

ERORI DE TRUNCHIERE

Erorile produse ca urmare a calculului în virgulă mobilă sunt erori de trunchiere. Acestea apar ca urmare a faptului că numerele sunt reprezentate cu un **număr finit de zecimale**.

Notății

x - valoarea exactă a unei mărimi;

\bar{x} - valoarea aproximativă a aceleiași mărimi;

e_x - eroarea de trunchiere.

Cu aceste notații, este valabilă relația: $x = \bar{x} + e_x$

Eroarea produsă la operația de **adunare**: $x + y = \bar{x} + \bar{y} + e_x + e_y \Rightarrow e_{x+y} = e_x + e_y$

Eroarea produsă la operația de **scădere**: $x - y = \bar{x} - \bar{y} + e_x - e_y \Rightarrow e_{x-y} = e_x - e_y$

Eroarea produsă la operația de **înmulțire**

$$x \cdot y = (\bar{x} + e_x) \cdot (\bar{y} + e_y) = \bar{x} \cdot \bar{y} + \bar{x} \cdot e_y + \bar{y} \cdot e_x + e_x \cdot e_y$$

neglijând termenul $e_x \cdot e_y$ obținem $e_{x \cdot y} = \bar{x} \cdot e_y + \bar{y} \cdot e_x$

Observație

La operația de înmulțire eroarea produsă depinde nu numai de erorile termenilor care intervin în operația respectivă dar și de valoarea acestora.

Eroarea produsă la operația de **împărțire**

$$\frac{x}{y} = \frac{\bar{x} + e_x}{\bar{y} + e_y} = \frac{\bar{x} + e_x}{\bar{y} \cdot \left(1 + \frac{e_y}{\bar{y}}\right)} = \frac{1}{1 + \frac{e_y}{\bar{y}}} \cdot \frac{\bar{x} + e_x}{\bar{y}} = \frac{1 - \frac{e_y}{\bar{y}}}{1 - \left(\frac{e_y}{\bar{y}}\right)^2} \cdot \frac{\bar{x} + e_x}{\bar{y}} \cong \frac{\bar{x} + e_x}{\bar{y}} \cdot \left(1 - \frac{e_y}{\bar{y}}\right)$$

În șirul precedent de egalități s-a neglijat termenul: $\left(\frac{e_y}{\bar{y}}\right)^2 \approx 0$

$$\frac{x}{y} = \frac{\bar{x}}{y} + \frac{e_x}{y} - \frac{\bar{x}}{y^2} \cdot e_y - \frac{e_x \cdot e_y}{y^2} \quad \text{neglijând termenul:} \quad \frac{e_x \cdot e_y}{y^2}$$

se obține următoarea expresie pentru eroarea de trunchiere produsă la operația de împărțire:

$$e_{\frac{x}{y}} = \frac{e_x}{y} - \frac{\bar{x}}{y^2} \cdot e_y$$

Observație

La operația de împărțire eroarea produsă depinde nu numai de erorile termenilor care intervin în operația respectivă dar și de valoarea acestora, **cu cât împărțitorul este mai mic cu atât valoarea erorii este mai mare.**

Ca urmare rezultă că în procesul de *inversare a unei matrici* este de preferat ca la fiecare etapă să se rearanjeze matricea, prin permutări de linii și coloane, astfel încât termenul de pe diagonală să fie cel cu valoarea absolută maximă la etapa respectivă. Acest proces poartă denumirea de pivotare și poate fi efectuat în două moduri:

→ *Pivotare parțială*, dacă la fiecare etapă se determină valoarea maximă de pe coloana corespunzătoare etapei respective

→ *Pivotare totală*, dacă la fiecare etapă se determină valoarea maximă din întreaga matrice.

Există o categorie de **matrici** numite **slab condiționate** care se caracterizează prin aceea că o mică modificare asupra valorii unui termen al matricii determină o modificare importantă asupra termenilor matricii inverse. Aceste matrici se caracterizează prin aceea că **valoarea determinantului matricii este mult mai mic decât termenii matricii**.

