j \leq n$, $j \neq n-1$ vom scădea coloana $n-1$ (obținută anterior) înmulțită cu a_{nj} .

Adică:

$$\begin{cases} a'_{i,n-1} = a_{i,n-1} / a_{n,n-1}, & 1 \leq i \leq n, \\ a'_{ij} = a_{ij} - a'_{i,n-1} \cdot a_{nj}, & 1 \leq j \leq n, j \neq n-1, 1 \leq i \leq n. \end{cases} \quad (21)$$

Observație. Linia n a matricei A se înlocuiește cu ultima linie din F , adică $0 \ 0 \ \dots \ 1 \ 0$ (vezi formulele (21) pentru $i = n$).

Observație. Transformările (21) sunt echivalente cu înmulțirea la dreapta a matricei A cu matricea:

$$M_{n-1} = \begin{pmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & 1 & a_{nn} \\ a_{n,n-1} & a_{n,n-1} & \dots & a_{n,n-1} & a_{n,n-1} \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix} \text{ linia „}n-1\text{”} \quad (22)$$

Deci, matricea M_{n-1} diferă de matricea unitate de ordinul n , doar în linia „ $n-1$ ”, cu componența din (22).

Observație. Dacă notăm cu $B_{n-1} = A \cdot M_{n-1}$, atunci ultima linie a matricei B_{n-1} este $0 \ 0 \ \dots \ 1 \ 0$ și elementele celorlalte linii i , $1 \leq i \leq n-1$ se calculează cu ajutorul relațiilor (21).

Observație. Matricea M_{n-1} este inversabilă și

$$M_{n-1}^{-1} = \begin{pmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{n,n-1} & a_{nn} \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix} \text{ linia „}n-1\text{”} \quad (23)$$

Deci, matricea M_{n-1}^{-1} diferă de matricea unitate de ordinul n doar în linia $n-1$ care este ultima linie din matricea A .

Cu toate aceste precizări făcute, în prima etapă se obține matricea $C_{n-1} = M_{n-1}^{-1} \cdot B_{n-1}$, adică $C_{n-1} = M_{n-1}^{-1} \cdot A \cdot M_{n-1}$, care este asemenea cu matricea A ($A \sim C_{n-1}$).

Observație. Cu excepția liniei „n-1“, matricele C_{n-1} și B_{n-1} au aceleași elemente. Elementele liniei „n-1“ din matricea C_{n-1} se obțin înmulțind elementele liniei „n-1“ din matricea M_{n-1}^{-1} , (23), cu elementele fiecărei coloane din B_{n-1} .

Obținem, astfel, relațiile:

$$a'_{n-1,j} = \sum_{i=1}^n a_{ni} \cdot a'_{ij}, \quad 1 \leq j \leq n \quad (24)$$

Relațiile (24) se referă la linia „n-1“ din matricea C_{n-1} , restul liniilor matricei C_{n-1} fiind preluate din B_{n-1} , fără nici-o modificare.

Etapa 2. În etapa a doua se reiau calculele bazate pe relațiile (21) + (24), presupunând $a'_{n-1,n-2} \neq 0$ ($a'_{n-1,n-2}$ este element din C_{n-1}). Cu alte cuvinte se construiește matricea:

$$C_{n-2} = M_{n-2}^{-1} \cdot C_{n-1} \cdot M_{n-2},$$

asemenea cu matricea C_{n-1} , deci asemenea cu A (relația de asemănare fiind relație de echivalență). În plus, ultimele două linii din C_{n-2} sunt identice cu ultimele două linii din F .

Observație. Putem scrie $C_{n-2} = M_{n-2}^{-1} (M_{n-1}^{-1} \cdot A \cdot M_{n-1}) \cdot M_{n-2} = (M_{n-2}^{-1} \cdot M_{n-1}^{-1}) \cdot A \cdot (M_{n-1} \cdot M_{n-2})$ – folosind asociativitatea produsului matricelor. Deci $A \sim C_{n-2}$.

Observație. Matricele M_{n-2}, M_{n-2}^{-1} se calculează asemănător matricelor M_{n-1}, M_{n-1}^{-1} din (22) respectiv (23), folosind linia „n-1“ din C_{n-1} pentru a obține linia „n-2“ din M_{n-2} , respectiv M_{n-2}^{-1} . Adică, linia „n-2“ din M_{n-2} este:

$$\frac{a'_{n-1,1}}{a'_{n-1,n-2}} \quad \frac{a'_{n-1,2}}{a'_{n-1,n-2}} \quad \dots \quad \frac{1}{a'_{n-1,n-2}} \quad \frac{a'_{n-1,n-1}}{a'_{n-1,n-2}} \quad \frac{a'_{n-1,n}}{a'_{n-1,n-2}}$$

