

B. Jora C. Popeea S. Barbulea

**METODE DE CALCUL
NUMERIC ÎN
AUTOMATICĂ**

SISTEME LINIARE

**EDITURA ENCICLOPEDICĂ
București 1996**

Cuprins

Cap. 1. Rezolvarea ecuațiilor matriciale liniare

- 1.1 Ecuații matriciale liniare 15
- 1.2 Rezolvarea ecuațiilor matriciale de tip $AX = C$ 17
- 1.3 Rezolvarea ecuațiilor Sylvester 20
- 1.4 Rezolvarea ecuațiilor Liapunov 27
- Exerciții 31
- Bibliografie 34

Cap. 2. Calculul funcțiilor de matrici. Exponențiala matricială

- 2.1 Funcții de matrici 35
- 2.2 Calculul funcțiilor de matrici 46
- 2.3 Calculul exponențialei matriciale 51
- Exerciții 60
- Bibliografie 63

Cap. 3. Tehnici de procesare a modelelor sistemice liniare

- 3.1 Modele sistemice liniare 65
- 3.2 Conexiuni 69
- 3.3 Realizări 74
- 3.4 Conversii de modele 84
- 3.5 Algoritmi de calcul polinomial 88
- Exerciții 97

Bibliografie 100

Cap. 4. Calculul răspunsului în timp al sistemelor liniare

- 4.1 Răspunsul liber al sistemelor liniare 101
- 4.2 Răspunsul la intrări liniar generate 103
- 4.3 Răspunsul la intrări etajate 109
- 4.4 Răspunsul staționar 112
- 4.5 Calculul caracteristicilor de frecvență 115
- 4.6 Răspunsul sistemelor liniare discrete 117
- Exerciții 118
- Bibliografie 122

Cap. 5. Proceduri de analiză sistemică

- 5.1 Stabilitatea sistemelor liniare 123
- 5.2 Descompunerea spectrală 125
- 5.3 Controlabilitate și observabilitate 130
- 5.4 Teste elementare de controlabilitate 134
- 5.5 Forma Hessenberg controlabilă 139
- 5.6 Descompunerea controlabilă 153
- 5.7 Stabilizabilitate și detectabilitate 155
- 5.8 Realizări minimale 157
- 5.9 Gramieni de controlabilitate și observabilitate 158
- 5.10 Echilibrarea sistemelor liniare 162
- Exerciții 168
- Bibliografie 174

Cap. 6. Proceduri de alocare a polilor

- 6.1 Formularea problemei de alocare 175
- 6.2 Proceduri de alocare pentru sisteme cu o singură intrare 184
- 6.3 Proceduri de alocare pentru sisteme cu mai multe intrări 193
- Proceduri de alocare suboptimală 194
- Proceduri de alocare robustă 214

Exerciții 224

Bibliografie 229

Cap. 7. Rezolvarea ecuațiilor matriciale Riccati

7.1 Problema de comandă optimală liniar-patrată 231

Ecuția matricială Riccati (EMR) 238

7.2 Problema de comandă optimală liniar-patrată discretă 245

Ecuția matricială Riccati discretă (DEMR) 249

7.3 Problema liniar-patrată cu orizont infinit 259

Ecuția matricială algebrică Riccati (EMAR) 260

Calculul soluției stabilizatoare a EMAR 261

7.4 Problema liniar-patrată discretă cu orizont infinit 274

Ecuția matricială algebrică Riccati discretă (DEMAR) 275

Calculul soluției stabilizatoare a DEMAR 277

Exerciții 290

Bibliografie 296

Cap. 8. Proceduri de sinteză optimală

8.1 Generalizări ale problemelor de comandă optimală 299

8.2 Problema de estimare optimală 319

Ecuția matricială Riccati de estimare 320

Ecuția matricială Riccati de estimare discretă 326

8.3 Proceduri de sinteză sistemică 335

Exerciții 347

Bibliografie 353

Anexa A. Proceduri de identificare sistemică 357

Bibliografie 370

Cuvânt introductiv

Lucrarea de față reprezintă o primă parte a cursului universitar **Metode Numerice**, destinat studenților anului III al facultății Automatică și Calculatoare din Universitatea POLITEHNICA din București și, evident, poate fi utilă tuturor studenților de la facultățile de același profil din țară. Desigur, ea poate interesa și pe alți utilizatori potențiali (studenți, cadre didactice, cercetători, ingineri, economiști, etc.) ai metodelor moderne de calcul în analiza și proiectarea sistemică.