Exemplu:

$$A = \begin{bmatrix} 100 & 10 \\ 9,5 & 1 \end{bmatrix} ; \quad |A| = 5 \quad ; \quad A^{-1} = \begin{bmatrix} 0,2 & -2 \\ -1,9 & 20 \end{bmatrix}$$

presupunem că se produce o mică modificare a unuia dintre termenii matricii:

$$A = \begin{bmatrix} 100 & 10 \\ 9,9 & 1 \end{bmatrix} ; \quad |A| = 1 \quad ; \quad A^{-1} = \begin{bmatrix} 1 & -10 \\ -9,9 & 100 \end{bmatrix}$$

se poate constata o modificare semnificativă a termenilor matricii inverse.

Erorile de trunchiere pot determina obținerea unor rezultate complet diferite față de valoarea reală a matricii inverse. Ca urmare pentru calculul matricii inverse se va utiliza un **algoritm iterativ** care permite efectuarea **corectării matricii inverse până când se obține rezultatul cu precizia dorită**.

ALGORITMUL LUI HOTTELING PENTRU CALCULUL INVERSEI UNEI MATRICI

Algoritmul se bazează pe efectuarea unei operații de **împărțire prin operația de înmulțire**. Vom exemplifica algoritmul pentru cazul numerelor reale:

Presupunem că se cunoaște o primă aproximare $(1/a)$ a inversului numărului real a pe care o notăm d_1 . Condiția de convergența a algoritmului este ca:

$$e_1 = |1 - a \cdot d_1| < 1$$

Algoritmul constă în determinarea unui șir de aproximări d_2, d_3, d_4, \dots a inversului numărului a astfel încât erorile să tindă spre zero.

Pentru a realiza condiția cerută este necesar să se determine o relație de recurență astfel încât eroarea la un moment dat să poată fi exprimată sub forma unei puteri a lui e_1 . Aceasta deoarece:

$$\lim_{k \rightarrow \infty} e_1^k = 0$$

Se pune condiția: $e_2 = 1 - a \cdot d_2 = e_1^2$

din care rezultă succesiv:

$$e_1^2 = e_1 \cdot e_1 = (1 - a \cdot d_1) \cdot (1 - a \cdot d_1) = 1 - 2 \cdot a \cdot d_1 + a^2 \cdot d_1^2$$

adică:

$$e_1^2 = 1 - a \cdot d_1 \cdot [1 + (1 - a \cdot d_1)] = 1 - a \cdot d_1 \cdot (1 + e_1)$$

și

$$d_2 = d_1 \cdot (1 + e_1)$$

Procesul de determinare al inversului unui număr va fi următorul:

$$d_2 = d_1 \cdot (1 + e_1) \Rightarrow e_2 = 1 - a \cdot d_2$$

$$d_3 = d_2 \cdot (1 + e_2) \Rightarrow e_3 = 1 - a \cdot d_3$$

.....

$$d_k = d_{k-1} \cdot (1 + e_{k-1}) \Rightarrow e_k = 1 - a \cdot d_k$$

care continuă până când e_k devine mai mic decât o valoare impusă.

Pentru aplicarea algoritmului în cazul matricilor se va utiliza noțiunea de *normă* a unei matrici, definită astfel:

$$N[A] = \max_i \left(\sum_{j=1}^n |a_{i,j}| \right) \quad \text{Dacă} \quad N[A] < 1 \quad \text{atunci} \quad \lim_{n \rightarrow \infty} N[A]^n = 0$$

Pentru a calcula inversa unei matrici utilizând algoritmul lui Hotteling se va proceda după cum urmează:

Presupunem:

A - Matricea pentru care dorim să calculăm inversa

D_1 - Prima aproximare a inversei matricii, determinată cu algoritmul descris anterior

Eroarea care se produce în acest caz se va calcula astfel:

$$E_1 = U - A \cdot D_1$$

în care s-a notat cu U matricea unitate.

Următoarele aproximări ale matricii inverse se vor obține astfel:

$$D_2 = D_1 \cdot (U + E_1) \Rightarrow E_2 = U - A \cdot D_2$$

$$D_3 = D_2 \cdot (U + E_2) \Rightarrow E_3 = U - A \cdot D_3$$

.....

$$D_k = D_{k-1} \cdot (U + E_{k-1}) \Rightarrow E_k = U - A \cdot D_k$$

Procesul se întrerupe în momentul în care este îndeplinită relația $N[E_k] < \varepsilon$

în care ε este o valoare impusă inițial.