Analog, linia „n-2“ din M_{n-2}^{-1} este:

$$a'_{n-1,1} \quad a'_{n-1,2} \quad \dots \quad a'_{n-1,n-2} \quad a'_{n-1,n-1} \quad a'_{n-1,n} \quad \text{ș.a.m.d.}$$

Etapa n-1. În etapa n-1 se obține matricea $C_1 = M_1^{-1} \cdot C_2 \cdot M_1$ (în ipoteza $c_{21} \neq 0$, c_{21} fiind element din C_2), $A \sim C_1$, C_1 având forma (19) deoarece se obțin ultimele n-1 linii din (19). Relația dintre C_1 și A este dată de:

$$C_1 = (M_1^{-1} M_2^{-1} \dots M_{n-1}^{-1}) A \cdot (M_{n-1} \cdot M_{n-2} \dots M_1)$$

Notând $S = M_{n-1} M_{n-2} \dots M_1 \Rightarrow C_1 = S^{-1} \cdot A \cdot S$.

4. APROXIMAREA FUNCȚIILOR

Introducere

Vom încerca în cele ce urmează să prezentăm pe scurt, procesul de aproximare a unei funcții cunoscută printr-un număr finit de valori.

Printr-un anumit procedeu, utilizarea unor aparate de măsură de exemplu, putem cunoaște valorile unei funcții (temperatură, tensiuni sau intensitățile electrice, magnetism, etc.) cu expresie necunoscută pentru anumite valori ale variabilei.

Vom nota cu x_0, x_1, \dots, x_n valorile pentru care au avut loc măsurătorile și cu f_0, f_1, \dots, f_n valorile măsurate, adică: $f_i = f(x_i)$, $i = \overline{0, n}$.

Există situații în care este necesară cunoașterea valorii funcției f într-un punct \bar{x} , $\bar{x} \in [x_0, x_n]$, $\bar{x} \neq x_i$, $i = \overline{0, n}$, care nu a făcut deci parte din setul de măsurători inițiale.

Dacă reluarea măsurătorilor pentru determinarea valorii $f(\bar{x})$ nu mai este posibilă din diverse motive, de exemplu acela că fenomenul nu se mai poate repeta, suntem nevoiți ca pe baza datelor acumulate să căutăm o aproximare pentru $f(\bar{x})$.

Având în vedere faptul că pot exista mai multe valori de genul \bar{x} , vom căuta un aproximant pentru funcția f pe intervalul studiat $[x_0, x_n]$, care să aproximeze de fiecare dată valorile $f(\bar{x})$, pentru \bar{x} dat.

Vom expune în cadrul acestui capitol două metode de construcție a aproximantului:

a) ca polinom de interpolare (când n este mic!);

b) ca element de cea mai bună aproximare obținut prin metoda celor mai mici pătrate (când n este mare!).

Aproximanții construiți pot de asemenea înlocui funcția în operații de integrare și derivare aproximativă așa cum se va prezenta în capitolele următoare.

Aproximarea funcțiilor prin interpolare

Fie „ $n+1$ ” puncte **distincte**, x_i , $i = \overline{0, n}$ (numite noduri) pe un interval $I \subset \mathbb{R}$, și $f(x_i) = f_i$, $i = \overline{0, n}$ valorile **date** ale funcției reale f în aceste puncte ($f : \mathbb{R} \rightarrow \mathbb{R}$).

Să determinăm un polinom de grad **cel mult** egal cu n , notat P_n care să treacă prin punctele (x_i, f_i) , $i = \overline{0, n}$.

Un asemenea polinom îl numim polinom de interpolare, iar despre f spunem că se aproximează polinomial pe I .

Se pun următoarele întrebări:

a) dacă există, care este expresia polinomului de interpolare în funcție de x_i și f_i , $i = \overline{0, n}$?

b) cât de bine aproximează polinomul de interpolare funcția f pe I ?

Să încercăm să răspundem la aceste întrebări.

4.1. Polinomul de interpolare Lagrange

4.1.1. Construirea polinomului

Calea directă de determinare a polinomului de interpolare ar fi folosirea expresiei generale a polinomului de gradul n , coeficienții lui (în număr de $n+1$) urmând să rezulte din cele $n+1$ condiții.

$$P_n(x_i) = f_i, \quad i = \overline{0, n} \quad (1)$$

Adică, fie $P_n(x) = a_0x^n + a_1x^{n-1} + \dots + a_n$. Astfel (1) reprezintă un sistem liniar cu $n+1$ ecuații (deci cele $n+1$ condiții (1)) și $n+1$ necunoscute a_0, a_1, \dots, a_n (coeficienții polinomului de interpolare). Determinantul matricei acestui sistem este determinantul Vandermonde $V(x_0, x_1, \dots, x_n) \neq 0$, deoarece nodurile sunt distincte. Rezultă soluția unică a_0, a_1, \dots, a_n .