Cititorul *ideal* al lucrării trebuie să posede cunoștințe de bază privind metodele de calcul numeric matricial – enumerate în finalul acestui cuvânt introductiv – precum și unele informații sigure asupra principalelor noțiuni de teoria sistemelor liniare. De asemenea, el trebuie să aibă obișnuința unei gândiri limpezi, ”algorithmice”, precum și pasiunea lucrului la calculator.

Desigur, cititorul *real*, dispunând într-un grad mai mare sau mai mic de calitățile menționate mai sus, și le poate perfecționa prin studiul lucrării, confirmând, și pe această cale, viabilitatea principiilor sistemice.

Lucrarea este structurată în 8 capitole. Ea este însoțită de o bibliografie de bază (cu lucrări referite prin utilizarea cifrelor romane), utilă în general pentru abordarea tuturor temelor propuse. Bibliografiile specifice selective, dar care pot fi extinse cu ușurință de către însuși cititorul interesat prin cumularea bibliografiilor lucrărilor citate, sunt atașate fiecărui capitol. De asemenea, lucrarea este însoțită de o listă de notații uzuale iar cuprinsul poate constitui, în același timp, indice de noțiuni.

Capitolul 1 completează cunoștințele cititorului privind metodele de rezolvare a sistemelor de ecuații liniare cu cazul matricial în care necunoscutele sunt structurate într-o matrice $X \in \mathcal{R}^{m \times n}$, $m > 1$, $n > 1$, fapt care permite dezvoltarea unor tehnici specifice. În particular, ecuațiile matriciale Liapunov vor fi întâlnite aproape în toate capitolele lucrării.

Capitolul 2 tratează problema importantă a calculului funcțiilor de matrici. Un accent special este pus pe calculul exponențialei matriciale, uti-

lizate intensiv în capitolele 4 și 7.

Capitolul 3 expune principalele tehnici de procesare a modelelor sistemice liniare (conexiuni, realizări, conversii de modele, algoritmi de calcul polinomial) fără de care nu este posibilă abordarea niciunei probleme de analiză sau sinteză sistemică asistate de calculator.

Capitolul 4 descrie metodele numerice de simulare, i.e. de calcul al răspunsului în timp al sistemelor liniare la stimuli externi caracteristici pentru diverse regimuri de funcționare. Aceste tehnici de simulare sunt absolut necesare pentru validarea demersului teoretic și calculatoriu de analiză și sinteză sistemică și reprezintă o punte de legătură esențială către lumea reală a experimentului fizic.

Capitolul 5 prezintă procedurile numerice de analiză a proprietăților sistemice fundamentale (stabilitate, controlabilitate, observabilitate, etc.) precum și principalele metode de construcție a realizărilor minimale și echilibrate. În acest fel conținutul acestui capitol se constituie într-o colecție de algoritmi fundamentali, componente de bază ale oricărei activități de concepție asistate de calculator a unor sisteme de conducere complexe.

Capitolul 6 realizează o trecere în revistă a principalelor proceduri numerice de alocare a polilor sistemelor liniare. În contextul general al tehnicilor de sinteză, alocarea polilor reprezintă un instrument necesar pentru asigurarea unei dinamici dorite sistemelor automate elaborate. Dacă în cazul sistemelor simple reacția după stare ce realizează alocarea este unic determinată și, deci, libertățile de calcul sunt strict procedurale, în cazul sistemelor multiple apar posibilități suplimentare de optimizare sau robustețiere a spectrului alocat.

Capitolul 7 abordează, din punctul de vedere al metodelor numerice de calcul, problema de sinteză liniar-patrată, care vizează optimizarea unui sistem liniar cu indice de calitate patrată. În aplicații problema liniar-patrată generală apare sub forma unor probleme cu finalități sau cu tehnici de tratare specifice. În acest context sunt prezentate metodele de calcul cele mai performante pentru rezolvarea ecuațiilor matriciale Riccati, însoțite de analize comparative ale eficienței și stabilității lor numerice.