Concluzie

Calculul inversei unei matrici se va face parcurgând următoarele etape:

- determinarea unei prime aproximări a inversei prin algoritmul de inversare direct, de preferat utilizând pivotarea parțială sau totală;
- corectarea valorii inversei prin parcurgerea algoritmului lui Hotteling.

REZOLVAREA SISTEMELOR DE ECUAȚII LINIARE

Există două metode de rezolvare a sistemelor de ecuații liniare

Metode directe

Constau în obținerea soluției printr-un **algoritm cu un număr finit de pași** (de regulă egal cu numărul de ecuații ale sistemului) prin aducerea matricii de coeficienți ai sistemului la o formă particulară (triunghiulară sau matrice unitate).

În cazul sistemelor slab condiționate metoda se aplică repetat până la obținerea soluției cu o precizie impusă inițial.

Metode iterative (indirecte)

Constau în determinarea soluției într-un **număr de pași care depinde de precizia cu care se efectuează calculul**.

METODE DIRECTE. *Metoda lui Gauss*

Metoda are drept scop rezolvarea sistemului: $A \cdot X = B$

Algoritmul constă în **aducerea matricii de coeficienți A la forma de matrice unitate modificând simultan matricea B , a termenilor liberi.**

Se poate determina o secvență de matrici elementare cu ajutorul căreia, prin **preamplificare**, matricea A să fie adusă la forma de matrice unitate:

$$E_n \cdot E_{n-1} \cdots E_2 \cdot E_1 \cdot A \cdot X = E_n \cdot E_{n-1} \cdots E_2 \cdot E_1 \cdot B$$

rezultă că:

$$X = E_n \cdot E_{n-1} \cdots E_2 \cdot E_1 \cdot B$$

Algoritmul de rezolvare a sistemului va produce următoarea transformare:

$$\left[\begin{array}{cccccc} a_{1,1} & \cdot & \cdot & \cdot & \cdot & a_{1,n} & b_1 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ a_{n,1} & \cdot & \cdot & \cdot & \cdot & a_{n,n} & b_n \end{array} \right] \Rightarrow \left[\begin{array}{cccccc} 1 & \cdot & \cdot & \cdot & \cdot & 0 & b_1^{(n)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & \cdot & 1 & b_n^{(n)} \end{array} \right]$$

Ca urmare **soluția sistemului** va fi: $x_i = b_i^{(n)}$, $i = \overline{1, n}$

Se poate reduce numărul de operații dacă aducem matricea A la forma triunghiulară.

Acest algoritm nu este suficient pentru **sistemele slab condiționate**. Este necesară aplicarea repetată a metodei directe. Succesiunea calculelor se va desfășura în etape, după cum se prezintă în continuare:

Etapa 0

Se rezolvă, cu metoda Gauss, sistemul $A \cdot X = B$ obținându-se soluția: X_1

Etapa 1

Se corectează soluția obținută la etapa precedentă, astfel:

Se înlocuiește soluția în sistem: $A \cdot X_1 - B = E_1$

în care: E_1 este matricea de eroare.

Se rescrie sistemul sub forma: $A \cdot X_1 - E_1 = B$

Prin rezolvarea ecuației (numită *ecuație de corecție*): $A \cdot \delta X_1 = -E_1$

se obține matricea δX_1 care conține corecțiile care trebuie aplicate soluțiilor. Ca urmare ecuația inițială se va putea scrie sub forma:

$$A \cdot X_1 + A \cdot \delta X_1 = B \quad \Rightarrow \quad A \cdot (X_1 + \delta X_1) = B$$

După ***Etapă 1***, soluțiile vor fi formate din soluțiile obținute la ***Etapă 0*** la care se adaugă corecțiile obținute la ***Etapă 1***, prin rezolvarea ***ecuației de corecție***, adică:

$$X = X_1 + \delta X_1$$

Dacă **erorile sunt mai mici** decât o valoare impusă inițial ϵ , atunci procedeul de corectare se încheie, în caz contrar se efectuează o nouă etapă de corecție.

Procedeul se reia până când erorile obținute devin mai mici decât valoarea ϵ impusă inițial. În final soluția sistemului se va obține sub forma:

$$X = X_1 + \delta X_1 + \dots + \delta X_k$$

în care ***k*** reprezintă numărul de etape de corectare efectuate.