Deci polinomul de interpolare există și este unic, expresia sa fiind dependentă de $x_i, f_i, i = \overline{0, n}$.

Fiind destul de laborioasă această cale, vom prefera o altă cale care conduce direct la expresia polinomului de interpolare.

Să notăm cu $l_k, k = \overline{0, n}$, un polinom de grad n pentru care $l_k(x_j) = \delta_{kj}$ (simbolul lui Kronecker), $(\forall) j = \overline{0, n}$.

Așadar

$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j}, \quad k = \overline{0, n} \quad (2)$$

$$\text{adică } l_k(x) = \frac{(x - x_0) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)}, \quad k = \overline{0, n}$$

Observație. Polinoamele de grad n , definite de (2) se numesc polinoame fundamentale Lagrange (justificarea notației).

Este valabilă următoarea proprietate:

Proprietatea 1. Polinomul $P_n(x) = \sum_{k=0}^n f_k l_k(x)$ este polinom de

interpolare (numit polinom de interpolare Lagrange).

Așadar, polinomul de interpolare Lagrange:

$$P_n(x) = \sum_{k=0}^n f_k \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j} \quad (3)$$

este un polinom de interpolare pentru datele $(x_i, f_i), i = \overline{0, n}$.

Proprietatea 2. Polinomul de interpolare corespunzător la „ $n+1$ ” puncte distincte, există și este unic (sau polinomul de interpolare Lagrange constituie **singurul** polinom de grad cel mult n ce interpolează datele $(x_i, f_i), i = \overline{0, n}$).

4.2. Polinomul de interpolare Newton

4.2.1. Diferențe divizate

Fie „n+1“ puncte **distincte** $x_i, i = \overline{0, n}$ (numite noduri) pe un interval $I \subset \mathbb{R}$ și $f(x_i) = f_i, i = \overline{0, n}$ valorile date ale funcției reale f în aceste puncte.

Definiția 1. Definim diferențele divizate de ordin $k, k = \overline{0, n}$, recursiv astfel:

– diferențele divizate de ordin **zero** coincid cu valorile f_i , ale funcției f în nodurile $x_i, i = \overline{0, n}$ (deci sunt „n+1“ diferențe divizate de ordin 0). Le notăm $f(x_i)$;

– diferențele divizate de ordin **unu** se definesc prin egalitatea

$$f(x_i; x_{i+1}) = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}, 0 \leq i \leq n-1,$$

adică $f(x_0; x_1), f(x_1; x_2), \dots, f(x_{n-1}, x_n)$ (deci sunt „n“ diferențe divizate de ordin 1);

– diferențele divizate de ordin k , se obțin din diferențe divizate de ordin $k-1$ după formula:

$$f(x_0; x_1; \dots; x_k) = \frac{f(x_1; x_2; \dots; x_k) - f(x_0; x_1; \dots; x_{k-1})}{x_k - x_0} \quad (7)$$

Observație. Există o singură diferență divizată de ordinul n și anume $f(x_0; x_1; \dots, x_n) = \frac{f(x_1; x_2; \dots; x_n) - f(x_0; x_1; \dots; x_{n-1})}{x_n - x_0}$.

Observație. Din definiția 1 rezultă că putem trece diferențele divizate (de toate ordinele) într-un tablou astfel:

x_i	dif.div. de ord.0	dif.div. de ord.1	ord.2	...	ord. n
x_0	$f(x_0) = f_0$	$f(x_0; x_1)$	$f(x_0; x_1; x_2)$		$f(x_0; x_1; \dots; x_n)$
x_1	$f(x_1) = f_1$	$f(x_1; x_2)$	$f(x_1; x_2; x_3)$		–
x_2	$f(x_2) = f_2$	$f(x_2; x_3)$	$f(x_2; x_3; x_4)$		–
\vdots	\vdots	\vdots	\vdots		\vdots
x_{n-2}	$f(x_{n-2}) = f_{n-2}$	$f(x_{n-2}; x_{n-1})$	$f(x_{n-2}; x_{n-1}; x_n)$		–
x_{n-1}	$f(x_{n-1}) = f_{n-1}$	$f(x_{n-1}; x_n)$	–		–
x_n	$f(x_n) = f_n$	–	–		–

folosind formulele (7).

