

METODE ȘI PROGRAME DE CALCUL NUMERIC

-NOTE DE CURS-

GRECU LUMINIȚA

I. CONCEPTE DE BAZĂ ȘI TIPURI DE ERORI

I.1. INTRODUCERE

Metodele numerice sunt acele tehnici care permit transformarea modelelor matematice în modele numerice (ce operează pe spații finite), și presupun algoritmi ce pot fi ușor transformați în coduri sursă, folosind diferite limbaje de programare, iar prin intermediul acestora rezolvarea problemelor cărora li se adresează cu ajutorul calculatorului. Pe scurt putem spune că ele permit rezolvarea problemelor matematice cu ajutorul calculatorului.

Trecerea de la modelul matematic la cel numeric se face, în general, pe baza unor aproximări și de aceea soluția oferită de algoritmi rezultați ca urmare a aplicării metodelor numerice este, de cele mai multe ori, una aproximativă. Ea poartă numele de soluție numerică și, în cele mai multe situații, este diferită de cea exactă.

I.2. TIPURI DE ERORI

Pentru a înțelege mai bine principalele surse de erori care pot să apară în rezolvarea numerică a problemelor este bine să înțelegem mai întâi modul în care numerele sunt înregistrate și depozitate în memoria unui calculator. Deoarece numai numere cu un număr finit de zecimale pot fi reprezentate într-un calculator, se folosește reprezentarea în virgulă mobilă a numerelor reale.

Reprezentarea în virgulă mobilă (sau flotantă) presupune că numerele sunt reprezentate în calculator sub forma:

$$y = \pm(0.x_1x_2 \dots x_t) \cdot b^e$$

$$y = \pm b^e \cdot \sum_{k=1}^t x_k b^{-k}$$

unde b reprezintă baza de numerație iar x_k sunt cifre din axa respectivă, adică valori din mulțimea: $\{0,1, \dots, b-1\}$.

Deci, reprezentarea în virgulă mobilă poate fi sintetizată astfel:

$$y = s \cdot M \cdot b^e,$$

unde s reprezintă bitul de semn (0 pentru + și 1 pentru -), M este mantisa, b este baza de numerație (de obicei 2) iar e este exponentul bazei.

Astfel pentru reprezentarea unui număr este necesar un cuvânt binar cu trei câmpuri: semnul, mantisa și exponentul.

Lungimea mantisei reprezintă precizia de reprezentare a numărului. Astfel numerele pot avea mai multe formate: formate cu precizie simplă (24 biți pentru mantisă) și precizie simplă extinsă (≥ 32 biți pentru mantisă), și formate cu precizie dublă (52 biți) și precizie dublă extinsă (≥ 64 biți).

Calculatorul operează cu numere normalizate, adică numere reprezentate astfel încât mantisa să fie un număr subunitar cu prima cifră după virgulă diferită de zero, deci astfel:

$$y = \pm 0.x_1x_2 \dots x_t, \quad y = \pm \sum_{k=1}^t x_k b^{-k}.$$

De exemplu numărul 1,5 admite mai multe reprezentări (folosind baza 2) printre care: $+1111$, $+111,1 \cdot 2^1$, $+11,11 \cdot 2^2$, $+1,111 \cdot 2^3$, $+0,1111 \cdot 2^4$ și $+0,001111 \cdot 2^6$.

Dintre aceste reprezentări ale numărului 1,5 penultima este cea normalizată.

Când se operează cu numere normalizate acestea trebuie să fie scrise cu ajutorul aceluiași exponent. Dacă numerele au exponenți diferiți se aduce numărul cu exponent mai mic la o formă corespunzătoare, ce utilizează exponentul mai mare. Se face operația propriu-zisă, și apoi se normalizează rezultatul.

Exemplul I.1. Să se adune următoarele două numere: $x = 4,156832$; $y = 246,548$, considerând că ele sunt reprezentate sub formă normalizată într-un sistem de calcul cu virgulă mobilă având baza de numerație $b = 10$, $t = 6$.

Soluție:

Forma normalizată a numerelor este:

$$x = 0,4156832 \cdot 10 = 0,415683 \cdot 10; y = 0,246548 \cdot 10^3.$$

Cifra 2 din scrierea lui x se pierde ținând cont că t , lungimea mantisei, este 6. Observăm că formele normalizate nu prezintă același exponent. Aducem numărul scris cu exponent mai mic la o formă în care apare exponentul mai mare, și astfel avem $x = 0,004156832 \cdot 10^3$. Păstrând doar 6 cifre semnificative obținem reprezentarea $x = 0,004156 \cdot 10^3$. Efectuăm suma numerelor adunând mantisele: $x + y = 0,250704 \cdot 10^3$.

Efectuând calculele în mod obișnuit obținem:

$$x + y = 250,704832.$$

Exemplul I.2. Se consideră un sistem de calcul cu virgulă mobilă având baza de numerație $b = 10$, $t = 4$ și o reprezentare normalizată a numerelor. Se consideră următoarele numere:

$$x_1 = 0.2146, x_2 = 3.175, x_3 = 15.421, x_4 = 176.86$$

Să se efectueze adunarea acestor numere :

a) în ordine crescătoare

b) în ordine descrescătoare

Soluție:

Formele normalizate ale numerelor sunt:

$$x_1 = 0.2146, x_2 = 0.3175 \cdot 10, x_3 = 0.1542 \cdot 10^2, x_4 = 0.1768 \cdot 10^3$$

a) Adunarea în ordine crescătoare

Se observă că numerele sunt date chiar în ordine crescătoare. Astfel algoritmul este: se adună x_1 cu x_2 , iar rezultatul cu x_3 . Noul rezultat obținut se adună cu x_4 . Notăm cu a_i ($i = 1, 2, 3$) rezultatele acestui proces.

$$a_1 = x_1 + x_2 = (0.0214 + 0.3175) \cdot 10 = 0.3389 \cdot 10,$$

$$a_2 = a_1 + x_3 = (0.0338 + 0.1542) \cdot 10^2 = 0.1880 \cdot 10^2,$$

$$a_3 = a_2 + x_4 = (0.0188 + 0.1768) \cdot 10^3 = 0.1956 \cdot 10^3,$$

b) Adunarea în ordine descrescătoare

Algoritmul de calcul este următorul: se adună x_4 cu x_3 , iar rezultatul cu x_2 . Noul rezultat obținut se adună cu x_1 . Rezultatele intermediare se notează cu b_i ($i = 1, 2, 3$)

$$b_1 = x_4 + x_3 = (0.0154 + 0.1768) \cdot 10^3 = 0.1922 \cdot 10^3$$

$$b_2 = b_1 + x_2 = (0.1922 + 0.0031) \cdot 10^3 = 0.1953 \cdot 10^3$$

$$b_3 = b_2 + x_1 = (0.1953 + 0.0002) \cdot 10^3 = 0.1955 \cdot 10^3$$

Concluzie: Adunarea numerelor în calculator nu mai are aceleași proprietăți ca operația cunoscută de adunare a numerelor reale. Soluția finală a calculelor generate de modelarea numerică, ce presupune efectuarea unor operații aritmetice, care se execută repetat și în număr finit, poate fi afectată de acestea.

Principalele clase de erori sunt următoarele: erori inerente, erori de metodă, erori de rotunjire și de trunchiere (sau reziduale).

Erorile inerente nu au legătură cu elaborarea modelului numeric. Ele se pot produce fie când se obține modelul fizic sau cel matematic, fie când se fac măsuratori, sau chiar când se introduc datele de intrare într-un program de calcul numeric.

Erorile de metodă sunt specifice folosirii metodelor numerice și algoritmilor de calcul.

Exemplul I.3. Să se evalueze numeric $\int_0^1 \frac{x^n}{2x+9} dx$, pentru $n=7$.

Soluție:

Notând cu I_n integrala precedentă se demonstrează cu ușurință că are loc relația de recurență: $I_n = \frac{1}{2n} - \frac{9}{2} I_{n-1}$.

$$\begin{aligned} \int_0^1 \frac{x^n}{2x+9} dx &= \int_0^1 \frac{\frac{1}{2} x^{n-1} (2x+9) - \frac{9}{2} x^{n-1}}{2x+9} dx = \frac{1}{2} \int_0^1 x^{n-1} dx - \frac{9}{2} \int_0^1 \frac{x^{n-1}}{2x+9} dx = \\ &= \frac{x^n}{2n} \Big|_0^1 - \frac{9}{2} I_{n-1} = \frac{1}{2n} - \frac{9}{2} I_{n-1} \end{aligned}$$

Calculăm I_7 .

Avem:

$$I_0 = \frac{1}{2} \ln(2x+9) \Big|_0^1 = \frac{1}{2} \ln \frac{11}{9} \cong 0,10034$$

$$I_1 = \frac{1}{2} - \frac{9}{2} \cdot 0,10034 = 0,04847$$

$$I_2 = 0,25 - \frac{9}{2} \cdot 0,04847 = 0,031885$$

$$I_3 = \frac{1}{6} - \frac{9}{2} \cdot 0,031885 = 0,023184$$

$$I_4 = \frac{1}{8} - \frac{9}{2} \cdot 0,023184 = 0,020672$$

$$I_5 = \frac{1}{10} - \frac{9}{2} \cdot 0,020672 = 0,006976$$

$$I_6 = \frac{1}{12} - \frac{9}{2} \cdot 0,006976 = 0,051941$$

$$I_7 = \frac{1}{14} - \frac{9}{2} \cdot 0,051941 = -0,162306 .$$

Rezultatul este incorect deoarece integrandul este un număr pozitiv astfel ca integrala este pozitivă.

Eroarea care apare se datorează faptului că eroarea obținută în calculul lui I_0 , prin aproximarea numărului $\ln \frac{11}{9}$ s-a amplificat la fiecare pas al algoritmului datorită produsului $\frac{9}{2} I_{n-1}$. La pasul n de calcul eroarea cumulată devine proporțională cu 4.5^n .

Erori de rotunjire și de trunchiere sunt erorile generate în principal de faptul că dispunem de un spațiu limitat de reprezentare a numerelor într-un calculator, dar nu numai.

Întâlnim erori de trunchiere atunci când aproximăm o mărime reprezentată în calcule printr-un număr infinit de termeni, printr-un număr finit de astfel de termeni. Eroarea apare datorită termenilor neglijați. În mare parte aceste erori sunt acceptate în calculul numeric deoarece nu pot fi evitate.

Un exemplu de astfel de eroare este cea care se face când aproximăm numărul e cu o parte finită din suma dată de dezvoltarea

lui în serie Taylor: $e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots$, deci prin înlăturarea unui

număr infinit de termeni din această sumă.

Erorile de rotunjire apar deoarece numerele sunt reprezentate în calculator cu un număr redus de cifre semnificative, depinzând de lungimea cuvântului (numărul de biți) utilizat de calculator.

I.3. EVALUAREA ERORILOR

În calculul numeric noțiunea de eroare are semnificația de precizie a aproximării. În calculele numerice sunt acceptate erorile mai mici decât eroarea admisibilă precizată.

Dacă ne raportăm la momentul în care se determină erorile în cadrul aplicării unei metode numerice întâlnim două tipuri de erori:

- apriori- se estimează înainte de începerea calculului propriu-zis;
- aposteriori- se estimează după începerea aplicării algoritmului, după una sau mai multe etape de rezolvare.

Indiferent de natura lor, erorile se evaluează în funcție de valoarea exactă (reală) x a mărimii de interes și de valoarea aproximată (calculată) \bar{x} a acestei mărimi. Evaluarea se poate face în mod absolut sau relativ.

Definiția I.1. Eroarea absolută se notează cu e_x și este dată de relația:

$$\text{a) } e_x = x - \bar{x}, \text{ sau b) } e_x = |x - \bar{x}| \quad (1)$$

Eroarea absolută nu ține cont de ordinul mărimii studiate și din acest motiv nu este totdeauna suficientă doar cunoașterea acesteia.

Definiția I.2. Eroarea relativă se notează cu ε_x și este dată de relația:

$$\text{a) } \varepsilon_x = \frac{e_x}{x}, \text{ sau b) } \varepsilon_x = \left| \frac{e_x}{x} \right|$$

Uneori eroarea relativă se exprimă, în procente, prin relația:

$$\varepsilon_x = \left| \frac{e_x}{x} \right| \cdot 100 [\%]. \quad (2)$$

Deoarece valoarea exactă nu se cunoaște de cele mai multe ori, se introduce în relațiile precedente valoarea calculată în locul valorii reale și se obțin relațiile:

$$\text{a) } \varepsilon_x = \frac{e_x}{\bar{x}}, \text{ sau b) } \varepsilon_x = \left| \frac{e_x}{\bar{x}} \right| \text{ sau } \varepsilon_x = \left| \frac{e_x}{\bar{x}} \right| \cdot 100 [\%], \text{ în procente.} \quad (3)$$

Exemplul 1.5. Să considerăm o mărime, $x = 13$, a cărei valoare calculată este $\bar{x} = 14$. Eroarea absolută este $e_x = 1$. Aceeași eroare absolută, $e_y = 1$, se obține și dacă de exemplu analizăm o altă mărime y , pentru care $y = 1386$, iar $\bar{y} = 1387$.

Folosirea valorii calculate drept aproximare a mărimii reale în prima situație convine mai mult decât în cazul doi, dacă ne raportăm la ordinul de mărime al celor două mărimi analizate. Erorile absolute, fiind egale, nu dovedesc acest lucru, însă erorile relative sugerează

asta: $\varepsilon_x = \frac{1}{14} \cong 0.0714$, $\varepsilon_y = \frac{1}{1387} = 7.2098 \cdot 10^{-4}$, și deci $\varepsilon_y < \varepsilon_x$.

Astfel se realizează o aproximare mai bună a mărimii y prin \bar{y} decât aproximarea mărimii x prin \bar{x} .

1.4. CONDIȚIONARE ȘI STABILITATE

Condiționarea unei probleme se referă la sensibilitatea datelor de ieșire, adică a soluției, față de variațiile (perturbațiile) datelor de intrare. Problemele pot fi bine condiționate, sau slab condiționate. Când variații mici ale datelor de intrare determină variații mici în soluție, problema este bine condiționată, iar când determină variații mari ale soluției ea este slab condiționată.

Astfel pentru diferite clase de probleme de calcul numeric se definesc numere de condiționare care exprimă factorul de amplificare a erorii.

O problemă caracterizată de un număr de condiționare mare va fi o problemă slab condiționată, iar una caracterizată de un număr de condiționare mic (în general subunitar) va fi bine condiționată.

Conceptele de stabilitate sau instabilitate numerică se referă la algoritmi ce însoțesc metodele numerice. Acei algoritmi care se dovedesc a nu amplifica erorile în timpul calculelor pe care le presupun se spune că sunt stabili din punct de vedere numeric, respectiv instabili, dacă produc un astfel de efect.

De exemplu în cazul algoritmilor care oferă soluții aproximative care converg către o limită, se poate întâmpla ca, deși condițiile de convergență să fie satisfăcute, datorită cumulării erorilor de rotunjire, soluția generată să se depărteze, uneori chiar foarte mult, de această limită. Spunem în acest caz că algoritmul este instabil.

Cele mai mari erori care ar putea să apară în rezolvarea pe cale numerică a problemelor s-ar obține deci în cazul în care s-ar aplica un algoritm instabil la rezolvarea unei probleme slab condiționate.

O astfel de situație trebuie evitată prin găsirea altor metode de soluționare. Pentru a obține o bună soluție numerică la o problemă aceasta trebuie să fie bine condiționată iar algoritmul de rezolvare stabil.

II. METODE DIRECTE PENTRU REZOLVAREA SISTEMELOR DE ECUAȚII LINIARE

II.1. INTRODUCERE

În ceea ce privește un sistem de n ecuații cu n necunoscute, rezolvarea însăși a unui astfel de sistem, și deci posibilitatea de a vedea dacă avem sau nu soluție, este la fel de costisitoare (din punct de vedere al timpului și al memoriei) ca și calcularea determinantului sistemului sau a rangului matricei acestuia. De aceea de multe ori se trece direct la soluționarea sistemelor fără ca în prealabil să se studieze existența soluției.

Pentru rezolvarea sistemelor liniare există două clase de metode :

a) Metodele directe sau exacte, cum sunt metodele de tip Gauss, ce furnizează soluția exactă într-un număr finit de pași, abstracție făcând de erorile de rotunjire sau trunchiere care pot să apară.

b) Metodele iterative, precum metodele Jacobi și Gauss-Seidel, ce aproximează soluția generând un șir ce converge către aceasta.

Metodele exacte necesită multe operații aritmetice, ce duc la acumulări mari ale erorilor de rotunjire, și presupun timp mare de utilizare a calculatorului și spații mari de memorie.

Metodele iterative sunt mai rapide decât cele exacte, necesită un număr mai mic de operații aritmetice, dar e mai greu de precizat

de la început numărul de iterații necesare obținerii unei soluții suficient de bune.

II.2. REZOLVAREA SISTEMELOR TRIUNGHILARE

Un sistem de n ecuații cu n necunoscute se numește triunghiular dacă matricea atașată acestuia este superior sau inferior triunghiulară.

Definiția II.1. Matricea $T = (t_{ij}) \in M(C)$ se numește superior (inferior) triunghiulară dacă:

$$t_{ij} = 0 \text{ pentru } i > j \quad (t_{ij} = 0 \text{ pentru } i < j).$$

Un sistem de n ecuații cu n necunoscute a cărei matrice este superior, respectiv inferior triunghiulară se rezolvă ușor prin substituție inversă (sau înapoi), respectiv substituție directă (substituție înainte).

Vom considera cazul sistemelor ce admit o soluție unică, deci cazul în care elementele $t_{ii} \neq 0, \forall i = \overline{1, n}$.

Să considerăm un sistem caracterizat de o matrice inferior triunghiulară. Sistemul are în acest caz forma:

$$\begin{cases} t_{11}x_1 & = b_1 \\ t_{21}x_1 + t_{22}x_2 & = b_2 \\ \dots & \\ t_{n1}x_1 + t_{n2}x_2 + \dots + t_{nn}x_n & = b_n \end{cases}$$

Metoda substituției directe presupune obținerea necunoscutelor succesiv, începând cu prima din ecuațiile sistemului. Astfel obținem:

$$\begin{aligned}
 x_1 &= \frac{b_1}{t_{11}}, \\
 x_2 &= \frac{b_2 - t_{21}x_1}{t_{22}}, \\
 &\dots \\
 x_n &= \frac{b_n - \sum_{k=1}^{n-1} t_{nk}x_k}{t_{nn}}
 \end{aligned}$$

Algoritmul de rezolvare a sistemelor inferior triunghiulare

Date de intrare: n , numărul de ecuații și necunoscute; matricea triunghiulară a sistemului, T , care îndeplinește condiția $t_{ii} \neq 0, \forall i = \overline{1, n}$; termenii liberi $b_i, i = \overline{1, n}$.

Date de ieșire: soluția sistemului

1. Citește datele de intrare;
2. Inițializează variabila $i \leftarrow 1$;
3. Atribuie lui $x_1 \leftarrow \frac{b_1}{t_{11}}$;
4. Pentru i de la 2 la n , cu pasul 1, execută
 - 4.1. Inițializează S cu 0 $S \leftarrow 0$;
 - 4.2. Pentru $k = \overline{1, i-1}$, cu pasul 1, execută

$$s \leftarrow s + t_{ik}x_k;$$

- 4.3. Atribuie valoare necunoscutei $x_i, x_i \leftarrow \frac{b_i - S}{t_{ii}}$;

5. Afișează valorile necunoscutelor;
6. Stop.

Analog se obține și algoritmul pentru rezolvarea sistemelor superior triunghiulare, denumit substituție inversă, sau înapoi. Prima necunoscută care se află este x_n , apoi se determină x_{n-1}, \dots, x_1 . Sistemul se rezolvă practic începând cu ultima ecuație:

$$x_n = \frac{b_n}{t_{nn}},$$

$$x_i = \frac{b_i - \sum_{j=i+1}^n t_{ij} x_j}{t_{ii}}, \quad i = n-1, n-2, \dots, 1.$$

Algoritmul de rezolvare a sistemelor superior triunghiulare

Date de intrare: n , numărul de ecuații și necunoscute; matricea triunghiulară a sistemului, T , care îndeplinește condiția $t_{ii} \neq 0, \forall i = \overline{1, n}$; termenii liberi $b_i, i = \overline{1, n}$.

Date de ieșire: soluția sistemului

1. Citește datele de intrare;
2. Inițializează variabila $i, i \leftarrow n$;
3. Atribuie lui $x_n \leftarrow \frac{b_n}{t_{nn}}$;
4. Pentru i de la $n-1$ la 1 , cu pasul -1 , execută
 - 4.1. Inițializează S cu $0, S \leftarrow 0$;
 - 4.2. Pentru $k = \overline{n, i-1}$, cu pasul -1 , execută

$$S \leftarrow S + t_{ik} x_k;$$

4.3. Atribuie valoare necunoscutei x_i , $x_i \leftarrow \frac{b_i - S}{t_{ii}}$;

5. Afișează valorile necunoscutele;

6. Stop.

Exemplul II.1.

1) Să se rezolve sistemul de ecuații de mai jos prin metoda substituției inverse:

$$4x_1 + x_2 - x_3 + x_4 = 5$$

$$3x_2 - x_3 + x_4 = 3$$

$$5x_3 + x_4 = 6$$

$$2x_4 = 2$$

Soluție:

Matricea coeficienților sistemului dat este superior triunghiulară și deci aplicăm substituția înapoi. Se începe cu calculul lui x_4 , apoi se calculează și celelalte necunoscute: x_3, x_2, x_1 .

Obținem astfel: $x_4 = 1, x_3 = 1, x_2 = 1, x_1 = 1$.

Folosind limbajul C de programare s-au realizat următoarele programe pentru rezolvarea sistemelor superior, respectiv inferior triunghiulare.

/*Program Rezolvare Sistem Superior Triunghiular*/

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
float s,x[30],T[30][30],B[30];
```

```
int i,k,j,n;
```

```
void main()
```

```
{
```

```

printf("Dati nr de ecuatii n:\n");
scanf("%d",&n);
printf("Dati coeficientii necunoscutelor\n");
for (i=1;i<=n;i++)
    for (j=i;j<=n;j++)
        {
            printf("T[%d][%d]: ",i,j);
            scanf("%f",&T[i][j]);
        }
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
    {
        printf("B[%d]: ",i);
        scanf("%f",&B[i]);
    }
x[n]=B[n]/T[n][n];
for (i=n-1;i>=1;i--)
    {
        s=0;
        for (k=n;k>=i-1;k--)
            s=s+T[i][k]*x[k];
        x[i]=(B[i]-s)/T[i][i];
    }
printf("Solutia este: \n");
for (i=1;i<=n;i++)
    printf("x[%d]=%f\n",i,x[i]);
getch();
}

```

/*Program Rezolvare Sistem Inferior Triunghiular*/

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
float s,x[30],T[30][30],B[30];
```

```
int i,k,j,n;
```

```
void main()
```

```
{
```

```
printf("Dati nr de ecuatii n:\n");
```

```
scanf("%d",&n);
```

```

printf("Dati coeficientii necunoscutelor\n");
for (i=1;i<=n;i++)
    for (j=1;j<=i;j++)
    {
        printf("T[%d][%d]: ",i,j);
        scanf("%f",&T[i][j]);
    }
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
    {
        printf("B[%d]: ",i);
        scanf("%f",&B[i]);
    }
x[1]=B[1]/T[1][1];
for (i=2;i<=n;i++)
    { s=0;
      for (k=1;k<=i-1;k++)
        s=s+T[i][k]*x[k];
      x[i]=(B[i]-s)/T[i][i];
    }
printf("Solutia este: \n");
for (i=1;i<=n;i++)
    printf("x[%d]=%f \n",i,x[i]);
getch();
}

```

În continuare sunt prezentate câteva rezulte obținute pe baza rulării acestor programe:

```

C:\cc\bin\rundos.exe
Dati nr de ecuatii n:
4
Dati coeficientii necunoscutelor
T[1][1]: 2
T[1][2]: -3
T[1][3]: 1
T[1][4]: -1
T[2][2]: 11
T[2][3]: -3
T[2][4]: 1
T[3][3]: 16
T[3][4]: -9
T[4][4]: 143
Dati termenii liberi:
B[1]: -4
B[2]: 0
B[3]: -11
B[4]: 429
Solutia este:
x[1]=-1.000000
x[2]=0.000000
x[3]=1.000000
x[4]=3.000000

C:\cc\bin\rundos.exe
Dati nr de ecuatii n:
4
Dati coeficientii necunoscutelor
T[1][1]: 3
T[2][1]: -1
T[2][2]: 2
T[3][1]: 1
T[3][2]: -1
T[3][3]: 3
T[4][1]: 4
T[4][2]: -5
T[4][3]: 1
T[4][4]: -1
Dati termenii liberi:
B[1]: 6
B[2]: 0
B[3]: 1
B[4]: 4
Solutia este:
x[1]=2.000000
x[2]=1.000000
x[3]=0.000000
x[4]=-1.000000

```

II.3. ALGORITMUL GAUSS PENTRU REZOLVAREA SISTEMELOR LINIARE. VARIANTA GAUSS-JORDAN.

Metoda lui Gauss este una din cele mai cunoscute metode de rezolvare a sistemelor de ecuații liniare, de calcul al determinanților și de determinare a inversei unei matrici. Se poate utiliza inclusiv la aflarea rangului unei matrici. Considerăm un sistem linear de n ecuații cu n necunoscute:

$$Ax = B, \text{ unde } A \in M_n(R), B \in M_{n \times 1}(R).$$

Algoritmul Gauss, denumit și metoda eliminării parțiale sau tehnica de pivotare, constă în a aduce sistemul, prin $n-1$ transformări elementare, la o formă echivalentă, caracterizată de o matrice superior triunghiulară, notată A' , prin eliminarea succesivă a necunoscutelor din ecuațiile sistemului inițial.

Pentru aflarea soluției sistemului considerat se rezolvă apoi sistemul triunghiular obținut, $A'x = B'$, prin substituție inversă.

Transformările succesive la fiecare etapă se aplică matricii extinse a sistemului.

Pentru a prezenta aceste transformări succesive notăm printr-un indice superior pasul curent, astfel că $a_{ij}^{(k)}$ reprezintă elementul de pe linia i coloana j al matricii de la pasul k . Notăm elementele matricii A , cu indicele 0 sus, adică $a_{ij} = a_{ij}^0 \forall i, j \in \{1, \dots, n\}$.

Transformările efectuate la etapa k pot fi rezumate în formule:

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

$$b_i^{(k)} = \begin{cases} b_i^{(k-1)} & \text{pentru } i \leq k, \\ b_i^{(k-1)} - \frac{a_{i,k}^{(k-1)}}{a_{k,k}^{(k-1)}} b_k^{(k-1)} & \text{pentru } i \geq k+1 \end{cases}$$

Elementul $a_{k,k}^{(k-1)}$ poartă numele de pivot al transformării de la pasul k . Evident el trebuie să fie nenul pentru ca algoritmul să funcționeze.

Se obține în final un sistem echivalent caracterizat de o matrice superior triunghiulară a cărei soluție se obține prin substituție inversă.

Dacă se întâmplă ca pe parcursul aplicării algoritmului să întâlnim un element nul pe poziția pivotului se permută linia curentă cu o linie situată sub ea, astfel încât noul element, ce ajunge pe poziția pivotului, să fie nenul.

Algoritmul de rezolvare a sistemului $Ax = B$ prin metoda Gauss este descris în continuare.

Date de intrare: n , matricea A și vectorul coloană format cu termenii liberi, B

Date de ieșire: soluția (x_1, x_2, \dots, x_n) sau un mesaj de eroare

1. Pentru i de la 1 la $n-1$ execută:

- a. Caută p cel mai mic întreg $p = \overline{i, n}$ astfel ca $a_{pi} \neq 0$
- b. Dacă nu există, atunci scrie „sistemul nu are soluție unică” și du-te la pasul 4
- c. Dacă $p \neq i$, atunci permută liniile p și i
- d. Pentru j de la $i+1$ la n execută

i.
$$m_{ji} := \frac{a_{ji}}{a_{ii}},$$

ii.
$$\text{linia } j := \text{linia } j - m_{ji} \text{linia } i$$

2. Dacă $a_{nn} \neq 0$, atunci

2.1
$$x_n := \frac{b_n}{a_{nn}}$$

2.2 pentru i de la $n-1$ la 1 execută

$$x_i := \left[b_i - \sum_{j=i+1}^n a_{ij} x_j \right] / a_{ii}$$

altfel

2.3 Dacă $b_n = 0$, afisează sistemul nu are soluție unică și du-te la pasul 4

altfel

2.4 afisează sistemul nu are soluție și du-te la pasul 4

3. Afişează x_1, x_2, \dots, x_n

4. Stop.

Metoda Gauss poate fi folosită și la calculul determinantului, unei matrici, precum și la aflarea inversei unei matrici.

Valoarea determinantului unei matrici superior triunghiulare, de forma matricii A' se obține simplu cu ajutorul relației:

$$\det(A) = \Delta = a_{11} a_{22}^{(1)} a_{33}^{(2)} \dots a_{nn}^{(n-1)}.$$

Reamintim că pentru a putea folosi această metodă este necesar ca elementele alese drept pivoți să fie nenule. În plus, pot apărea erori numerice dacă ele au valori prea mici. Pentru a evita această situație se pot inversa ecuațiile sistemului între ele, fapt care duce însă la schimbarea semnului determinantului.

II.4. TEHNICI DE PIVOTARE

Pentru a micșora erorile ce pot să apară în cazul aplicării metodei Gauss s-au conturat două tehnici de pivotare: pivotarea parțială și pivotarea totală.

În cazul pivotării parțiale se alege drept pivot elementul de sub diagonală, de pe aceeași coloană, maxim în modul. Acesta este adus pe poziția pivotului printr-o permutare convenabilă de linii, după care algoritmul se continuă în modul obișnuit.

În cazul pivotării totale, la pasul i se alege pivot elementul de pe linia i coloana i doar dacă acesta este elementul maxim în modul între elementele submatricei obținute la pasul precedent $A_i^{(i-1)} = (a_{kl}^{(i-1)})_{i \leq k, l \leq n}$ (formate cu liniile și coloanele cu indici mai mari sau egali cu i). Dacă nu el nu respectă această condiție se caută un astfel de element în submatricea precizată și se aduce acesta pe poziția (i, i) prin permutări de linii dar și de coloane.

Observația II.1. Permutările de linii nu afectează structura soluției sistemului, pe când cele de coloane, da. De aceea permutările de coloane trebuie reținute, în ordinea efectuării lor, pentru ca, parcurgând în sens invers șirul lor ordonat, și permutând corespunzător componentele soluției sistemului rezultat, să obținem soluția sistemului inițial.

Exemplul II.2.

Aplicând metoda Gauss să se rezolve sistemul:

$$\begin{cases} -x + 2y + z = 2 \\ 2x + y - 4z = 8 \\ 3x - y - z = 2 \end{cases}$$

Soluție:

$$\begin{array}{ccc|c}
 [-1] & 2 & 1 & 2 \\
 2 & 1 & -4 & 8 \\
 3 & -1 & -1 & 2 \\
 \hline
 -1 & 2 & 1 & 2 \\
 0 & [5] & -2 & 12 \\
 0 & 5 & 2 & 8 \\
 \hline
 -1 & 2 & 1 & 2 \\
 0 & 5 & -2 & 12 \\
 0 & 0 & 4 & -4
 \end{array}$$

Astfel se obține forma echivalentă, caracterizată de o matrice superior triunghiulară:

$$-x + 2y + z = 2$$

$$5y - 2z = 12$$

$$4z = -4 \Rightarrow z = -1$$

$$y = \frac{12 + 2 \cdot (-1)}{5} = \frac{10}{5} = 2$$

$$x = -2 + 2 \cdot 2 + (-1) = 1$$

Soluția sistemului este deci tripletul: $(1, 2, -1)$

Exemplul II.3

Să se rezolve același sistem aplicând tehnica pivotării parțiale.

Soluție:

În cadrul tehnicii de pivotare parțială, se permută între ele liniile sistemului astfel încât la fiecare pas pivotul ales să fie maxim în modul între elementele aflate pe coloana sa, sub linia pe care este situat.

La primul pas se observă că $\max(|-1|, 2, 3) = 3$.

Elementul maxim în modul de pe prima coloană, adică 3, trebuie adus, printr-o permutare de linii, pe prima poziție. Astfel permutăm liniile 1 și 3.