Capitolul 8 consideră unele generalizări ale problemei de sinteză liniar-patrată și ilustrează aplicarea algoritmilor de calcul stabiliți la construcția compensatoarelor H_2 și H_∞ (sub)optimale.

Principalele rezultate ale expunerii sunt concretizate sub formă de algoritmi de calcul direct implementabili iar fiecare capitol este însoțit de un set de probleme.

Ne exprimăm opinia că pentru însușirea materialului prezentat este absolut necesară *rezolvarea* problemelor propuse și, mai ales, *implementarea* algoritmilor, urmată de *experimentarea* lor pe exemple numerice concrete și semnificative.

Autorii mulțumesc colegului lor prof. Paul Flondor pentru comentariile constructive făcute pe marginea lucrării. De asemenea, mulțumesc domnului Marcel Popa, directorul Editurii Enciclopedice, pentru atenția plină de solicitudine cu care a urmărit apariția promptă și în cele mai bune condiții a lucrării. Autorii apreciază în mod deosebit interesul manifestat de domnul ing. Sabin Stamatescu, directorul firmei ASTI București, față de conținutul lucrării și posibilele aplicații ale procedurilor de calcul prezentate în conducerea cu calculator a proceselor industriale. În sfârșit, dar, desigur, nu în ultimul rând, mulțumiri se cuvin studenților Siana Petropol, Simona Chirvase, Ștefan Murgu, Cătălin Petrescu, Simona Ruxandu, Adrian Sandu, Laura Stan, Florin Roman, Mădălina Zamfir, Ion Vasile care au realizat redactarea computerizată a lucrării cu mențiunea expresă că fără munca lor plină de dăruire nimic din ce am realizat împreună nu ar fi fost posibil.

*
* *

Enumerăm în continuare principalele probleme de calcul numeric matricial, împreună cu succinte referiri la modul de rezolvare al acestora, pe care cititorul interesat de cuprinsul lucrării de față este bine să le aibă în vedere.

1. Pentru rezolvarea unui sistem de ecuații liniare $Ax = b$, în care matricea $A \in \mathcal{R}^{n \times n}$ este *inversabilă*, se utilizează *procedura de triangularizare prin eliminare gaussiană cu pivotare parțială*

$$A \leftarrow R = MA, \quad b \leftarrow Mb,$$

unde R rezultă superior triunghiulară inversabilă iar M este o secvență de transformări elementare stabilizate, urmată de rezolvarea – prin substituție înapoi – a sistemului superior triunghiular rezultat. Efortul de calcul necesar este $N_{op} \approx n^3/3$.

În cazul *simetric* $A = A^T$, se recomandă conservarea acestei proprietăți prin utilizarea schemei de eliminare simetrizate $A \leftarrow D = MAM^T$, propuse de Parlett [III], în care D rezultă cvasidiagonală.

În cazul *simetric pozitiv definit* $A = A^T > 0$, se impune aplicarea factorizării Cholesky $A = LL^T$.

2. Pentru rezolvarea, în sensul celor mai mici patrate (CMMP), a sistemului supradeterminat $Ax = b$, unde matricea $A \in \mathcal{R}^{m \times n}$ este *monică*,

i.e. are coloanele liniar independente ($\text{rang}A = n$), se utilizează *procedura de triangularizare ortogonală*

$$A \leftarrow R = UA, \quad b \leftarrow Ub,$$

unde R rezultă superior triunghiulară iar U este o secvență de transformări ortogonale elementare (reflectorii sau rotații). Triangularizarea ortogonală este, evident, echivalentă cu factorizarea (sau descompunerea) $A = QR$, unde factorul ortogonal $Q = U^T$ se obține în formă factorizată. Efortul de calcul necesar este $N_{op} \approx mn^2 - n^3/3$ iar eventuala formare explicită a matricii Q , prin acumularea transformărilor elementare parțiale, necesită încă $\approx 2n^3/3$ operații în cazul cel mai defavorabil ($m = n$).

Dacă matricea A nu este neapărat monică, atunci, pentru asigurarea stabilității numerice este necesară utilizarea unei strategii adecvate de *pivotare a coloanelor*, i.e. $A \leftarrow R = UAP$, unde P este o matrice de permutare iar R rezultă superior trapezoidală, vezi [III].