Deci notând d_{ij} , elementele acestui tablou
 $\Rightarrow d_{i0} = f_i, i = \overline{0, n}$

$$d_{ij} = \frac{d_{i+1,j-1} - d_{i,j-1}}{x_{j+i} - x_i}, \quad j = \overline{1, n}; \quad i = \overline{0, n-j}$$

Observație. Se poate demonstra, prin inducție după m , că

$$f(x_0; \dots; x_m) = \sum_{j=0}^m \frac{f(x_j)}{\prod_{\substack{i=0 \\ i \neq j}}^m (x_j - x_i)} \quad (8)$$

ceea ce arată că diferențele divizate sunt simetrice în argumente.

Polinomul notat

$$N_n(x) = f(x_0) + \sum_{k=0}^{n-1} f(x_0; \dots; x_{k+1}) \prod_{j=0}^k (x - x_j) \quad (12)$$

se numește **polinomul lui Newton de interpolare cu diferențe divizate**. Din (9), combinând cu (5), deducem:

$$f(x; x_0; \dots; x_n) = \frac{f^{(n+1)}(\xi)}{(n+1)!},$$

ξ este intermediar punctelor x, x_0, \dots, x_n .

Observații:

1. Polinomul de interpolare Lagrange, cu exprimarea (3) este identic cu polinomul de interpolare al lui Newton cu exprimarea (12).

2. Din tabloul diferențelor divizate, se observă că ne interesează doar prima linie.

4.4. Aproximarea cu funcții spline cubice

4.4.1. Construcția aproximantului spline cubic

Dacă gradul polinomului este 3 vom spune că aproximăm cu **spline cubic**, definit prin condițiile de mai jos:

1. $S(x_i) = f_i, 0 \leq i \leq n$;
2. S, S', S'' continue pe $[a, b]$;

3. S polinom de gradul trei pe fiecare subinterval $[x_{i-1}, x_i]$, $1 \leq i \leq n$.

Notăm:

S_i – restricția lui S pe $[x_{i-1}, x_i]$ adică $S_i(x) = S(x)$, $x \in [x_{i-1}, x_i]$.

$u_i = S''(x_i)$, $0 \leq i \leq n$, deci $u_i = S''_i(x_i) = S''_{i+1}(x_i)$.

Forma generala a ramurii „i” a spline-ului cubic este

$$S_i(x) = \frac{u_i(x-x_{i-1})^3 + u_{i-1}(x_i-x)^3}{6h_i} + \left(f_i - u_i \frac{h_i^2}{6} \right) \frac{x-x_{i-1}}{h_i} + \left(f_{i-1} - u_{i-1} \frac{h_i^2}{6} \right) \frac{x_i-x}{h_i}, \quad (22)$$

Se observă continuitatea lui S.

În continuare, vom impune condițiile de continuitate și pentru S', adică

$$S'_i(x_i) = S'_{i+1}(x_i), \quad 1 \leq i \leq n-1$$

de unde avem:

$$u_{i-1} \frac{h_i}{6} + u_i \frac{h_i + h_{i+1}}{3} + u_{i+1} \frac{h_{i+1}}{6} = \frac{f_{i+1} - f_i}{h_{i+1}} - \frac{f_i - f_{i-1}}{h_i} \quad (23)$$

$$1 \leq i \leq n-1$$

Relațiile (23) reprezintă un sistem de n-1 ecuații cu n+1 necunoscute u_0, u_1, \dots, u_n , pentru a căror determinare există două cazuri practice:

I. $S''(a) = S''(b) = 0$, adică $u_0 = u_n = 0$. Rezultă n – 1 ecuații cu n – 1 necunoscute: u_1, u_2, \dots, u_{n-1} , adică am redus numărul de necunoscute.

$$\text{II. } \left. \begin{array}{l} S'(a) = f'(a) = f'_0 \\ S'(b) = f'(b) = f'_n \end{array} \right\} \Rightarrow u_0 \frac{h_1}{3} + u_1 \frac{h_1}{6} = \frac{f_1 - f_0}{h_1} - f'_0 \quad (24)$$

și

$$u_{n-1} \frac{h_n}{6} + u_n \frac{h_n}{3} = f'_n - \frac{f_n - f_{n-1}}{h_n}, \quad (25)$$

adică mărime numărul de ecuații și deci (24) + (23) + (25) reprezintă un sistem de $n+1$ ecuații cu $n+1$ necunoscute: u_0, u_1, \dots, u_n .

În ambele cazuri I și II sistemele obținute sunt tridiagonale simetrice cu diagonala principală dominantă.

4.5. Metoda celor mai mici pătrate – cazul funcțiilor tabelate

4.5.1. Construcția elementului de cea mai bună aproximare

Fiind dată o funcție $f : [a, b] \rightarrow \mathbb{R}$ vom spune că este tabelată atunci când se cunosc valorile $f(x_i) = f_i$ ale funcției f pe un sistem de noduri echidistante sau nu: $a = x_1 < x_2 < \dots < x_n = b$.