Obținem:

$$\begin{cases} 3x - y - z = 2 \\ 2x + y - 4z = 8 \\ -x + 2y + z = 2 \end{cases} \Rightarrow$$

$$\begin{array}{ccc|c} [3] & -1 & -1 & 2 \\ 2 & 1 & -4 & 8 \\ -1 & 2 & 1 & 2 \\ \hline 3 & -1 & -1 & 2 \\ 0 & [5/3] & -10/3 & 20/3 \\ 0 & 5/3 & 2/3 & 8/3 \\ \hline 3 & -1 & -1 & 2 \\ 0 & 5/3 & -10/3 & 20/3 \\ 0 & 0 & 4 & -4 \end{array}$$

Forma echivalentă a sistemului, dată de ultimul tabel este :

$$\begin{aligned} 3x - y - z &= 2 \\ \frac{5}{3}y - \frac{10}{3}z &= \frac{20}{3} \\ 4z &= -4 \Rightarrow z = -1 \\ y &= \frac{20 + 10(-1)}{5} = 2 \\ x &= \frac{2 + 2 + (-1)}{3} = 1 \end{aligned}$$

Obținem astfel soluția unică: $(1, 2, -1)$.

Exemplul II.4

Să se rezolve același sistem aplicând tehnica pivotării totale.

Soluție:

La primul pas se caută elementul maxim în modul din toată matricea sistemului. Se observă că elementul maxim în modul este 4. Acesta trebuie adus pe linia 1 coloana 1. Astfel se permută liniile 1 și 2, rezultând:

$$\begin{cases} 2x + y - 4z = 8 \\ 3x - y - z = 2 \\ -x + 2y + z = 2 \end{cases}$$

Apoi se permută coloanele 1 și 3. Reținem această permutare de coloane ($C_1 \leftrightarrow C_3$).

Sistemul devine:

$$\begin{cases} -4x_1 + x_2 + 2x_3 = 8 \\ -x_1 - x_2 + 3x_3 = 2 \\ x_1 + 2x_2 - x_3 = 2 \end{cases}$$

$$\begin{array}{ccc|c} [-4] & 1 & 2 & 8 \\ -1 & -1 & 3 & 2 \\ 1 & 2 & -1 & 2 \\ \hline -4 & 1 & 2 & 8 \\ 0 & [-5/4] & 10/4 & 0 \\ 0 & 9/4 & -2/4 & 4 \end{array}$$

La o aplicare a tehnicii de pivotare obișnuite, la pasul doi ar trebui ales drept pivot elementul $-5/4$.

Se observă însă că elementul maxim în modul al submatricei

$$\begin{pmatrix} -5/4 & 10/4 \\ 9/4 & -2/4 \end{pmatrix} \text{ este } 10/4$$

Acesta va fi pivotul următoarei transformări. El trebuie adus pe linia 2, coloana 2. Permutăm pentru aceasta coloana 2 cu coloana 3, și fiind vorba de o permutare de coloane o reținem și pe aceasta ($C_2 \leftrightarrow C_3$). Obținem:

$$\begin{array}{ccc|c} -4 & 2 & 1 & 8 \\ 0 & [10/4] & -5/4 & 0 \\ 0 & -2/4 & 9/4 & 4 \\ \hline -4 & 2 & 1 & 8 \\ 0 & 10/4 & -5/4 & 0 \\ 0 & 0 & 2 & 4 \end{array}$$

Forma echivalentă (scrisă pentru alte necunoscute decât cele inițiale, deoarece au avut loc permutări de coloane) este :

$$(*) \begin{cases} -4y_1 + 2y_2 + y_3 = 8 \\ 10/4 y_2 - 5/4 y_3 = 0 \Rightarrow \\ 2y_3 = 4 \\ \\ y_3 = 2 \\ y_2 = \frac{5 \cdot (2)}{+10} = 1 \\ y_1 = \frac{2 + 2 - 8}{4} = \frac{-4}{4} = -1 \end{cases}$$

Soluția sistemului (*) este tripletul: $(-1, 1, 2)$.

Revenim la permutările de coloane pe care le-am făcut.

Parcurgându-le în ordine inversă, și efectuându-le asupra componentelor soluției, obținem în final soluția sistemului inițial.

Astfel avem :

$$(-1, 1, 2) \xrightarrow{C_2 \rightarrow C_3} (-1, 2, 1) \xrightarrow{C_1 \rightarrow C_3} (1, 2, -1).$$

Astfel soluția sistemului inițial este tripletul $(1, 2, -1)$.

Exemplul II.5.

Folosind tehnica pivotării să se rezolve sistemul :

$$\begin{cases} 2x - 3y + z - t = -4 \\ x + 4y - z = -2 \\ -2x - y + 3z - t = 2 \\ -x + 2y + t = 4 \end{cases}$$

Soluție:

[2]	-3	1	-1	-4
1	4	-1	0	-2
-2	-1	3	-1	2
-1	2	0	1	4
2	-3	1	-1	-4
0	$\left[\frac{11}{2}\right]$	$-\frac{3}{2}$	$\frac{1}{2}$	0
0	-4	4	-2	-2
0	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	2
2	-3	1	-1	-4
0	11	-3	1	0
0	0	$\left[\frac{32}{11}\right]$	$-\frac{18}{11}$	$-\frac{22}{11}$
0	0	$\frac{7}{11}$	$\frac{5}{11}$	$\frac{22}{11}$
2	-3	1	-1	-4
0	11	-3	1	0
0	0	$\frac{32}{11}$	$-\frac{18}{11}$	$-\frac{22}{11}$
0	0	0	$\frac{143}{176}$	$\frac{429}{176}$

Se obține astfel următoarea formă echivalentă pentru sistem:

$$\begin{cases} 2x - 3y + z - t = -4 \\ 11y - 3z + t = 0 \\ \frac{32}{11}z - \frac{18}{11}t = -2 \\ \frac{143}{176}t = \frac{429}{176} \end{cases}$$

Prin substituție inversă obținem soluția: $(-1, 0, 1, 3)$.

Un cod sursă (în C) pentru algoritmul eliminării gaussiene-varianta pivotării parțiale (până la obținerea matricei superior triunghiulare) este prezentat în continuare.

```
/*Program Eliminare Gaussiana*/
#include<stdio.h>
#include<conio.h>
#include<math.h>
double t,m;
int i,j,k,l,n;
float s, a[10][10], b[10];
void main()
{
printf("Introduceti nr. de linii si de coloane ale matricei");
scanf("%d",&n);
printf("Dati matricea\n");
for(i=1;i<=n;i++)
for(j=1;j<=n;j++)
{
printf("a[%d][%d]: ",i,j);
scanf("%f",&a[i][j]);
}
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
{
printf("b[%d]: ",i);
```

```

        scanf("%f",&b[i]);
    }
for(k=1;k<=n;k++)
{
l=k;
for (i=k+1;i<=n;i++)
if (abs(a[i][k])>abs(a[k][k]))
l=i;
if(l!=k)
{
for(j=1;j<=n;j++)
{
t=a[k][j];a[k][j]=a[l][j];a[l][j]=t;
}
t=b[k];b[k]=b[l];b[l]=t;
}
for(i=k+1;i<=n;i++)
{
m=a[i][k]/a[k][k];
a[i][k]=0;
for(j=k+1;j<=n;j++)
a[i][j]=a[i][j]-m*a[k][j];
b[i]=b[i]-m*b[k];
}
}
printf("Dupa eliminarea gaussiana matricea este: \n");
for (i=1;i<=n;i++)
{
for (j=1;j<=n;j++)
printf("a[%d][%d]=%f ",i,j,a[i][j]);
printf("b[%d]=%f ",i,b[i]);
printf("\n");
}
getch();
}

```

Soluționarea problemei de la exemplul II.3 cu ajutorul programului de calcul numeric anterior este prezentată în continuare.

```

C:\>dir\andros.exe
Introduceti nr. de linii si de coloane ale matricei:3
Dati matricea
a[1][1]: -1
a[1][2]: 2
a[1][3]: 1
a[2][1]: 2
a[2][2]: 1
a[2][3]: -4
a[3][1]: 3
a[3][2]: -1
a[3][3]: -1
Dati termenii liberi:
b[1]: 2
b[2]: 3
b[3]: 2
Dupa eliminarea gaussiana matricea este:
a[1][1]:3.000000 a[1][2]:-1.000000 a[1][3]:-1.000000 b[1]:2.000000
a[2][1]:0.400000 a[2][2]:1.666667 a[2][3]:-3.333333 b[2]:6.666667
a[3][1]:0.000000 a[3][2]:-0.000000 a[3][3]:4.000000 b[3]:-4.000000

```

II.5. VARIANTA GAUSS-JORDAN.

Transformările efectuate sunt asemănătoare cu cele din cadrul metodei Gauss, însă afectează toată matricea de la pasul curent, în final obținându-se un o matrice diagonală.

Formulele de transformare iterativă a elementelor matricei A la etapa k sunt :

$$\begin{cases}
 a_{kj}^{(k)} = a_{kj}^{(k-1)}, \text{ pentru } j \geq 1 \\
 a_{ij}^{(k)} = a_{ij}^{(k-1)} - \frac{a_{i,k}^{(k-1)}}{a_{k,k}^{(k-1)}} a_{k,j}^{(k-1)}, \text{ pentru } 1 \leq i \leq n, i \neq k, j \geq k+1 \\
 a_{ij}^{(k)} = a_{ij}^{(k-1)}, \text{ pentru } j < k \\
 a_{ik}^{(k)} = 0, \text{ pentru } 1 \leq i \leq n, i \neq k
 \end{cases}$$

iar membrul drept se modifică după relațiile :

$$\begin{cases}
 b_k^{(k)} = b_k^{(k-1)} \\
 b_i^{(k)} = b_i^{(k-1)} - \frac{a_{i,k}^{(k-1)}}{a_{k,k}^{(k-1)}} b_k^{(k-1)} \quad \text{pentru } i \neq k
 \end{cases}$$

În acest caz soluția sistemului va fi : $x_i = \frac{b_i^{(n-1)}}{a_{ii}^{(n-1)}}$, $i = \overline{1, n}$.

Metoda Gauss-Jordan mai poartă denumirea de metoda eliminării complete. Se poate aplica și ușor modificată în sensul că se poate obține în final un sistem caracterizat de matricea unitate. Pentru aceasta linia pivotului se împarte la pivot. Se fac într-adevăr mai multe operații algebrice asupra matricei sistemului, față de varianta prezentată, dar se găsește direct soluția acestuia ($x_i = b_i^{(n-1)}, i = \overline{1, n}$).

Exemplul II.6.

Să se rezolve următoarele sisteme folosind metoda Gauss-Jordan.

$$\text{a) } \begin{cases} x_1 - x_2 + x_3 = 4 \\ 2x_1 + x_2 - x_3 = 1 \\ -x_1 + 2x_2 + x_3 = 6 \end{cases}, \quad \text{b) } \begin{cases} x_1 - x_2 + x_3 + 2x_4 = 3 \\ x_2 - 2x_3 + x_4 = 2 \\ 2x_1 + 2x_2 - x_3 + x_4 = 1 \end{cases}.$$

Soluție:

a) Obținem succesiv următoarele tablouri:

$$\begin{array}{ccc|c}
 1 & -1 & 1 & 4 \\
 2 & 1 & -1 & 1 \\
 -1 & 2 & 1 & 6 \\
 \hline
 1 & -1 & 1 & 4 \\
 0 & 3 & -3 & -7 \\
 0 & 1 & 2 & 10 \\
 \hline
 1 & 0 & 0 & \frac{5}{3} \\
 0 & 1 & -1 & -\frac{7}{3} \\
 0 & 0 & 3 & \frac{37}{3} \\
 \hline
 1 & 0 & 0 & \frac{15}{9} \\
 0 & 1 & 0 & \frac{16}{9} \\
 0 & 0 & 1 & \frac{37}{9}
 \end{array}$$

În acest caz ultima formă echivalentă pentru sistem este una diagonală, cu elemente unitate, deci soluția se obține direct. Astfel,

$$x_1 = \frac{15}{9}, \quad x_2 = \frac{16}{9}, \quad x_3 = \frac{37}{9}$$

reprezintă soluția unică a sistemului.

b) În acest caz avem un sistem cu 3 ecuații și 4 necunoscute. Procedând analog în cazul sistemului de la punctul b) găsim:

$$\begin{array}{cccc|c}
 1 & -1 & 1 & 2 & 3 \\
 0 & 1 & -2 & 1 & 2 \\
 2 & 2 & -1 & 1 & 1 \\
 \hline
 1 & -1 & 1 & 2 & 3 \\
 0 & 1 & -2 & 1 & 2 \\
 0 & 4 & -3 & -3 & -5 \\
 \hline
 1 & 0 & -1 & 3 & 5 \\
 0 & 1 & -2 & 1 & 2 \\
 0 & 0 & 5 & -7 & -13 \\
 \hline
 1 & 0 & 0 & \frac{8}{5} & \frac{12}{5} \\
 0 & 1 & 0 & \frac{1}{5} & -\frac{16}{5} \\
 0 & 0 & 1 & -\frac{7}{5} & -\frac{13}{5}
 \end{array}$$

O formă echivalentă pentru sistem este:

$$\begin{cases}
 x_1 + \frac{8}{5}x_4 = \frac{12}{5} \\
 x_2 - \frac{7}{5}x_4 = -\frac{16}{5} \Rightarrow \text{rang} A = 3, \text{ rang } \bar{A} = 3, \\
 x_3 - \frac{7}{5}x_4 = -\frac{13}{5}
 \end{cases}$$

Sistemul este deci compatibil simplu nedeterminat. Necunoscutele principale sunt x_1, x_2, x_3 iar necunoscuta secundară x_4 .

Alegem $x_4 = \alpha$.

$$x_4 = \alpha \Rightarrow S: \begin{cases}
 x_1 = \frac{12}{5} - \frac{8}{5}\alpha \\
 x_2 = \frac{-16}{5} + \frac{9}{5}\alpha \\
 x_3 = -\frac{13}{5} + \frac{7}{5}\alpha \\
 x_4 = \alpha
 \end{cases}$$

III. METODE NUMERICE ÎN CALCULUL MATRICEAL

III.1. FACTORIZAREA MATRICELEOR

Prin factorizarea unei matrici înțelegem scrierea acesteia ca un produs de două matrice: una inferior triunghiulară cealaltă superior triunghiulară (denumită factorizare LR).

Definiția III.1. Fie $A \in M_n(R)$, reprezentarea lui A sub forma:

$$A = L \cdot R,$$

cu L - matrice inferior triunghiulară și R - matrice superior triunghiulară, reprezintă factorizarea LR a lui A .

Factorizarea unei matrici se folosește în soluționarea sistemelor liniare cu n ecuații și n necunoscute.

Considerând sistemul $Ax=B$, prin înlocuirea lui A în relația precedentă cu expresia obținută prin factorizare, obținem:

$$LRx=B$$

Astfel rezolvarea sistemului $Ax = B$ presupune descompunerea acestuia în rezolvarea a doua sisteme:

1) $Ly = B$ (sistem caracterizat de o matrice inferior triunghiulară, care se rezolvă printr-o “substituție înainte”);

2) $Rx = Y$ (sistem caracterizat de o matrice superior triunghiulară, care se rezolvă printr-o “substituție înapoi”).

Teorema III.1. Dacă matricea A are toți minorii diagonali principali nenuli, atunci există o factorizare LR a lui A în care L este nesingulară.

Pentru determinarea matricelor L și R există două tipuri importante de algoritmi asemănători, care diferă doar datorită aspectului matricelor rezultate în urma descompunerii:

- 1) Metoda Crout – caz în care diagonala lui R este unitară.
- 2) Metoda Doolittle –caz în care diagonala lui L este unitară

1) **Factorizarea Crout (metoda Crout)**

În acest caz avem următoarea descompunere $A = L \cdot R$ cu

$$L = \begin{pmatrix} l_{11} & 0 & 0 & \cdot & \cdot & 0 \\ l_{21} & l_{22} & 0 & \cdot & \cdot & 0 \\ l_{31} & l_{32} & l_{33} & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ l_{n1} & l_{n2} & l_{n3} & \cdot & \cdot & l_{nn} \end{pmatrix} \quad R = \begin{pmatrix} 1 & r_{12} & r_{13} & \cdot & \cdot & r_{1n} \\ 0 & 1 & r_{23} & \cdot & \cdot & r_{2n} \\ 0 & 0 & 1 & \cdot & \cdot & r_{3n} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \cdot & \cdot & 1 \end{pmatrix}$$

Determinarea elementelor se face impunând condiția ca rezultatul produsului celor două matrice L și R , în această ordine, să fie egal cu matricea A , adică:

$$l_{i1}r_{1j} + l_{i2}r_{2j} + \dots + l_{in}r_{nj} = a_{ij}, \quad 1 \leq i, j \leq n$$

Deoarece se știe că $l_{ij} = 0$ pentru $j \geq i + 1$, $r_{ij} = 0$ pentru $j \leq i - 1$ și

$$r_{ii} = 1, \quad 1 \leq i \leq n \quad , \quad \text{obținem: } a_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik}r_{kj} \quad \text{și de aici relațiile ce}$$

determină elementele necunoscute ale celor două matrice:

- Pentru matricea inferior triunghiulară:

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} r_{kj}, \quad j \leq i,$$

unde pentru $j=1$ se consideră $\sum_{k=1}^0 = 0$ și $l_{i1} = a_{i1}$

- Pentru matricea superior triunghiulară:

$$r_{ij} = \left[a_{ij} - \sum_{k=1}^{i-1} l_{ik} r_{kj} \right] \cdot l_{ii}^{-1}, \quad j \geq i+1,$$

unde pentru $i=1$ se consideră $\sum_{k=1}^0 = 0$ și $r_{1j} = \frac{a_{1j}}{a_{11}}$

Pentru a face mai accesibilă aplicarea algoritmului de factorizare prezentat în vom detalia, insistând asupra unor reguli practice, ușor de reținut.

Observația III.1. Determinarea practică a matricelor se face în paralel, astfel că la fiecare pas se determină o coloană a lui L și o linie a lui R , în modul următor:

Algoritmul factorizării Crout

Pasul 1:

- Se determină prima coloană a lui L , care coincide cu cea a lui A .
- Se determină prima linie a lui R , folosind formulele date. Înmulțim practic prima linie a lui L cu toate coloanele lui R ce au indicele mai

mare strict decât 1 și egalăm elementele corespunzătoare cu cele ale lui A .

Pasul 2:

- Se determină a doua coloană a lui L , înmulțind toate liniile lui L , de la 2 la n , cu a doua coloană a lui R și egalând elementele corespunzătoare cu cele ale lui A .

- Se determină a doua linie a lui R , înmulțind a doua linie a lui L cu toate coloanele lui R de la 3 la n , și egalând elementele corespunzătoare cu cele ale lui A .

...,etc.

Pasul n :

- Se determină ultima coloană a lui L (practic elementul l_{nn}).

Observația III.2. Numărul de elemente calculate, și implicit numărul de operații ce se efectuează la fiecare pas, scade pe măsura creșterii numărului de ordine al pasului, astfel că la ultimul pas se calculează de fiecare dată doar un singur element.

Exemplul III.1. Folosind metoda Crout să se factorizeze matricea

$$A = \begin{pmatrix} 2 & -1 & 2 \\ 1 & 3 & 1 \\ -1 & 2 & 1 \end{pmatrix}$$

Soluție:

Calculând minorii diagonali principali observăm că ei respectă condiția din teoremă.

Fie cele două matrice L și R .

$$L = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \quad R = \begin{pmatrix} 1 & r_{12} & r_{13} \\ 0 & 1 & r_{23} \\ 0 & 0 & 1 \end{pmatrix}$$

Pasul 1:

-Se determină prima coloană a lui L . Avem relațiile:

$$l_{11} = a_{11} = 2 \Rightarrow l_{11} = 2$$

$$l_{21} = a_{21} = 1 \Rightarrow l_{21} = 1$$

$$l_{31} = a_{31} = -1 \Rightarrow l_{31} = -1$$

-Se determină prima linie a lui R :

$$l_{11} \cdot r_{12} = a_{12} = -1 \Rightarrow r_{12} = \frac{-1}{l_{11}} = \frac{-1}{2}, \Rightarrow r_{12} = -\frac{1}{2}$$

$$l_{11} \cdot r_{13} = a_{13} = 2 \Rightarrow r_{13} = \frac{2}{2} = 1$$

Pasul 2:

-Se determină a doua coloană a lui L .

$$l_{21} \cdot r_{12} + l_{22} = a_{22} = 3, \Rightarrow l_{22} = 3 + \frac{1}{2} = \frac{7}{2}$$

$$l_{31} \cdot r_{12} + l_{32} = a_{32} = 2 \Rightarrow l_{32} = 2 - \frac{1}{2} = \frac{3}{2}$$

-Se determină a doua linie a lui R .

$$l_{21} \cdot r_{13} + l_{22} \cdot r_{23} = a_{23} = 1 \Rightarrow r_{23} = \frac{1 - 1 \cdot 1}{\frac{7}{2}} = 0.$$

Pasul 3:

-Se determină a treia coloană a lui L , adică elementul l_{33} .

$$l_{31} \cdot r_{13} + l_{32} \cdot r_{23} + l_{33} = a_{33} = 1 \Rightarrow l_{33} = 1 + 1 = 2.$$

Am obținut următoarea factorizare pentru matricea A :

$$A = L \cdot R, \text{ unde } L = \begin{pmatrix} 2 & 0 & 0 \\ 1 & \frac{7}{2} & 0 \\ -1 & \frac{3}{2} & 2 \end{pmatrix} \text{ și } R = \begin{pmatrix} 1 & -\frac{1}{2} & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

/*Program Factorizare Crout*/

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
float s, A[30][30], L[30][30], R[30][30];
```

```
int i,k,j,n;
```

```
void main()
```

```
{
printf("Dati ordinul matricei A, n:\n");
scanf("%d",&n);
printf("Dati A\n");
for (i=1;i<=n;i++)
for (j=1;j<=n;j++)
{
printf("A[%d][%d]: ",i,j);
scanf("%f",&A[i][j]);
}
for (i=1;i<=n;i++)
L[i][1]=A[i][1];
for (i=1;i<=n;i++)
{
R[1][i]=A[1][i]/L[1][1];
R[i][i]=1;
}
for (j=2;j<=n;j++)
```

```

{
    for (i=j;i<=n;i++)
    {
        s=0;
        for (k=1;k<=j-1;k++)
            s=s+L[i][k]*R[k][j];
        L[i][j]=A[i][j]-s;
    }
    for (i=j+1;i<=n;i++)
    {
        s=0;
        for (k=1;k<=i-1;k++)
            s=s+L[j][k]*R[k][i];
        R[j][i]=(A[j][i]-s)/L[j][j];
    }
}
printf("Solutia este: \n");
for (i=1;i<=n;i++)
{
    for (j=1;j<=n;j++)
        printf("L[%d][%d]=%f ",i,j,L[i][j]);
    printf("\n");
}
for (i=1;i<=n;i++)
{
    for (j=1;j<=n;j++)
        printf("R[%d][%d]=%f ",i,j,R[i][j]);
    printf("\n");
}
getch();
}

```

Folosirea programului anterior pentru factorizarea matricii

$$A = \begin{pmatrix} 2 & -1 & 2 \\ 1 & 3 & 1 \\ -1 & 2 & 1 \end{pmatrix}, \text{ conduce la următoarele rezultate:}$$

```

C:\cc\bin\rundos.exe
Dati ordinul matricei A, n:
3
Dati A
A[1][1]: 2
A[1][2]: -1
A[1][3]: 3
A[2][1]: 1
A[2][2]: 1
A[2][3]: -1
A[3][1]: 1
A[3][2]: 2
A[3][3]: -1
Solutia este:
L[1][1]=2.000000 L[1][2]=0.000000 L[1][3]=0.000000
L[2][1]=1.000000 L[2][2]=1.500000 L[2][3]=0.000000
L[3][1]=1.000000 L[3][2]=2.500000 L[3][3]=1.666667
R[1][1]=1.000000 R[1][2]=-0.500000 R[1][3]=1.500000
R[2][1]=0.000000 R[2][2]=1.000000 R[2][3]=-1.666667
R[3][1]=0.000000 R[3][2]=0.000000 R[3][3]=1.000000

```

2) Factorizarea Doolittle (metoda Doolittle) :

În acest caz matricele L și R implicate sunt următoarele:

$$L = \begin{pmatrix} 1 & 0 & 0 & - & - & 0 \\ l_{21} & 1 & 0 & - & - & 0 \\ l_{31} & l_{32} & 1 & - & - & 0 \\ - & - & - & - & - & - \\ l_{n1} & l_{n2} & - & - & - & 1 \end{pmatrix} \quad R = \begin{pmatrix} r_{11} & r_{12} & r_{13} & - & - & r_{1n} \\ 0 & r_{22} & r_{23} & - & - & r_{2n} \\ 0 & 0 & r_{33} & - & - & r_{3n} \\ - & - & - & - & - & - \\ 0 & 0 & - & - & - & r_{nn} \end{pmatrix}$$

Determinarea elementelor lor se face analog, diferă doar ordinea de determinare a acestora.

Algoritmul factorizării Doolittle

Pasul 1:

- Se determină prima linie a lui R .
- Se determină prima coloană a lui L .

Pasul 2:

- Se determină a doua linie a lui R .
- Se determină a doua coloană a lui L .

.....

Pasul n :

- Se determină linia n a lui R , practic elementul r_{nn} al lui R .

Exemplul III.2. Fie A cea din exemplul III.1:

$$A = \begin{pmatrix} 2 & -1 & 2 \\ 1 & 3 & 1 \\ -1 & 2 & 1 \end{pmatrix}$$

$$\text{Fie } L = \begin{pmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{pmatrix} \quad R = \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{pmatrix}$$

Pasul 1:

Se determina prima linie a lui R

$$r_{11} = 2, \quad r_{12} = -1, \quad r_{13} = 2$$

Se determină prima coloană a lui L

$$l_{21}r_{11} = a_{21} \Rightarrow l_{21} = \frac{1}{2}, \quad l_{31}r_{11} = a_{31} \Rightarrow l_{31} = \frac{-1}{2}$$

Pasul 2:

Se determina a doua linie a lui R .

$$l_{21}r_{12} + r_{22} = a_{22} \Rightarrow r_{22} = 3 + \frac{1}{2} = \frac{7}{2}$$

$$l_{21}r_{13} + r_{23} = a_{23} \Rightarrow r_{23} = 1 - 1 = 0$$

Se determina a doua coloană a lui L

$$l_{31}r_{12} + l_{32}r_{22} = a_{32} \Rightarrow l_{32} = \frac{\left(2 - \frac{1}{2}\right)}{\frac{7}{2}} = \frac{3}{7}$$

Pasul 3

Se determină r_{33}

$$l_{31}r_{13} + l_{32}r_{23} + r_{33} = a_{33} \Rightarrow r_{33} = 1 + 1 - 0 = 2$$

Obținem astfel matricele:

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 1/2 & 1 & 0 \\ -1/2 & 3/7 & 1 \end{pmatrix} \quad R = \begin{pmatrix} 2 & -1 & 2 \\ 0 & 7/2 & 0 \\ 0 & 0 & 2 \end{pmatrix}.$$

/*Program Factorizare Doolittle*/

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
float s, A[30][30], L[30][30], R[30][30];
```

```
int i,k,j,n;
```

```
void main()
```

```
{
```

```
    printf("Dati ordinul matricei A, n:\n");
```

```
    scanf("%d",&n);
```

```
    printf("Dati A\n");
```

```
    for (i=1;i<=n;i++)
```

```
        for (j=1;j<=n;j++)
```

```
            {
```

```
                printf("A[%d][%d]: ",i,j);
```

```
                scanf("%f",&A[i][j]);
```

```
            }
```

```
    for (j=1;j<=n;j++)
```

```
        R[1][j]=A[1][j];
```

```

for (j=1;j<=n;j++)
    {
        L[j][1]=A[j][1]/R[1][1];
        L[j][j]=1;
    }
for (i=2;i<=n;i++)
{
    for (j=i;j<=n;j++)
        {
            s=0;
            for (k=1;k<=i-1;k++)
                s=s+L[i][k]*R[k][j];
            R[i][j]=A[i][j]-s;
        }
    for (j=i+1;j<=n;j++)
        {
            s=0;
            for (k=1;k<=j-1;k++)
                s=s+L[j][k]*R[k][i];
            L[j][i]=(A[j][i]-s)/R[i][i];
        }
}
printf("Solutia este: \n");
for (i=1;i<=n;i++)
    {
        for (j=1;j<=n;j++)
            printf("L[%d][%d]=%f ",i,j,L[i][j]);
        printf("\n");
    }
for (i=1;i<=n;i++)
    {
        for (j=1;j<=n;j++)
            printf("R[%d][%d]=%f ",i,j,R[i][j]);
        printf("\n");
    }
getch();
}

```

Rularea programului anterior pentru factorizarea Doolittle a matricii

$$A = \begin{pmatrix} 2 & -1 & 3 \\ 1 & 1 & -1 \\ 1 & 2 & -1 \end{pmatrix}, \text{ conduce la rezultatele de mai jos.}$$

```

C:\cc\bin\rundos.exe
Dati ordinul matricii A, n:
3
Dati A
A[1][1]: 2
A[1][2]: -1
A[1][3]: 3
A[2][1]: 1
A[2][2]: 1
A[2][3]: -1
A[3][1]: 1
A[3][2]: 2
A[3][3]: -1
Solutia este:
L[1][1]=1.000000 L[1][2]=0.000000 L[1][3]=0.000000
L[2][1]=0.500000 L[2][2]=1.000000 L[2][3]=0.000000
L[3][1]=0.500000 L[3][2]=1.666667 L[3][3]=1.000000
R[1][1]=2.000000 R[1][2]=-1.000000 R[1][3]=3.000000
R[2][1]=0.000000 R[2][2]=1.500000 R[2][3]=-2.500000
R[3][1]=0.000000 R[3][2]=0.000000 R[3][3]=1.666667
    
```

3) Factorizarea Cholesky (sau metoda rădăcinii pătrate)

Teorema III.2. Dacă $A \in M_n(\mathbb{R})$ e simetrică și pozitiv definită, adică dacă:

$$A = A^t, \quad x^t A x \geq 0 \quad (\forall) x \in \mathbb{R}^n \quad \text{și} \quad x^t A x = 0 \Leftrightarrow x = 0,$$

atunci există o factorizare LR a lui A cu $R = L^t$, deci $A = L \cdot L^t$

Astfel, în aceste condiții avem:

$$A = \begin{pmatrix} l_{11} & 0 & - & - & 0 \\ l_{21} & l_{22} & - & - & 0 \\ - & - & - & - & - \\ l_{n1} & l_{n2} & - & - & l_{nn} \end{pmatrix}, \begin{pmatrix} l_{11} & l_{21} & - & - & l_{n1} \\ 0 & l_{22} & - & - & l_{n2} \\ - & - & - & - & - \\ 0 & 0 & - & - & l_{nn} \end{pmatrix}.$$

$$l_{11} = \sqrt{a_{11}}, \dots, l_{ii} = \sqrt{a_{ii} - \sum_{j=1}^{i-1} l_{ij}^2}, \quad i = \overline{2, n}.$$

$$l_{11}l_{21} = a_{12} \Rightarrow l_{21} = \frac{a_{12}}{l_{11}} \Rightarrow l_{i1} = \frac{a_{1i}}{l_{11}}, \quad (\forall) i = \overline{2, n} \Rightarrow$$

Prima coloană și prima linie se calculează simultan. La pasul următor se determină mai întâi l_{22} , apoi și restul elementelor de pe coloana 2 a lui L . Se continuă analog pentru toți pașii până la pasul n .

În continuare este prezentat un program în C realizat pe baza algoritmului factorizării Cholesky.