3. Pentru rezolvarea, în sens CMMP, a sistemului subdeterminat $Ax = b$, unde matricea $A \in \mathcal{R}^{m \times n}$ este *epică*, i.e. are liniile liniar independente ($\text{rang}A = m$), se utilizează *procedura de triangularizare ortogonală la dreapta*

$$A \leftarrow L = AV,$$

unde L rezultă inferior triunghiulară. (Această procedură este echivalentă cu procedura precedentă reformulată prin dualitate, i.e. aplicată lui A^T .) Efortul de calcul necesar este $N_{op} \approx nm^2 - m^3/3$.

4. Pentru calculul valorilor proprii ale unei matrici $A \in \mathcal{R}^{n \times n}$ se utilizează *algoritmul QR* care, în esență, construiește iterativ forma Schur S a lui A , i.e.

$$A \leftarrow S = Q^T A Q, \quad (N_{op} \approx 6n^3),$$

unde S rezultă cvasisuperior triunghiulară iar Q este o secvență de transformări ortogonale. În caz de necesitate valorile proprii pot fi ordonate pe diagonala principală a lui S , conform oricărui criteriu de ordonare impus, cu ajutorul unei secvențe suplimentare de transformări ortogonale de asemănare. Calculul matricii Q prin acumularea transformărilor parțiale (necesară, e.g. pentru calculul vectorilor proprii) implică aproape dublarea numărului de operații menționat.

În cazul *simetric* $A = A^T$, se utilizează avantajos versiunea simetrică a algoritmului **QR**, i.e. $A \leftarrow \Lambda = Q^T A Q$, unde Λ rezultă diagonală.

5. Pentru rezolvarea problemei CMMP generale în care matricea $A \in \mathcal{R}^{m \times n}$ este de rang nu neapărat maximal $r \stackrel{\text{def}}{=} \text{rang}A \leq \min(m, n)$, precum

și a altor probleme importante de calcul numeric matricial (determinarea rangului, a pseudoinversei, operații cu subspații liniare, etc.) se utilizează *descompunerea valorilor singulare DVS*

$$U^T AV = \Sigma,$$

unde matricea Σ conține pe diagonala principală valorile singulare nenule $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ ale lui A (restul elementelor fiind nule) iar matricile U, V sunt ortogonale. (Algoritmul **DVS** utilizat pentru calculul descompunerii de mai sus constituie o adaptare ingenioasă a algoritmului **QR** simetric). Din nou, formarea factorilor ortogonali U, V necesită acumularea (relativ costisitoare) a transformărilor parțiale.

6. Pentru calculul valorilor proprii generalizate ale unui fascicol $\lambda B - A$, unde $A, B \in \mathcal{R}^{n \times n}$, se utilizează algoritmul **QZ** care, în esență, aduce perechea (A, B) la *forma Schur generalizată*

$$A \leftarrow S = Q^T AZ, \quad B \leftarrow T = Q^T AZ,$$

unde S rezultă cvasisuperior triunghiulară, T rezultă superior triunghiulară iar Q și Z sunt secvențe de transformări ortogonale.

În cazul *simetric* $A = A^T, B = B^T > 0$ se utilizează factorizarea Cholesky $B = LL^T$ și se aplică algoritmul **QR** simetric matricii $\tilde{A} = L^{-1}AL^{-T}$, ceea ce conduce la diagonalizarea simultană a ambelor matrici A, B .

Pentru detalii privind aspectele teoretice și procedurale ale problemelor de calcul matricial enumerate mai sus recomandăm consultarea referințelor bibliografice [V], [VI] sau [IX], [X].

Bibliografie

- *Pentru programe de calcul și indicații de utilizare:*

- [I] **Smith B.T., Boyle J.M., Ikebe Y., Klema V.C., Moler C.B.** MATRIX EIGENSYSTEM ROUTINES: EISPACK GUIDE, 2-nd ed., Springer-Verlag, New York, 1970.
- [II] **Garbow B.S., Boyle J.M., Dongarra J.J., Moler C.B.** MATRIX EIGENSYSTEM ROUTINES: EISPACK GUIDE EXTENSION, Springer - Verlag, New York, 1972.
- [III] **Dongarra J.J., Bunch J.R., Moler C.B., Stewart G.W.** LINPACK USER'S GUIDE, SIAM Publications, Philadelphia, PA, 1978.
- [IV] **Moler C.B., Little J.N., Bangert S.** PC-MATLAB USER'S GUIDE, The Math Works Inc., 20 N. Main St., Sherborn, Mass., 1987.