Metoda celor mai mici pătrate (m.c.m.p.) își propune ca pentru funcțiile tabelate să determine un aproximant care să aproximeze suficient de bine funcția în orice punct din intervalul $[a, b]$, fără să coincidă cu funcția în mod expres în vreun punct.

Alegerea aproximantului se va face dintr-o familie de funcții liniar independentă, numită sistem de generatori sau funcții standard și evident forma lui va depinde de familia din care face parte.

Cea mai comodă alegere pentru aproximant va fi de formă polinomială.

Elementul de cea mai bună aproximare în sensul celor mai mici pătrate se va determina din condiția de minimizare a abaterii dintre funcția f și aproximantul g :

$$\|f - g\|^2 = \sum_{i=1}^n |f(x_i) - g(x_i)|^2 \quad (33)$$

unde:

$$\|f\|^2 = \sum_{i=1}^n f_i^2 .$$

Aproximantul g va fi generat de o familie finită de funcții liniar independente: $\varphi_1, \varphi_2, \dots, \varphi_m$ adică g va avea forma:

$$g(x) = \sum_{i=1}^m c_i \varphi_i(x) \quad (34)$$

Coefficienții $c_i \in \mathbb{R}$ îi vom determina din condiția de minimizare a erorii de aproximare:

$$\min R(c_1, c_2, \dots, c_m) = \min \sum_{j=1}^n \left(f(x_j) - \sum_{i=1}^m c_i \varphi_i(x_j) \right)^2 \quad (35)$$

echivalentă cu rezolvarea sistemului liniar:

$$\frac{\partial R}{\partial c_1} = 0; \quad \frac{\partial R}{\partial c_2} = 0; \quad \dots \quad \frac{\partial R}{\partial c_m} = 0 \quad (36)$$

sau

$$\sum_{i=1}^m c_i \sum_{j=1}^n \varphi_k(x_j) \varphi_i(x_j) = \sum_{j=1}^n \varphi_k(x_j) f(x_j), \quad k = \overline{1, m} \quad (37)$$

elementului de cea mai bună aproximare pentru o funcție dată.

4.5.2. Aproximarea funcțiilor prin metoda celor mai mici pătrate utilizând polinoame algebrice

Vom alege funcțiile standard ca polinoame algebrice, o primă posibilitate fiind ca:

$$\varphi_1(x) = 1; \quad \varphi_i(x) = x^{i-1}, \quad i \in \overline{2, m} \quad (41)$$

Ecuțiile sistemului (37) (care se mai numesc și ecuații normale) vor deveni:

$$\left\{ \begin{array}{l} nc_1 + \sum_{i=2}^m c_i \sum_{j=1}^n x_j^{i-1} = \sum_{j=1}^n f_j \\ \sum_{i=1}^m c_i \sum_{j=1}^n x_j^i = \sum_{j=1}^n x_j f_j \\ \dots\dots\dots \\ \sum_{i=1}^m c_i \sum_{j=1}^n x_j^{m+i-2} = \sum_{j=1}^n x_j^{m-1} f_j \end{array} \right. \quad (42)$$

unde x_1, x_2, \dots, x_n sunt punctele în care se cunoaște fenomenul fizic a cărui modelare se studiază.

5. EVALUAREA NUMERICĂ A INTEGRALELOR

Introducere

Vom studia în continuare metode de evaluare aproximativă a integralelor definite $\int_a^b f(x)dx$ pentru o funcție f integrabilă a cărei primitivă este dificil de evaluat.

Aproximarea integralei definite se va face cu formule de forma:

$$\int_a^b f(x)dx = \sum_{i=0}^n c_i f(x_i) + E(f),$$

(1)

unde: c_i sunt coeficienți, x_i sunt noduri echidistante alcătuind o diviziune a intervalului $[a, b]$ iar $E(f)$ este restul metodei de aproximare numerică.

Pentru fiecare metodă restul (eroarea) este specific, iar coeficienții c_i și nodurile x_i se aleg astfel încât restul să fie nul atunci când funcția se înlocuiește cu un polinom de un anumit grad m .

Există o varietate de metode de aproximare a integralelor cunoscute și sub numele de formule de cvadratură, clasificarea acestora făcându-se în funcție de modul de obținere.

a) Dacă pentru obținerea formulei de cvadratură (1) nodurile echidistante x_0, x_1, \dots, x_n sunt fixate și se determină coeficienții c_i , astfel încât $E(f)$ să fie nul pentru orice f – polinom cu gradul cel mult egal cu m , ordinul de exactitate va fi m .