/*Program Factorizare Cholesky*/

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
#include<math.h>
```

```
float s, sum, A[30][30], L[30][30];
```

```
int i, k, j, n;
```

```
void main()
```

```
{
```

```
printf("Dati ordinul matricei A, n:\n");
```

```
scanf("%d", &n);
```

```
printf("Introduceti Asimetrice si pozitiv definite\n");
```

```
for (i=1; i<=n; i++)
```



```

for (j=1;j<=n;j++)
{
    printf("A[%d][%d]: ",i,j);
    scanf("%f",&A[i][j]);
}
j=1;
L[j][j]=sqrt(A[1][1]);
for (i=2;i<=n;i++)
    L[i][j]=A[1][i]/L[1][1];

for (j=2;j<=n;j++)
{
    s=0;
    for (k=1;k<=j-1;k++)
        s=s+L[j][k]*L[j][k];
    L[j][j]=sqrt(A[j][j]-s);
    for (i=j+1;i<=n; i++)
    {
        sum=0;
        for (k=1;k<=j-1;k++)
            sum=sum+L[i][k]*L[j][k];
        L[i][j]=(A[i][j]-sum)/L[j][j];
    }
}
printf(" Matricea L este: \n");
for (i=1;i<=n;i++)
{
    for (j=1;j<=n;j++)
        printf("L[%d][%d]=%f ",i,j,L[i][j]);
    printf("\n");
}
getch();
}

```

III. 2. METODE NUMERICE PENTRU AFLAREA INVERSEI UNEI MATRICI

III.2.1. Matrice elementare

Dacă permutăm liniile sau coloanele matricii unitate, sau dacă facem o serie de combinații liniare cu acestea, obținem o serie de matrice ce se numesc matrice elementare. Aceste matrice elementare, notate în general cu E sunt folosite pentru a introduce o serie de schimbări în structura unei matrici date, fără a-i schimba rangul, sau determinantul.

De exemplu, pentru a înmulți linia a treia dintr-o matrice $A = (a_{ij})_{1 \leq i, j \leq 3}$ cu un scalar λ , construim matricea elementară care să aiba pe linia a treia în loc de 1 valoarea λ .

$$E = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \lambda \end{pmatrix}$$

Făcând produsul dintre E și A , în această ordine, obținem schimbarea propusă :

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \lambda \end{pmatrix} \cdot \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ \lambda a_{31} & \lambda a_{32} & \lambda a_{33} \end{pmatrix}$$

Să construim acum o serie de matrice elementare care să inducă în matricea data A următoarele schimbări :

E_1 să înmulțească linia a doua cu 3 ;

E_2 să adune linia a doua la linia a treia ;

E_3 să permute linia a treia cu a doua.

$$E_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix}, E_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}, E_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

Aplicăm succesiv aceste transformări elementare matriciei A . Se obține următorul rezultat :

$$E_3 \cdot E_2 \cdot E_1 \cdot A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ 3a_{21} + a_{31} & 3a_{22} + a_{32} & 3a_{23} + a_{33} \\ 3a_{21} & 3a_{22} & 3a_{23} \end{pmatrix}$$

Matricele elementare pot fi folosite pentru aflarea inversei unei matrici date, în anumite condiții restrictive.

III.2.2. Metoda transformărilor elementare succesive pentru aflarea inversei unei matrici

Această metodă constă în construirea a n matrice elementare E_1, E_2, \dots, E_n astfel ca:

1. $E_1 A$ să aibă prima coloană identică cu prima coloană a lui I_n .
2. $E_2 E_1 A$ să aibă primele două coloane identice cu cele ale lui I_n , s.a.m.d..

3. Procedeu se repetă recursiv până când se obține matricea I_n , adică relația: $E_n \dots E_2 E_1 A = I_n$.

4. Dacă A este inversabilă, atunci $E_n \dots E_2 E_1 = A^{-1}$, deci inversa ei se află făcând produsul matricelor astfel construite, considerate în ordinea inversă obținerii lor.

Fie A matricea pătratică a cărei inversă o vom afla:

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

Pasul 1

Construim

$$E_1 = \begin{bmatrix} 1 & 0 & \dots & 0 \\ a_{11} & & & \\ -\frac{a_{21}}{a_{11}} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ -\frac{a_{n1}}{a_{11}} & 0 & \dots & 1 \\ a_{11} & & & \end{bmatrix}$$

Apoi calculăm $A_1 = E_1 A$, matrice ale cărei elemente le vom nota cu a_{ij}^1

Deci:

$$A_1 = E_1 A = \begin{bmatrix} 1 & a_{12}^1 & \dots & a_{1n}^1 \\ 0 & a_{22}^1 & \dots & a_{2n}^1 \\ \dots & \dots & \dots & \dots \\ 0 & a_{n2}^1 & \dots & a_{nn}^1 \end{bmatrix}$$

Pasul 2

Construim acum matricea elementară E_2 .

$$E_2 = \begin{bmatrix} 1 & -\frac{a_{12}^1}{a_{22}^1} & \dots & 0 \\ 0 & \frac{1}{a_{22}^1} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & -\frac{a_{n2}^1}{a_{22}^1} & \dots & 1 \end{bmatrix}$$

Calculăm $A_2 = E_2 A_1 = E_2 E_1 A$, matrice ale cărei elemente le vom nota cu a_{ij}^2

$$A_2 = E_2 E_1 A = \begin{bmatrix} 1 & 0 & \dots & a_{1n}^2 \\ 0 & 1 & \dots & a_{2n}^2 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & a_{nn}^2 \end{bmatrix}$$

Continuând în mod analog se construiesc matricele E_3, \dots, E_n și apoi cu relația (2) matricea A^{-1} .

Observația III.4. Relațiile precedente arată că la pasul i trebuie să avem $a_{ii}^{i-1} \neq 0$ pentru ca algoritmul să poată fi aplicat, unde $a_{11}^0 = a_{11}$.

Exemplul III.3 Să se determine inversa matricei A , unde

$$A = \begin{pmatrix} 3 & -4 \\ 1 & 2 \end{pmatrix}.$$

Soluție:

Determinăm $E_1 = \begin{pmatrix} 1/3 & 0 \\ -1/3 & 1 \end{pmatrix}$ și calculăm $A_1 = E_1 \cdot A$:

$$E_1 \cdot A = \begin{pmatrix} 1/3 & 0 \\ -1/3 & 1 \end{pmatrix} \begin{pmatrix} 3 & -4 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & -4/3 \\ 0 & 10/3 \end{pmatrix} = A_1$$

Determinăm $E_2 = \begin{pmatrix} 1 & 4/10 \\ 0 & 3/10 \end{pmatrix}$ și calculăm $E_2 \cdot A_1$

$$E_2 \cdot A_1 = \begin{pmatrix} 1 & 4/10 \\ 0 & 3/10 \end{pmatrix} \cdot \begin{pmatrix} 1 & -4/3 \\ 0 & 10/3 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

În concluzie inversa lui A este dată de:

$$A^{-1} = E_2 \cdot E_1 = \begin{pmatrix} 1 & 4/10 \\ 0 & 3/10 \end{pmatrix} \begin{pmatrix} 1/3 & 0 \\ -1/3 & 1 \end{pmatrix}$$

$$A^{-1} = \begin{pmatrix} 1/5 & 2/5 \\ -1/10 & 3/10 \end{pmatrix}$$

Procedeu prezentat mai sus nu este singurul procedeu numeric ce poate fi utilizat pentru calculul inversei unei matrici când aceasta există.

Alte două procedee larg utilizate sunt:

-metoda lui Gauss-Jordan

-metoda iterativă.

Metoda Gauss-Jordan se aplică tabloului format din matricea A , a cărei inversă dorim să o aflăm, așezată în partea stângă, urmată de matricea unitate corespunzătoare ordinului lui A .

În tabloul final se obține în stânga matricea unitate iar în dreapta inverse matricii A .

În cazul aflării inversei unei matrici pivoții se aleg pe diagonala principală. Dacă la un anumit pas elementul ce urmează a fi ales pivot este nul, și totuși există inversa lui A , se permută două linii (coloane) astfel încât să se obțină un pivot nenul și se continuă algoritmul. Matricea inversă se obține prin permutarea coloanelor (liniilor) corespunzătoare în matricea din partea dreaptă a ultimului tablou.

Exemplul III.4. Să se afle cu ajutorul metodei Gauss-Jordan inversa matricei:

$$A = \begin{pmatrix} 1 & -3 \\ 1 & 7 \end{pmatrix}$$

Construim primul tablou ce conține în dreapta matricea unitate de ordinul trei și în partea stângă matricea A .

[1]	-3	1	0
1	7	0	1
1	-3	1	0
0	[10]	-1	1
1	0	7/10	3/10
0	1	-1/10	1/10

Din ultimul tablou deducem inversa lui A , A^{-1} .

$$A^{-1} = \begin{pmatrix} 7/10 & 3/10 \\ -1/10 & 1/10 \end{pmatrix}.$$

Exemplul III.5. Să se afle cu ajutorul metodei Gauss-Jordan inversa matricei:

$$A = \begin{pmatrix} 1 & 2 & 0 \\ -1 & 3 & 1 \\ 2 & 0 & 1 \end{pmatrix}$$

Soluție:

1	2	0	1	0	0
-1	3	1	0	1	0
2	0	1	0	0	1
1	2	0	1	0	0
0	5	1	-2	0	1
0	-4	1	-2	0	1
1	0	-2/5	3/5	-2/5	0
0	1	1/5	1/5	1/5	0
0	0	9/5	-6/5	4/5	1
1	0	0	1/3	-2/9	2/9
0	1	0	1/3	1/9	-1/9
0	0	1	-6/9	4/9	5/9

Din ultimul tablou deducem:

$$A^{-1} = \begin{pmatrix} \frac{1}{3} & \frac{-2}{9} & \frac{2}{9} \\ \frac{1}{3} & \frac{1}{9} & \frac{-1}{9} \\ \frac{-2}{3} & \frac{4}{9} & \frac{5}{9} \end{pmatrix}.$$

```
/*Program Inversa -varianta G-J cu pivotare partiala*/
```

```
#include<stdio.h>
```

```
#include<conio.h>
```

```
#include<math.h>
```

```
float t,m;
```

```
int i,j,k,l,n;
```

```
float s, A[10][20], u[10][10];
```

```
void main()
```

```
{
```

```
printf("Introduceti nr. de linii si de coloane ale matricei");
```

```
scanf("%d",&n);
```

```
printf("Dati matricea\n");
```

```
for(i=1;i<=n;i++)
```

```
for(j=1;j<=n;j++)
```

```
{
```

```
printf("A[%d][%d]: ",i,j);
```

```
scanf("%f",&A[i][j]);
```

```
}
```

```
for (i=1;i<=n;i++)
```

```
for(j=n+1;j<=n+n;j++)
```

```
{
```

```
if (j==i+n) A[i][j]=1;
```

```
else A[i][j]=0;
```

```
}
```

```
for(k=1;k<=n;k++)
```

```
{
```

```
l=k;
```

```
for (i=k;i<=n;i++)
```

```
{u[i][k]=sqrt(A[i][k]*A[i][k]);
```

```
if (u[i][k]>u[k][k])
```

```
l=i;}
```

```
if(l!=k)
```

```
for(j=1;j<=n+n;j++)
```

```
{
```

```
t=A[k][j];A[k][j]=A[l][j];A[l][j]=t;
```

```
}
```

```
for(i=1;i<=n;i++)
```

```
{
```

```

        m=A[i][k]/A[k][k];
        A[i][k]=0;
        for(j=k+1;j<=n+n;j++)
            A[i][j]=A[i][j]-m*A[k][j];
    }
    A[k][k]=1;
}
printf("Matricea inversa este: \n");
for (i=1;i<=n;i++)
{
    for (j=n+1;j<=n+n;j++)
        printf("A[%d][%d]=%f",i,j-n,A[i][j]);
    printf("\n");
}
getch();
}

```

Exemplul III.6. Să se determine, folosind metoda

transformărilor elementare, inversa matricei $A = \begin{pmatrix} 4 & -1 & 3 \\ 2 & 0 & -1 \\ 2 & 1 & -2 \end{pmatrix}$.

Soluție:

Pasul 1: Determinăm $E_1 = \begin{pmatrix} 1/4 & 0 & 0 \\ -1/2 & 1 & 0 \\ -1/2 & 0 & 1 \end{pmatrix}$

Calculăm

$$E_1 \cdot A = \begin{pmatrix} 1/4 & 0 & 0 \\ -1/2 & 1 & 0 \\ -1/2 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 4 & -1 & 3 \\ 2 & 0 & -1 \\ 2 & 1 & -2 \end{pmatrix} = \begin{pmatrix} 1 & -1/4 & 3/4 \\ 0 & 1/2 & -5/2 \\ 0 & 3/2 & -7/2 \end{pmatrix} \stackrel{not}{=} A_1$$

Pasul 2: Determinăm $E_2 = \begin{pmatrix} 1 & 1/2 & 0 \\ 0 & 2 & 0 \\ 0 & -3 & 1 \end{pmatrix}$

Calculăm

$$E_2 \cdot A_1 = \begin{pmatrix} 1 & 1/2 & 0 \\ 0 & 2 & 0 \\ 0 & -3 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & -1/4 & 3/4 \\ 0 & 1/2 & -5/2 \\ 0 & 3/2 & -7/2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & -1/2 \\ 0 & 1 & -5 \\ 0 & 0 & 4 \end{pmatrix} \stackrel{not}{=} A_2$$

Pasul 3: Determinăm $E_3 = \begin{pmatrix} 1 & 0 & 1/8 \\ 0 & 1 & 5/4 \\ 0 & 0 & 1/4 \end{pmatrix} \Rightarrow$

$$E_3 \cdot A^2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \Rightarrow A^{-1} = E_3 E_2 E_1.$$

Efectuând calculele obținem:

$$A^{-1} = \begin{pmatrix} 1/8 & 1/8 & 1/8 \\ 1/4 & -7/4 & 5/4 \\ 1/4 & -3/4 & 1/4 \end{pmatrix}.$$

IV. METODA APROXIMAȚILOR SUCCESIVE ÎN REZOLVAREA SISTEMELOR LINIARE

IV.1. METODA APROXIMAȚILOR SUCCESIVE

Definiția IV.1. Spunem că (X, d) este un spațiu metric complet dacă orice șir Cauchy cu elemente din X este șir convergent.

Definiția IV.2. Să considerăm un spațiu metric (X, d) și o aplicație $f : X \rightarrow X$. Vom spune că $x^* \in X$ este punct fix pentru f dacă $f(x^*) = x^*$.

Definiția IV.3. Aplicația f se numește contracție dacă există $a \in \mathbb{R}$, $0 < a < 1$, astfel încât $d(f(x), f(y)) \leq ad(x, y)$ pentru oricare ar fi $x, y \in X$. a se numește constanta contracției.

Teorema IV.1. (Teorema de punct fix a lui Banach) Fie (X, d) un spațiu metric complet și $f : X \rightarrow X$ o contracție. Atunci f are un unic punct fix $x^* \in X$. El este limita șirului

$$x_{n+1} = f(x_n), n = 0, 1, \dots$$

unde $x_0 \in X$ este arbitrar.

Din demonstrația teoremei se poate deduce că:

$$d(x_n, x_{n+1}) \leq a^n d(x_0, x_1), n \geq 0$$

$$d(x_n, x_{n+p}) \leq \frac{a^n}{1-a} d(x_0, x_1), n, p \in \mathbb{N}$$

Relația precedentă spune că șirul (x_n) este șir Cauchy. Întrucât (X, d) este spațiu metric complet, (x_n) va fi șir convergent. Notăm cu $x^* \in X$ limita sa. Avem:

$$\lim_{n \rightarrow \infty} d(x_n, x^*) = 0 \text{ și } x^* = f(x^*)$$

Observația IV.1. Elementele x_0, x_1, x_2, \dots se numesc aproximațiile succesive ale lui x^* . Aproximația inițială x_0 poate fi luată arbitrar în X .

Observația IV.2.

Se poate deduce că:

$$1. d(x_n, x^*) \leq \frac{a^n}{1-a} d(x_0, x_1), n \geq 0$$

Această inegalitate ne arată cât de mare poate fi eroarea pe care o obținem când aproximăm pe x^* cu x_n .

$$2. d(x_n, x^*) \leq \frac{a}{1-a} d(x_{n-1}, x_n), n \geq 1$$

IV.2. METODE ITERATIVE PENTRU REZOLVAREA SISTEMELOR DE ECUAȚII LINIARE

Metodele iterative sunt mai simple însă prezintă dezavantajul că nu se poate stabili de la început numărul de pași necesari rezolvării. Astfel sunt costisitoare, mai ales în ceea ce privește timpul de calcul. Se folosesc în general pentru

sisteme de dimensiuni mari și sunt de asemenea potrivite rezolvării sistemelor ce prezintă mulți coeficienți nuli.

Metodele iterative constau în construirea unui șir $(x^{(k)} \in R^n)$, $k = 0, 1, \dots$, convergent către soluția exactă x a sistemului: $Ax = b$, unde $A \in R^{n \times n}$, $b \in R^n$.

Oprirea procesului iterativ este influențată de eroarea admisibilă dată, notată cu ε sau cu ε_{adm} .

Observația IV.3. Pentru rezolvarea sistemului $Ax = b$ definim $f: R^n \rightarrow R^n$ astfel: $f(x) = y$, unde componentele lui y sunt date de:

$$\begin{aligned}
 y_1 &= \frac{1}{a_{11}} \left(b_1 - \sum_{j=2}^n a_{1j} x_j \right) \\
 y_2 &= \frac{1}{a_{22}} \left(b_2 - \sum_{j=3}^n a_{2j} x_j - a_{21} y_1 \right) \\
 &\vdots \\
 y_i &= \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} y_j - \sum_{j=i+1}^n a_{ij} x_j \right) \\
 &\vdots \\
 y_n &= \frac{1}{a_{nn}} \left(b_n - \sum_{j=1}^{n-1} a_{nj} y_j \right)
 \end{aligned}$$

Sunt adevărate următoarele enunțuri:

1. Această aplicație se poate defini doar atunci când diagonală principală a matricei A nu conține elemente nule ($a_{ii} \neq 0$).

2. Aplicația $d : R^n \times R^n \rightarrow R$ $d(x, y) = \max_i |x_i - y_i|$ este o metrică pe R^n iar (R^n, d) reprezintă un spațiu metric complet.

Teorema IV.2. Dacă $A \in M_n(R)$ este astfel încât

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

aplicația f este o contracție de coeficient q dat de relația:

$$q = \max_{i=1, n} \frac{\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|}{|a_{ii}|}$$

De aici obținem un algoritm de soluționare a unui sistem algebric caracterizat de o matrice ce satisface cerințele teoremei.

Algoritm. Alegem $x^{(0)} \in R^n$ și apoi construim șirul de aproximații succesive (iterații successive) folosind relația:

$$x^{(n+1)} = f(x^{(n)}) \quad n = 0, 1, 2, \dots$$

Atunci $\Rightarrow x^* = \lim_{n \rightarrow \infty} x^{(n)}$ este singurul punct fix al lui f .

Mai mult $x^* = f(x^*) \Rightarrow x^*$ este soluția sistemului $Ax = b$.

Există în principal două metode iterative ce se folosesc la rezolvarea sistemelor liniare: metoda Jacobi și metoda Gauss-Seidel.

I. Metoda Jacobi (metoda iterațiilor simultane)

Această metodă iterativă constă în construcția șirului $x^{(k)}$, $k = 0, 1, \dots$ astfel:

$x^{(0)}$ se alege arbitrar,

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

Observația IV.4. Se observă că la o iterație în formule de calcul din membrul drept se folosesc toate componentele iterației precedente.

II. Metoda Seidel-Gauss (metoda iterațiilor succesive)

Ea constă în construcția șirului $x^{(k)}, k = 0, 1, \dots$ după modelul de construcție a funcției f din paragrafele anterioare (astfel $y \rightarrow x^{(k+1)}$).

Observația IV.5. La o iterație, în formule, se înlocuiesc valorile obținute la iterația precedentă cu valorile de la iterația curentă, dacă acestea sunt calculate.

Există și cazuri în care condiția (*) nu este îndeplinită și totuși algoritmul Seidel-Gauss converge.

Metoda Jacobi este mai lent convergentă decât metoda Seidel–Gauss în aceleași condiții inițiale date $x^{(0)}$ și ε .

Observația IV.6.

Condiția (*) din teorema IV.2 poate fi înlocuită cu condiția:

A să fie diagonal dominantă pe coloane, adică

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

Observația IV.7.

În general iterațiile se opresc atunci când, pentru un ε dat, adică pentru o eroare admisibilă dată, este îndeplinită condiția:

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

Eroarea admisibilă dată reflectă gradul de precizie dorit.

Observația IV.8. Algoritmul se poate opri și dacă nu s-a ajuns la gradul de precizie dorit dar s-au efectuat prea multe iterații, adică dacă numărul de iterații atinge o valoare precizată, denumită număr maxim admis de iterații. Uneori se folosește chiar o combinație de condiții în care apar atât eroarea admisibilă cât și numărul maxim de iterații.

În continuare descriem pe larg cele două metode iterative amintite.

I. Metoda Jacobi

Metoda Jacobi este cea mai simplă dintre metodele iterative.

Considerăm un sistem cu n ecuații și n ce satisface condiția de convergență a metodei.

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n = b_2 \\ \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \\ a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + a_{nn}x_n = b_n \end{cases}$$

Rearanjăm sistemul scriindu-l sub o nouă formă, denumită formă iterativă, astfel:

$$\begin{cases} x_1 = \frac{b_1}{a_{11}} - \left(\frac{a_{12}}{a_{11}} x_2 + \frac{a_{13}}{a_{11}} x_3 + \cdots + \frac{a_{1n}}{a_{11}} x_n \right) \\ x_2 = \frac{b_2}{a_{22}} - \left(\frac{a_{21}}{a_{22}} x_1 + \frac{a_{23}}{a_{22}} x_3 + \cdots + \frac{a_{2n}}{a_{22}} x_n \right) \\ \dots \\ x_n = \frac{b_n}{a_{nn}} - \left(\frac{a_{n1}}{a_{nn}} x_1 + \frac{a_{n2}}{a_{nn}} x_2 + \cdots + \frac{a_{n,n-1}}{a_{nn}} x_{n-1} \right) \end{cases}$$

Această formă ne permite ca la o nouă iterație să evaluăm noi valori ale soluției, folosind în calcul, în membrii dreپți, valorile determinate la iterația anterioară. Notând cu x_i^0 componentele vectorului soluție inițială calculăm la pasul k , $k \geq 1$ valorile soluțiilor astfel:

$$\begin{cases} x_1^k = \frac{b_1}{a_{11}} - \left(\frac{a_{12}}{a_{11}} x_2^{k-1} + \frac{a_{13}}{a_{11}} x_3^{k-1} + \cdots + \frac{a_{1n}}{a_{11}} x_n^{k-1} \right) \\ x_2^k = \frac{b_2}{a_{22}} - \left(\frac{a_{21}}{a_{22}} x_1^{k-1} + \frac{a_{23}}{a_{22}} x_3^{k-1} + \cdots + \frac{a_{2n}}{a_{22}} x_n^{k-1} \right) \\ \dots \\ x_n^k = \frac{b_n}{a_{nn}} - \left(\frac{a_{n1}}{a_{nn}} x_1^{k-1} + \frac{a_{n2}}{a_{nn}} x_2^{k-1} + \cdots + \frac{a_{n,n-1}}{a_{nn}} x_{n-1}^{k-1} \right) \end{cases}$$

Șirul de valori găsite converge spre soluția sistemului. Se consideră de obicei $x_i^0 = 0$.

Algoritmul de rezolvare a sistemelor de ecuații liniare prin metoda Jacobi

Date de intrare: rangul n , matricea coeficienților A , vectorul termenilor liberi b , eroarea admisibilă ε_{adm} și numărul maxim de iterații admise N_{max} ;

Date de ieșire: soluția aproximativă

1. Se introduc datele de intrare;
2. Se stabilește aproximația inițială, se alege ca fiind cea nulă:

$$x_i \leftarrow 0, i = \overline{1, n};$$

3. Se inițializează procesului iterativ: $k \leftarrow 1$;
4. Se execută calculele:
 - a. Se calculează noua aproximație în vectorul y :

$$y_i \leftarrow \frac{b_i}{a_{ii}} - \sum_{\substack{j=1 \\ j \neq i}}^n \frac{a_{ij}}{a_{ii}} x_j, i = \overline{1, n}$$

- b. Se calculează eroarea absolută maximă în iterația curentă:

$$D = \max_i |y_i - x_i|$$

- c. Se consemnează execuția iterației $k \leftarrow k + 1$

5. Atâta vreme cât $D \geq \varepsilon_{adm}$ și $k \leq N_{max}$
6. Dacă $k \leq N_{max}$ se afișează soluția aproximativă

$$y_i, i = \overline{1, n} \text{ și numărul de iterații efectuate } k;$$

7. Altfel se afișează mesajul "S-a depășit numărul maxim de iterații" .
8. Stop.

II. Metoda Gauss-Seidel (G-S)

Fie sistemul liniar de n ecuații cu n necunoscute scris sub formă iterativă:

$$\left\{ \begin{array}{l} x_1 = \frac{b_1}{a_{11}} - \left(\frac{a_{12}}{a_{11}} x_2 + \frac{a_{13}}{a_{11}} x_3 + \dots + \frac{a_{1n}}{a_{11}} x_n \right) \\ x_2 = \frac{b_2}{a_{22}} - \left(\frac{a_{21}}{a_{22}} x_1 + \frac{a_{23}}{a_{22}} x_3 + \dots + \frac{a_{2n}}{a_{22}} x_n \right) \\ \vdots \\ x_n = \frac{b_n}{a_{nn}} - \left(\frac{a_{n1}}{a_{nn}} x_1 + \frac{a_{n2}}{a_{nn}} x_2 + \dots + \frac{a_{n,n-1}}{a_{nn}} x_{n-1} \right) \end{array} \right.$$

Amintim că starea inițială a soluției este definită prin valorile x_i^0 , $i = 1, n$.

Determinăm valoarea lui x_1 la pasul 1 (notată cu x_1^1), ca la metoda Jacobi:

$$x_1^1 = \frac{b_1}{a_{11}} - \left(\frac{a_{12}}{a_{11}} x_2^0 + \frac{a_{13}}{a_{11}} x_3^0 + \dots + \frac{a_{1n}}{a_{11}} x_n^0 \right)$$

În evaluarea lui x_2^1 ținem cont că pentru x_1 cunoaștem două valori: x_1^0 și x_1^1 (calculată anterior). Cum x_1^1 reprezintă o valoare mai recentă o vom prefera pentru accelerarea convergenței. Deci:

$$x_2^1 = \frac{b_2}{a_{22}} - \left(\frac{a_{21}}{a_{22}} x_1^1 + \frac{a_{23}}{a_{22}} x_3^0 + \dots + \frac{a_{2n}}{a_{22}} x_n^0 \right), \text{ etc.}$$

Pentru pasul k putem scrie:

$$x_1^k = \frac{b_1}{a_{11}} - \left(\frac{a_{12}}{a_{11}} x_2^{k-1} + \frac{a_{13}}{a_{11}} x_3^{k-1} + \dots + \frac{a_{1n}}{a_{11}} x_n^{k-1} \right)$$

⋮

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

⋮

$$x_n^k = \frac{b_n}{a_{nn}} - \left(\frac{a_{n1}}{a_{nn}} x_1^k + \frac{a_{n2}}{a_{nn}} x_2^k + \dots + \frac{a_{nn-1}}{a_{nn}} x_{n-1}^k \right)$$

Observația IV.9. Algoritmul Gauss-Seidel are deci la bază algoritmul Jacobi, dar se distinge de acesta prin folosirea celor mai recente valori ale necunoscutelor. Astfel se obține pe ansamblu o accelerare a procesului de iterație.

Algoritm de rezolvare a sistemelor de ecuații liniare prin metoda Gauss-Seidel

Date de intrare: rangul n , matricea coeficienților A , vectorul termenilor liberi b , eroarea admisibilă ε_{adm} și numărul maxim de iterații admise N_{max} ;

Date de ieșire: soluția aproximativă

1. Se introduc datele de intrare;
2. Se alege aproximația inițială, aceasta se alege a fi identic nulă:

$$x_i \leftarrow 0, i = \overline{1, n}$$

3. Se inițializează procesului iterativ: $k \leftarrow 0$;
4. Se inițializează eroarea absolută maximă în iterația curentă cu o valoare superioară lui ε_{adm} : $D \leftarrow \varepsilon_{adm} + 1$;

5. Se trece la o nouă iterație: $k \leftarrow k + 1$;
6. Se calculează noua aproximație în vectorul y :

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

7. Se calculează eroarea absolută maximă în iterația curentă:

$$D = \max_i |y_i - x_i|$$

8. Dacă $D \leq \varepsilon_{adm}$ (metoda converge) se trece la pasul 11, altfel (dacă $D > \varepsilon_{adm}$), se compară k și N_{max} ; dacă $k \geq N_{max}$ (metoda nu converge), se afișează mesajul "Depasire nr. maxim iteratii" și se trece la pasul 12; altfel se trece la pasul următor;

9. Se actualizează vectorul aproximațiilor din ultima iterație:

$$x_i \leftarrow y_i, i = \overline{1, n}$$

10. Se revine la pasul 5;
11. Se afișează soluția aproximativă $y_i, i = \overline{1, n}$ și numărul de iterații efectuate k ;
12. Stop.

Exemplul IV.1. Să se rezolve următorul sistem de ecuații prin metoda Jacobi :

$$\begin{cases} -x_1 + 10x_2 + x_3 = 20 \\ 2x_1 - x_2 + 10x_3 = 22 \\ 10x_1 - x_2 + 2x_3 = 12 \end{cases}$$

Se va accepta o soluție pentru care $\max_i |x_i^k - x_i^{k-1}| < \varepsilon_{adm}$, unde $\varepsilon_{adm} = 0,002$. Se va lucra cu patru zecimale.

Soluție:

Pentru ca sistemul de ecuații trebuie să fie astfel ca elementele diagonale din matricea coeficienților să fie dominante în valoare absolută, schimbăm poziția ecuațiilor astfel: ecuația 3 devine ecuația 1, ecuația 1 devine ecuația 2 iar ecuația 2 devine ecuația 3 :

$$\begin{cases} 10x_1 - x_2 + 2x_3 = 12 \\ -x_1 + 10x_2 + x_3 = 20 \\ 2x_1 - x_2 + 10x_3 = 22 \end{cases}$$

Rescriem sistemul în forma sa iterativă :

$$\begin{cases} x_1 = 1,2 + 0,1x_2 - 0,2x_3 \\ x_2 = 2 + 0,1x_1 - 0,1x_3 \\ x_3 = 2,2 - 0,2x_1 + 0,1x_2 \end{cases}$$

- Pasul 1.

Considerăm ca soluție inițială : $x_1^0 = x_2^0 = x_3^0 = 0$.

$$\begin{cases} x_1^1 = 1,2 + 0,1 \cdot 0 - 0,2 \cdot 0 = 1,2 & ; & |x_1^1 - x_1^0| = 1,2 \\ x_2^1 = 2 + 0,1 \cdot 0 - 0,1 \cdot 0 = 2 & ; & |x_2^1 - x_2^0| = 2 \\ x_3^1 = 2,2 - 0,2 \cdot 0 + 0,2 \cdot 0 = 2,2 & ; & |x_3^1 - x_3^0| = 2,2 \end{cases}$$

$$\max_i |x_i^1 - x_i^0| = 2,2 > \varepsilon_{adm}$$

- Pasul 2

Valorile determinate în pasul precedent x_1^1, x_2^1, x_3^1 se introduc în membrul drept al sistemului:

$$\begin{cases} x_1^2 = 1,2 + 0,1 \cdot 2 - 0,2 \cdot 2,2 = 0,96 & ; \quad |x_1^2 - x_1^1| = 0,24 \\ x_2^2 = 2 + 0,1 \cdot 1,2 - 0,1 \cdot 2,2 = 1,9 & ; \quad |x_2^2 - x_2^1| = 0,1 \\ x_3^2 = 2,2 - 0,2 \cdot 1,2 + 0,2 \cdot 2 = 2,16 & ; \quad |x_3^2 - x_3^1| = 0,16 \end{cases}$$

$$\max_i |x_i^2 - x_i^1| = 0,24 > \varepsilon_{adm}$$

- Pasul 3

Valorile determinate la pasul precedent $x_1^2; x_2^2; x_3^2$ se introduc în membrul drept al sistenului și se obține:

$$\begin{cases} x_1^3 = 0,958 & ; \quad |x_1^3 - x_1^2| = 0,002 \\ x_2^3 = 1,88 & ; \quad |x_2^3 - x_2^2| = 0,02 \\ x_3^3 = 2,198 & ; \quad |x_3^3 - x_3^2| = 0,038 \end{cases}$$

$$\max_i |x_i^3 - x_i^2| = 0,038 > \varepsilon_{adm}$$

- Pasul 4

Valorile determinate la pasul precedent $x_1^3; x_2^3; x_3^3$ se introduc în membrul drept al sistemului de ecuații și se obține:

$$\begin{cases} x_1^4 = 0,9484 & ; \quad |x_1^4 - x_1^3| = 0,0096 \\ x_2^4 = 1,8760 & ; \quad |x_2^4 - x_2^3| = 0,004 \\ x_3^4 = 2,1964 & ; \quad |x_3^4 - x_3^3| = 0,0016 \end{cases}$$

$$\max_i |x_i^4 - x_i^3| = 0,0096 > \varepsilon_{adm}$$

- Pasul 5

Valorile determinate la pasul precedent $x_1^4; x_2^4; x_3^4$ se introduc în membrul drept al sistemului de ecuații și se obține:

$$\begin{cases} x_1^5 = 0,9483 & ; & |x_1^5 - x_1^4| = 0,0001 \\ x_2^5 = 1,8752 & ; & |x_2^5 - x_2^4| = 0,0008 \\ x_3^5 = 2,1979 & ; & |x_3^5 - x_3^4| = 0,0015 \end{cases}$$

$$\max_i |x_i^5 - x_i^4| = 0,0015 < 0,002 = \varepsilon_{adm}$$

Deoarece s-a obținut precizia de calcul solicitată, procesul de iterație se oprește.

Soluția sistemului este :

$$\begin{cases} x_1 = x_1^5 = 0,9483 \\ x_2 = x_2^5 = 1,8752 \\ x_3 = x_3^5 = 2,1979 \end{cases}$$

Un program în C, care furnizează soluția aproximativă a unui sistem liniar, obținută după un număr dat de iterații, prin metoda Jacobi este prezentat în continuare.