METODE ITERATIVE

Aceste metode sunt **convergente** numai pentru sistemele de ecuații a căror matrice de coeficienți este *diagonal dominantă*.

Prin matrice diagonal dominantă se înțelege matricea la care **termenii de pe diagonală au valorile absolute mai mari sau cel mult egale cu suma valorilor absolute a termenilor aflați pe aceeași linie cu ei**. Adică este îndeplinită condiția:

$$\left| a_{i,i} \right| \geq \sum_{\substack{j=1 \\ j \neq i}}^{j=n} \left| a_{i,j} \right| \quad , \quad i = \overline{1, n}$$

În cazul tuturor metodelor iterative algoritmul pornește de la o **soluție a sistemului aleasă în mod arbitrar**, de obicei soluția nulă:

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

Scopul metodelor este acela de a **corecta succesiv soluția inițială până la obținerea soluției reale** a sistemului.

Metoda Jacobi

Aplicarea metodei presupune scrierea sistemului sub o formă în care **din fiecare ecuație se explicitază succesiv câte o necunoscută**. Ca urmare, dacă sistemul inițial este:

$$\sum_{j=1}^n a_{i,j} \cdot x_j = b_i, \quad i = \overline{1, n}$$

acesta se rescrie sub forma:

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

O iterație, k , a metodei Jacobi constă în **determinarea soluției sistemului, corespunzătoare acelei iterații, prin înlocuirea în sistem a soluțiilor determinate la iterația precedentă**.

La prima iterație se utilizează soluția de start. Ca urmare pentru o iterație oarecare k , se poate scrie:

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

Algoritmul se încheie atunci când **diferența dintre două soluții determinate succesiv este mai mică decât o valoare impusă inițial**, adică este îndeplinită condiția:

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

Metoda Gauss-Seidel

Metoda reprezintă o îmbunătățire a metodei Jacobi în sensul că la iterația ***k***, valorile necunoscutelor sunt calculate nu numai funcție de valorile determinate la **iterația precedentă** dar și de cele calculate la **iterația în curs**.

Ca urmare, relația de calcul a valorilor necunoscutelor devine:

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

Metoda Southwell

Aplicarea metodei presupune scrierea sistemului de ecuații $A \cdot X = B$ sub o formă în care **din fiecare ecuație a fost separată o necunoscută**. Ca urmare forma inițială a sistemului:

$$\sum_{j=1}^n a_{i,j} \cdot x_j = b_i, \quad i = \overline{1, n}$$

va suferi următoarele transformări:

- se împarte fiecare ecuație la coeficientul de pe **diagonala principală**:

$$x_i + \sum_{\substack{j=1 \\ j \neq i}}^n \frac{a_{i,j}}{a_{i,i}} \cdot x_j - \frac{b_i}{a_{i,i}} = 0$$

- după care, sistemul se va scrie sub forma:

$$-x_i + \sum_{\substack{j=1 \\ j \neq i}}^n b_{i,j} \cdot x_j + c_i = 0$$

în care s-a notat:

$$b_{i,j} = -\frac{a_{i,j}}{a_{i,i}} \quad c_i = \frac{b_i}{a_{i,i}}$$

Algoritmul constă în parcurgerea mai multor etape plecând de la o **soluție inițială aleasă în mod arbitrar**. De regulă se aleg ca **valori inițiale valorile nule**:

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

Etapa 1:

Se înlocuiesc valorile inițiale în sistem, ca urmare din fiecare ecuație va rezulta o valoare numită **rest**. Valorile acestora se vor obține cu relațiile:

$$r_i^{(1)} = -x_i^{(0)} + \sum_{\substack{j=1 \\ j \neq i}}^n b_{i,j} \cdot x_j^{(0)} + c_i$$

Se determină **valoarea absolută maximă dintre resturile** calculate anterior:

$$\left| r_m^{(1)} \right| = \max_i \left| r_i^{(1)} \right|, \quad i = \overline{1, n}$$

valoarea astfel determinată **va corecta valoarea necunoscutei având indicele corespunzător**

$$x_m^{(1)} = x_m^{(0)} + r_m^{(1)}$$

După ce a fost obținută valoarea $x_m^{(1)}$ **se recalculează resturile utilizând această valoare a necunoscutei**, se obțin astfel **valori noi pentru celelalte $n-1$ resturi**. **Din nou se alege valoarea maximă dintre cele $n-1$ resturi corectându-se necunoscuta corespunzătoare.**

Procedeul se repetă, în același mod, **până la corectarea tuturor necunoscutelor**, după care se trece la etapa următoare.