Câteva exemple de metode pe care le vom prezenta în continuare din cadrul acestei familii sunt: metoda trapezului, metoda Simpson, metoda lui Newton, cunoscute în general sub denumirea de formule de tip Newton-Côtes.

Aceste metode permit și evaluarea integralelor din funcții a căror expresie nu este cunoscută dar avem acces prin măsurători la valorile care ne interesează.

b) Dacă în formula (1) toți coeficienții c_i sunt egali: $c_i = c$, $i = \overline{0, n}$ și se determină nodurile x_i astfel încât restul $E(f)$ să fie nul pentru orice f – polinom de gradul $\leq m$, atunci se obține o formulă de tip Cebîșev:

$$\int_a^b f(x) dx = c \sum_{i=0}^n f(x_i) + E(f)$$

de ordinul m .

c) Dacă în formula (1) se determină atât coeficienții c_i cât și valorile x_i iar restul va fi nul pentru orice f – polinom de grad maxim $2m$ obținem formule de cvadratură de tip Gauss având ordinul de exactitate $2m$.

Formulele de tip Newton-Côtes vor fi expuse în continuare iar cele de tip Cebîșev și Gauss vor fi expuse în capitolul 7.

5.1. Metoda trapezului

Se obține aproximând funcția de integrat cu un polinom de interpolare Lagrange construit pe nodurile a și b .

$$\text{Deci } \int_a^b f(x) dx \approx \frac{b-a}{2} [f(a) + f(b)]$$

Geometric, membrul drept reprezintă aria trapezului cu bazele $f(a)$, $f(b)$ și înălțimea $b - a$ – care se obține aproximând curba dată pe $[a, b]$ printr-o dreaptă (de aici denumirea metodei).

De la interpolarea Lagrange știm că

$$f(x) = L(x) + \frac{f''(\xi)}{2} (x-a)(x-b), \text{ în condiția } f \in C^2[a, b],$$

$$\xi \in (a, b)$$

$$\text{Deci } \int_a^b f(x) dx = \frac{b-a}{2} [f(a) + f(b)] + E_T(f),$$

$$\text{cu } E_T(f) = \frac{1}{2} \int_a^b f''(\xi) (x-a)(x-b) dx$$

eroarea de aproximare în formula trapezului.

De aceea vom considera o diviziune a intervalului $[a, b]$, prin punctele echidistante (pentru comoditatea formulei pe care o vom stabili) $x_i = x_0 + ih, i = \overline{0, n}$, unde $x_0 = a, x_n = b$, iar pasul $h = \frac{b-a}{n}$.

$$\text{Obținem: } \int_a^b f(x) dx = \frac{h}{2} \left[f(x_0) + 2 \sum_{i=1}^{n-1} f(x_i) + f(x_n) \right] + E_T(f),$$

$$\text{unde } E_T(f) = -\frac{h^3}{12} [f''(\xi_1) + \dots + f''(\xi_n)].$$

Cum f'' este continuă rezultă că $(\exists)\xi \in (a, b)$ astfel încât

$$f''(\xi) = \frac{f''(\xi_1) + \dots + f''(\xi_n)}{n} \text{ și rezultă că}$$

$$E_T(f) = -\frac{h^3}{12} n f''(\xi) = -\frac{h^2}{12} (b-a) f''(\xi) = -\frac{(b-a)^3}{12n^2} f''(\xi)$$

Notând cu $M_2 = \max_{x \in [a, b]} |f''(x)|$, obținem:

$$|E_T(f)| \leq \frac{(b-a)^3 M_2}{12n^2},$$

deci putem mărgini eroarea în formula trapezului.

6. APROXIMAREA NUMERICĂ A SOLUȚIILOR ECUAȚIILOR DIFERENȚIALE

Introducere

Ecuatiile diferențiale ordinare sau cu derivate parțiale constituie modelele matematice pentru majoritatea problemelor ingineresti: studiul eforturilor la care sunt supuse elementele de rezistență: bare, grinzi, plăci subțiri, groase, conducte; studiul problemelor de câmp electric în dielectrics, câmp magnetic, câmp termic, propagarea undelor de toate felurile și lista poate continua.

Odată stabilit fenomenul fizico-tehnic și ecuațiile diferențiale care îl guvernează, ca formă, coeficienți, condiții la limită (pe frontieră) rămâne de rezolvat ultima problemă: rezolvarea acestui

model matematic. Din diverse motive: neomogenitățile fizice, frontiere cu geometrie dificilă, număr de necunoscute, etc., rezolvarea o vom face căutând o soluție aproximativă cu ajutorul unui calculator.

Vom expune în cadrul acestui capitol principalele metode aproximative care se pretează la implementarea pe calculator pentru rezolvarea ecuațiilor diferențiale.