```
/*Program - Metoda JACOBI-nr maxim iteratii*/
#include<stdio.h>
#include<conio.h>
#include<math.h>
float s, A[30][30], x[30],B[30], y[30];
int n,k,j,i,Nmax;
void main()
{
printf ("Introduceti dimensiunea matricei sistemului ");
scanf("%d",&n);
printf("Dati matricea coeficientilor\n");
for(i=1;i<=n;i++)
for(j=1;j<=n;j++)
{
printf("A[%d][%d]: ",i,j);
```

```

scanf("%f",&A[i][j]);
}
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
    {
        printf("B[%d]: ",i);
        scanf("%f",&B[i]);
    }
    printf ("Introduceti nr. maxim de iteratii Nmax= ");
scanf("%d",&Nmax);
for(i=1;i<=n;i++)
    x[i]=0;
for(k=1;k<=Nmax;k++)
    {
    for(i=1;i<=n;i++)
        {
        s=0;
        for(j=1;j<=n;j++)
            s=s+A[i][j] * x[j];
        s=s- A[i][i] * x[i];
        y[i]= (B[i]-s)/A[i][i];
        }
    for(i=1;i<=n;i++)
        x[i]=y[i];
    }
printf("Solutia dupa %d iteratii este:\n", Nmax);
for(i=1;i<=n;i++)
    printf("x[%d]=%f\n",i,x[i]);
    getch();
}

```

Prin rularea programului precedent, obținem după 5 iterații, pentru problema anterioară, rezultatele următoare (abstracție făcând de erorile de rotunjire, putem spune că sunt identice cu cele anterioare).

```
C:\cc\bin\rundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 10
A[1][2]: -1
A[1][3]: 2
A[2][1]: -1
A[2][2]: 10
A[2][3]: 1
A[3][1]: 2
A[3][2]: -1
A[3][3]: 10
Dati termenii liberi:
B[1]: 12
B[2]: 20
B[3]: 22
Introduceti nr. maxim de iteratii Nmax= 5
Solutia dupa 5 iteratii este:
x[1]=0.948320
x[2]=1.875200
x[3]=2.197920
```

Exemplul IV.2. Să se determine, folosind metoda Jacobi, soluția aproximativă, după 10 iterații, pentru sistemul următor:

$$\begin{cases} 7x_1 + x_2 - x_3 = -4 \\ 2x_1 + 10x_2 - 2x_3 = 1 \\ x_1 + 3x_2 + 5x_3 = 2 \end{cases}$$

Prin rularea programului sursă ce folosește numărul maxim de iterații obținem, după 10 iterații pentru problema precedentă rezultatele următoare:

```
C:\cc\bin\rundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 7
A[1][2]: 1
A[1][3]: -1
A[2][1]: 2
A[2][2]: 10
A[2][3]: -2
A[3][1]: 1
A[3][2]: 3
A[3][3]: 5
Dati termenii liberi:
B[1]: -4
B[2]: 1
B[3]: 2
Introduceti nr. maxim de iteratii Nmax= 10
Solutia dupa 10 iteratii este:
x[1]=-0.562485
x[2]=0.281270
x[3]=0.343736
```

Următorul cod sursă (în C) rezolvă un sistem de ecuații prin algoritmul metodei Jacobi dar folosește pentru oprire atât un număr maxim de iterații admis, cât și o eroare admisibilă dată.

```
/*Program Metoda JACOBI-epsilon si Nmax*/
#include<stdio.h>
#include<conio.h>
#include<math.h>
float s, max, eps,A[30][30],x[30],B[30],y[30];
int n,Nmax,j,i,k;
void main()
{
printf ("Introduceti dimensiunea matricei sistemului ");
scanf("%d",&n);
printf("Dati matricea coeficientilor\n");
for(i=1;i<=n;i++)
for(j=1;j<=n;j++)
{
printf("A[%d][%d]: ",i,j);
scanf("%f",&A[i][j]);
}
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
{
printf("B[%d]: ",i);
scanf("%f",&B[i]);
}
printf ("Introduceti nr. maxim de iteratii m= ");
scanf("%d",&Nmax);
printf ("Introduceti precizia dorita eps= ");
scanf("%f",&eps);
for(i=1;i<=n;i++) x[i]=0;
k=0;
```

```

do
{
for(i=1;i<=n;i++)
{
s=0;
for(j=1;j<=n;j++)
s=s+A[i][j] * x[j];
s=s- A[i][i] * x[i];
y[i]= (B[i]-s)/A[i][i];
}
max=fabs(y[1]-x[1]);
for(i=2;i<=n;i++)
if (max<fabs(y[i]-x[i]))
max=fabs((y[i]-x[i]));
for(i=1;i<=n;i++) x[i]=y[i];
k=k+1;
}
while (max>=eps && k<=Nmax);
if (k<=Nmax)
{
printf("Dupa %d iteratii solutia este:\n",k);
for(i=1;i<=n;i++)
printf("x[%d]=%f\n",i,x[i]);
}
else printf("S-a depasit numarul maxim de iteratii");
getch();
}

```

Exemplul IV.3. Să se rezolve sistemul de ecuații din exemplul anterior utilizând sursa de mai sus, considerând eroarea admisibilă 10^{-3} .

```
C:\cd\bin\rundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 7
A[1][2]: 1
A[1][3]: -1
A[2][1]: 2
A[2][2]: 10
A[2][3]: -2
A[3][1]: 1
A[3][2]: 3
A[3][3]: 5
Dati termenii liberi:
B[1]: -4
B[2]: 1
B[3]: 2
Introduceti nr. maxim de iteratii n= 10
Introduceti precizia dorita eps= 0.001
Dupa 8 iteratii solutia este:
x[1]=-0.562620
x[2]=0.281091
x[3]=0.343715
```

Se observă că soluția aproximativă a sistemului s-a obținut după un număr de 8 iterații, în limitele preciziei dorite, iar numărul maxim de iterații admis nu a fost depășit.

Vom căuta soluția acestui sistem măbind precizia, adică micșorând eroarea admisibilă dată de 10 ori (aducând-o la valoarea 10^{-4}), dar menținând același timp de execuție, adică același număr maxim de iterații admis.

Se observă din rezultatele următoare că nu putem atinge precizia dorită folosind același număr maxim de iterații admis.

```
C:\cc\bin\yundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 7
A[1][2]: 1
A[1][3]: -1
A[2][1]: 2
A[2][2]: 10
A[2][3]: -2
A[3][1]: 1
A[3][2]: 3
A[3][3]: 5
Dati termenii liberi:
B[1]: -4
B[2]: 1
B[3]: 2
Introduceti nr. maxim de iteratii n= 10
Introduceti precizia dorita eps= 0.0001
S-a depasit numarul maxim de iteratii
```

Pentru rezolvarea sistemelor de ecuatii liniare prin metoda Seidel Gauss sunt concepute în continuare două programe simple, în limbajul C. Primul se bazează pe afișarea soluției numerice ce se obține după un număr maxim de iterații admis, Nmax, ce se cere a fi introdus de către utilizator.

```
/*Program - Metoda Seidel-Gauss-nr maxim iteratii*/
#include<stdio.h>
#include<conio.h>
float s, A[30][30],x[30],B[30], y[30];
int n,k,j,i,Nmax;
void main()
{
printf ("Introduceti dimensiunea matricei sistemului ");
scanf("%d",&n);
printf("Dati matricea coeficientilor\n");
for(i=1;i<=n;i++)
```

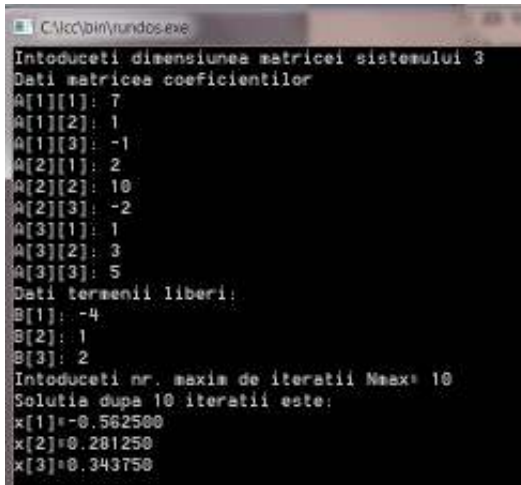


```

for(j=1;j<=n;j++)
{
printf("A[%d][%d]: ",i,j);
scanf("%f",&A[i][j]);
}
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
{
printf("B[%d]: ",i);
scanf("%f",&B[i]);
}
printf ("Introduceti nr. maxim de iteratii Nmax= ");
scanf("%d",&Nmax);
for(i=1;i<=n;i++)
x[i]=0;
for(k=1;k<=Nmax;k++)
{
for(i=1;i<=n;i++)
{
s=0;
for(j=1;j<=i-1;j++)
s=s+A[i][j] * y[j];
for(j=i+1;j<=n;j++)
s=s+A[i][j] * x[j];
y[i]= (B[i]-s)/A[i][i];
x[i]=y[i];
}
}
printf("Solutia dupa %d iteratii este:\n", Nmax);
for(i=1;i<=n;i++)
printf("x[%d]=%f\n",i,x[i]);
getch();
}

```

Soluția obținută cu varianta algoritmului ce folosește doar numărul maxim de iterații pentru oprire, pentru sistemul de la exemplul IV.2., în cazul în care acest număr maxim este 10 este:



```
C:\cc\bin\yundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 7
A[1][2]: 1
A[1][3]: -1
A[2][1]: 2
A[2][2]: 10
A[2][3]: -2
A[3][1]: 1
A[3][2]: 3
A[3][3]: 5
Dati termenii liberi:
B[1]: -4
B[2]: 1
B[3]: 2
Introduceti nr. maxia de iteratii Nmax: 10
Solutia dupa 10 iteratii este:
x[1]=0.562500
x[2]=0.281250
x[3]=0.343750
```

```
/*Program Metoda Seidel Gauss-epsilon si Nmax*/
#include<stdio.h>
#include<conio.h>
#include<math.h>
float s, max, eps,A[30][30],x[30],B[30],y[30];
int n,Nmax,j,i,k;
void main()
{
printf ("Introduceti dimensiunea matricei sistemului ");
scanf("%d",&n);
printf("Dati matricea coeficientilor\n");
for(i=1;i<=n;i++)
for(j=1;j<=n;j++)
```

```

{
printf("A[%d][%d]: ",i,j);
scanf("%f",&A[i][j]);
}
printf("Dati termenii liberi:\n");
for (i=1;i<=n;i++)
    {
    printf("B[%d]: ",i);
    scanf("%f",&B[i]);
    }
    printf ("Introduceti nr. maxim de iteratii m= ");
scanf("%d",&Nmax);
printf ("Introduceti precizia dorita eps= ");
scanf("%f",&eps);
for(i=1;i<=n;i++) x[i]=0;
k=0;
do
    {
    for(i=1;i<=n;i++)
        {
        s=0;
        for(j=1;j<=i-1;j++)
            s=s+A[i][j] * y[j];
        for(j=i+1;j<=n;j++)
            s=s+A[i][j] * x[j];
        y[i]= (B[i]-s)/A[i][i];
        }
    max=fabs(y[1]-x[1]);
    for(i=2;i<=n;i++)
        if (max<fabs(y[i]-x[i]))
            max=fabs((y[i]-x[i]));
    for(i=1;i<=n;i++) x[i]=y[i];
    k=k+1;
    }
while (max>=eps && k<=Nmax);

```

```

if (k<=Nmax)
{
    printf("Dupa %d iteratii solutia este:\n",k);
    for(i=1;i<=n;i++)
        printf("x[%d]=%f\n",i,x[i]);
}
else printf("S-a depasit numarul maxim de iteratii");
getch();
}

```

Exemplul IV.4. Să se rezolve sistemul de ecuații din exemplul IV.3. prin metoda Gauss-Seidel, cu aceeași eroare admisibilă pentru soluția finală.

```

C:\vc\bin\rundos.exe
Introduceti dimensiunea matricei sistemului 3
Dati matricea coeficientilor
A[1][1]: 7
A[1][2]: 1
A[1][3]: -1
A[2][1]: 2
A[2][2]: 10
A[2][3]: 2
A[3][1]: 1
A[3][2]: 3
A[3][3]: 5
Dati termenii liberi:
B[1]: -4
B[2]: 1
B[3]: 2
Introduceti nr. maxim de iteratii Nmax 10
Introduceti precizia dorita eps: 0.001
Dupa 5 iteratii solutia este:
x[1]=-0.526333
x[2]=0.118428
x[3]=0.434210

```

Remarcăm faptul că se ajunge la precizia dorită după doar 5 iterații prin folosirea metodei Seidel-Gauss. Același sistem de ecuații a fost rezolvat cu o aceeași eroare admisibilă în 8 pași prin metoda Jacobi. Se deduce de aici accelerarea convergenței în cazul metodei Gauss-Seidel.

V. METODE NUMERICE PENTRU REZOLVAREA ECUAȚIILOR ȘI SISTEMELOR DE ECUAȚII NELINIARE

V.1. METODE NUMERICE PENTRU REZOLVAREA ECUAȚIILOR NELINIARE

Pentru o funcție : $f: [a,b] \rightarrow \mathbb{R}$, continuă și derivabilă (sau doar când $f \in C([a,b])$) dorim să găsim rădăcinile ecuației:

$$f(x) = 0. \quad (*)$$

Dacă f este o funcție polinomială de gradul 1, 2, sau 3 există formule generale de rezolvare, dar dacă gradul lui f este mai mare decât 4 nu mai dispunem de astfel de formule.

Definiția V.1 Spunem că o rădăcină reală α este separată într-un interval $[a,b]$ dacă acest interval conține o singură rădăcină a ecuației, adică doar rădăcina α .

Definiția V.2 Spunem că o rădăcină α , separată într-un interval $[a,b]$, este localizată în limitele unei precizii ε , prin x_n , într-un interval $[a_n, b_n]$, dacă este îndeplinită una din condițiile:

$$|f(x_n)| < \varepsilon, \quad (1)$$

$$|b_n - a_n| < \varepsilon. \quad (2)$$

Relațiile (1), respectiv (2), sunt folosite la oprirea algoritmilor de determinare a rădăcinilor ecuațiilor neliniare.

Cele mai cunoscute metode numerice pentru aflarea rădăcinilor ecuațiilor neliniare sunt: metoda biseecției, metoda coardei, metoda secantei, metoda aproximațiilor successive (sau a contracției), metoda lui Newton (sau a tangentei), etc.

V.2. METODA BISECȚIEI

Această metodă constă în reducerea intervalului de separare, prin înjumătățiri repetate și selectarea subintervalului în care se află rădăcina. Procesul se încheie în momentul satisfacerii uneia dintre condițiile (1) sau (2).

Considerăm deci ecuația: $f(x) = 0$, $x \in [a, b]$ cu f continuă, astfel încât $f(a)f(b) < 0$, iar în $[a, b]$ a fost separată o soluție a ecuației.

Introducem notația $a_0 = a$, $b_0 = b$. Determinăm mijlocul intervalului $[a_0, b_0]$, pe care îl notăm cu $x_0 = \frac{a_0 + b_0}{2}$.

Acesta separă intervalul inițial în două subintervale $[a_0, x_0]$, $[x_0, b_0]$, rădăcina α găsindu-se într-unul din ele, și anume în acela la capetele căruia funcția prezintă semne contrare. Vom nota acel interval cu $[a_1, b_1]$. În mod evident:

$$|\alpha - x_0| < \frac{b_0 - a_0}{2}$$

Urmând același procedeu se obțin intervalele $[a_1, b_1]$, $[a_2, b_2]$, $[a_3, b_3]$, ..., $[a_n, b_n]$ și mijloacele acestora:

$$x_n = \frac{a_n + b_n}{2}, \quad n \geq 0 \quad (3)$$

La fiecare pas selectarea intervalului următor se face astfel:

Dacă $f(a_n)f(x_n) < 0$ atunci:

$$a_{n+1} = a_n, \quad b_{n+1} = x_n,$$

altfel

$$a_{n+1} = x_n, \quad b_{n+1} = b_n. \quad (4)$$

Se observă că:

$$b_{n+1} - a_{n+1} = \frac{b_n - a_n}{2} = \dots = \frac{b_0 - a_0}{2^{n+1}} \quad (5)$$

$$\text{și } |\alpha - x_n| < b_{n+1} - a_{n+1} = \frac{b_0 - a_0}{2^{n+1}} \quad (6)$$

Relația (6) ne arată că șirurile $\{a_n\}$ și $\{b_n\}$ sunt convergente către α .

Algoritmul metodei bisecției

Date de intrare:

a, b = capetele intervalului de separare,

f = funcția căreia i se localizează rădăcina.

ε = precizia determinării.

Date de ieșire:

x = valoarea aproximativă a rădăcinii din intervalul considerat.

1. Se inițializează valoarea aproximativă a rădăcinii (valoarea corespunzătoare mijlocului intervalului dat)

$$x \leftarrow (a+b)/2$$

2. Atâta timp cât $|f(x)| \geq \varepsilon$ repetă:

- dacă $f(a) \cdot f(x) < 0$ atunci
 - $b \leftarrow x$ altfel
 - $a \leftarrow x$
- $x \leftarrow (a+b)/2$.

3. Dacă $f(x) = 0$ atunci

- se afișează x precizându-se că este soluția exactă a ecuației neliniare, altfel
- se afișează x precizându-se că este valoarea aproximativă a rădăcinii

4. Stop

Observația V.1. Oprirea algoritmului de mai sus se bazează pe o condiție de tip (1). Ea poate fi înlocuită de o condiție de tip (2), caz în care algoritmul devine:

1. Se inițializează valoarea aproximativă a rădăcinii (valoarea corespunzătoare mijlocului intervalului dat)
 $x \leftarrow (a+b)/2$
2. Atâta timp cât $|b - a| \geq \varepsilon$ repetă
 - dacă $f(a), f(x) < 0$ atunci
 - $b \leftarrow x$ altfel
 - $a \leftarrow x$
 - $x \leftarrow (a+b)/2$.
3. Se afișează valoarea lui x
4. Stop

Se observă că în acest caz nu se mai precizează dacă este vorba de soluția exactă sau aproximativă a ecuației.

Observația V.2. Se poate introduce o condiție de oprire și cu ajutorul numărului maxim admis de iterații, respectiv cu ajutorul unor combinații între condițiile amintite.

Exemplul V.1. Să se găsească valoarea aproximativă a lui $\sqrt{3}$ folosind metoda biseției și executând un număr maxim de 5 iterații.

Soluție:

Atașăm ecuația corespunzătoare: $x^2 - 3 = 0 \Rightarrow f(x) = x^2 - 3$

Considerăm $f : [1, 2] \rightarrow \mathbb{R} \Rightarrow [a, b] = [1, 2]$. În acest interval este separată o singură rădăcină a acestei ecuații. Acest lucru poate fi justificat simplu deoarece se știe că rădăcinile acestei ecuații sunt $\pm \sqrt{3}$, iar $f(1)f(2) < 0$.

$$\text{Calculăm } x_0 = \frac{a+b}{2} = \frac{1+2}{2} = \frac{3}{2}$$

$$f(2) \cdot f\left(\frac{3}{2}\right) = 1 \cdot \left(-\frac{3}{4}\right) = -\frac{3}{4} < 0 \Rightarrow [a, b] = \left[\frac{3}{2}, 2\right]$$

$$\text{Calculăm } x_1 = \frac{a_1+b_1}{2} = \frac{7}{4} = 1,75$$

$$f\left(\frac{3}{2}\right) \cdot f\left(\frac{7}{4}\right) = -\frac{3}{4} \cdot \frac{1}{16} < 0 \Rightarrow [a_2, b_2] = \left[\frac{3}{2}, \frac{7}{4}\right]$$

$$\text{Calculăm } x_2 = \frac{a_2+b_2}{2} = \frac{\frac{3}{2} + \frac{7}{4}}{2} = \frac{13}{8}$$

$$\Rightarrow f\left(\frac{13}{8}\right) \cdot f\left(\frac{7}{4}\right) = \left(-\frac{169}{64} - 3\right) \cdot \left(\frac{49}{16} - 3\right) < 0 \Rightarrow [a_3, b_3] = \left[\frac{13}{8}, \frac{7}{4}\right]$$

$$\text{Calculăm } x_3 = \frac{a_3+b_3}{2} = \frac{\frac{13}{8} + \frac{7}{4}}{2} = \frac{27}{16}$$

$$f\left(\frac{27}{16}\right) \cdot f\left(\frac{7}{4}\right) = (-) \cdot (+) < 0 \Rightarrow [a_4, b_4] = \left[\frac{27}{16}, \frac{7}{4}\right]$$

$$\text{Calculăm } x_4 = \frac{a_4+b_4}{2} = \frac{\frac{27}{16} + \frac{7}{4}}{2} = \frac{55}{32} = 1.71875$$

$$f\left(\frac{55}{32}\right) \cdot f\left(\frac{7}{4}\right) = (-) \cdot (+) < 0 \Rightarrow [a_5, b_5] = \left[\frac{55}{32}, \frac{7}{4}\right] \Rightarrow \frac{55}{32} < \sqrt{3} < \frac{7}{4}$$

Astfel $x_5 = \frac{a_5+b_5}{2} = 1.734375$, este valoarea aproximativă cerută.

Programul următor, realizat în limbajul de programare C, se referă la calculul valorii aproximative a lui $\sqrt{3}$ cu ajutorul metodei biseției. El poate fi folosit și pentru aflarea valorilor aproximative ale rădăcinilor reale ale altor ecuații, modificând expresia funcției cu cea corespunzătoare ecuației în cauză).

```

/* Program Bisectia1-condiție tip (1) */
#include<conio.h>
#include<stdio.h>
#include<math.h>
float a,b,eps,x,C;
float F(float u)
{return u*u-3;}
void main ()
{
printf("introduceti intervalul in care este separata radacina \n");
printf("introduceti capatul din stanga a:");
scanf("%f",&a);
printf("\nintroduceti capatul din dreapta b:");
scanf("%f",&b);
printf("eroarea:");
scanf("%f",&eps);
x=(a+b)/2;
A=F(x);
while (fabs(A)>=eps)
{
C=F(a)*F(x);
if (C<0)
b=x;
else a=x;
x=(a+b)/2;
printf("\n%f",x); /*afiseaza valorile intermediare*/
}
}

```

```

        A=F(x);
    }
if (A==0)
    printf ("\n%f este solutia exacta a ecuatiei", x);
else
    printf ("\n%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

Rezultatele obținute prin rularea acestui program pentru datele de intrare de mai jos sunt:

```

C:\oc\bin\runDOS.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1
introduceti capatul din dreapta b:2
eroarea:0.001
1.750000
1.625000
1.687500
1.718750
1.734375
1.726563
1.730469
1.732422
1.731445
1.731934
1.731934 este solutia aproximativa a ecuatiei_

```

Se observă că $x_5 = 1.734375$, adică cea obținută anterior.

Modificând intervalul de separație dar menținând aceeași eroare admisibilă obținem soluția dorită după un număr mai mare de iterații, așa cum se observă din datele următoare:

```

C:\cc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:3
eroarea:0.001

1.500000
1.750000
1.625000
1.687500
1.718750
1.734375
1.726563
1.730469
1.732422
1.731445
1.731934
1.731934 este solutia aproximativa a ecuatiei_

```

V.3. METODA COARDEI

Fie f continuă și derivabilă de două ori pe $[a, b]$, astfel încât $f(a)f(b) < 0$ și $f''(x)$ păstrează semn constant pe intervalul considerat. De asemenea presupunem că în $[a, b]$ a fost separată o soluție a ecuației: $f(x) = 0$.

Metoda coardei sau a falsei poziții constă în construirea unui șir de aproximații care să convergă către soluția ecuației separată în acest interval, în care termenii șirului reprezintă abscisele punctelor de intersecție ale axei Ox cu diverse drepte determinate de două puncte situate pe graficul funcției f , convenabil alese.

Prima aproximație, notată x_1 , reprezintă abscisa punctului în care Ox taie dreapta ce trece prin punctele $(a, f(a))$ și $(b, f(b))$:

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

Alegând $y = 0$ obținem:

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

Se construiește apoi o nouă dreaptă ce trece prin $(x_1, f(x_1))$ și un al doilea punct în funcție de intervalul în care se află rădăcina ecuației: (a, x_1) sau (x_1, b) , continuându-se apoi în același mod.

Dacă $f(a)f(x_1) < 0$ atunci al doilea punct al dreptei va fi $(a, f(a))$, și obținem un șir descrescător și mărginit, deci convergent, dat de relația de recurență:

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

Dacă $f(x_1)f(b) < 0$ atunci al doilea punct al dreptei va fi $(b, f(b))$, și obținem un șir crescător și mărginit, deci convergent dat prin relația de recurență:

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

Oprirea algoritmului se face în mod analog cazului precedent, adică prin folosirea condițiilor de localizare de tip (1) sau (2), sau prin folosirea unui număr maxim de iterații admis sau a oricărei combinații între acestea.

Următorul program în limbajul C implementează algoritmul acestei metode pentru aflarea rădăcinilor ecuației $x^2 - 2 = 0$, în cazul în care oprirea algoritmului folosește o condiție de tip (1).

```

/* Program Metoda Coardei 1-condiție de tip(1) */
#include<conio.h>
#include<stdio.h>
#include <math.h>
float a,b,eps,x,A, C;
float F(float u)
{return u*u-2;}
void main ()
{
printf("introduceti intervalul in care este separata radacina\n");
printf("introduceti capatul din stanga a:");
scanf("%f",&a);
printf("\nintroduceti capatul din dreapta b:");
scanf("%f",&b);
printf("eroarea:");
scanf("%f",&eps);
x=a-F(a)*(b-a)/(F(b)-F(a));
A=F(x);
while (fabs(A)>=eps)
{
C=F(x)*F(a);
if (C<0)
x=x-F(x)*(x-a)/(F(x)-F(a));
else
{
x=x-F(x)*(b-x)/(F(b)-F(x));
a=b;
}
A=F(x);
}
if (A==0)
printf ("%f este solutia exacta a ecuatiei", x);
else
printf ("%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

Rezultatele obținute la rulare sunt:

```
C:\cc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:3
eroarea:0.00001
1.414211 este solutia aproximativa a ecuatiei
```

Varianta care folosește o condiție de oprire de tip (2) este:

```
/* Program Metoda Coardei 2 –conditie de tip (2) */
#include<conio.h>
#include<stdio.h>
#include <math.h>
float a,b,eps,x,y,A,C,ER;
float F(float u)
{return u*u-2;}
void main ()
{
printf("introduceti intervalul in care este separata radacina\n");
printf("introduceti capatul din stanga a:");
scanf("%f",&a);
printf("\nintroduceti capatul din dreapta b:");
scanf("%f",&b);
printf("eroarea:");
scanf("%f",&eps);
x=a-F(a)*(b-a)/(F(b)-F(a));
ER=b-a;
while (fabs(ER)>=eps)
{
C=F(x)*F(a);
if (C<0)
y=x-F(x)*(x-a)/(F(x)-F(a));
else
{
y=x-F(x)*(b-x)/(F(b)-F(x));
a=b;
}
```

```

}
ER=y-x;
x=y;
}
A=F(x);
if (A==0)
    printf ("%f este solutia exacta a ecuatiei", x);
else
    printf ("%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

Rezultate obținute la rulare:

```

C:\jcc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:3
eroarea:0.00001
1.414211 este solutia aproximativa a ecuatiei_

```

```

C:\jcc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:4
eroarea:0.001
1.413586 este solutia aproximativa a ecuatiei

```

V.4. METODA SECANTEI

Considerăm din nou ecuația: $f(x) = 0, x \in [a, b]$ cu f continuă și derivabilă de două ori pe $[a, b]$, astfel încât în $[a, b]$ a fost separată o soluție a ecuației, pe care o notăm cu α .

Presupunând că $f'(x) \neq 0$ pe $[a, b]$, se poate demonstra că șirul de numere reale $(x_k)_{k \geq 0}$ definit de relația de recurență:

$$x_{k+1} = \frac{x_{k-1}f(x_k) - x_k f(x_{k-1})}{f(x_k) - f(x_{k-1})},$$
 cu $x_0, x_1 \in [a, b]$ este un șir convergent către α .

Metoda secantei constă în construcția acestui șir de aproximații, ce se poate demonstra că este un șir convergent la α .

Algoritmul folosește ca date de intrare valorile inițiale ale șirului aproximațiilor, adică valorile x_0, x_1 și expresia lui f , la care se adaugă fie numărul maxim admis de iterații, fie o eroare admisibilă dată, ε_{adm} , fie ambele, în funcție de criteriul ales pentru oprirea acestuia (se poate folosi oricare din criteriile amintite la celelalte metode). Valorile inițiale se pot alege chiar capetele intervalului de separație.

Observație: valorile din șirul recursiv ce se construiește nu reprezintă nimic altceva decât punctele în care secanta la graficul funcției f , dusă prin punctele de pe graficul ei, de coordonate $(x_{k-1}, f(x_{k-1}))$ și $(x_k, f(x_k))$ intersectează axa Ox .

Pentru varianta algoritmului metodei secantei ce folosește eroarea admisibilă dată și o condiție de oprire de tip (1) s-a realizat următorul program în C, care permite calculul valorii aproximative a lui $\sqrt{5}$ cu eroarea ε_{adm} , introdusă de la tastatură. El poate fi ușor modificat pentru a permite rezolvarea oricărei ecuații neliniare analog ca în cazul celorlalte metode.

```

/* Program Metoda Secantei- conditie de tip(1) */
#include<conio.h>
#include<stdio.h>
#include<math.h>
float a,b,eps,x, A;
float F(float u)
{return u*u-5;}
void main ()
{
printf("introduceti intervalul in care este separata radacina\n");
printf("introduceti capatul din stanga a:");
scanf("%lf",&a);
printf("\nintroduceti capatul din dreapta b:");
scanf("%lf",&b);
printf("eroarea:");
scanf("%lf",&eps);
do
{
x=(a*F(b)-b*F(a))/(F(b)-F(a));
a=b;
b=x;
A=F(x);
}
while (fabs(A)>=eps);
if (fabs(A)==0)
printf ("%f este solutia exacta a ecuatiei", x);
else
printf ("%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

Considerând intervalul $[1 ; 4]$ și eroarea admisibilă dată 0.001, s-
 a obținut pentru $\sqrt{5}$ valoarea aproximativă: 2.236066.