Etapa k: Se înlocuiesc valorile necunoscutele obținute la etapa precedentă în sistem, se obțin resturile:

$$r_i^{(k)} = -x_i^{(k-1)} + \sum_{\substack{j=1 \\ j \neq i}}^n b_{i,j} \cdot x_j^{(k-1)} + c_i$$

Se determină **restul cu valoarea absolută cea mai mare**:

$$\left| r_m^{(k)} \right| = \max_i \left| r_i^{(k)} \right|, \quad i = \overline{1, n}$$

valoarea astfel determinată va **corecta valoarea necunoscutei având indicele corespunzător**

$$x_m^{(k)} = x_m^{(0)} + r_m^{(1)} + r_m^{(2)} + \dots + r_m^{(k)}$$

După ce a fost obținută valoarea $x_m^{(k)}$ **se recalculează resturile utilizând această valoare a necunoscutei, se obțin astfel valori noi pentru celelalte $n-1$ resturi și se alege valoarea maximă corectându-se necunoscuta corespunzătoare**. Procedul se repetă, în același mod, până la corectarea tuturor necunoscutelelor. Algoritmul se încheie atunci când diferența dintre două soluții determinate succesiv este mai mică decât o valoare impusă inițial, adică este îndeplinită condiția:

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

Exemplu

Se consideră sistemul:

$$\begin{cases} 10 \cdot x_1 - 2 \cdot x_2 = 8 \\ -x_1 + 10 \cdot x_2 = 9 \end{cases}$$

Sistemul are, evident, soluțiile: $x_1 = 1$, $x_2 = 1$.

Rezolvarea sistemului prin metoda Southwell presupune scrierea acestuia, succesiv, în următoarele forme:

$$\begin{cases} x_1 - 0,2 \cdot x_2 = 0,8 \\ -0,1 \cdot x_1 + x_2 = 0,9 \end{cases}$$

Respectiv:

$$\begin{cases} -x_1 + 0,2 \cdot x_2 + 0,8 = 0 \\ -x_2 + 0,1 \cdot x_1 + 0,9 = 0 \end{cases}$$

Se aleg ca valori inițiale ale soluțiilor valorile: $x_1 = 0$, $x_2 = 0$.

Etapa 1.

Se înlocuiesc **valorile inițiale** în sistem,
se obțin următoarele valori pentru resturi:

$$\begin{cases} r_1^{(1)} = 0,8 \\ r_2^{(1)} = 0,9 \end{cases}$$

$$\begin{cases} -x_1 + 0,2 \cdot x_2 + 0,8 = 0 \\ -x_2 + 0,1 \cdot x_1 + 0,9 = 0 \end{cases}$$

Se determină **valoarea absolută maximă** dintre *resturile* calculate anterior, care va fi:

$$r_2^{(1)} = 0,9$$

Valoarea astfel determinată va corecta **valoarea necunoscută având indicele corespunzător**:

$$x_2^{(1)} = x_2^{(0)} + r_2^{(1)} = 0 + 0,9 = 0,9$$

Se **recalculează** restul $r_1^{(1)}$ utilizând această valoare a necunoscutei x_2 , și se obține:

$$r_1^{(1)} = 0 + 0,2 \cdot 0,9 + 0,8 = 0,98$$

Deoarece este singurul rest rămas, nu mai este necesară determinarea unei **valori maxime**. Ca urmare, va fi corectată valoarea necunoscutei corespunzătoare, adică x_1 , se obține:

$$x_1^{(1)} = 0 + 0,98$$

Ca urmare la **încheierea Etapei 1** se obțin următoarele valori ale soluțiilor sistemului:

$$\begin{cases} x_1^{(1)} = 0,98 \\ x_2^{(1)} = 0,9 \end{cases}$$

Se constată că valorile obținute s-au **apropiat semnificativ de soluțiile sistemului**. Se pot parcurge etape suplimentare până la obținerea soluțiilor cu precizia dorită.