- *Pentru algoritmi de calcul matricial:*

- [V] **Stewart G. W.** INTRODUCTION TO MATRIX COMPUTATIONS, Academic Press, New York and London, 1973.
- [VI] **Golub G. H., Van Loan Ch. F.** MATRIX COMPUTATIONS, Second edition, The Johns Hopkins University Press, Baltimore and London, 1988.

- *Pentru chestiuni teoretice de calcul matricial:*

- [VII] **Gantmaher F.R.** TEORIA MATRIȚ (ediția a 2-a), Ed. Nauka, Moscova, 1966. (THE THEORY OF MATRICES, vols. 1-2, Chelsea, New York, 1959).

- *Lucrări în limba română:*

- [VIII] **Ionescu V., Lupas L.** TEHNICI DE CALCUL ÎN TEORIA SISTEMELOR, vol.1. SISTEME LINIARE, vol.2. SISTEME OPTIMALE, E.T., București, 1973.
- [IX] **Bucur C.M., Popeea C.A., Simion Gh.Gh.** MATEMATICI SPECIALE. CALCUL NUMERIC, E.D.P., București, 1983.
- [X] **Ionescu V., Varga A.** TEORIA SISTEMELOR. SINTEZA ROBUSTĂ. METODE NUMERICE DE CALCUL., Ed. ALL, București, 1994.
- [XI] **Iorga V., Jora B., Nicolescu C., Lopătan I., Fătu I.** PROGRAMARE NUMERICĂ, Ed. Teora, București, 1996.

- *Alte titluri de uz general:*

- [XII] **Åström K.J., Wittenmark B.** COMPUTER CONTROLLED SYSTEMS, Prentice Hall Inc., Englewood Cliffs, NJ, 1984.
- [XIII] **Jamshidi M., Herget C.J. (ed.)**, COMPUTER-AIDED CONTROL SYSTEMS ENGINEERING, North-Holland, Amsterdam, 1985. (Maşinostroenie, Moscova, 1989.)
- [XIV] **Solodovnicov V.V. (ed.)**, AVTOMATIZIROVANNOE PROEKTIROVANIE SISTEM UPRAVLENIIA, Maşinostroenie, Moscova, 1990.
- [XV] **Jamshidi M., Tarokh M., Shafai B.** COMPUTER-AIDED ANALYSIS AND DESIGN OF LINEAR CONTROL SYSTEMS, Prentice Hall Inc., Englewood Cliffs, NJ, 1992.

Lista de notații

• Notații generale

\mathcal{N} – mulțimea numerelor naturale.

\mathcal{Z} – mulțimea numerelor întregi.

\mathcal{R} – mulțimea numerelor reale.

\mathcal{C} – mulțimea numerelor complexe.

\mathcal{C}^- – semiplanul stâng deschis al planului variabilei complexe s .

$\mathcal{D}_1(0)$ – discul unitate deschis al planului variabilei complexe z .

\mathcal{R}^n – spațiul liniar n -dimensional al vectorilor (coloană) x cu n componente reale $x_i \in \mathcal{R}$, $i = 1 : n$.

\mathcal{C}^n – spațiul liniar n -dimensional al vectorilor (coloană) x cu n componente complexe $x_i \in \mathcal{C}$, $i = 1 : n$.

e_k , $k = 1 : n$ – baza standard a spațiului liniar \mathcal{R}^n , respectiv \mathcal{C}^n .

$\mathcal{R}^{m \times n}$ – spațiul liniar al matricilor cu m linii și n coloane cu elemente reale $a_{ij} \in \mathcal{R}$, $i = 1 : m$, $j = 1 : n$.

$\mathcal{C}^{m \times n}$ – spațiul liniar al matricilor cu m linii și n coloane cu elemente complexe $a_{ij} \in \mathcal{C}$, $i = 1 : m$, $j = 1 : n$.¹

¹În calcule, vectorii se identifică cu matricile cu o singură coloană iar scalarii se identifică cu matricile (sau vectorii) cu un singur element.

A^T – transpusa matricii (reale sau complexe) A .