```
C:\lcc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:4
eroarea:0.001
2.236066 este solutia aproximativa a ecuatiei
```

Pentru alte date de intrare se obțin rezultatele următoare:

```
C:\lcc\bin\rundos.exe
introduceti intervalul in care este separata radacina
introduceti capatul din stanga a:1

introduceti capatul din dreapta b:3
eroarea:0.000001
2.236068 este solutia aproximativa a ecuatiei
```

V.5. METODA APROXIMAȚIILOR SUCCESIVE PENTRU REZOLVAREA ECUAȚIILOR ȘI SISTEMELOR NELINIARE

Fie ecuația: $f(x) = 0$ (*), în care $f: R \rightarrow R$ este o funcție dată, căreia ne propunem să-i găsim rădăcina dintr-un anumit interval $[a, b]$.

Dacă rescriem ecuația sub forma echivalentă:

$$x = \varphi(x)$$

cu $\varphi: [a, b] \rightarrow [a, b]$ contracție, cu coeficientul de contracție $q < 1$, putem s-o rezolvăm cu ajutorul metodei aproximațiilor succesive.

Rădăcina din intervalul $[a, b]$ reprezintă unicul punct fix al aplicației φ . Acesta se obține ca limită a șirului aproximațiilor succesive, date de:

$x_n = \varphi(x_{n-1}), n \geq 1$, unde $x_0 \in [a, b]$ arbitrar.

Relațiile:

$$|x_n - x| \leq \frac{q^n}{1-q} |x_1 - x_0|, n \geq 0$$

$$|x_n - x| \leq \frac{q}{1-q} |x_n - x_{n-1}|, n \geq 1$$

ne indică cât de bună este valoarea aproximativă calculată la pasul n .

Cel mai des întâlnit caz în care putem aplica metoda aproximațiilor succesive este:

1. Dacă $\varphi : [a, b] \rightarrow [a, b]$ este continuă pe $[a, b]$, derivabilă pe (a, b) , și $|\varphi'(x)| \leq q < 1$ oricare ar fi $x \in (a, b)$, atunci φ este o contracție.

Exemplul V.2. Să se găsesască valoarea aproximativă a rădăcinii din intervalul $[2.2; 4]$ ecuației $x^3 - 2x^2 = 5$ folosind metoda aproximațiilor succesive și executând un număr maxim de 5 iterații.

Soluție:

Evident $f(x) = x^3 - 2x^2 - 5$ are proprietatea că

$f(2.2)f(4) < 0$, iar $f'(x) = 3x^2 - 4x = x(3x - 4 > 0)$ pe intervalul

considerat, deci f admite o singură rădăcină în $[2.2; 4]$.

Scriem ecuația astfel: $x = \frac{5}{x^2} + 2$

Observăm că dacă $x \in [2.2; 4] \Rightarrow \frac{5}{16} + 2 \leq \frac{5}{x^2} + 2 \leq \frac{5}{4.4} + 2$.

Astfel considerând $\varphi : [2.2;4] \rightarrow [2.2;4]$, $\varphi(x) = \frac{5}{x^2} + 2$,
 aceasta este o contracție, deoarece $\varphi'(x) = -\frac{10}{x^3}$, iar
 $|\varphi'(x)| \leq \frac{10}{2.2^3} = \frac{10}{10.648} < 1$.

Alegând $x_0 = 2.5$ găsim după cele cinci iterații valorile:

$$x_1 = \varphi(2.5) = 2.8$$

$$x_2 = \varphi(2.8) = 2.637755$$

$$x_3 = \varphi(2.637755) = 2.718623$$

$$x_4 = \varphi(2.718623) = 2.676507$$

$$x_5 = \varphi(2.676507) = 2.697965.$$

Un program în C bazat pe metoda aproximațiilor succesive, ce folosește o condiție de oprire de tip (1), pentru problema anterioară este:

```

/* Program Contractia 1 */
#include<conio.h>
#include<stdio.h>
#include <math.h>
float x,eps, A;
float F(float u)
{return 5/ (u*u)+2;}

void main ()
{
printf("introduceti aproximatia initiala \n");
printf("x=");
scanf("%f",&x);
printf("introduceti eroarea admisibila eps:");
scanf("%f",&eps);
A=x-F(x);
while (fabs(A)>=eps)

```

```

    {
        x=F(x);
printf("\n%f",x); /*afiseaza valorile intermediare*/
        A=x-F(x);
    }
    if (A==0)
        printf("\n%f este solutia exacta a ecuatiei", x);
    else
        printf("\n%f este solutia aproximativa a ecuatiei", x);
    getch();
}

```

Linia evidențiată din cadrul lui permite ca, după rularea sa, să obținem și valorile intermediare calculate la fiecare iterație. Rezultatele rulării codului sursă de mai sus sunt prezentate în continuare.

```

C:\oc\bin\runidos.exe
introduceti aproximatia initiala
x=:2.5
introduceti eroarea admisibila eps:0.00001

2.800000
2.637755
2.718623
2.676507
2.697965
2.686906
2.692572
2.689661
2.691154
2.690387
2.690781
2.690579
2.690683
2.690629
2.690657
2.690643
2.690643 este solutia aproximativa a ecuatiei

```

Un program bazat pe metoda aproximațiilor succesive, ce folosește o condiție de oprire de tip (1) combinată cu cea care folosește un număr maxim de iterații admis introdus de la tastatură , notat Nmax, pentru problema anterioară, este prezentat în continuare.

Variabila întregă Nit înregistrează numărul iterațiilor efectuate până la atingerea preciziei dorite. Dacă nu este atinsă precizia dorită programul afișează acest lucru.

```
/* Program Contractia 2 */
#include<conio.h>
#include<stdio.h>
#include <math.h>
float x,eps,A;
int Nmax, Nit;
float F(float u)
{return 5/ (u*u)+2;}
void main ()
{
printf("introduceti aproximatia initiala \n");
printf("x=");
scanf("%f",&x);
printf("introduceti eroarea admisibila eps:");
scanf("%f",&eps);
printf("introducenumarul maxim de iteratii admis Nmax:");
scanf("%d",&Nmax);
A=x-F(x);
Nit=0;
while ((fabs(A)>=eps) && (Nit<=Nmax))
{
x=F(x);
printf("\n%f",x); /*afiseaza valorile intermediare*/
A=x-F(x);
Nit=Nit+1;
}
if (Nit>Nmax)
printf("\nnu s-a atins precizia dorita in %d iteratii", Nmax);
else
{
if (A==0)
printf ("\n%f este solutia exacta a ecuatiei", x);
```

```

else
    printf ("\n%f este solutia aproximativa a ecuatiei, obtinuta
dupa %d iteratii", x, Nit);
}
getch();
}

```

Rularea acestui program pentru diverse date de intrare a condus la rezultatele următoare.

```

C:\vc\bin\rundos.exe
introduceti aproximatia initiala
x:- 2.5
introduceti eroarea admisibila eps:0.000001
introducenumarul maxim de iteratii admis Nmax:10

2.800000
2.637755
2.718623
2.676507
2.697965
2.686906
2.692572
2.689661
2.691154
2.690387
2.690781
nu s-a atins precizia dorita dupa 10 iteratii_

```

```

C:\vc\bin\rundos.exe
introduceti aproximatia initiala
x:- 2.5
introduceti eroarea admisibila eps:0.0001
introducenumarul maxim de iteratii admis Nmax:30

2.800000
2.637755
2.718623
2.676507
2.697965
2.686906
2.692572
2.689661
2.691154
2.690387
2.690781
2.690579
2.690683
2.690683 este solutia aproximativa a ecuatiei, obtinuta dupa 13 iteratii_

```


V.6. METODA APROXIMAȚIILOR SUCCESIVE PENTRU REZOLVAREA SISTEMELOR DE ECUAȚII NELINIARE

Vom considera sistemul de ecuații neliniare:

$$\begin{cases} F(x, y) = 0 \\ G(x, y) = 0 \end{cases}$$

transcris sub forma echivalentă:

$$\begin{cases} x = f(x, y) \\ y = g(x, y) \end{cases}$$

Fie $D = \{(x, y) \in \mathbb{R}^2 : a \leq x \leq b, c \leq y \leq d\}$. Dacă $(f(x, y), g(x, y)) \in D$

oricare ar fi $(x, y) \in D$, putem defini aplicația:

$$\varphi : D \rightarrow D, \varphi(x, y) = (f(x, y), g(x, y)), x, y \in D.$$

Sistemul dat este echivalent cu ecuația:

$$\varphi(x, y) = (x, y) \quad (**)$$

Cele mai importante situații în care φ este o contracție.

1. Să presupunem că există $q < 1$ astfel încât în raport cu norma $\|h\| = \max\{|h(x, y)| : (x, y) \in D\}$ să avem:

$$\left\| \frac{\partial f}{\partial x} \right\| + \left\| \frac{\partial g}{\partial x} \right\| \leq q; \quad \left\| \frac{\partial f}{\partial y} \right\| + \left\| \frac{\partial g}{\partial y} \right\| \leq q$$

În aceste condiții φ este o contracție, dacă ne raportăm la metrica definită prin: $\rho[(x, y), (z, t)] = |x - z| + |y - t|$, care face ca (D, ρ) să fie un spațiu metric complet.

2. Dacă există $q < 1$ astfel încât:

$$\left\| \frac{\partial f}{\partial x} \right\| + \left\| \frac{\partial f}{\partial y} \right\| \leq q; \quad \left\| \frac{\partial g}{\partial x} \right\| + \left\| \frac{\partial g}{\partial y} \right\| \leq q$$

atunci φ este o contracție dacă ne raportăm la metrica definită prin:

$$\rho[(x, y), (z, t)] = \max\{|x - z|, |y - t|\},$$

care face ca (D, ρ) să fie un spațiu metric complet.

Astfel, în oricare din aceste situații, ecuația (***) are o singură soluție în D , soluție care poate fi obținută prin metoda aproximațiilor succesive; aceasta va fi și unica soluție în D a sistemului dat;

Soluția sistemului nelinier poate fi obținută prin metoda aproximațiilor succesive pornind de la $(x_0, y_0) \in D$ și luând $(x_n, y_n) = \varphi(x_{n-1}, y_{n-1}), n \geq 1$, adică:

$$\begin{cases} x_n = f(x_{n-1}, y_{n-1}) \\ y_n = g(x_{n-1}, y_{n-1}) \end{cases}, n \geq 1$$

V.7. METODA NEWTON

Considerăm ecuația: $f(x)=0$, unde $f : R \rightarrow R$ derivabilă, cu derivata nenulă într-un interval $[a, b]$ și presupunem că am separat o rădăcină α în acest interval.

Fie $x_0 \in [a, b]$ o aproximație inițială a rădăcinii α .

Notăm cu x_1 intersecția tangentei la graficul funcției în x_0 cu axa Ox . Ecuația tangentei este:

$$y - f(x_0) = f'(x_0)(x - x_0) \Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Notăm cu x_2 intersecția tangentei la graficul funcției în x_1 cu axa Ox și continuăm analog obținând procedeul ce poartă numele de *metoda Newton Raphson* sau *metoda tangentei*.

Metoda Newton constă astfel în a găsi un șir de aproximații $(x_n)_{n \geq 1}$ pentru rădăcină, în felul următor: pentru fiecare $n \geq 1$, x_n este

punctul de intersecție dintre axa Ox și tangenta la graficul lui f în punctul de abscisă x_{n-1} .

Ecuția acestei tangente este:

$$y - f(x_{n-1}) = f'(x_{n-1})(x - x_{n-1})$$

Intersecția ei cu axa Ox este punctul de abscisă:

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

Prin urmare procesul iterativ al metodei lui Newton poate fi descris astfel:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, n \geq 0$$

Observație: Dacă f'' are semn constant pe $[a, b]$ se poate alege x_0 un punct în care are loc: $f(x_0)f''(x_0) > 0$.

Un program în C, ce folosește o condiție de oprire de tip (1), pentru metoda Newton este prezentat în continuare.

```

/* Program Newton*/
#include<conio.h>
#include<stdio.h>
#include<math.h>
float x0, x, eps, A;
float f(float u)
{return f(u);}
float Df(float u)
{return fderivat(u);}
void main ()
{
printf("introduceti aproximatia initiala a radacinii cautate
x0:");
scanf("%f",&x0);
printf("introduceti eroarea:");
scanf("%f",&eps);
x=x0-f(x0)/Df(x0);

```

```

A=f(x);
while (fabs(A)>=eps)
{
x0=x;
x=x0-f(x0)/Df(x0);
A=f(x);
}
if (A==0)
    printf ("%f este solutia exacta a ecuatiei", x);
else
    printf ("%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

Pentru rezolvarea unei probleme concrete trebuie înlocuite expresiile funcțiilor f , respectiv Df , $Df(x) = f'(x)$, în punctele corespunzătoare. Astfel pentru a găsi valoarea lui $\sqrt{2}$ putem rula programul următor:

```

/* Program Newton*/
#include<conio.h>
#include<stdio.h>
#include<math.h>
float x0,x,eps,A;
float f(float u)
{return u*u-2;}
float Df(float u)
{return 2*u;}
void main ()
{
printf("introduceti aproximatia initiala a radacinii cautate
x0:");
scanf("%f",&x0);
printf("introduceti eroarea:");
scanf("%f",&eps);

```

```

x=x0-f(x0)/Df(x0);
A=f(x);
while (fabs(A)>=eps)
{
    x0=x;
    x=x0-f(x0)/Df(x0);
    A=f(x);
}
if (A==0)
    printf ("%f este solutia exacta a ecuatiei", x);
else
    printf ("%f este solutia aproximativa a ecuatiei", x);
getch();
}

```

```

C:\lcc\bin\rundos.exe
introduceti aproximatia initiala a radacinii cautate x0:2
introduceti eroarea:0.000001
1.414214 este solutia aproximativa a ecuatiei

```

Observația V.7. Metoda poate fi aplicată și la rezolvarea sistemelor neliniare de forma:

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

În acest caz locul funcției reale f este luat de funcția vectorială de variabilă vectorială:

$$F : R^m \rightarrow R^m, \quad F(\bar{x}) = \begin{pmatrix} f_1 \\ \dots \\ f_m \end{pmatrix}, \quad \bar{x} \in R^m.$$

Dacă $f_i : R^m \rightarrow R$ au derivate parțiale de ordinul întâi continue atunci există $F'(\bar{x})$ și acesta este operator liniar de la R^m în R^m , dat de matricea (jacobianul):

$$F'(\bar{x}) = \left(\frac{\partial f_i}{\partial x_j} \right)_{1 \leq i, j \leq m}.$$

Dacă $(\exists) [F'(\bar{x})]^{-1} \Rightarrow \bar{x}^{(n+1)} = \bar{x}^{(n)} - [F'(\bar{x}^{(n)})]^{-1} F(\bar{x}^{(n)})$, aceasta fiind formula recursivă pentru calculul elementelor din șirul aproximațiilor. Indicele (n) așezat sus desemnează elementul calculat la pasul n .

Observația V.8.

Metoda Newton modificată înlocuiește pe $[F'(\bar{x}^{(n)})]^{-1}$ cu $[F'(\bar{x}^{(0)})]^{-1}$, ceea ce implică un volum de calcule mult mai mic deoarece la fiecare pas se folosește aceeași inversă – cea de la primul pas. Procesul iterativ astfel rezultat este convergent, dar converge mai lent decât cel din cazul metodei Newton.

Exemplul V.5. Să considerăm ecuația:

$$x^3 + x + 1 = 0.$$

Să se aplice metoda lui Newton pentru aflarea aproximativă a unei rădăcini din intervalul $\left[-1, -\frac{1}{2}\right]$. Se folosește o condiție de oprire de tip (1), cu $\varepsilon = 10^{-3}$

Soluție:

Notând $f(x) = x^3 + x + 1$ avem:

$$f(-1) = -1 \text{ și } f\left(\frac{-1}{2}\right) = -\frac{1}{8} - \frac{1}{2} + 1 = \frac{3}{8}.$$

De asemenea avem:

$$f'(x) = 3x^2 + 1 > 0 \text{ și } f''(x) = 6x.$$

Deoarece $f'(x) > 0 \ (\forall) x \in \left[-1, -\frac{1}{2}\right]$, f este strict crescătoare pe intervalul considerat. Cum $f(-1)f\left(-\frac{1}{2}\right) < 0$ atunci ecuația dată are o singură rădăcină în intervalul considerat.

Deoarece $(\forall) x \in \left[-1, -\frac{1}{2}\right]$ avem: $f''(x) < 0$, putem aplica metoda Newton alegând o aproximație inițială x_0 care satisface condiția: $f(x_0)f''(x_0) > 0$, deci alegem $x_0 = -1$.

Obținem succesiv:

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} = -1 + \frac{-1}{4} = -0.75.$$

Evaluăm $f(x_1) = -0.171875$.

Cum $|f(x_1)| = 0.171875 > \varepsilon$, trecem la pasul următor

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)} = -0.75 + \frac{-0.171875}{2.6875} = -0.686047.$$

Evaluăm $f(x_2) = -0.008949$.

Cum $|f(x_2)| = 0.008949 > \varepsilon$, trecem la pasul următor

$$x_3 = x_2 - \frac{f(x_2)}{f'(x_2)} = -0.686047 + \frac{-0.008949}{2.411994} = -0.682340.$$

Avem $f(x_3) = f(0.682340) = -0.000029$ și

$$|f(x_3)| = 0.000029 < \varepsilon = 0.0001.$$

Condiția (1) fiind îndeplinită algoritmul se oprește. Astfel valoarea aproximativă a rădăcinii este: $\alpha \cong -0.682340$

V.8. Metoda șirului Sturm pentru separarea rădăcinilor reale ale unei ecuații algebrice

Această metodă se aplică pentru determinarea numărului rădăcinilor reale ale unei ecuații algebrice, a semnelor acestora și chiar pentru separarea lor.

$$\text{Fie } P \in R[X], \quad P(X) = a_0 X^n + a_1 X^{n-1} + \dots + a_{n-1} X + a_n.$$

Considerăm mai întâi cazul în care polinomul P nu are rădăcini reale multiple, adică are rădăcini reale distincte, (fiecare cu ordinul de multiplicitate unu). Pentru orice astfel de polinom $P \in R[X]$ construim un șir de polinoame atașate astfel:

$$P_0(X) = P(X), \quad P_1(X) = P'(X), \quad P_k(X) = -R_k(X), \quad 2 \leq k \leq s,$$

unde $P_{k-2}(X) = P_{k-1}(X) \cdot Q_{k-1}(X) + R_k(X)$.

Observația V.9. Acest șir de polinoame este cel pe care-l obținem în algoritmul lui Euclid pentru aflarea unui cel mai mare divizor comun al polinoamelor P, P' notat (P, P') , cu observația că resturile se înmulțesc cu (-1) . Cum P nu are rădăcini multiple rezultă că P și P' sunt prime între ele, deci rezultă că $(P, P') = \text{const} \neq 0$. Această constantă este chiar P_s .

Observația V.10. Acest șir se numește șir Sturm asociat polinomului P . El nu este unicul șir Sturm asociat polinomului P .

Observația V.11. Polinoamele din șirul Sturm construit pe baza algoritmului lui Euclid pot fi înlocuite de polinoame cu coeficienții mai simpli, asociate în divizibilitate cu primele dar care să păstreze semnele coeficienților.

Definiția V.4. Pentru polinomul P și șirul Sturm asociat lui P definim funcția $S: R \rightarrow N$, care asociază oricărui număr real x numărul de variații de semn din șirul de numere

$$P_0(x), P_1(x), \dots, P_s(x).$$

Teorema V.2. Fie $a, b \in R$, $a < b$, astfel încât $P(a) \neq 0$ și $P(b) \neq 0$. Avem:

1) $S(a) \geq S(b)$;

2) Numărul rădăcinilor reale ale ecuației $P(x) = 0$, situate în (a, b) este egal cu $S(a) - S(b)$.

Observația V.12. Dacă P are rădăcini multiple atunci P și P' au rădăcini comune și astfel rezultă că $(P, P') = P_s$, cu grad $P_s \geq 1$. În acest caz șirul Sturm asociat lui P se poate alege ca fiind format din:

$$F_0 = \frac{P_0}{P_s}, F_1 = \frac{P_1}{P_s}, \dots, F_s = \frac{P_s}{P_s} = 1.$$

Consecința V.1. Dacă $P(0) \neq 0$ se poate determina numărul rădăcinilor pozitive și negative, N_+ , respectiv N_- astfel:

$N_+ = S(0) - S(\infty)$, reprezintă numărul rădăcinilor pozitive

$N_- = S(-\infty) - S(0)$, reprezintă numărul rădăcinilor negative.

Consecința V.2. Fie $\alpha, \beta \in R$, $\alpha < \beta$, astfel încât $P(\alpha) \neq 0$, $P(\beta) \neq 0$ și, $S(\alpha) - S(\beta) = 1$, atunci în intervalul (α, β) polinomul P are o singură rădăcină reală.

Exemplul V.7. Fie ecuația $x^3 - 6x + 1 = 0$. Să se determine câte rădăcini pozitive și câte negative are această ecuație și să se determine intervale de separație pentru acestea.

Soluție:

Construim un șir Sturm asociat.

$$P_0(X) = X^3 - 6X + 1, \quad P_1(X) = 3X^2 - 6.$$

Observăm că putem înlocui polinomul $P_1(X)$ cu un altul, mai simplu, asociat în divizibilitate cu el: $X^2 - 2$

$$\text{Obținem: } P_2(X) = 4X - 1, \quad P_3(X) = 31/16.$$

Calculăm $S(-\infty)$.

Șirul de semne este: $-, +, -, +$.

Observăm că apar trei variații de semn, și deci:

$$S(-\infty) = 3.$$

Șirul de semne corespunzător lui $+\infty$ este: $+, +, +, +$.

Astfel $S(\infty) = 0$.

Analog obținem: $S(0) = 2$.

Deducem de aici că: $N_- = 1, \quad N_+ = 2$

Calculând alte valori ale funcției S găsim:

$$S(-2) = 2, \quad S(1) = 1, \quad S(3) = 0$$

Astfel $S(-2) - S(0) = 3 - 2 = 1$, deci există o unică rădăcină în intervalul $(-2, 0)$.

Cum $S(0) - S(1) = 2 - 1 = 1$, deci există o unică rădăcină și în intervalul $(0, 1)$.

Analog $S(1) - S(3) = 1 - 0 = 1$, deci există o unică rădăcină în intervalul $(1, 2)$.

Astfel au fost separate cele trei rădăcini reale ale polinomului.

Pentru a putea separa rădăcinile unei ecuații în care apare o funcție derivabilă se poate utiliza și metoda șirului Rolle pe care o vom aminti în continuare.

Fie $f : I \rightarrow \mathbb{R}$, o funcție derivabilă pe I . Metoda șirului Rolle se bazează pe o proprietate importantă a funcțiilor derivabile pe care

o amintim în continuare: între două zerouri consecutive ale derivatei, există cel mult un zerou al funcției.

Observația V.13. Cu ajutorul șirului lui Rolle, determinăm numărul rădăcinilor reale ale ecuației $f(x)=0$, indicând și intervalele în care se află aceste rădăcini.

Etapele formării șirului lui Rolle.

- Se stabilește intervalul I de studiu, al ecuației $f(x)=0$, funcția $f: I \rightarrow R$, fiind presupusă derivabilă pe I .

- Se rezolvă ecuația: $f'(x)=0$ și se ordonează (crescător), rădăcinile reale, din I , ale ecuației: $x_m < \dots < x_1 < x_2 < \dots < x_M$.

- Se calculează, valorile funcției în aceste puncte, la care se adaugă limitele funcției, notate l_1, l_2 la capetele intervalului I ; obținem șirul de valori:

$$l_1, f(x_m), \dots, f(x_1), f(x_2), \dots, f(x_M), l_2.$$

- Datele obținute se trec într-un tabel/tablou, pentru x , $f(x)$; șirul lui Rolle este șirul semnelor acestor valori (poate apărea și zero) :

- Comentarii:

- Dacă $f(x_1)f(x_2) < 0 \Rightarrow$ ecuația: $f(x)=0$, are în intervalul (x_1, x_2) , o singură rădăcină reală.

- Dacă $f(x_1)f(x_2) > 0 \Rightarrow$ ecuația: $f(x)=0$, nu are în intervalul (x_1, x_2) , nici o soluție reală.

- Dacă $f(x_i)=0$, atunci x_i este rădăcină multiplă, a ecuației $f(x)=0$, și în intervalul (x_{i-1}, x_i) , (x_i, x_{i+1}) , ecuația $f(x)=0$, nu mai are soluții reale.

Observația V.14. Metoda șirului lui Rolle, este eficientă în cazul în care, există posibilitatea, rezolvării efective a ecuației $f'(x) = 0$.

Observația V.15. Numărând schimbările de semn și zerourile, se determină, numărul de soluții reale ale ecuației date, și intervalele în care se află acestea.

Exemplul V.9. Să se determine folosind șirul Rolle numărul de rădăcini reale ale ecuației $x^5 + 2x - 1 = 0$, și să se separe acestea. Să se determine, folosind metoda biseției, o valoare aproximativă pentru cea mai mare rădăcina a sa, efectuând 4 iterații.

Soluție:

Se consideră $f : R \rightarrow R$, $f(x) = x^5 + 2x - 1$

Calculăm limitele funcției la $\pm \infty$ și obținem:

$$\lim_{x \rightarrow \infty} f(x) = \infty, \quad \lim_{x \rightarrow -\infty} f(x) = -\infty.$$

Derivata acestei funcții este: $f' : R \rightarrow R$, $f'(x) = 5x^4 + 2 > 0$.

Deducem că funcția este strict crescătoare pe R .

Calculăm valorile funcției în două puncte:

$$f(0) = -1, \quad f(1) = 2$$

Realizăm următorul tablou:

x	$-\infty$		0		1		$+\infty$
$f(x)$	$-\infty$	\nearrow	-1	\nearrow	2	\nearrow	$+\infty$
$f'(x)$	+	+	+	+	+	+	+

Se observă astfel că ecuația admite o singură rădăcină reală situată în intervalul $(0,1)$.

Aplicând metoda biseecției găsim :

$$a = 0, b = 1 \Rightarrow x = \frac{1}{2} = 0.5, \text{ valoarea găsită la prima iterație}$$

$$f(0.5) = 0.031250$$

$$f(0.5) \cdot f(0) < 0 \Rightarrow b = 0.5$$

$$a = 0, b = 0.5 \Rightarrow x = \frac{0.5}{2} = 0.25, \text{ valoarea găsită la a doua iterație}$$

$$f(0.25) = -0.499023$$

$$f(0.75) \cdot f(0) > 0 \Rightarrow a = 0.25$$

$$a = 0.25, b = 0.5 \Rightarrow x = \frac{0.25 + 0.5}{2} = 0.375, \text{ valoarea găsită la a treia iterație.}$$

$$f(0.375) = -0.242584$$

$$f(0.375) \cdot f(0.25) > 0 \Rightarrow a = 0.375$$

$$a = 0.375, b = 0.5 \Rightarrow x = \frac{0.375 + 0.5}{2} = 0.4375$$

Astfel soluția aproximativă cerută este 0.4375.

VI. METODE NUMERICE PENTRU DETERMINAREA POLINOMULUI CARACTERISTIC, A VECTORILOR ȘI VALORILOR PROPRII

VI.1. POLINOM CARACTERISTIC, VECTORI ȘI VALORI PROPRII

Definiția VI.1. Fie $A \in M_n(\mathbb{R})$. Polinomul definit prin:
 $P_A(\lambda) = \det(\lambda I - A)$ se numește polinomul caracteristic al lui A .

Observația VI.1. Polinomul caracteristic este un polinom de grad n , cu coeficienți reali, și cu coeficientul dominant egal cu unu.

Definiția VI.2. Ecuația: $P_A(\lambda) = 0$ se numește ecuația caracteristică a lui A .

$$\lambda^n + c_1 \lambda^{n-1} + \dots + c_{n-1} \lambda + c_n = 0 \quad (*)$$

Polinomul caracteristic are cel mult „ n ” rădăcini distincte (reale sau complexe).

(**) $(A - \lambda I) \cdot \{x\} = 0$ (vectorul nul din \mathbb{R}^n); are și soluții nebanale.

Definiția VI.3. Rădăcinile lui P_A se numesc valori proprii (sau valori caracteristice) ale matricei A .

Definiția VI.4. Un vector nenul, $\{x\} \neq 0$, se numește vector propriu asociat valorii proprii λ dacă satisface relația:
 $(\lambda I - A)\{x\} = 0$.

Teorema VI.1. (Gerschgorin- localizarea valorilor proprii)
Fie A o matrice pătratică de ordinul n . Dacă λ este o valoare proprie, arbitrară, a lui A atunci:

$$|\lambda - a_{ii}| \leq r_i, \quad r_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad 1 \leq i \leq n.$$

VI.2. METODELE NUMERICE PENTRU CALCULUL VALORILOR ȘI VECTORILOR PROPRII

1. Metoda minorilor diagonali

Dacă $A \in M_n(\mathbb{R})$ este de forma $A = (a_{ij})$, $1 \leq i, j \leq n$, atunci $P_A(\lambda)$ se poate scrie astfel:

$$P_A(\lambda) = \lambda^n - \tau_1 \lambda^{n-1} + \tau_2 \lambda^{n-2} + \dots + (-1)^n \tau_n,$$

unde: $\tau_1 = \sum_{i=1}^n a_{ii}$ (τ_1 = suma minorilor diagonali de ordin 1)

$$\tau_2 = \sum_{1 \leq i < j \leq n} \begin{vmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{vmatrix} \quad (\tau_2 = \text{suma minorilor diagonali de ordin 2})$$

...

...

$$\tau_n = \det A$$

Exemplul VI.1. Să se calculeze, cu ajutorul minorilor diagonali, polinomul caracteristic al matricei:

$$A = \begin{pmatrix} 1 & 0 & -1 & 2 \\ 1 & -1 & 0 & 0 \\ 3 & -1 & 2 & 1 \\ 2 & -1 & 1 & 2 \end{pmatrix};$$

Soluție: Calculăm minorii diagonali și obținem:

$$\tau_1 = 1 - 1 + 2 + 2 = 4$$

$$\tau_2 = \begin{vmatrix} 1 & 0 \\ 1 & -1 \end{vmatrix} + \begin{vmatrix} 1 & -1 \\ 3 & 2 \end{vmatrix} + \begin{vmatrix} 1 & 2 \\ 2 & 2 \end{vmatrix} + \begin{vmatrix} -1 & 0 \\ -1 & 2 \end{vmatrix} + \begin{vmatrix} -1 & 0 \\ -1 & 2 \end{vmatrix} + \begin{vmatrix} 2 & 1 \\ 1 & 2 \end{vmatrix} = 1$$

$$\tau_3 = \begin{vmatrix} 1 & 0 & -1 \\ 1 & -1 & 0 \\ 3 & -1 & 2 \end{vmatrix} + \begin{vmatrix} 1 & 0 & 2 \\ 1 & -1 & 0 \\ 2 & -1 & 2 \end{vmatrix} + \begin{vmatrix} 1 & -1 & 2 \\ 3 & 2 & 1 \\ 2 & 1 & 2 \end{vmatrix} + \begin{vmatrix} -1 & 0 & 0 \\ -1 & 2 & 1 \\ -1 & 1 & 2 \end{vmatrix} = -2$$

$$\tau_4 = \det A = -6$$

$$P_A(\lambda) = \lambda^4 - 4\lambda^3 + \lambda^2 + 2\lambda - 6$$

Exemplul VI.2. Să se calculeze, cu ajutorul minorilor diagonali, polinomul caracteristic al matricei:

$$A = \begin{pmatrix} 0 & -1 & 1 & 0 \\ 3 & 1 & 0 & 0 \\ 1 & -1 & 3 & 4 \\ 3 & -1 & 1 & 2 \end{pmatrix};$$

Soluție: Calculăm minorii diagonali și obținem:

$$\tau_1 = 0 + 1 + 3 + 2 = 6$$

$$\tau_2 = \begin{vmatrix} 0 & -1 \\ 3 & 1 \end{vmatrix} + \begin{vmatrix} 0 & 1 \\ 1 & 3 \end{vmatrix} + \begin{vmatrix} 0 & 0 \\ 3 & 2 \end{vmatrix} + \begin{vmatrix} 1 & 0 \\ -1 & 3 \end{vmatrix} + \begin{vmatrix} 1 & 0 \\ -1 & 2 \end{vmatrix} + \begin{vmatrix} 3 & 4 \\ 1 & 2 \end{vmatrix} = 9$$

$$\tau_3 = \begin{vmatrix} 0 & -1 & 1 \\ 3 & 1 & 0 \\ 1 & -1 & 3 \end{vmatrix} + \begin{vmatrix} 0 & -1 & 0 \\ 3 & 1 & 0 \\ 3 & -1 & 2 \end{vmatrix} + \begin{vmatrix} 0 & 1 & 0 \\ 1 & 3 & 4 \\ 3 & 1 & 2 \end{vmatrix} + \begin{vmatrix} 1 & 0 & 0 \\ -1 & 3 & 4 \\ -1 & 1 & 2 \end{vmatrix} = 23$$

$$\tau_4 = \det A = 22$$

$$\Rightarrow P_A(\lambda) = \lambda^4 - 6\lambda^3 + 9\lambda^2 - 23\lambda + 22$$

2. Metoda Leverrier

Fie $P_A(\lambda) = \lambda^n - \tau_1 \lambda^{n-1} + \tau_2 \lambda^{n-2} + \dots + (-1)^n \tau_n$.