Pentru ecuații diferențiale ordinare acestea se pot clasifica în două mari tipuri:

- metode unipas (Euler, Runge-Kutta) în care determinarea soluției în fiecare punct se va obține direct;
- metode multipas, sau predictor-corector în care valoarea soluției în fiecare punct se va „predicte” și apoi se va corecta iterativ.

Evident este vorba de soluții aproximative pe care nu avem cum să le comparăm cu o soluție exactă, deoarece practic aceasta este imposibil de găsit.

De aceea în practică trebuie să procedăm cu atenție pentru alegerea algoritmilor cei mai potriviți pentru problema concretă de rezolvat.

6.1. METODE DE TP EULER

6.1.1. Metoda Euler

Se consideră ecuația diferențială:

$$y' = f(x, y) \quad (1)$$

cu condiția inițială:

$$y(x_0) = y_0, \quad (2)$$

unde funcția f este definită într-un domeniu D din planul xOy .

Metoda lui Euler propune aproximarea soluției printr-o linie poligonală în care fiecare segment are direcția data de capatul sau stang. Astfel se consideră nodurile echidistante $x_i = x_0 + ih$, $i = \overline{0, n}$.

Considerând cunoscută aproximarea y_i a soluției problemei (1)+(2) în x_i , procedeul de aproximare Euler, poate fi acum rezumat astfel:

$$\begin{cases} f_i = f(x_i, y_i); \\ x_{i+1} = x_i + h; \\ y_{i+1} = y_i + hf_i; \end{cases} \quad (4)$$

Neglijarea termenilor de ordin superior în (4) face ca metoda să fie comodă în calcul, dar puțin precisă, erorile cumulându-se la fiecare pas.

Metoda se poate aplica și dacă nodurile x_i nu sunt echidistante, având la fiecare iterație alt pas, h în acest caz.

6.1.2. Metoda Euler modificată

Considerăm pe $[x_i, x_{i+1}]$ ca direcție a segmentului $M_i M_{i+1}$ direcția definită de punctul de la mijlocul segmentului (nu de extremitatea stângă ca în formula inițială) se obține metoda Euler modificată. Dacă x_i, y_i sunt valori calculate, procesul iterativ este următorul:

$$\begin{cases} x_i = x_0 + ih; f_i = f(x_i, y_i); x_{i+\frac{1}{2}} = x_i + \frac{h}{2}; \\ y_{i+\frac{1}{2}} = y_i + \frac{h}{2} f_i; f_{i+\frac{1}{2}} = f\left(x_{i+\frac{1}{2}}, y_{i+\frac{1}{2}}\right); y_{i+1} = y_i + hf_{i+\frac{1}{2}}; \end{cases}$$

6.1.4. Metoda Euler–Heun

Vom încheia trecerea în revistă a variantelor Euler pentru rezolvarea ecuațiilor de tipul (1) prezentând o ultimă variantă de ordinul doi a metodei lui Euler și anume aceea propusă de Heun, [26].

Presupunând calculată valoarea y_i la pasul x_i , se propune pentru calcularea soluției la pasul x_{i+1} expresia:

$$y_{i+1} = y_i + \frac{h}{4} \left[f(x_i, y_i) + 3f\left(x_i + \frac{2}{3}h, y_i + \frac{2}{3}hf(x_i, y_i)\right) \right].$$

6.2. Metode de tip Runge-Kutta (R-K) pentru problema Cauchy

Introducere

Dacă metodele de tip Euler prezentate anterior au mai mult un caracter didactic pentru familiarizarea cititorului cu problematica, metodele de tip Runge Kutta pe care le vom expune în continuare pot fi utilizate cu succes în rezolvarea problemelor concrete din practica ingierească.

Metodele Runge-Kutta au trei proprietăți distincte [30]:

1. Sunt metode directe, adică pentru determinarea aproximării soluției la pasul $i+1$ avem nevoie de informațiile existente în punctul precedent x_i, y_i .

2. Sunt identice cu seriile Taylor până la termenii h^n , unde h este pasul curent iar n este diferit pentru metode diferite din această familie și definește ordinul metodei.

Metodele de tip Euler pot fi și ele incluse în familia Runge-Kutta și putem astfel observa că metoda Euler este o metodă R-K de ordinul întâi iar metodele Euler-Cauchy și Euler-Heun sunt metode R-K de ordinul 2.

3. Metodele de tip R-K necesită în procesul de calcul doar evaluarea funcției din membrul drept în diferite puncte nu și a derivatelor ei. Aceasta este, alături de precizia bună a metodelor de ordin 3, 4, 5 motivul principal al utilizării lor frecvente în practică.

Oricum și pentru metodele din această familie, formulele de evaluare a erorii păstrează doar un caracter teoretic neputând fi aplicate în practică decât pentru exemple simple cu caracter didactic.