A^H – conjugata hermitică a matricii (complexe) A , i.e. $A^H = \bar{A}^T$, unde \bar{A} este conjugata complexă a lui A .

A^+ – pseudoinversa normală (Moore-Penrose) a matricii A ; dacă A este monică $A^+ = (A^T A)^{-1} A^T$, dacă A este epică atunci $A^+ = A^T (A A^T)^{-1}$.

$\sigma_i(A)$, $i = 1 : p$, $p = \min(m, n)$ – valorile singulare ale matricii A ordonate astfel încât $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$.

$\sigma(A)$ – mulțimea $\{\sigma_1(A), \sigma_2(A), \dots, \sigma_p(A)\}$ a valorilor singulare ale matricii A .

$r = \text{rang} A$ – rangul matricii A , i.e. numărul valorilor singulare nenule.

I_n – matricea unitate de ordinul n .

A^{-1} – inversa matricii patrate nesingulare A , i.e. $A A^{-1} = A^{-1} A = I_n$.

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

$$A^{-H} = (A^{-1})^H = (A^H)^{-1}$$

$\text{tr} A$ – urma matricii patrate A , i.e. suma elementelor diagonale.

$\det A$ – determinantul matricii patrate A .

$\lambda_i(A)$, $i = 1 : n$ – valorile proprii ale matricii patrate A de ordin n .

$\lambda(A)$ – spectrul (de valori proprii) $\{\lambda_1(A), \lambda_2(A), \dots, \lambda_n(A)\}$ al matricii A .

$\rho(A) = \max_{i=1:n} |\lambda_i(A)|$ – raza spectrală a matricii A .

$\text{cond}(A) = \|A\| \|A^{-1}\|$ – numărul de condiție la inversare al matricii A ($\|\cdot\|$ este o normă matricială consistentă, vezi mai jos).

$(x, y) = y^T x$ – produsul scalar standard a doi vectori reali; în cazul complex produsul scalar este $(x, y) = y^H x$.

$\|x\| = (x, x)^{1/2}$ – norma euclidiană a vectorului x ; se notează $\|x\|_Q^2 = x^T Q x$ unde Q este o matrice simetrică ($Q = Q^T$).

$\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ – p -normele vectorului n -dimensional x , $p \geq 1$; în calcule se utilizează în special $\|x\|_1, \|x\|_2 = \|x\|$ și $\|x\|_\infty = \max_{i=1:n} |x_i|$.

$(A, B) = \text{tr}(B^T A)$ ($\text{tr}(B^H A)$) – produsul scalar a două matrici reale (complexe).

$\|A\|_F = (A, B)^{1/2}$ – norma Frobenius a matricii A ,
 $\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2$ sau $\|A\|_F^2 = \sum_{i=1}^r \sigma_i^2$.

$|A|_p = (\sum_{i=1}^r \sigma_i^p)^{1/p}$ – p -normele Schatten, $p \geq 1$; în calcule se utilizează în special norma-urmă $|A|_1 = \sum_{i=1}^r \sigma_i$, norma Frobenius $|A|_2 = \|A\|_F$ și norma spectrală $|A|_\infty = \sigma_1(A)$.

$\|A\|_p = \max_{\|x\|_p=1} \|Ax\|_p$ – p -normele induse; în calcule se utilizează în special norma $\|A\|_1 = \max_{j=1:n} \sum_{i=1}^m |a_{ij}|$, norma spectrală $\|A\|_2 = \sigma_1(A)$ și norma $\|A\|_\infty = \max_{i=1:m} \sum_{j=1}^n |a_{ij}|$.

• Notății sistemice generale

$x \in \mathcal{R}^n$ – vectorul (mărimilor) de stare.

$u \in \mathcal{R}^m$ – vectorul (mărimilor) de intrare.

$y \in \mathcal{R}^l$ – vectorul (mărimilor) de ieșire.

$S = (A, B, C, D)$ – reprezentarea de stare a unui sistem liniar ($A \in \mathcal{R}^{n \times n}$, $B \in \mathcal{R}^{n \times m}$, $C \in \mathcal{R}^{l \times n}$ și $D \in \mathcal{R}^{l \times m}$).