În cadrul acestei metode, pentru a determina coeficienții polinomului caracteristic parcurgem două etape:

- 1) Determinăm $s_k = \text{Tr}(A^k)$, $1 \leq k \leq n$
- 2) Coeficienții τ_k , $1 \leq k \leq n$, se calculează recursiv cu

relațiile:

$$\tau_1 = s_1; \quad \tau_2 = \frac{s_1 \tau_1 - s_2}{2};$$

$$\tau_k = \frac{s_1 \cdot \tau_{k-1} - s_2 \cdot \tau_{k-2} + \dots + (-1)^{k+1} \cdot s_k}{k}; \quad 2 \leq k \leq n.$$

Exemplul VI.3. Să se determine polinomul caracteristic al matricei următoare, folosind metoda Leverrier.

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & 2 \\ 1 & 0 & 0 \end{pmatrix};$$

- 1) Determinăm mai întâi elementele $s_k = \text{Tr}(A^k)$.

$$s_1 = \text{Tr}(A) = 1 - 1 + 0 = 0$$

Calculăm puterile matricei A și apoi urmele acestora. Avem:

$$A^2 = \begin{pmatrix} 2 & 0 & 3 \\ 2 & 1 & -2 \\ 1 & 1 & 1 \end{pmatrix} \quad \Rightarrow s_2 = \text{Tr}(A^2) = 2 + 1 + 1 = 4$$

$$A^3 = \begin{pmatrix} 5 & 2 & 2 \\ 0 & 1 & 4 \\ 2 & 0 & 3 \end{pmatrix} \quad \Rightarrow s_3 = \text{Tr}(A^3) = 5 + 1 + 3 = 9$$

2) Determinăm coeficienții τ_k .

$$\tau_1 = s_1 = 0; \tau_2 = \frac{(s_1 \cdot \tau_1 - s_2)}{2} = \frac{(0 \cdot 0 - 4)}{2} = -2$$

$$\tau_3 = \frac{(s_1 \tau_2 - s_2 \tau_1 + s_3)}{3} = \frac{(0 \cdot (-2) - 4 \cdot 0 + 9)}{3} = 3$$

$$\Rightarrow P_A(\lambda) = \lambda^3 - 2\lambda - 3.$$

Exemplul VI.4. Să se determine polinomul caracteristic al matricei următoare folosind metoda Leverrier.

$$A = \begin{pmatrix} 1 & 1 & 1 \\ -1 & 2 & 1 \\ 1 & 0 & 1 \end{pmatrix};$$

1) Determinăm mai întâi elementele $s_k = \text{Tr}(A^k)$.

$$s_1 = \text{Tr}(A) = 2 + 1 + 1 = 4$$

Calculăm puterile matricei A și apoi urmele acestora. Avem:

$$A^2 = \begin{pmatrix} 1 & 3 & 3 \\ -2 & 3 & 2 \\ 2 & 1 & 2 \end{pmatrix} \quad \Rightarrow s_2 = \text{Tr}(A^2) = 3 + 2 + 1 = 6$$

$$A^3 = \begin{pmatrix} 1 & 7 & 7 \\ -3 & 4 & 3 \\ 3 & 4 & 5 \end{pmatrix} \quad \Rightarrow s_3 = \text{Tr}(A^3) = 1 + 4 + 5 = 10$$

2) Determinăm coeficienții τ_k .

$$\tau_1 = s_1 = 4; \tau_2 = \frac{(s_1 \cdot \tau_1 - s_2)}{2} = \frac{(4 \cdot 4 - 6)}{2} = 5$$

$$\tau_3 = \frac{(s_1 \tau_2 - s_2 \tau_1 + s_3)}{3} = \frac{(4 \cdot 5 - 6 \cdot 4 + 10)}{3} = 2$$

$$\Rightarrow P_A(\lambda) = \lambda^3 - 4\lambda^2 + 5\lambda - 2. P_A(\lambda) = (\lambda - 1)^2(\lambda - 2)$$

Rădăcinile sale sunt: $\lambda_1 = \lambda_2 = 1$, $\lambda_3 = 2$.

Astfel au fost găsite două valori proprii distincte ale matricei considerate, una având ordinul de multiplicitate doi iar cealaltă unu.

3. Metoda Krylov

Fie $A \in M_n(R)$ și $P_A(\lambda) = \lambda^n + c_1 \lambda^{n-1} + \dots + c_{n-1} \lambda + c_n$

polinomul ei caracteristic. Folosind teorema Cayley-Hamilton care afirmă că: matricea A verifică ecuația caracteristică (forma matriceală), putem scrie: $P_A(A) = 0_n$.

$$\text{Deci: } A^n + c_1 \cdot A^{n-1} + \dots + c_{n-1} \cdot A + c_n \cdot I_n = 0_n.$$

Fie $y^{(0)} \in R^n$, oarecare, nenul. Prin înmulțirea relației precedente la dreapta cu $y^{(0)}$ obținem relația:

$$A^n \cdot y^{(0)} + c_1 A^{n-1} \cdot y^{(0)} + \dots + c_x \cdot A \cdot y^{(0)} + c_n \cdot y^{(0)} = \bar{0}.$$

Relația de mai sus reprezintă un sistem de n ecuații liniare cu n necunoscute c_1, c_2, \dots, c_n .

Dacă alegerea lui $y^{(0)} \in R^n$ conduce la obținerea unui determinant nenul soluția unică a acestuia reprezintă coeficienții polinomului caracteristic. Dacă determinantul este nul, se reiau calculele cu un alt vector inițial $y^{(0)}$.

Pentru simplitate se notează:

$$A^k \cdot y^{(0)} = y^{(k)} = A(A^{k-1} y^{(0)}) = A y^{(k-1)}, \quad k = \overline{1, n} \Rightarrow$$

sistemul se poate scrie astfel:

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

Determinantul său este $|y^{(n-1)} y^{(n-2)} \dots y^{(0)}|$ și trebuie să fie nenul pentru a continua calculele.

Exemplul VI.3. Folosind metoda Krylov să se determine $P_A(\lambda)$, unde:

$$A = \begin{pmatrix} 1 & 1 & 1 \\ -1 & 2 & 1 \\ 1 & 0 & 1 \end{pmatrix}.$$

$$\text{Alegem } y^{(0)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}.$$

Găsim:

$$y^{(1)} = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \quad y^{(2)} = \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix}, \quad y^{(3)} = \begin{pmatrix} 1 \\ -1 \\ 9 \end{pmatrix}.$$

Sistemul pe care îl obținem în acest caz are determinantul nul, deci alegem alt vector inițial $y^{(0)}$.

$$\text{Fie } y^{(0)} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}. \text{ Găsim în acest caz:}$$

$$y^{(1)} = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \quad y^{(2)} = \begin{pmatrix} 3 \\ 3 \\ 1 \end{pmatrix}, \quad y^{(3)} = \begin{pmatrix} 7 \\ 4 \\ 4 \end{pmatrix}.$$

$$\Rightarrow \begin{cases} 3 \cdot c_1 + c_2 = -7 \\ 3 \cdot c_1 + 2c_2 + c_3 = -4 \\ c_1 = -4 \end{cases} \Rightarrow$$

Soluția unică a acestui sistem este:

$$c_1 = -4, c_2 = 5, c_3 = -2$$

Astfel polinomul caracteristic ce se obține este:

$$P_A(\lambda) = \lambda^3 - 4 \cdot \lambda^2 + 5 \cdot \lambda - 2.$$

$$P_A(\lambda) = (\lambda - 1)^2(\lambda - 2)$$

Rădăcinile sale sunt: $\lambda_1 = \lambda_2 = 1, \lambda_3 = 2.$

Observația VI.7. Dacă polinomul caracteristic are rădăcini distincte atunci, cu ajutorul metodei Krylov se pot determina vectorii proprii.

Notând

$$q_j(\lambda) = \frac{P_A(\lambda)}{\lambda - \lambda_j} = \lambda^{n-1} + q_{1j}\lambda^{n-2} + \dots + q_{n-1j}, j = \overline{1, n},$$

atunci un vector propriu corespunzător valorii caracteristice λ_j este dat de relația:

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

Exemplu VI.7. Să se calculeze, folosind metoda Krylov, polinomul caracteristic și vectorii proprii ai matricei A , unde:

$$A = \begin{pmatrix} 2 & 1 & 1 \\ -1 & 2 & -1 \\ 1 & -1 & 2 \end{pmatrix}$$

Soluție:

$$\text{Fie } y^{(0)} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \Rightarrow y^{(1)} = \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}, y^{(2)} = \begin{pmatrix} 3 \\ -5 \\ 6 \end{pmatrix}, y^{(3)} = \begin{pmatrix} 7 \\ -19 \\ 20 \end{pmatrix}.$$

$$\text{Cum } \begin{vmatrix} 3 & 1 & 0 \\ -5 & -1 & 0 \\ 6 & 2 & 1 \end{vmatrix} = -3 + 5 = 2 \neq 0, \text{ sistemul ce trebuie rezolvat}$$

este:

$$c_1 \begin{pmatrix} 3 \\ -5 \\ 6 \end{pmatrix} + c_2 \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix} + c_3 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = - \begin{pmatrix} 7 \\ 19 \\ 20 \end{pmatrix}$$

Soluția acestuia este: $c_1 = -6, c_2 = 11, c_3 = -6$

Ecuția caracteristică a matricei A este:

$$\lambda^3 - 6 \cdot \lambda^2 + 11 \cdot \lambda - 6 = 0$$

Rădăcinile ei, adică valorile proprii, sunt: $\lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$.

Fiind distincte două câte două vom determina vectorii proprii.

Pentru $\lambda_1 = 1$ avem:

$$q_1(\lambda) = \frac{P_A(\lambda)}{\lambda - \lambda_1} = \lambda^2 - 5 \cdot \lambda + 6$$

$$y^{(2)} - 5 \cdot y^{(1)} + 3 \cdot y^{(0)} = \begin{pmatrix} 3 \\ -5 \\ 6 \end{pmatrix} - 5 \cdot \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix} + 6 \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \\ 2 \end{pmatrix} = 2 \cdot \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$$

este un vector propriu corespunzător valorii proprii $\lambda_1 = 1$.

$$\text{Pentru } \lambda_2 = 2 \Rightarrow q_2(\lambda) = \frac{P_A(\lambda)}{\lambda - \lambda_2} = \lambda^2 - 4 \cdot \lambda + 3.$$

Astfel un vector propriu corespunzător ei este:

$$y^{(2)} - 4 \cdot y^{(1)} + 3 \cdot y^{(0)} = \begin{pmatrix} 3 \\ -5 \\ 6 \end{pmatrix} - 4 \cdot \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix} + 3 \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix}$$

$$\text{Pentru } \lambda_3 = 3 \Rightarrow q_3(\lambda) = \frac{P_A(\lambda)}{\lambda - \lambda_3} = \lambda^2 - 3 \cdot \lambda + 2, \text{ deci un}$$

vector propriu corespunzător ei este:

$$y^{(2)} - 3 \cdot y^{(1)} + 2 \cdot y^{(0)} = \begin{pmatrix} 0 \\ -2 \\ 2 \end{pmatrix} = 2 \cdot \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}.$$

4. Metoda Fadeev

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

Coefficienții polinomului caracteristic, cât și inversa A^{-1} se pot obține din relațiile ce caracterizează algoritmul acestei metode:

$$A_1 = A, \quad c_1 = -\text{Tr}(A_1), \quad B_1 = A_1 + c_1 I_n$$

$$A_2 = AB_1, \quad c_2 = -\frac{1}{2} \text{Tr}(A_2), \quad B_2 = A_2 + c_2 I_n$$

.....

$$A_k = AB_{k-1}, \quad c_k = -\frac{1}{k} \text{Tr}(A_k), \quad B_k = A_k + c_k I_n$$

.....

$$\Rightarrow \begin{cases} B_n = O_n \\ A^{-1} = -\frac{1}{c_n} \cdot B_{n-1}, \text{ dacă } c_n \neq 0 \end{cases}$$

Exemplul VI.5. Folosind metoda Fadeev să se determine polinomul caracteristic, $P_A(\lambda)$, pentru matricea:

$$A = \begin{pmatrix} -1 & 2 & 0 \\ 1 & -1 & 2 \\ 3 & 1 & -3 \end{pmatrix},$$

și dacă este inversibilă și inversa ei A^{-1} .

Soluție:

Urmăm pașii din algoritm și avem:

$$A_1 = A = \begin{pmatrix} -1 & 2 & 0 \\ 1 & -1 & 2 \\ 3 & 1 & -3 \end{pmatrix} \Rightarrow c_1 = 5, \quad B_1 = \begin{pmatrix} 4 & 2 & 0 \\ 1 & 4 & 2 \\ 3 & 1 & 2 \end{pmatrix}$$

$$A_2 = A \cdot B_1 = \begin{pmatrix} -2 & 6 & 4 \\ 9 & 0 & 2 \\ 4 & 13 & -4 \end{pmatrix} \Rightarrow c_2 = 3, \quad B_2 = \begin{pmatrix} 1 & 6 & 4 \\ 9 & 3 & 2 \\ 4 & 7 & -1 \end{pmatrix}$$

$$A_3 = A \cdot B_2 = \begin{pmatrix} 17 & 0 & 0 \\ 0 & 17 & 0 \\ 0 & 0 & 17 \end{pmatrix} \Rightarrow c_3 = -17, \quad B_3 = O_3.$$

Astfel obținem următoarea expresie pentru polinomul caracteristic

$$P_A(\lambda) = \lambda^3 + 5\lambda^2 + 3\lambda - 17.$$

Deoarece $c_3 = -17 \neq 0$, se poate deduce și inversa lui A .

Obținem:

$$A^{-1} = \frac{1}{17} \cdot B_2 = \frac{1}{17} \begin{pmatrix} 1 & 6 & 4 \\ 9 & 3 & 2 \\ 4 & 7 & -1 \end{pmatrix}.$$

VII APROXIMAREA FUNCȚIILOR

VII.1. APROXIMAREA PRIN INTERPOLARE

Fie o funcție $f : [a, b] \rightarrow R$, cunoscută numai într-un număr limitat de puncte (sau noduri): x_0, x_1, \dots, x_n , ce poartă numele de *suportul interpolării*, prin valorile: $f(x_0), f(x_1), \dots, f(x_n)$, notate în general: f_0, f_1, \dots, f_n .

De obicei aceste date sunt prezentate într-un *tabel de interpolare*:

x	x_0	x_1	..	x_n
$f(x)$	$f(x_0)$	$f(x_1)$..	$f(x_n)$

Dacă punctele x_0, x_1, \dots, x_n sunt echidistante, se spune că dispunem de o *discretizare uniformă*, iar dacă nu sunt echidistante *discretizarea este neuniformă*.

Dacă presupunem că cele $n+1$ noduri reprezintă abscisele unor puncte din plan, iar f_0, f_1, \dots, f_n reprezintă ordonatele lor se pune problema determinării unei curbe, imaginea geometrică a unei funcții reale, care trece prin acestea.

Vom aproxima comportarea funcției printr-un *polinom generalizat de interpolare*, de forma:

$$P_n(x) = a_0 u_0(x) + a_1 u_1(x) + \dots + a_n u_n(x)$$

în care funcțiile *liniar independente*: $u_0(x), u_1(x), \dots, u_n(x)$ sunt cunoscute și constituie *baza interpolării*. Aceasta poate fi formată din funcții simple: polinoame, funcții trigonometrice, exponențiale, etc.

Determinarea polinomului generalizat de interpolare presupune practic determinarea coeficienților a_0, a_1, \dots, a_n . Se impune practic ca, pe suportul interpolării, polinomul de interpolare să coincidă cu funcția f . Astfel se obțin relațiile:

$$P_n(x_i) = f(x_i), i = \overline{0, n} \quad (*)$$

Relațiile (*) sunt cunoscute sub numele de *condiții de interpolare*.

Se poate considera o abordare mai generală a aproximării unei funcții, impunând condiții de interpolare de forma:

$$P_n(x_i) = f(x_i), P_n'(x_i) = f'(x_i), \dots, P_n^{(k)}(x_i) = f^{(k)}(x_i).$$

Condițiile de interpolare (*) conduc la sistemul de ecuații liniare:

$$\sum_{k=0}^n a_k u_k(x_i) = f(x_i), i = \overline{0, n} \quad .$$

1.1. INTERPOLARE POLINOMIALĂ. POLINOMUL DE INTERPOLARE LAGRANGE

Dacă se alege baza polinomială:

$$u_0(x)=1 \quad u_1(x)=x \quad \dots \quad u_n(x)=x^n$$

sistemul determinat de condițiile de interpolare devine:

$$\sum_{k=0}^n a_k \cdot x_i^k = f(x_i) \quad .$$

El reprezintă un sistem linear de $n+1$ ecuații cu $n+1$ necunoscute caracterizat de un determinant Vandermonde nenul, dacă punctele x_i sunt distincte, caz în care sistemul este compatibil determinat, având soluție unică, ceea ce implică un polinom de interpolare unic.

Teorema 1. Date fiind $n+1$ puncte distincte x_0, x_1, \dots, x_n și $n+1$ valori reale arbitrare f_0, f_1, \dots, f_n , există și este unic un polinom de interpolare corespunzător tabelului de interpolare construit pe baza lor.

Polinomul general de interpolare, de grad cel mult n , se poate obține sub forma:

$$P(x) = \sum_{k=0}^n f(x_k) \cdot \prod_{\substack{i=0 \\ i \neq k}}^n \frac{x - x_i}{x_k - x_i} = \sum_{k=0}^n f_k l_k(x)$$

(**)

și poartă numele de polinom de interpolare Lagrange asociat tabelului de interpoale dat, notat $L_n(x)$, iar polinoamele

$$l_k(x) = \prod_{i=0, i \neq k}^n \frac{x - x_i}{x_k - x_i}$$

ce apar în formulă se numesc polinoame

fundamentale de interpolare Lagrange.

Ele au proprietatea:

$$l_k(x_i) = \begin{cases} 1, & i = k \\ 0, & i \neq k \end{cases}$$

Exemplul 1. Să se determine polinomul de interpolare Lagrange pentru o funcție f știind că suportul interpolării este: $x_0 = -1, x_1 = 0, x_2 = 1, x_3 = 3$, iar valorile funcției în aceste puncte sunt: 7, 1, 2, 1. Să se determine apoi valoarea acestuia în punctul 2.

Soluție:

Polinomul de interpolare Lagrange este dat de:

$$P(x) = \sum_{k=0}^3 f(x_k) \cdot l_k(x).$$

$$\text{Notăm } \omega(x) = (x - x_0)(x - x_1)(x - x_2)(x - x_3).$$

$$\text{Avem: } \omega(x) = (x + 1)x(x - 1)(x - 3).$$

Notăm: $\omega_k(x) = \frac{\omega(x)}{x - x_k}$, $k = \overline{0,3}$. Obținem următoarele

relații:

$$\omega_0(x) = x(x-1)(x-3), \quad \omega_1(x) = (x+1)(x-1)(x-3),$$

$$\omega_2(x) = (x+1)x(x-3), \quad \omega_3(x) = (x+1)x(x-1).$$

Cu ajutorul lor construim polinoamele fundamentale

Lagrange date de relația: $l_k(x) = \frac{\omega_k(x)}{\omega_k(x_k)}$, $k = \overline{0,3}$.

Obținem deci:

$$l_0(x) = \frac{x(x-1)(x-3)}{-8},$$

$$l_1(x) = \frac{(x+1)(x-1)(x-3)}{3},$$

$$l_2(x) = \frac{(x+1)x(x-3)}{-4},$$

$$l_3(x) = \frac{(x+1)x(x-1)}{24}.$$

Deducem deci următoarea formă pentru polinomul de interpolare Lagrange:

$$L_3(x) = 7 \frac{x(x-1)(x-3)}{-8} + \frac{(x+1)(x-1)(x-3)}{3} + 2 \frac{(x+1)x(x-3)}{-4} + \frac{(x+1)x(x-1)}{24}$$

Se obține apoi valoarea polinomului în punctul 2: $L_3(2) = 4$.

Exemplul 2. Considerând variația densității unui material în intervalul de temperatura $t = 60^{\circ}C \dots 120^{\circ}C$ din tabelul următor să se determine densitatea acestuia la $70^{\circ}C$, folosindu-se un polinom Lagrange de gradul 2.

$t [^{\circ}C]$	60	80	100	120
$\rho [Kg/m^3]$	435,60	413,40	402,65	398.50

Soluție:

Folosind toate datele din tabelul de interpolare dat s-ar obține un polinom Lagrange de grad 3. Polinomul cerut trebuie să aibă gradul doi, deci selectăm un suport de interpolare mai restrâns: ca suport de interpolare punctele $t_1 = 60$, $t_2 = 80$, $t_3 = 100$ și valorile corespunzătoare ale funcției densitate.

Se obțin polinoamele fundamentale Lagrange:

$$l_1(t) = \frac{(t-t_2)(t-t_3)}{(t_1-t_2)(t_1-t_3)} = \frac{(t-80)(t-100)}{(60-80)(60-100)} = \frac{t^2 - 180t + 8000}{800}$$

$$l_2(t) = \frac{(t-t_1)(t-t_3)}{(t_2-t_1)(t_2-t_3)} = \frac{(t-60)(t-100)}{(80-60)(80-100)} = \frac{t^2 - 160t + 6000}{-400}$$

$$l_3(t) = \frac{(t-t_1)(t-t_2)}{(t_3-t_1)(t_3-t_2)} = \frac{(t-60)(t-80)}{(100-60)(100-80)} = \frac{t^2 - 140t + 4800}{800}$$

Astfel polinomul de interpolare Lagrange este :

$$L(t) = 435,60 \cdot l_1(t) + 413,40 \cdot l_2(t) + 402,65 \cdot l_3(t).$$

Pentru $t = 70^\circ C$ se obține următoarea
valoare : $L(70) = 423,06875$.

1.2. DIFERENȚE DIVIZATE

O altă metodă de a obține polinomul de interpolare corespunzător unui tabel de interpolare dat este cea bazată pe așa-numitele diferențe divizate, care reprezintă în fapt niște notații asociate unor rapoarte construite după anumite reguli cu datele din tabelul de interpolare.

Diferențele divizate pot fi introduse prin recurență, cu ajutorul definițiilor următoare:

- Diferențele divizate de ordin zero, notate $F_0[]$ coincid cu valorile funcției în nodurile date. Există astfel $n+1$ diferențe divizate de ordinul zero pentru datele din tabloul de interpolare considerat;

$$F_0[x_i] = f(x_i)$$

- Diferențele divizate de ordinul unu, notate cu $F_1[]$, se definesc prin egalitatea:

$$F_1[x_i, x_{i+1}] = \frac{F_0[x_i] - F_0[x_{i+1}]}{x_i - x_{i+1}}$$

și sunt în număr de n : $F_1[x_0, x_1], F_1[x_1, x_2], \dots, F_1[x_{n-1}, x_n]$;

.....

• Diferențele divizate de ordinal p se definesc cu ajutorul diferențelor divizate de ordin $p-1$ prin relația:

$$F_p[x_0, x_1, \dots, x_p] = \frac{F_{p-1}[x_0, \dots, x_{p-1}] - F_{p-1}[x_1, \dots, x_p]}{x_0 - x_p};$$

.....

• Există o singură diferență divizată de ordinul n dată de:

$$F_n[x_0, x_1, \dots, x_n] = \frac{F_{n-1}[x_0, \dots, x_{n-1}] - F_{n-1}[x_1, \dots, x_n]}{x_0 - x_n}$$

Pentru a simplifica evaluarea diferențelor divizate se poate utiliza așezarea acestora într-un tabel de forma următoare:

x_i	F_0	F_1	F_2	...	F_n
x_0	$F_0[x_0] = f_0$	$F_1[x_0, x_1]$	$F_2[x_0, x_1, x_2]$		$F_n[x_0, \dots, x_n]$
x_1	$F_0[x_1] = f_1$	$F_1[x_1, x_2]$	$F_2[x_1, x_2, x_3]$		-
x_2	$F_0[x_2] = f_2$	$F_1[x_2, x_3]$	$F_2[x_2, x_3, x_4]$		-
...		-
x_{n-2}	$F_0[x_{n-2}] = f_{n-2}$	$F_1[x_{n-2}, x_{n-1}]$	$F_2[x_{n-2}, x_{n-1}, x_n]$		-

x_{n-1}	$F_0[x_{n-1}] = f_{n-1}$	$F_1[x_{n-1}, x_1]$	-		-
x_n	$F_0[x_n] = f_n$	-	-		-

Notând cu d_{ij} elementele acestui tablou observăm că sunt valabile relațiile:

$$d_{i0} = f_i, i = \overline{0, n},$$

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

O altă formă a polinomului de interpolare de grad cel mult n , se poate obține cu ajutorul ediferențelor divizate și poartă numele de polinom Newton de interpolare, notat cu $N_n(x)$:

$$N_n(x) = f(x_0) + (x-x_0)F_1[x_0, x_1] + (x-x_0)(x-x_1)F_2[x_0, x_1, x_2] \dots + (x-x_0) \dots (x-x_{n-1})F_n[x_0, \dots, x_n]$$

Exemplul 3 Folosind diferențele divizate să se determine polinomul Newton de interpolare ce interpoalează datele din tabelul următor:

x	-1	1	2	3
f	4	2	1	-2

Soluție:

Construim tabelule diferențelor divizate și apoi cu ajutorul elementelor de pe prima linie a tabelului polinomul Newton de interpolare.

x	F_0	F_1	F_2	F_3
-1	4	-1	0	-1/4
1	2	-1	-1	
2	1	-3		
3	-2			

Utilizând relația (14) obținem:

$$N_3(x) = F_0(x_0) + (x - x_0)F_1[x_0, x_1] + (x - x_0)(x - x_1)F_2[x_0, x_1, x_2] + (x - x_0)(x - x_1)(x - x_2)F_3[x_0, \dots, x_3]$$

Astfel polinomul Newton de interpolare este:

$$N_3(x) = 4 + (x + 1)(-1) + (x + 1)(x - 1)(x - 2)\left(-\frac{1}{4}\right).$$

Exemplul 4.

Să se determine polinomul Newton de interpolare de grad cel mult 3, ce interpoolează funcția $f: R \rightarrow R$, $f(x) = 2x^2 + 3x - 1$ în nodurile $x = 0, 1, 2, 3$. Care este gradul polinomului găsit?

x	0	1	2	3
f	-1	4	13	26

Soluție:

Construim tabelule diferențelor divizate și apoi polinomul Newton de interpolare.

x	F_0	F_1	F_2	F_3
0	-1	5	2	0
1	4	9	2	
2	13	13		
3	26			

$$N_3(x) = F_0(x_0) + (x - x_0)F_1[x_0, x_1] + (x - x_0)(x - x_1)F_2[x_0, x_1, x_2] + (x - x_0)(x - x_1)(x - x_2)F_3[x_0, x_1, x_2, x_3]$$

Astfel obținem:

$$N_3(x) = -1 + 5x + 2x(x - 1).$$

Remarcăm că el are gradul doi și coincide cu funcția dată:

$$N_3(x) = f(x)$$

În ceea ce privește eroarea pe care o introducem aproximând o funcție prin polinomul de interpolare de grad cel mult n , se poate demonstra că în cazul unei funcții de clasă $C^{n+1}([a, b])$ aceasta este dată de:

$$E_n(x) = f(x) - P_n(x) = f^{(n+1)}(z) \frac{\omega(x)}{(n+1)!}, \text{ cu } z \in (a, b).$$

Dacă în plus nodurile din intervalul $[a, b]$ sunt echidistante, atunci

$$|E_n(x)| = |f(x) - P_n(x)| \leq \frac{1}{4(n+1)} M \left(\frac{b-a}{n} \right)^{n+1}, \quad \text{unde}$$

$$M = \sup_{[a,b]} |f^{(n+1)}(x)|.$$

Exemplul 5.

Să se determine cât de mare este eroarea pe care o introducem când aproximăm funcția $f(x) = \sin x$ pe intervalul $[0, 2]$ cu un polinom de interpolare corespunzător unui tabel de interpolare cu 21 noduri echidistante.

$$\text{Avem } n = 20, a = 0, b = 2, |f^{(21)}(x)| = |\cos x| \leq 1 = M.$$

Astfel obținem o eroare foarte mică:

$$|E_{20}(x)| = |f(x) - P_{20}(x)| \leq \frac{1}{4 \cdot 21} \left(\frac{2-0}{20} \right)^{21} \cong 0,012 \cdot 10^{-21}$$

1.3. DIFERENȚE FINITE

Diferențele finite sunt utile în evaluarea numerică a derivatelor și în formulele de interpolare. Ele sunt de mai multe feluri și se aplică unor funcții definite în puncte echidistante:

$$x_i = x_0 + ih :$$

-diferențe progresive:

$$\Delta f(x_i) = \frac{f(x_i + h) - f(x_i)}{h} = \frac{f(x_{i+1}) - f(x_i)}{h}$$

$$\Delta^k f(x_i) = \frac{\Delta^{k-1} f(x_{i+1}) - \Delta^{k-1} f(x_i)}{h}$$

-diferențe regresive:

$$\begin{aligned} \nabla f(x_i) &= \frac{f(x_i) - f(x_i - h)}{h} = \frac{f(x_i) - f(x_{i-1})}{h} \\ \nabla^k f(x_i) &= \frac{\nabla^{k-1} f(x_i) - \nabla^{k-1} f(x_{i-1})}{h} \end{aligned} \quad (15)$$

-diferențe centrate:

$$\begin{aligned} \delta f(x_i) &= \frac{f(x_{i+\frac{1}{2}}) - f(x_{i-\frac{1}{2}})}{h} \\ \delta^k f(x_i) &= \frac{\delta^{k-1} f(x_{i+\frac{1}{2}}) - \delta^{k-1} f(x_{i-\frac{1}{2}})}{h} \end{aligned}$$

Pentru a trece de la diferențele finite la diferențele divizate se folosește relația:

$$\Delta^n f(x_i) = \nabla^n f(x_{i+n}) = \delta^n f(x_{i+n/2}) = n! h^n F^n[x_i, x_{i+1}, \dots, x_{i+n}].$$

Astfel, cu ajutorul lor pot fi construite polinoamele de interpolare, în cazul funcțiilor definite pe un suport de interpolare cu noduri echidistante.

Fie deci funcția f cunoscută printr-un tabel de interpolare în care abscisele x_i sunt echidistante.

Dorim să evaluăm funcția în punctul intermediar x , ($x \neq x_i$). Vom considera în mod simplificator că acest punct poate fi situat:

-la începutul tabloului $x_0 < x < x_1$

-la sfârșitul tabloului $x_{n-1} < x < x_n$.

Dacă $x_i < x < x_{i+1}$ cu $i \neq 0$ și $i \neq n-1$ vom neglija un număr de puncte de la începutul sau sfârșitul tabloului, pentru a ne situa într-unul din cazurile de mai sus.

Formulele de interpolare construite cu ajutorul diferențelor finite poartă numele de formulele Newton – Gregory.

Prima formulă Newton-Gregory realizează interpolare la început de tablou, adică consideră: $x = x_0 + uh$ cu $0 < u < 1$.

Pornind de la expresia polinomului Newton de interpolare în punctul x , și exprimând diferențele divizate cu ajutorul diferențelor finite progresive: $F_k[x_0, x_1, \dots, x_k] =$

$\frac{\Delta^k f(x_0)}{k! \cdot h^k}$ se obține:

$$P_1(x) = P_1(x + u \cdot h) = p_1(u) = f(x_0) + \frac{x - x_0}{1! \cdot h} \cdot \Delta f(x_0) + \frac{(x - x_0) \cdot (x - x_1)}{2! \cdot h^2} \cdot \Delta^2 f(x_0) + \dots + \frac{(x - x_0) \cdot (x - x_1) \dots (x - x_{n-1})}{n! \cdot h^n} \cdot \Delta^n f(x_0)$$

Ținând cont de faptul că:

$x - x_k = x - x_0 - (x_k - x_0) = uh - kh = (u - k)h$, obținem:

$$p_1(u) = f_0 + \frac{u}{1!} \cdot \Delta f_0 + \frac{u \cdot (u - 1)}{2!} \cdot \Delta^2 f_0 + \dots + \frac{u \cdot (u - 1) \dots (u - n + 1)}{n!} \cdot \Delta^n f_0$$

Următoarea formulă de interpolare Newton-Gregory se referă la interpolarea la sfârșit de tablou și se obțin exprimând diferențele divizate prin diferențe regresive. Astfel dacă se ia:

$$x = x_n - u \cdot h, \quad 0 < u < 1, \text{ și se parcurg aceleași etape, se obține:}$$

$$p_2(u) = f_0 - \frac{u}{1!} \cdot \nabla f_n + \frac{u(u-1)}{2!} \cdot \nabla^2 f_n + \dots + (-1)^n \frac{u(u-1)\dots(u-n+1)}{n!} \cdot \nabla^n f_n$$

.