O formulă Runge-Kutta de ordinul 2:

$$\begin{cases} y(x+h) = y(x) + \frac{1}{4}(k_1 + 3k_2) \\ k_1 = hf(x, y) \\ k_2 = hf\left(x + \frac{2}{3}h, y + \frac{2}{3}k_1\right) \end{cases} \quad (10'')$$

cunoscută ca formula Euler-Heun.

O formulă de ordinul 3 este următoarea:

$$y(x+h) = y(x) + \frac{1}{4}(k_1 + 3k_3)$$

$$k_1 = hf(x, y)$$

$$k_2 = hf\left(x + \frac{h}{3}, y + \frac{k_1}{3}\right)$$

$$k_3 = hf\left(x + \frac{2h}{3}, y + \frac{2k_2}{3}\right)$$

Formule Runge-Kutta de ordinul 4:

$$y(x+h) = y(x) + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4)$$

$$k_1 = hf(x, y)$$

$$k_2 = hf\left(x + \frac{h}{2}, y + \frac{k_1}{2}\right)$$

$$k_3 = hf\left(x + \frac{h}{2}, y + \frac{k_2}{2}\right)$$

$$k_4 = hf(x+h, y+k_3)$$

formula *propriu-zisă* Runge-Kutta și

$$y(x+h) = y(x) + \frac{1}{8}(k_1 + 3k_2 + 3k_3 + k_4)$$

$$k_1 = hf(x, y)$$

$$k_2 = hf\left(x + \frac{h}{3}, y + \frac{k_1}{3}\right)$$

$$k_2^j = hf_j \left(x_i + \frac{h}{2}, y_{1,i} + \frac{k_1^1}{2}, \dots, y_{n,i} + \frac{k_1^n}{2} \right), j = \overline{1, n};$$

$$k_3^j = hf_j \left(x_i + \frac{h}{2}, y_{1,i} + \frac{k_2^1}{2}, \dots, y_{n,i} + \frac{k_2^n}{2} \right), j = \overline{1, n};$$

$$k_4^j = hf_j \left(x_i + h, y_{1,i} + k_3^1, \dots, y_{n,i} + k_3^n \right), j = \overline{1, n};$$

6.3. Metoda diferențelor finite pentru probleme Sturm-Liouville

Considerăm problema bilocală (Sturm-Liouville) dată de ecuația:

$$u(x)y''(x) + v(x)y'(x) + w(x)y(x) = f(x) \quad (24)$$

și condițiile de limită

$$\begin{cases} y(a) = \alpha \\ y(b) = \beta \end{cases} \quad (25)$$

unde $x \in [a, b]$ și presupunem că u, v, w, f sunt continue pe $[a, b]$, $u(x) > 0$, $w(x) < 0$ ($\forall x \in [a, b]$) iar (24)+(25) are soluție unică pe $[a, b]$.

Fie Δ o diviziune echidistantă de pas h a lui $[a, b]$, $x_0 = a$, $x_i = a + ih$, $0 \leq i \leq n$ cu $x_n = b$.

Dacă în ecuația (24) facem $x = x_i$, $1 \leq i \leq n - 1$ și folosim diferențele finite centrate pentru aproximarea operatorilor diferentiați avem

$$u(x_i) \frac{y(x_{i+1}) - 2y(x_i) + y(x_{i-1}))}{h^2} + v(x_i) \frac{y(x_{i+1}) - y(x_{i-1}))}{2h} + w(x_i)y(x_i) = f(x_i)$$

$$\text{Iar condițiile (25) devin } \begin{cases} y(x_0) = \alpha \\ y(x_n) = \beta \end{cases}$$

Notând simplificat $u_i = u(x_i)$, $v_i = v(x_i)$, $w_i = w(x_i)$, $f_i = f(x_i)$, $y_i = y(x_i)$, $0 \leq i \leq n$, avem $y_0 = \alpha$, $y_n = \beta$, iar

$$u_i \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} + v_i \frac{y_{i+1} - y_{i-1}}{2h} + w_i y_i = f_i, \quad 1 \leq i \leq n-1$$

sau

$$\left(\frac{u_i}{h^2} - \frac{v_i}{2h} \right) y_{i-1} + \left(w_i - \frac{2u_i}{h^2} \right) y_i + \left(\frac{u_i}{h^2} + \frac{v_i}{2h} \right) y_{i+1} = f_i, \quad (28)$$

$1 \leq i \leq n-1$, care reprezintă un sistem liniar de $n-1$ ecuații și $n-1$ necunoscute y_1, y_2, \dots, y_{n-1} , matricea sistemului fiind tridiagonală dominant diagonală.