$T = (N, p)$ – reprezentarea de transfer a unui sistem liniar (N este matricea coeficienților polinoamelor ce definesc numărătorul matricii de transfer iar p este vectorul coeficienților polinomului numitor comun i.e. $T(s) = C(sI - A)^{-1}B + D = N(s)/p(s)$).

$\delta(t)$ ($\delta(k)$) – impulsul unitate continuu (discret).

$1(t)$ ($1(k)$) – funcția treaptă unitară continuă (discretă).

• Transformări

(1) MAN (MAN^{-1} sau MAN^T) – transformare de echivalență (bilaterală) a matricii $A \in \mathcal{R}^{m \times n}$ (M și N sunt matrici patrute nesingulare; transformarea de echivalență conservă rangul iar dacă M , N sunt ortogonale atunci conservă și valorile singulare).

(2) NAN^{-1} – transformare de asemănare a matricii $A \in \mathcal{R}^{n \times n}$ (transformarea de asemănare conservă valorile proprii).

- (3) NAN^T – transformare de *congruență* a matricii $A \in \mathcal{R}^{n \times n}$ (N este nesingulară; aplicată unei matrici A simetrice, transformarea de congruență conservă rangul și inerția i.e. numerele de valori proprii negative, nule și, respectiv, pozitive).

Dacă N este ortogonală atunci transformările (2) și (3) coincid și definesc transformarea de *asemănare ortogonală*.

• Prescurtări

SISO – siglă pentru sisteme simple având o singură intrare și o singură ieșire (**S**ingle-**I**nput **S**ingle-**O**utput).

SIMO, (MISO, MIMO) – siglă pentru sistemele cu o singură intrare și mai multe ieșiri. Celelalte prescurtări au semnificații evidente.

FSR(G) – forma Schur reală (generalizată).

FSC(G) – forma Schur complexă (generalizată).

DVS – descompunerea valorilor singulare.

FSH – forma (bloc-)superior Hessenberg.

EM(A)L – ecuație matricială (algebrică) Liapunov continuă.

DEM(A)L – ecuație matricială (algebrică) Liapunov discretă.

EM(A)R – ecuație matricială (algebrică) Riccati continuă.

DEM(A)R – ecuație matricială (algebrică) Riccati discretă.

Capitolul 1

Rezolvarea ecuațiilor matriceale liniare

Acest capitol este consacrat tehnicilor de rezolvare a unor sisteme liniare, în general de mari dimensiuni, structurate în exprimări matriceale care permit dezvoltarea unor metode specifice de calcul.

1.1 Ecuații matriceale liniare

Ecuațiile matriceale liniare sunt sisteme de ecuații liniare care pot fi scrise compact într-o formă matriceală de tipul

$$f(X) = C \quad (1.1)$$

unde $X \in \mathcal{R}^{m \times n}$ este matricea necunoscutelor, $C \in \mathcal{R}^{p \times q}$ o matrice dată, iar $f : \mathcal{R}^{m \times n} \rightarrow \mathcal{R}^{p \times q}$ o aplicație liniară (i.e. f satisface $f(X_1 + X_2) = f(X_1) + f(X_2)$ și $f(\alpha X_1) = \alpha f(X_1)$ pentru orice matrice $X_1, X_2 \in \mathcal{R}^{m \times n}$ și orice scalar α). Se poate arăta că orice aplicație liniară de argument matriceal poate fi scrisă sub forma

$$f(X) = \sum_{i=1}^k A_i X B_i, \quad A_i \in \mathcal{R}^{p \times m}, B_i \in \mathcal{R}^{n \times q} \quad (1.2)$$

pentru un anumit k și, în consecință, (1.1) devine

$$\sum_{i=1}^k A_i X B_i = C. \quad (1.3)$$

Evident, ecuația (1.3) poate fi scrisă într-o formă ”desfășurată” ca sistem de pq ecuații liniare cu mn necunoscute care poate servi ca bază pentru exprimarea condițiilor de existență/unicitate ca și pentru elaborarea procedurilor de rezolvare într-un sens sau altul. O astfel de abordare eludează ”structura matriceală” a sistemului iar aplicarea tehnicilor clasice de rezolvare pe sistemul desfășurat, neexploatând structura internă a datelor de intrare, este, de cele mai multe ori, ineficientă (vezi exercițiul 1.9).