1.3. INTERPOLARE CU FUNCȚII SPLINE

Polinomul de interpolare Lagrange prezintă un grad înalt pentru un număr mare de puncte. Aceasta ne determină uneori să alegem polinoame de interpolare de grad mic, valabile pe subintervale ale intervalului $[a, b]$. În astfel de situații folosim funcțiile spline pentru aproximare.

Definiția 1.

O funcție S se numește funcție spline de grad k dacă ea îndeplinește condițiile:

Este definită pe $[a, b]$,

Derivatele ei de ordin r sunt continue pe $[a, b]$, pentru $0 \leq r \leq k - 1$

Restricțiile ei la subintervalele $[x_i, x_{i+1}]$, $i = \overline{1, n-1}$ sunt polinoame de grad cel mult k .

Funcții spline liniare

Fie $n + 1$ puncte: $x_0 < x_1 < \dots < x_n$ în care se cunosc valorile funcției: $f(x_0), f(x_1), \dots, f(x_n)$, vom considera în acest caz funcții de interpolare liniare locale, pe subintervalele:

$$[x_0, x_1], [x_1, x_2], \dots, [x_i, x_{i+1}], \dots, [x_{n-1}, x_n]$$

de forma:

$$s_i(x) = a_i x + b_i, \quad i = 0 \dots n - 1,$$

în care cei $2n$ parametri a_i și b_i se determină impunând satisfacerea condițiilor de interpolare:

$$s_i(x_i) = f(x_i) \quad i = 0 \dots n - 1 \qquad s_{n-1}(x_n) = f(x_n)$$

și a condițiilor de racordare în punctele interioare:

$$s_i(x_{i+1}) = s_{i+1}(x_{i+1}) \quad i = 0 \dots n - 2$$

Interpolarea liniară ne oferă o funcție continuă și derivabilă pe porțiuni, dar care prezintă totuși dezavantajul discontinuității derivatelor în punctele interioare.

Condițiile de interpolare ne dau coeficienții:

$$a_i = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}, \quad i = 0 \dots n - 1$$

$$b_i = \frac{x_{i+1} f(x_i) - x_i f(x_{i+1})}{x_{i+1} - x_i}.$$

Observăm că pe fiecare din subintervalele formate se construiesc practic segmente de drepte ce unesc punctele de

coordonate $(x_i, f(x_i))$ și $(x_{i+1}, f(x_{i+1}))$ și astfel se realizează o aproximare a curbei reale printr-o linie poligonală.

Funcția spline se poate prelungi și în afara intervalului $[a, b]$ astfel:

$$s(x) = \begin{cases} s_1(x), & x < a \\ s_{n-1}(x), & x > b \end{cases}$$

Funcții spline pătratice

Dacă în cazul funcțiilor spline liniare pe fiecare subinterval format am construit practic segmente, în cazul funcțiilor spline pătratice se vor construi arce de parabole prin intermediul funcțiilor de gradul doi. Astfel pe fiecare subinterval $[x_i, x_{i+1}]$ se alege

$s_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2$, cu a_i, b_i, c_i coeficienți ce urmează a fi determinați.

Astfel funcțiile spline trebuie să îndeplinească relațiile:

$$s_i(x) = f_i$$

$$s_i(x_{i+1}) = f_{i+1}$$

$$s'_i(x_i) = d_i$$

$$s'_i(x_{i+1}) = d_{i+1}$$

cu d_i construite recursiv cu ajutorul relației de ordinul unu:

$$d_{i+1} = -d_i + 2 \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}, \quad i = \overline{1, n-1}, \text{ iar } d_1 \text{ arbitrar.}$$

Se obțin relațiile:

$$s_i(x) = f(x_i) + d_i(x - x_i) + \frac{d_{i+1} - d_i}{2(x_{i+1} - x_i)}(x - x_i)^2$$

Cu ajutorul funcțiilor spline pătratice de mai sus se obține o funcție de aproximare care are derivatele continue pe $[a, b]$.

Exemple:

Să se determine funcția spline pătratică ce interpolează datele din tabelul:

x	0	1	2	3
f	1	4	8	10

Construim funcțiile spline calculând mai întâi coeficienții d_i pornind de la $d_1 = 0$.

Astfel obținem:

$$d_2 = 2 \frac{4-1}{1} = 6, \quad d_3 = -6 + 2 \frac{8-4}{1} = 2, \quad d_4 = -2 + 2 \frac{10-8}{1} = 2.$$

Funcțiile spline pătratice atașate subintervalelor formate sunt:

$$s_1(x) = 1 + \frac{6}{2}x^2 = 1 + 3x^2,$$

$$s_2(x) = 4 + 6(x-1) + \frac{2-6}{2}(x-1)^2 = -2x^2 + 10x - 4$$

$$s_3(x) = 8 + 2(x-2) = 2x + 4.$$

Astfel funcția spline corespunzătoare tabelului dat are următoarea expresie:

$$S(x) = \begin{cases} 1 + 3x^2, & x \in [0,1] \\ -2x^2 + 10x - 4, & x \in [1,2] \\ 2x + 4, & x \in [2,3] \end{cases}$$

Se observă ca ea este nu numai continuă, ci și derivabilă, iar derivata sa este continuă pe $[0,3]$.

Funcții spline cubice

Prin alegerea unor funcții de interpolare de gradul 3 se poate realiza o interpolare care presupune și fixarea valorii derivatelor pe suportul interpolării:

$$f'(x_0), f'(x_1), \dots, f'(x_n)$$

O funcție spline cubică se poate exprima astfel:

$$S_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3$$

Cei $4n$ coeficienți a_i, b_i, c_i, d_i ai funcțiilor spline se determină impunând $2n + 2$ condiții de interpolare:

$$S_i(x_i) = f(x_i), \quad S'_i(x_i) = f'(x_i) \text{ cu } i = 0 \dots n$$

și $2n - 2$ condiții de racordare, asigurând continuitatea și derivabilitatea funcțiilor de interpolare vecine în punctele interioare:

$$S_i(x_{i+1}) = S_{i+1}(x_{i+1}), \quad S'_i(x_{i+1}) = S'_{i+1}(x_{i+1}) \text{ cu } i = \overline{0, n-2}$$

2. APROXIMAREA ÎN SENSUL CELOR MAI MICI PĂTRATE

În paragrafele precedente s-au construit funcții care aproximează datele dintr-un tabel prin interpolare, adică funcții ce trec prin toate punctele din plan corespunzătoare tabelului dat.

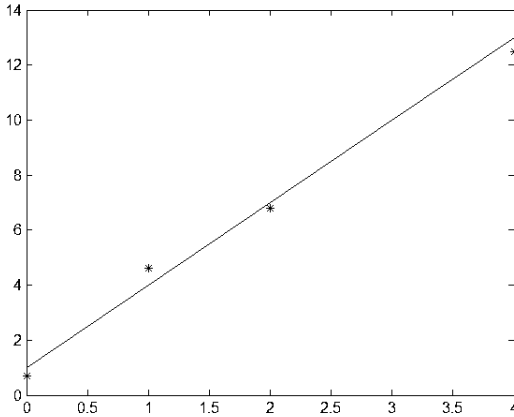
Uneori se pune problema găsirii unei funcții care să reprezinte cel mai bine datele din tabel, dar care să nu treacă neapărat prin toate punctele date.

Primul tip de aproximare se folosește când datele din tabel sunt date cu o mare acuratețe, iar cel de-al doilea în cazul în care în determinarea lor pot să apară mici erori (de exemplu de măsurare).

Cea mai des utilizată tehnică pentru a găsi o aproximantă care se potrivește cel mai bine datelor dintr-un tabel este metoda celor mai mici pătrate.

Primul pas ce trebuie făcut pentru a aplica această metodă este reprezentarea punctelor în plan, pentru a putea vizualiza forma norului de puncte și pentru a intui astfel ce model de funcție se potrivește pentru aproximantă.

În cazul în care se observă pe grafic că norul de puncte urmează forma unei drepte, precum în figura următoare, se folosește pentru a aproxima un model liniar: $f(x) = ax + b$.



Erorile care apar în fiecare nod, denumite și reziduuri, $e_i = f(x_i) - y_i$, trebuie minimize.

În cazul metodei celor mai mici pătrate coeficienții a și b se determină impunând condiția ca suma pătratelor erorilor ce apar să fie minimă. De aceea metoda folosită se numește metoda celor mai mici pătrate.

Astfel se formează suma: $E(a, b) = \sum_{i=1}^N ((ax_i + b) - y_i)^2$ căreia

vrem să-i determinăm minimul. Vom calcula mai întâi punctele ei staționare.

$$\left\{ \begin{array}{l} \frac{\partial E}{\partial a} = 0 \\ \frac{\partial E}{\partial b} = 0 \end{array} \right. , \text{ cu } \left\{ \begin{array}{l} \frac{\partial E}{\partial a} = 2 \sum_{i=1}^N ((ax_i + b) - y_i)x_i \\ \frac{\partial E}{\partial b} = 2 \sum_{i=1}^N ((ax_i + b) - y_i) \end{array} \right. .$$

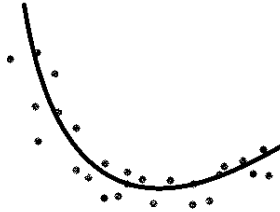
Se obține astfel sistemul:

$$\left\{ \begin{array}{l} a \sum_{i=1}^N x_i^2 + b \sum_{i=1}^N x_i - \sum_{i=1}^N x_i y_i = 0 \\ a \sum_{i=1}^N x_i + bN - \sum_{i=1}^N y_i = 0 \end{array} \right. .$$

Coeficienții din ecuația funcției de gradul întâi care aproximează cel mai bine datele în sensul celor mai mici pătrate sunt:

$$a = \frac{\left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right) - N \sum_{i=1}^N x_i y_i}{\left(\sum_{i=1}^N x_i \right)^2 - N \sum_{i=1}^N x_i^2}, \quad b = \frac{\sum_{i=1}^N y_i - a \sum_{i=1}^N x_i}{N} .$$

În cadrul aproximărilor pe care le facem locul funcției de gradul întâi poate fi luat de orice altă funcție polinomială și nu numai. În cazul în care se observă pe grafic că norul de puncte urmează forma unei parabole, precum în figura următoare, se folosește o funcție de gradul doi.



Pentru exprimarea unui model parabolic de forma: $y = a + bx + cx^2$, cu ajutorul metodei celor mai mici pătrate, construim mai întâi funcția E , ce însumează pătratele reziduurilor care apar : $E(a, b, c) = \sum_{i=1}^N ((cx_i^2 + bx_i + a) - y_i)^2$.

Coefficienții a, b, c se obțin din condiția ca funcția E să admită un minimum, adică:

$$\begin{cases} \frac{\partial E}{\partial a} = 0, \\ \frac{\partial E}{\partial b} = 0, \\ \frac{\partial E}{\partial c} = 0. \end{cases}$$

Calculând derivatele parțiale obținem expresiile:

$$\begin{cases} \frac{\partial E}{\partial a} = 2 \sum_{i=1}^N ((cx_i^2 + bx_i + a) - y_i), \\ \frac{\partial E}{\partial b} = 2 \sum_{i=1}^N ((cx_i^2 + bx_i + a) - y_i) x_i, \\ \frac{\partial E}{\partial c} = 2 \sum_{i=1}^N ((cx_i^2 + bx_i + a) - y_i) x_i^2. \end{cases}$$

Simplificând scrierea sumelor obținem sistemul următor:

$$\begin{cases} na + b\sum x_i + c\sum x_i^2 = \sum y_i \\ a\sum x_i + b\sum x_i^2 + c\sum x_i^3 = \sum x_i y_i \\ a\sum x_i^2 + b\sum x_i^3 + c\sum x_i^4 = \sum x_i^2 y_i \end{cases}$$

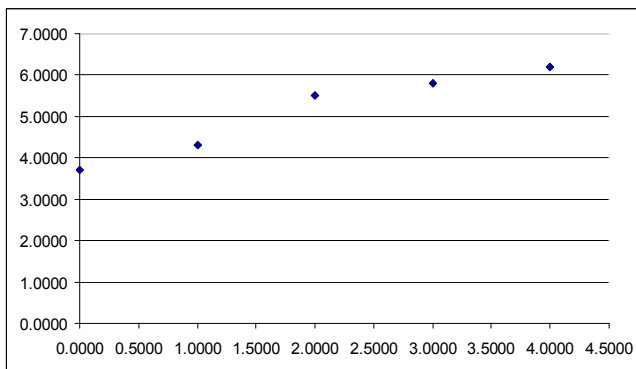
Rezolvând acest sistem obținem practic parabla căutăată.

Exemple:

1. Să se determine curba care aproximează cel mai bine în sensul celor mai mici pătrate datele din tabelul următor:

x	0	1	2	3	4
f	3,7	4,3	5,5	5,8	6,2

Reprezentând grafic punctele observăm că ele se așează în jurul unei drepte



Astfel vom determina o aproximantă de tipul unei funcții de gradul întâi: $f(x) = ax + b$, cu a și b dați de relațiile:

$$a = \frac{\left(\sum_{i=1}^N x_i\right)\left(\sum_{i=1}^N y_i\right) - N \sum_{i=1}^N x_i y_i}{\left(\sum_{i=1}^N x_i\right)^2 - N \sum_{i=1}^N x_i^2}, \quad b = \frac{\sum_{i=1}^N y_i - a \sum_{i=1}^N x_i}{N},$$

unde $N=5$.

Pentru ușurință calculele sunt trecute într-un tabel:

i	x_i	y_i	$x_i y_i$	x_i^2
1	0.0000	3.7000	0.0000	0.0000
2	1.0000	4.3000	4.3000	1.0000
3	2.0000	5.5000	11.0000	4.0000
4	3.0000	5.8000	17.4000	9.0000
5	4.0000	6.2000	24.8000	16.0000
Σ	10.0000	25.5000	57.5000	30.0000

Obținem: $a = \frac{10 \cdot 25.5 - 5 \cdot 57.5}{10^2 - 5 \cdot 30} = 0.75$, $b = \frac{25.5 - 0.75 \cdot 10}{5} = 3.7$

Astfel funcția care aproximează cel mai bine datele din tabel are expresia: $f(x) = 0.75x + 3.7$.

VIII. EVALUAREA NUMERICĂ A INTEGRALELOR

Prin evaluarea numerică a integralelor înțelegem calculul aproximativ al integralelor definite, $\int_a^b f(x)dx$, cu ajutorul unor formule care folosesc valorile funcției în anumite noduri din intervalul $[a,b]$. Necesitatea evaluării numerice a integralelor apare în special atunci când primitivele funcțiilor de integrat sunt dificil de aflat sau când acestea sunt date tabelar. Metodele cel mai des utilizate pentru evaluarea numerică a integralelor sunt: metoda trapezului, metoda Simpson și metoda Newton.

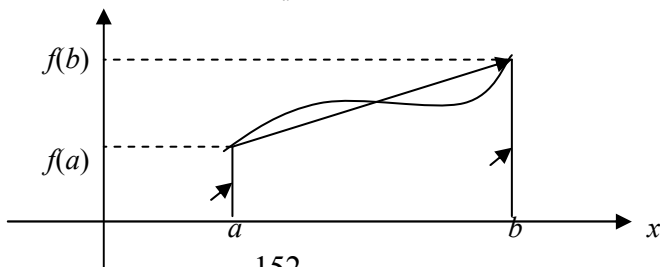
VIII.1. METODA TRAPEZULUI

Această metodă se obține aproximând funcția de integrat cu un polinom de interpolare Lagrange construit pe nodurile a și b , adică printr-un polinom Lagrange de gradul întâi.

x	a	b
f	$f(a)=f_1$	$f(b)=f_2$

Pe baza tabelului de mai sus obținem polinomul Lagrange:

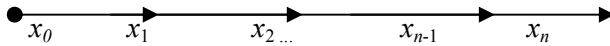
$$L_1(x) = \frac{x-b}{a-b} f(a) + \frac{x-a}{b-a} f(b) \Rightarrow \int_a^b f(x)dx \approx \frac{b-a}{2} [f(a) + f(b)]$$



Polinomul Lagrange de gradul 1, ne oferă o funcție care, restricționată la $[a, b]$, are ca imagine geometrică un segment de dreaptă care aproximează curba pe acest interval. Aria subgraficului acesteia este tocmai aria trapezului cu bazele $f(a)$, $f(b)$ și înălțimea $b - a$. De aici rezultă și denumirea metodei drept *metoda trapezului*.

Observația VIII.1. Dacă $f \in C^2[a, b]$, eroarea care apare aproximând integrala definită cu formula precedentă este: $-\frac{1}{12} f''(\xi)(b-a)^3$, $\xi \in (a, b)$.

Se observă că eroarea este mare dacă $[a, b]$ este de lungime mare. Pentru a micșora această eroare se consideră o diviziune a intervalului $[a, b]$ prin puncte echidistante. Se determină pasul h dat de: $h = \frac{b-a}{n}$ și apoi nodurile echidistante: $x_i = x_0 + ih$, $i = \overline{0, n}$, unde $x_0 = a$, $x_n = b$.



Aplicând proprietatea de aditivitate (față de domeniu) a integralei definite și formula precedentă pe fiecare subinterval format, obținem:

$$\int_a^b f(x) dx \cong \frac{h}{2} \left[f(x_0) + f(x_n) + 2 \sum_{i=1}^{n-1} f(x_i) \right] \quad (*)$$

Aproximând $\int_a^b f(x) dx$ cu formula precedentă eroarea care se obține, notată $E(f)$, verifică inegalitatea:

$$\text{cu } |E(f)| \leq \frac{(b-a)^3 M}{12n^2}, \text{ unde } M = \max_{x \in [a, b]} |f^{(2)}(x)|.$$

Observația VIII.2. Mărginirea erorii ne permite să evaluăm numeric integrala cu o anumită precizie, stabilită prin eroarea admisibilă dată, ε_{adm} .

Impunând condiția:

$$\frac{(b-a)^3}{12n^2}M < \varepsilon_{adm} \Rightarrow n^2 > \frac{(b-a)^3}{12\varepsilon_{adm}}M \Rightarrow$$

putem determina $n^* = \left\lceil \frac{(b-a)^3}{12\varepsilon_{adm}}M \right\rceil + 1$.

Putem evalua integrala pe baza formulei (*) considerând $n \geq n^*$ și vom obține o eroare mai mică decât eroarea admisibilă dată.

În continuare este prezentat un cod sursă în C, realizat pe baza metodei trapezelor, pentru evaluarea numerică a integralei funcției $f: R \rightarrow R$, $f(x) = x^2 + 3x$, pe un interval introdus de la tastatură.

```

/* Metoda Trapezelor */
#include<conio.h>
#include<stdio.h>
float a,b,x[1000],h, s;
float f(float);
int i,n;
void main()
{
printf("introduceti limita din stanga a:");
scanf("%f",&a);
printf("\nintroduceti limita din dreapta b:");
scanf("%f",&b);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
h=(b-a)/n;
for(i=1;i<=n-1;i++)

```

```

        x[i]=a+i*h;
s=f(a)+f(b);
for(i=1;i<=n-1;i++)
    s=s+2*f(x[i]);
s=s*h/2;
printf(" valoarea aproximativa a integralei este I=%f",s);
}
float f(float x)
{return x*x+3*x;}

```

```

C:\cc\bin\rundos.exe
introduceti limita din stanga a:2
introduceti limita din dreapta b:4
introduceti numarul de subintervale n:10
valoarea aproximativa a integralei este I:36.679996

```

```

C:\cc\bin\rundos.exe
introduceti limita din stanga a:2
introduceti limita din dreapta b:4
introduceti numarul de subintervale n:100
valoarea aproximativa a integralei este I:36.666801

```

Valoarea exactă a integralei este:

$$\int_2^4 (x^2 + 3x) dx = \frac{110}{3} = 36, (6)$$

Valoarea numerică obținută cu ajutorul sursei precedente în cazul n=10 este 36.379996

Cea obținută în cazul n=100 este 36.666801

VIII.3. METODA SIMPSON

În cadrul acestei metode, formula de calcul a valorii aproximative a integralei, ce se obține aproximând funcția f printr-un

polinom Lagrange de gradul cel mult doi, care interpolează funcția f

în a , $c = \frac{a+b}{2}$ și b , este:

$$\int_a^b f(x)dx = \frac{b-a}{6}(f(a)+4f(c)+f(b))+E(f).$$

Eroarea este dată de:

$$E(f) = -\frac{1}{2880}(b-a)^5 f^{(4)}(\xi), \xi \in (a, b)$$

Fiind o eroare mare, în cazul unui interval de lungime supraunitară, se recurge și în acest caz la discretizarea intervalului prin noduri echidistante ($x_i = x_0 + ih$, $i = \overline{0, n}$, unde $x_0 = a$, $x_n = b$,

$h = \frac{b-a}{n}$) și aplicarea formulei de mai sus pe fiecare subinterval

format. Însurându-se acestea se obține formula lui Simpson generalizată:

$$\int_a^b f(x)dx \cong \frac{h}{6} \left[f(x_0) + f(x_n) + 2 \sum_{i=1}^{n-1} f(x_i) + 4 \sum_{i=1}^n f\left(\frac{x_{i-1} + x_i}{2}\right) \right]$$

Observația VIII.3. Aproximarea integralei $\int_a^b f(x)dx$ prin

formula lui Simpson generalizată, în cazul unei funcții de clasă C^4 , se face cu eroarea $E(f)$, cu

$$|E(f)| \leq \frac{(b-a)^5}{2880n^4} M', \text{ unde } M' = \max_{x \in [a, b]} |f^{(4)}(x)|.$$

Utilizând limbajul de programare C s-a realizat următorul program de calcul.

```

/* Metoda Simpson */
#include<conio.h>
#include<stdio.h>
float f(float);
float a,b,x[1000],y[1000], p, h, s;
int i,n;
void main ()
{
printf("introduceti limita din stanga a:");
scanf("%f",&a);
printf("\nintroduceti limita din dreapta b:");
scanf("%f",&b);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
h=(b-a)/n;
for(i=0;i<=n;i++)
    x[i]=a+i*h;
s=f(a)+f(b);
for(i=1;i<=n-1;i++)
    s=s+2*f(x[i]);
p=0;
for(i=1;i<=n;i++)
    y[i]=(x[i-1]+x[i])/2;
for(i=1;i<=n;i++)
    p=p+4*f(y[i]);
s=(s+p)*h/6;
printf(" valoarea aproximativa a integralei este I=%f",s);
}
float f(float x)
{return x*x+3*x;}

```

VIII. 4. METODA NEWTON

Fie $f \in C^4[a, b]$ și diviziunea:

$$a < x_1 < x_2 < b, \quad x_i = a + i \frac{b-a}{3}, \quad i = 1, 2.$$

Folosind pentru aproximarea funcției un polinom de interpolare construit pe baza suportului de interpolare corespunzător diviziunii precedente se obține, în mod analog, o nouă formulă de evaluare numerică a unei integrale definite, ce poartă numele de formula lui

$$\text{Newton: } \int_a^b f(x)dx = \frac{b-a}{8} [f(a) + 3f(x_1) + 3f(x_2) + f(b)] + E(f)$$

$$E(f) = -\frac{1}{6480} (b-a)^5 f^{(4)}(\xi), \xi \in (a, b).$$

În cazul formulei generalizate se realizează o diviziune a intervalului $[a, b]$ în n părți egale prin nodurile echidistante $x_i, i = \overline{0, n}, x_0 = a, x_n = b$ și se aplică formula lui Newton pe fiecare interval $[x_{i-1}, x_i]$. Însușind relațiile obținem formula lui Newton generalizată:

$$\int_a^b f(x)dx \cong \frac{h}{8} \left\{ f(a) + f(b) + 2 \sum_{i=1}^{n-1} f(x_i) + 3 \sum_{i=0}^{n-1} \left[f(x_i + \frac{h}{3}) + f(x_i + \frac{2h}{3}) \right] \right\}$$

O margine superioară pentru valoarea absolută a erorii ce apare este dată de relația:

$$|E(f)| \leq \frac{(b-a)^5}{6480n^4} M'', \text{ unde } M'' = \sup_{\xi \in [a, b]} |f^{(4)}(\xi)|$$

Observația VIII.4. M' din formula lui Simpson generalizată și M'' din formula lui Newton generalizată coincid.

Deoarece eroarea obținută este mai mică atunci când se utilizează formula lui Newton

$$\left(\frac{(b-a)^5}{6480n^4} M'' \leq \frac{(b-a)^5}{2880n^4} M' \right),$$

pentru un același n , este preferabil să se folosească în aplicații formula lui Newton. Programul în C corespunzător pentru cazul funcției $f : R \rightarrow R, f(x) = x^2 + 3x$ este prezentat în continuare.

```

/* Metoda Newton */
#include<conio.h>
#include<stdio.h>
#include<math.h>
float f(float);
float a,b,x[1000],y[1000],z[1000], p, h, s;
int i,n;
void main ()
{
printf("introduceti limita din stanga a:");
scanf("%f",&a);
printf("\nintroduceti limita din dreapta b:");
scanf("%f",&b);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
h=(b-a)/n;
for(i=0;i<=n;i++)
    x[i]=a+i*h;
s=f(a)+f(b);
for(i=1;i<=n-1;i++)
    s=s+2*f(x[i]);
p=0;
for(i=1;i<=n;i++)
    y[i]=x[i-1]+h/3;
for(i=1;i<=n;i++)
    z[i]=x[i-1]+2*h/3;
for(i=1;i<=n;i++)
    p=p+3*(f(y[i])+f(z[i]));
s=(s+p)*h/8;
printf(" valoarea aproximativa a integralei este I=%f",s);
}float f(float x)
{return x*x+3*x;}

```

IX. METODE NUMERICE PENTRU REZOLVAREA ECUAȚIILOR DIFERENȚIALE

IX.1. INTRODUCERE

În acest capitol sunt prezentate câteva metode numerice pentru găsirea soluției unei probleme de forma:

$$\begin{cases} y' = f(x, y) \\ y(x_0) = y_0 \end{cases} \quad (1)$$

Ea reprezintă de fapt o ecuație diferențială însoțită de o condiție inițială și poartă numele de problemă Cauchy (sau cu date inițiale). Ele sunt foarte des întâlnite în practică deoarece majoritatea fenomenelor fizice studiate sunt scrise sub forma unor ecuații diferențiale, căci este mai ușor de determinat cum variază o mărime decât mărimea în sine.

Aplicarea metodelor numerice pentru rezolvarea acestor tipuri de probleme se face în cazul în care condiția de unicitate a soluției este asigurată.

Teorema IX.1. (stabilește condiții suficiente pentru ca problema precedentă să admită soluție unică)

Fie U o vecinătate a punctului (x_0, y_0) de forma: $U = \{(x, y) / x \in [x_0 - a, x_0 + a], y \in [y_0 - b, y_0 + b]\}$. Dacă f este mărginită și lipschitziană în raport cu variabila y (adică $|f(x, y_1) - f(x, y_2)| \leq L|y_1 - y_2|$), atunci există un număr $h > 0$ și o

mulțime $U_1 \subset U$ de forma: $U_1 = \{(x, y) / x \in [x_0, x_0 + h], |y - y_0| \leq b\}$ astfel încât, pe această mulțime, soluția problemei (1) există și este unică.

Observația IX.1. Cele mai frecvente cazuri în care f satisface cerințele teoremei precedente sunt:

1) Dacă f este astfel încât $\left| \frac{\partial f}{\partial y} \right| \leq L$ pe $D \Rightarrow f$ este lipschitziană în raport cu y pe D .

2) Dacă $f \in C'(D) \Rightarrow$ derivatele parțiale sunt mărginite și deci f este lipschitziană.

Există două direcții principale de abordare a rezolvării numerice a ecuației (1):

- se caută o funcție elementară $\tilde{y}(x)$, care să aproximeze suficient de bine soluția problemei (1)

$$\tilde{y}(x) = y_0 \text{ și } |y(x) - \tilde{y}(x)| < \varepsilon \quad (\forall) x \in [x_0, x_0 + h]$$

- se determină mărimile y_i , care aproximează valorile funcției $y(x)$ în punctele x_i , $i=0,1,2,\dots$, adică $|y(x) - y_i| < \varepsilon \quad i = 0,1,2,\dots$.

Cele mai simple metode numerice ce se pot utiliza în rezolvarea ecuațiilor diferențiale sunt: metoda integrării iterative, metoda dezvoltării în serie Taylor, metoda Euler și Runge Kutta. Metoda de integrare iterativă și cea de dezvoltare în serie fac parte din prima categorie iar metodele Euler și Runge-Kutta fac parte din cea de a doua categorie.

IX.2. METODA INTEGRĂRII ITERATIVE

Fie $f : U \rightarrow R$, continuă pe U și lipschitziană în raport cu y , de coeficient L . Se notează $M = \sup_U |f(x, y)|$. Dacă se alege

$h = \min\left(a, \frac{1}{L}, \frac{b}{M}\right)$, atunci șirul de aproximații succesive:

$$y_{n+1} = y_0 + \int_{x_0}^x f(t, y_n(t)) dt \quad n=0, 1, 2, \dots, \quad (2)$$

converge uniform pentru $x \in [x_0, x_0 + h]$ către soluția problemei (1).

Astfel $y(x) = \lim_{n \rightarrow \infty} y_n(x)$ este soluția problemei (1).

Are loc inegalitatea:

$$|y(x) - y_n(x)| < ML^n \frac{(x - x_0)^{n+1}}{(n+1)!}$$

Metoda integrării iterative mai poartă numele de metoda Picard și, deoarece construiește la fiecare iterație o aproximare a soluției problemei, se mai întâlnește și sub denumirea de metoda aproximațiilor succesive.

Algoritmul de aproximare este următorul:

Pasul 1: Se determină intervalul pe care procesul iterativ este convergent, adică valoarea lui h și intervalul $[x_0, x_0 + h]$.

Pasul 2: Se determină succesiv funcțiile $y_1(x), y_2(x), \dots$ cu ajutorul formulelor (2).

Pasul 3: Se determină indicele n pentru care:

$$|y(x) - y_n(x)| < \varepsilon \quad (\forall) x \in [x_0, x_0 + h],$$

unde ε este eroarea admisibilă, folosind inegalitatea:

$$|y(x) - y_n(x)| < M L^n \frac{(x - x_0)^{n+1}}{(n+1)!}.$$

Concluzie: y_n este aproximarea dorită.

Exemplul IX.1. Utilizând metoda aproximărilor succesive, să se determine o soluție aproximativă pentru problema:

$$\begin{cases} y' = 3x + \frac{1}{2}y \\ y(0) = 1 \end{cases} \Rightarrow x_0 = 0, y_0 = 1.$$

P1: Funcția $f(x, y) = 3x + \frac{1}{2}y \in C^1(D)$ pentru orice $D \subset \mathbb{R}^2$. Fie $U = \{(x, y) \mid |x| \leq 1, |y - 1| \leq 1\}$ ($a=1, b=1$). Determinăm

$$M = \sup_{(x, y) \in U} |f(x, y)| = \sup_{(x, y) \in U} \left| 3x + \frac{1}{2}y \right| = 4,$$

$$\text{și } L = \sup_{(x, y) \in U} \left| \frac{\partial f}{\partial y}(x, y) \right| = \sup_{(x, y) \in U} |1/2| = 1/2.$$

Alegem $h = \min\left(1, 2, \frac{1}{4}\right) = \frac{1}{4}$.

P2: Aproximațiile sunt:

$$Y_1(x) = 1 + \int_0^x \left(3t + \frac{1}{2}\right) dt = 1 + \frac{3x^2}{2} + \frac{x}{2}.$$

$$Y_2(x) = 1 + \int_0^x \left(3t + \frac{1}{2}\left(1 + \frac{t}{2} + \frac{3t^2}{2}\right)\right) dt = 1 + \frac{x}{2} + \frac{13}{8}x^2 + \frac{x^3}{4}.$$

$$Y_3(x) = 1 + \int_0^x \left(3t + \frac{1}{2}\left(1 + \frac{t}{2} + \frac{13}{8}t^2 + \frac{t^3}{4}\right)\right) dx = 1 + \frac{x}{2} + \frac{13x^2}{8} + \frac{13x^3}{48} + \frac{x^4}{32}$$

P3: Pentru $x \in [0, 0 + 1/4] = [0, 0.25]$ deducem că:

$$|y(x) - y_3(x)| <$$

$$4 \cdot \left(\frac{1}{2}\right)^3 \cdot \frac{(0.25)^4}{4!} = \frac{1}{3} \cdot \left(\frac{1}{2}\right)^{12} \approx 8.138 \cdot 10^{-5} \approx 0.8 \cdot 10^{-4} < 10^{-4}$$

De exemplu pentru

$\varepsilon = 10^{-4} \Rightarrow y(x) = y_3(x) = 1 + \frac{x}{2} + \frac{13x^2}{8} + \frac{13x^3}{48} + \frac{x^4}{32}$ este soluția problemei date. Ea aproximează soluția exactă cu o eroare de 10^{-4} .

Exemplu IX.2. Să considerăm problema Cauchy

$$\begin{cases} y' = x^2 + 2y \\ y(1) = 0 \end{cases}$$

Să se determine soluția aproximativă a ei, obținută după trei iterații, folosind metoda integrării iterative și să se evalueze un majorant pentru eroarea ce se realizează prin această aproximare.

Soluție:

Alegem $a=b=1$ și mulțimea $U = [0,2] \times [-1,1]$.

$$\text{Determinăm } M = \sup_{(x,y) \in U} |f(x,y)| = \sup_{(x,y) \in U} |x^2 + 2y| = 6$$

$$\text{Determinăm } L = \sup_{(x,y) \in U} \left| \frac{\partial f}{\partial y}(x,y) \right| = \sup_{(x,y) \in U} |2| = 2$$

Deci vom alege un număr $h < \frac{1}{6}$, atunci problema admite o singură soluție aproximativă y^* definită pe intervalul $\left[\frac{5}{6}, \frac{7}{4}\right]$, ca limită a șirului de aproximări (2).

Pentru $x \in \left[\frac{5}{6}, \frac{7}{4}\right]$ vom avea:

$$y_1(x) = Ty_0(x) = \int_0^x t^2 dt = \frac{x^3}{3}$$

$$y_2(x) = \int_0^x \left(t^2 + \frac{2t^3}{3}\right) dt = \frac{x^3}{3} + \frac{2x^4}{3 \cdot 4} = \frac{x^3}{3} + \frac{x^4}{6}$$

$$y_3(x) = \int_0^x \left(t^2 + 2\left(\frac{t^3}{3} + \frac{t^4}{6}\right)\right) dt = \\ = \frac{x^3}{3} + \frac{x^4}{6} + \frac{x^5}{15}$$

Observăm că avem:

$$\|y_3 - y^*\| \leq 6 \cdot 2^3 \left(\frac{1}{6}\right)^4 \cdot \frac{1}{4!} = \frac{8}{6^3} \frac{1}{3 \cdot 8} = \frac{1}{3 \cdot 6^3} \approx 0,00154$$

IX.2. METODA DEZVOLTĂRII ÎN SERIE (METODA TAYLOR)

Dacă se presupune că funcția f este derivabilă de o infinitate de ori într-o vecinătate a punctului (x_0, y_0) de rază r , atunci putem atașa soluției problemei (1) seria Taylor corespunzătoare :

$$y(x) = y_0 + \frac{x - x_0}{1!} y'(x_0) + \frac{(x - x_0)^2}{2!} y''(x_0) + \dots$$

Coefficienții se calculează consecutiv și sunt valori cunoscute:

$$y_0 = y(x_0), y'(x_0) = f(x_0, y(x_0)) = f_0.$$

Ei se obțin prin derivări succesive. De exemplu, $y''(x) = f'_x + y'f'_y + f''_x + f'_y \cdot f \Rightarrow$

$$y''(x_0) = f'_x(x_0, y_0) + f'_y(x_0, y_0) \cdot f(x_0, y_0) = f'_x(x_0, y_0) + f'_y(x_0, y_0) \cdot f_0.$$

Derivând funcția precedentă, se poate calcula $y'''(x)$ și apoi $y''''(x_0)$, etc. Astfel se pot obține valorile derivatelor de orice ordin în punctul x_0 .

Practic cu această metodă se aproximează soluția $y(x)$ a problemei (1) în vecinătatea punctului (x_0, y_0) printr-un polinom Taylor de grad n ($n=4$ sau 5 de cele mai multe ori).

Această metodă poate să ofere valorile soluției în nodurile echidistante ale unui interval.

Înlocuind x cu $x_0 + h$ vom deduce y_1 , adică valoarea din punctul x_1 . Avem relația:

$$y_1 = y_0 + h \left(f + \frac{h}{2} (f'_x + f \cdot f'_y) \right),$$

expresiile din dreapta egalității fiind evaluate în punctul (x_0, y_0)

Pentru celelalte valori deducem relația recursivă:

$$y_{i+1} = y_i + h \left(f + \frac{h}{2} (f'_x + f \cdot f'_y) \right), \text{ cu termenul drept}$$

evaluat în punctul (x_i, y_i) , $i = \overline{1, n-1}$.

Formule analoage se pot obține dacă reținem mai mulți termeni din seria Taylor corespunzătoare.

Exemplul IX.3. Fie problema Cauchy următoare:

$$\begin{cases} y' = 2x - y \\ y(0) = -1 \end{cases}.$$

Să se aproximeze pe $[0,1]$ soluția problemei folosind metoda Taylor, cu $n = 3$. Pasul pe axa Ox se alege $h = 0.1$. Să se compare rezultatele cu cele date de soluția exactă:

$$\hat{y}(x) = e^{-x} + 2x - 2.$$

Soluție:

Din datele problemei avem: $y(0) = -1$ $y'(0) = 1$.

Derivând în raport cu x ecuația diferențială din problemă deducem că:

$$y'' = 2 - y' = 2 - 2x + y.$$

Considerând în relația precedentă $x = 0$ și folosind valorile anterior calculate obținem:

$$y''(0) = 2 - 1 = 1.$$

Procedând analog și în continuare obținem;

$$y''' = -y'' = -2 + y' = -2 + 2x - y, \quad y'''(0) = -1$$

Astfel soluția aproximativă a problemei date obținută cu metoda Taylor este:

$$\begin{aligned} y(x) &\approx y(0) + \frac{x}{1!} y'(0) + \frac{x^2}{2!} y''(0) + \frac{x^3}{3!} y'''(0) = \\ &= -1 + x + \frac{x^2}{2} - \frac{x^3}{6}. \end{aligned}$$

Astfel, considerând $h = 0.1$ putem deduce valoarea soluției în punctul

$$x_1 = 0 + 0,1 = 0,1.$$

Obținem: $y_1 = -1 + 0,1 + \frac{0,1^2}{2} - \frac{0,1^3}{6} = -0,89517$.

Procedând analog în celelalte puncte găsim datele din tabelul următor. Sunt prezentate și erorile care apar între valorile calculate și cele exacte.

x	y	Soluti exacta	Eroarea absoluta
0.000000	-1.000000	-1.000000	0.000000
0.100000	-0.895163	-0.895167	0.000004
0.200000	-0.781269	-0.781277	0.000008
0.300000	-0.659182	-0.659192	0.000010
0.400000	-0.529680	-0.529692	0.000012
0.500000	-0.393469	-0.393483	0.000014
0.600000	-0.251188	-0.251203	0.000015
0.700000	-0.103415	-0.103430	0.000015
0.800000	0.049329	0.049313	0.000016
0.900000	0.206570	0.206553	0.000017
1.000000	0.367879	0.367863	0.000016

Eraoarea absolută maximă este deci 0.000017.

IX.3. Metoda EULER

Pornind de la faptul că soluția unei ecuații diferențiale de ordinul întâi reprezintă o familie de curbe din plan deducem că soluția problemei (1) reprezintă și ea o astfel de curbă.

Metoda Euler propune aproximarea acesteia, adică a soluției problemei, printr-o linie poligonală în care fiecare segment este coliniar cu direcția tangentei la această curbă în punctul care coincide cu extremitatea sa stângă (vezi Fig. IX.1.). Punctele alese în lungul axei Ox sunt echidistante, iar distanța dintre oricare două alăturate se notează cu h .

Fie astfel nodurile echidistante $x_i = x_0 + ih$.

În punctul $M_0(x_0, y_0)$ se trasează segmentul ce trece prin M_0 și are panta $y'(x_0) = f(x_0, y_0)$.

Ecuția este:

$$y - y_0 = y(x_0)(x - x_0) \text{ sau } y - y_0 = f'(x_0, y_0)(x - x_0).$$

Se aproximează soluția în punctul x_1 cu ordonata punctului M_1 de intersecție dintre dreapta precedentă și $x = x_1$.

Obținem deci: $y_1 = y_0 + f_0 h$.

Procedăm analog în continuare considerând că M_1 joacă rolul lui M_0 , etc. (vezi Fig. IX.1).

Obținem astfel formulele următoare pentru aflarea valorilor necunoscutei problemei în nodurile alese:

$$y_{i+1} = y_i + hf_i, \text{ unde } f_i = f(x_i, y_i).$$

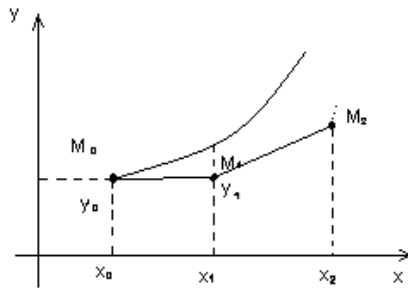


Fig. IX.1.

Astfel soluția problemei (1) este aproximată prin funcția tabel:

x	x ₀	x ₁	x _n
y(x)	y ₀	y ₁	y _n

Cunoscând valoarea lui h , care poartă numele de pasul metodei Euler, precum și numărul de valori dorite ale soluției problemei, adică n , se obține practic algoritmul de calcul numeric caracteristic metodei.

Algoritmul Metodei Euler:

Date de intrare: n -numărul de subintervale, h -pasul metodei, valorile inițiale x_0, y_0 și expresia funcției f .

Date de ieșire: vectorul ce conține valorile numerice: y_0, y_1, \dots, y_n .

P1: Se introduc datele de intrare;

P2: Se realizează diviziunea intervalului $[x_0, x_0 + nh]$ în n intervale prin nodurile x_0, x_1, \dots, x_n .

$$(x_i = x_0 + ih);$$

P3: Se determină valorile $y_i, i = \overline{0, n}$ cu relația:

$$y_{i+1} = y_i + hf_i, f_i = f(x_i, y_i);$$

P4: Se afișează valorile funcției și nodurile în care sunt calculate acestea, după care se oprește algoritmul.

Observația IX.3. Pentru micșorarea erorilor ce intervin în metoda Euler se înlocuiește direcția segmentului $\overline{M_i M_{i+1}}$ cu direcția corespunzătoare mijlocului segmentului $[x_i, x_{i+1}]$, rezultând metoda Euler modificată:

Formulele de calcul corespunzătoare acestei metode pot fi sintetizate astfel:

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

Exemplul IX.4. Fie problema Cauchy:

$$\begin{cases} y' = 2x - y \\ y(0) = -1 \end{cases} \quad x \in [0,1].$$

Să se determine o soluție aproximativă a ei folosind metoda Euler cu pasul $h = 0.1$.

Să se compare valorile găsite cu cele ale soluției exacte $\hat{y}(x) = e^{-x} + 2x - 2$.

Soluție:

Avem următoarele relații:

$$x_0 = 0; x_1 = 0,1; x_2 = 0,2; x_3 = 0,3; x_4 = 0,4, \dots, x_{10} = 1.$$

$$y_0 = -1 \Rightarrow y_1 = y_0 + hf(x_0, y_0) = -1 + 0,1(2 - (-1)) = -0,9.$$

$$y_2 = y_1 + hf(x_1, y_1) \cong -0,79, \text{ etc.}$$

Obținem următorul tabel:

x	$y(x)$	$\hat{y}(x)$	eroarea absolută
0	-1.000000	-1	0
0.1	-0.900000	-0.89516	0.004837
0.2	-0.790000	-0.78127	0.008731
0.3	-0.671000	-0.65918	0.011818
0.4	-0.543900	-0.52968	0.01422
0.5	-0.409510	-0.39347	0.016041
0.6	-0.268559	-0.25119	0.017371
0.7	-0.121703	-0.10341	0.018288
0.8	0.030467	0.049329	0.018862
0.9	0.187420	0.20657	0.01915
1	0.348678	0.367879	0.019201

Se observă că cea mai mare eroare absolută obținută este: 0.019201

Exemplul IX.4. Fie problema Cauchy:

$$\begin{cases} y' = y - \frac{y}{x} \\ y(1) = \frac{1}{2} \end{cases} \quad x \in [1,2].$$

Să se determine o soluție aproximativă a ei folosind metoda Euler cu pasul $h = 0.2$. Să se compare valorile găsite cu

valorile soluției exacte: $\hat{y}(x) = \frac{e^{x-1}}{2x}$

Soluție:

Avem următoarele relații:

$$x_0 = 1; x_1 = 1,2; x_2 = 1,4; x_3 = 1,6; x_4 = 1,8; x_5 = 2.$$

$$y_0 = \frac{1}{2} \Rightarrow y_1 = y_0 + hf(x_0, y_0) = \frac{1}{2} + 0,2\left(\frac{1}{2} - \frac{1}{2}\right) = 0,5.$$

$$y_2 = y_1 + hf(x_1, y_1) = 0,5 + 0,2(0,5 - 0,5/1,2) = 0,516667 \text{ etc.}$$

Pe baza metodei Euler s-a realizat în C următorul program:

```
/* Metoda Euler pentru rezolvare a ecuatiilor diferentiale */
```

```
#include<conio.h>
```

```
#include<stdio.h>
```

```
float fun(float, float);
```

```
float x0,y0,x[100],y[100],h;
```

```
int i,n;
```

```
float fun(float x,float y)
```

```
{return y-y/x;}
```

```
void main ()
```

```
{
```

```
printf("introduceti abscisa punctului initial x0:");
```

```
scanf("%f",&x0);
```

```
printf("\nintroduceti ordonata initiala y0:");
```

```
scanf("%f",&y0);
```

```

printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
printf("introduceti pasul metodei h:");
scanf("%f",&h);
x[0]=x0;
y[0]=y0;
for(i=1;i<=n;i++)
{
    x[i]=x0+i*h;
    y[i]=y[i-1]+fun(x[i-1],y[i-1])*h;
}
printf("valorile aproximative ale solutiei sunt:\n");
for(i=0;i<=n;i++)
    printf("y[%d]=%f",i,y[i]);
getch();
}

```

Rezultatele obținute prin rularea lui sunt:

```

C:\vc\bin\vrundose.exe
introduceti abscisa punctului initial x0:1
introduceti ordonata initiala y0:0.5
introduceti numarul de subintervale n:5
introduceti pasul metodei h:0.2
valorile aproximative ale solutiei sunt:
y[0]=0.500000y[1]=0.500000y[2]=0.516667y[3]=0.546190y[4]=0.587155y[5]=0.639346

```

Obținem următorul tabel:

x	$y(x)$	$\hat{y}(x)$	eroarea absolută
1	0.5	0.5	0
1.2	0.5	0.508918	0.008918
1.4	0.516667	0.532795	0.016128
1.6	0.546190	0.569412	0.023222
1.8	0.587155	0.618206	0.031051
2	0.639346	0.67957	0.040224

Se observă că cea mai mare eroare absolută este: 0.040224

Exemplul IX.4. Să se rezolve aceeași problemă folosind metoda Euler modificată.

Soluție:

Nodurile de pe axa Ox rămân aceleași:

$$x_0 = 1; x_1 = 1,2; x_2 = 1,4; x_3 = 1,6; x_4 = 1,8; x_5 = 2.$$

$$\text{Avem, din datele problemei: } y_0 = \frac{1}{2}.$$

Determinăm mijlocul primului subinterval format, valoarea ordonatei corespunzătoare și apoi valoarea funcției în acest nod:

$$x_{1/2} = x_0 + h/2 = 1 + 0,2/2 = 1,1;$$

$$y_{1/2} = y_0 + (h/2) * f(x_0, y_0) = \frac{1}{2} + 0,1 * 0 = 0,5;$$

$$f(x_{1/2}, y_{1/2}) = f(1,1; 0,5) = 0,5 - (0,5)/1,1 = 0,045455 = f_{1/2}.$$

Cu aceste valori calculăm prima valoare numerică a necunoscutei din problemă, adică valoarea y_1 :

$$y_1 = y_0 + f_{1/2} * h \Rightarrow y_1 = 0,5 + 0,2 * 0,045455 = 0,509091;$$

Rezultatele au fost obținute cu ajutorul următorului program realizat în C:

```
/* Metoda Euler modificata */
#include<conio.h>
#include<stdio.h>
#include<math.h>
float fun(float, float);
float x0,y0,x[100],y[100],Y[100],a[100],b[100],h;
int i,n;
float fun(float x,float y)
{return y-y/x;}
void main ()
```



```

{
printf("introduceti abscisa punctului initial x0:");
scanf("%f",&x0);
printf("\nintroduceti ordonata initiala y0:");
scanf("%f",&y0);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
printf("introduceti pasul metodei h:");
scanf("%f",&h);
x[0]=x0;
y[0]=y0;
for(i=0;i<=n;i++)
    {x[i+1]=x0+(i+1)*h;
      a[i]=x[i]+h/2;
      b[i]=y[i]+h*fun(x[i],y[i])/2;
      y[i+1]=y[i]+fun (a[i], b[i])*h;
    }
printf(" valorile aproximative ale solutiei sunt:\n");
for(i=0;i<=n;i++)
    printf("y[%d]=%f",i,y[i]);
getch();
}

```

```

C:\cc\bin\rundos.exe
introduceti abscisa punctului initial x0:1
introduceti ordonata initiala y0:0.5
introduceti numarul de subintervale n:5
introduceti pasul metodei h:0.2
 valorile aproximative ale solutiei sunt:
y[0]=0.500000y[1]=0.509091y[2]=0.532979y[3]=0.569526y[4]=0.618187y[5]=0.679355

```

Datele numerice obținute sunt trecute în tabelul următor, tabel în care se compară acestea cu valorile reale.

x	$y(x)$	$\hat{y}(x)$	eroarea absolută
	Euler modificata		
1	0.5	0.5	0
1.2	0.509091	0.508918	0.009091
1.4	0.532979	0.532795	0.016312
1.6	0.569526	0.569412	0.023336
1.8	0.618187	0.618206	0.031032
2	0.679355	0.67957	0.040009

Se observă că cea mai mare eroare absolută este: 0.040009 și ea este mai mică decât în cazul aplicării metodei Euler. Se observă o îmbunătățire a rezultatelor în cazul aplicării metodei Euler-modificată.

IX.4. Metode Runge-Kutta (R-K)

Aceste metode constau în aproximarea soluției $y(x)$ a problemei (1) în punctul $x_i + h$, dacă se cunoaște soluția în punctul x_i , cu ajutorul unor formule de tipul:

$$y_{i+1} \approx y_i + \Delta y_i,$$

unde $\Delta y_i = c_1 k_1 + c_2 k_2 + \dots + c_m k_m$, coeficienții k_i , $i = \overline{1, m}$ fiind valorile funcției $f(x, y)$ în anumite puncte din intervalul (x_i, x_{i+1}) .

Cea mai simplă dintre aceste metode este metoda Runge-Kutta de ordinul doi (RK2). În acest caz expresia lui Δy este:

$$\Delta y = c_1 k_1 + c_2 k_2, \text{ cu } k_1, k_2 \text{ dați de relațiile:}$$

$$k_1 = hf(x_i, y_i), \quad k_2 = hf(x_i + h, y_i + \beta k_1),$$

unde c_1, c_2, α, β sunt constante ce urmează a fi determinate.

Pentru a determina acești coeficienți evaluăm y_{i+1} cu formula lui Taylor corespunzătoare funcției y în punctul x_i .

Avem:

$$y(x_{i+1}) = y(x_i) + hy'(x_i) + \frac{h^2}{2} y''(x_i) + \dots = y(x_i) + hf(x_i, y_i) + \frac{h^2}{2} f''(x_i, y(x_i)) + \dots$$

Dar

$$f''(x_i, y(x_i)) = \frac{\partial^2 f}{\partial x^2}(x_i, y(x_i)) + \frac{\partial^2 f}{\partial y^2}(x_i, y(x_i))y'(x_i) = \frac{\partial^2 f}{\partial x^2}(x_i, y(x_i)) + \frac{\partial^2 f}{\partial y^2}(x_i, y(x_i))f(x_i, y(x_i))$$

Astfel obținem:

$$y(x_{i+1}) = y(x_i) + hf(x_i) + \frac{h^2}{2} \left(\frac{\partial^2 f}{\partial x^2}(x_i, y(x_i)) + \frac{\partial^2 f}{\partial y^2}(x_i, y(x_i))f(x_i, y(x_i)) \right) + O(h^3)$$

(*)

Dezvoltând funcția $f(x_i + h, y_i + \beta k_1)$ în jurul lui (x_i, y_i) avem relația:

$$f(x_i + \alpha h, y_i + \beta k_1) = f(x_i, y_i) + \alpha h \frac{\partial f}{\partial x} + \beta k_1 \frac{\partial f}{\partial y} + O(h^2)$$

Dar $y_{i+1} = y_i + c_1 k_1 + c_2 k_2$ și astfel obținem:

$$y(x_{i+1}) = y(x_i) + h(c_1 + c_2)f + c_2 h^2 \left(\alpha \frac{\partial f}{\partial x} + \beta f \frac{\partial f}{\partial y} \right) + O(h^3), \text{ funcțiile}$$

din membrul drept fiind evaluate în $(x_i, y(x_i))$.

Comparând relația găsită cu relația (*) deducem că c_1, c_2, α, β verifică sistemul compatibil nedeterminat:

$$\begin{cases} c_1 + c_2 = 1 \\ \alpha = \beta = \frac{1}{2c_2} \end{cases},$$

Deducem de aici că există mai multe formule corespunzătoare metodei RK2.

1. Una din soluțiile acestui sistem este: $c_1 = c_2 = \frac{1}{2}$, $\alpha = \beta = 1$.

Astfel formula corespunzătoare ei este:

$$y_{i+1} = y_i + \frac{h}{2}(f(x_i, y_i) + f(x_i + h, y_i + hf(x_i, y_i))),$$

$$i = 0, \dots, N-1$$

2. O altă alegere a coeficienților c_1, c_2, α, β , și anume:

$c_1 = 0, c_2 = 1$, $\alpha = \beta = \frac{1}{2}$, conduce la obținerea formulei corespunzătoare metode Euler modificată:

$$y_{i+1} = y_i + hf\left(x_i + \frac{h}{2}, y_i + \frac{h}{2}f(x_i, y_i)\right), \quad i = 0, \dots, N-1.$$

3. Dacă $c_1 = 1/4$, $c_2 = 3/4$, $\alpha = \beta = 2/3$, se obține formula:

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

Pentru cazul formulei Runge-Kutta de ordinul 2 de mai sus s-a realizat un program în C, pentru rezolvarea problemei considerate în exemplul IX.4.

/* Metoda Runge-Kutta de ordinul 2 */

```
#include<conio.h>
```

```
#include<stdio.h>
```

```
#include<math.h>
```

```
float fun(float, float);
```

```
float x0,y0,x[100],y[100],c1[100],c2[100],h;
```

```
int i,n;
```

```
float fun(float x,float y)
```

```
{return y-y/x;}
```

```

void main ()
{
printf("introduceti abscisa punctului initial x0:");
scanf("%f",&x0);
printf("\nintroduceti ordonata initiala y0:");
scanf("%f",&y0);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
printf("introduceti pasul metodei h:");
scanf("%f",&h);
x[0]=x0;
y[0]=y0;
for(i=0;i<=n;i++)
    {x[i+1]=x0+(i+1)*h;
      c1[i]=fun(x[i],y[i]);
      c2[i]=fun(x[i]+2*h/3,y[i]+2*h*c1[i]/3);
      y[i+1]=y[i]+h*(c1[i]+3*c2[i])/4;
    }
printf(" valorile aproximative ale solutiei sunt:\n");
for(i=0;i<=n;i++)
    printf("y[%d]=%f",i,y[i]);
getch();
}

```

```

C:\cc1\bin\rundos.exe
introduceti abscisa punctului initial x0:1
introduceti ordonata initiala y0:0.5
introduceti numarul de subintervale n:5
introduceti pasul metodei h:0.2
 valorile aproximative ale solutiei sunt:
y[0]=0.500000y[1]=0.508824y[2]=0.532569y[3]=0.569021y[4]=0.617607y[5]=0.678705

```

În urma rulării acestui program pentru pasul $h = 0.2$ și $n = 5$ s-au obținut datele din tabelul următor. Tot în acest tabel se face o comparație între valorile numerice și cele exacte evaluându-se erorile absolute ce apar.

x	$y(x)$ RK2	$\hat{y}(x)$	eroarea absolută
1	0.5	0.5	0.000000
1.2	0.508824	0.508918	0.000094
1.4	0.532569	0.532795	0.000226
1.6	0.569021	0.569412	0.000391
1.8	0.617607	0.618206	0.000599
2	0.678705	0.67957	0.000865

Se observă că erorile ce apar sunt într-adevăr mai mici decât $h^3 = 0.008$.

Programul se poate adapta foarte ușor pentru cazul altor formule Runge-Kutta de ordinul doi, și pentru alte ecuații diferențiale de ordinul întâi cu soluție unică.

Pentru a obține însă aproximări mai bune se folosesc de cele mai multe ori formule de ordin superior, cum sunt formule Runge-Kutta de ordinul 3 și 4. Și în acest caz există mai multe formule deoarece sistemul rezultat este tot compatibil nedeterminat.

Dintre cele mai importante formule Runge-Kutta de ordinul 3 amintim în continuare câteva:

- $y(x+h) = y(x) + 1/6(k_1 + 4k_2 + k_3)$,

unde $k_1 = hf(x, y)$, $k_2 = hf(x + h/2, y + k_1/2)$,
 $k_3 = hf(x + h, y - k_1 + 2k_2)$

- $y(x+h) = y(x) + 1/4(k_1 + 3k_3)$

unde $k_1 = hf(x, y)$, $k_2 = hf(x + h/3, y + k_1/3)$,
 $k_3 = hf(x + 2h/3, y_1 + 2k_2/3)$

În cazul folosirii acestor formule se obține o eroare de ordinul lui h^4 .

Pentru cazul primei formule Runge-Kutta de ordinul 3 avem:

$$y(x+h) = y(x) + h/6(c_1 + 4c_2 + c_3),$$

cu $c_1 = f(x, y)$, $c_2 = f(x + h/2, y + k_1/2)$,
 $c_3 = f(x + h, y - k_1 + 2k_2)$

Un program realizat în C pe baza acestei formule, pentru rezolvarea problemei considerate în exemplul IX.4, este următorul:

```

/* Metoda Runge-Kutta de ordinul 3 */
#include<conio.h>
#include<stdio.h>
#include<math.h>
float fun(float, float);
float x0,y0,x[1000],y[1000],c1[100],c2[100],c3[100],h;
int i,n;
float fun(float x,float y)
{return y-y/x;}
void main ()
{
printf("introduceti abscisa punctului initial x0:");
scanf("%f",&x0);
printf("\nintroduceti ordonata initiala y0:");

```

```

scanf("%f",&y0);
printf("introduceti numarul de subintervale n:");
scanf("%d",&n);
printf("introduceti pasul metodei h:");
scanf("%f",&h);
x[0]=x0;
y[0]=y0;
for(i=0;i<=n;i++)
    {x[i+1]=x0+(i+1)*h;
      c1[i]=fun(x[i],y[i]);
      c2[i]=fun(x[i]+h/2,y[i]+c1[i]*h/2);
      c3[i]=fun(x[i]+h,y[i]-c1[i]*h+2*c2[i]*h);
      y[i+1]=y[i]+(c1[i]+4*c2[i]+c3[i])*h/6;
    }
printf(" valorile aproximative ale solutiei sunt:\n");
for(i=0;i<=n;i++)
    printf("y[%d]=%f",i,y[i]);
getch();
}

```

```

C:\oc\bin\rundos.exe
introduceti abscisa punctului initial x0:1
introduceti ordonata initiala y0:0.5
introduceti numarul de subintervale n:5
introduceti pasul metodei h:0.2
valorile aproximative ale solutiei sunt:
y[0]=0.500000y[1]=0.508939y[2]=0.532828y[3]=0.569453y[4]=0.618253y[5]=0.679623

```

În urma rulării acestui program pentru pasul $h = 0.2$ și $n = 5$ s-au obținut datele de mai sus. În tabelul următor sunt comparate acestea cu soluția exactă și cu cea numerică obținută cu metoda RK2.

x	$y(x)$ RK3	$\hat{y}(x)$	eroarea absolută	$y(x)$ RK2	eroarea absolută
1	0.5	0.5	0.000000	0.5	0.000000
1.2	0.508824	0.508918	0.000094	0.508939	0.000021
1.4	0.532569	0.532795	0.000226	0.532828	0.000033
1.6	0.569021	0.569412	0.000391	0.569453	0.000041
1.8	0.617607	0.618206	0.000599	0.618253	0.000047
2	0.678705	0.67957	0.000865	0.679623	0.000053

Comparând aceste date cu cele din tabelul ce conține valorile soluției exacte a problemei și cu cele obținute ca urmare a aplicării metodei RK2, se observă că erorile ce apar sunt evident mai mici în cazul aplicării metodei RK3.

Programul se poate adapta foarte ușor pentru cazul altor formule Runge-Kutta, și evident pentru cazul altor ecuații diferențiale de ordinul întâi.

Cele mai cunoscute formule Runge-Kutta de ordinul 4 sunt:

- $$y(x+h) = y(x) + 1/6(k_1 + 4k_2 + k_3 + k_4)$$

unde
$$k_1 = hf(x, y), \quad k_2 = hf(x+h/2, y+k_1/2),$$

$$k_3 = hf(x+h/2, y+k_2/2), k_4 = hf(x+h, y+k_3)$$

- $$y(x+h) = y(x) + 1/8(k_1 + 4k_2 + k_3 + k_4),$$

unde
$$k_1 = hf(x, y), \quad k_2 = hf(x+h/3, y+k_1/3),$$

$$k_3 = hf(x+2h/3, y-k_1/3+k_2),$$

$$k_4 = hf(x+h, y+k_1-k_2+k_3).$$

BIBLIOGRAFIE

1. Abdelwahab Kharab and Ronald B Guenther, An introduction to numerical methods: A MATLAB Approach, 2002
2. Antia H. M. Numerical methods for scientists and engineers, Birkhaussen, 2002
3. Brătianu C., Bostan V., Cojocia L, Negreanu G., Metode numerice, Editura Tehnică, Bucuresti 1996.
4. Ebâncă D., Metode de calcul numeric, Editura Sitech, Craiova 1994.
5. Iorga V., Jora B., etc. Programare Numerică, Editura Teora București 1996.
6. Pavel L, Rizzoli I. Elemente de Analiză Funcțională și Analiză Numerică, Editura Universității din Bucuresti, 2002.
7. Pitea A., Postolache M., Modelare numerică pentru ecuații diferențiale si ecuații cu derivate parțiale, Editura Fair Partners, București 2007.
8. Popa M., Popa A., Militaru R., Noțiuni de Analiză Numerică, Editura Sitech, Craiova, 2001.
9. Postolache M., Modelare Numerică. Teorie si Aplicații, Editura Fair Partners, București 2008.
10. Rasa I., Vladislav T., Analiză Numerică. Curbe spline, operatori Bernstein, algoritmul lui Casteljaou, Editura Tehnică, București 1998.
11. Rinderiu P., Gruionu L., Metode Numerice-Elemente teoretice și aplicative, Editura Universitaria Craiova 2003.
12. Simionescu I., Dranga M., Moise V., Metode Numerice în Tehnică. Editura Tehnică, București 1995.
13. Vladislav T., Rasa I., Analiză Numerică, Editura Tehnică, București 1997.
14. Vraciu G., Popa M., Efreem R., Micu S., Analiză Numerică. Culegere de exerciții și probleme, (vol.I), Editura Sitech, Craiova 1996.
15. Vraciu G., Popa A., Metode numerice cu aplicații în tehnica de calcul, Editura Scrisul Românesc, Craiova, 1982.