Din acest motiv, pentru diferite cazuri particulare ale ecuației matriceale (1.3), condițiile de existență/unicitate ale soluțiilor, într-un sens precizat, se exprimă în raport cu matricele $A_i, B_i, i = 1 : k$, iar metodele de rezolvare fac apel la tehnici speciale.

Ecuatiile matriceale liniare (1.3) cele mai întâlnite se obțin pentru $k = 1$

$$AXB = C \quad (1.4)$$

a cărei rezolvare se reduce imediat la rezolvarea unor ecuații matriceale având forma particulară

$$AX = C \quad (1.5)$$

$$XB = C \quad (1.6)$$

respectiv, pentru $k = 2$,

$$A_1XB_1 + A_2XB_2 = C \quad (1.7)$$

cu particularizările

$$AX + XB = C, \quad (1.8)$$

$$AXB + X = C. \quad (1.9)$$

Ecuatiile (1.8), (1.9) sunt cunoscute sub denumirea de *ecuații matriceale tip Sylvester*. Datorită frecvenței utilizării a ecuațiilor matriceale de forma (1.8) în teoria sistemelor dinamice continue, respectiv a ecuațiilor de forma (1.9) în domeniul sistemelor dinamice discrete, în continuare ecuația (1.8) va fi referită ca *ecuație Sylvester continuă*, respectiv (1.9) ca *ecuație Sylvester discretă*.

În sfârșit, considerând în (1.8) $B \leftarrow A, A \leftarrow A^T$ și în (1.9) $B \leftarrow A, A \leftarrow -A^T, C \leftarrow -C$, obținem ecuațiile matriceale liniare cunoscute sub denumirile de *ecuație Liapunov continuă* pentru

$$A^T X + XA = C \quad (1.10)$$

respectiv *ecuație Liapunov discretă* pentru

$$A^T XA - X = C. \quad (1.11)$$

În cele mai multe aplicații matricele de intrare A , B , C sunt reale, matricea soluție X rezultând, la rândul său, reală. Totuși, algoritmi prezentați rămân aplicabili și în cazul matricelor de intrare complexe. Mai mult, chiar în cazul unor date inițiale reale, unele din metode de calcul, apelând la determinarea valorilor și vectorilor proprii, necesită, tranzitoriu, utilizarea unei aritmetici complexe. Atunci când pentru date de intrare reale se poate utiliza exclusiv o aritmetică reală (vezi exercițiul 1.14) vor fi făcute precizările cuvenite.

În cele ce urmează vom fi interesați în calculul soluțiilor ecuațiilor matriceale liniare cu semnificație pentru domeniul teoriei sistemelor în general, și al sistemelor automate în special, respectiv al sistemelor de forma (1.5), (1.6), (1.8), (1.9) și (1.10), (1.11). Condițiile de existență și unicitate ale soluțiilor se exprimă în mod specific pentru fiecare tip de ecuație, singurul rezultat comun fiind dat de

Propoziția 1.1 *Dacă ecuația matriceală liniară (1.3) admite o soluție $X \in \mathcal{R}^{m \times n}$, atunci această soluție este unică dacă și numai dacă ecuația matriceală omogenă (obținută din ecuația inițială pentru $C = 0$) admite drept unică soluție $X = 0$.*

Demonstrația este propusă ca exercițiu pentru cititor sau poate fi găsită în [VII]. \diamond

În cazul ecuațiilor matriceale liniare (1.8), (1.9) și (1.10), (1.11) o analiză simplă arată că matricea necunoscută X și membrul drept C au aceleași dimensiuni ($m \times n$). În consecință, enunțul de mai sus capătă următoarea formă mai precisă.

Propoziția 1.2 *Ecuațiile matriceale liniare (1.8), (1.9) și (1.10), (1.11) sunt global solubile, i.e. admit o soluție $X \in \mathcal{R}^{m \times n}$ oricare ar fi membrul drept C , dacă și numai dacă ecuațiile matriceale omogene corespunzătoare admit drept unică soluție $X = 0$.*

Altfel spus, are loc următoarea alternativă (a lui Fredholm): fie ecuațiile considerate sunt global solubile, fie ecuațiile omogene admit o soluție netrivială $X \neq 0$.

1.2 Rezolvarea ecuațiilor matriceale de tip $AX = C$

Ecuațiile matriceale de tipul $AX = C$, unde $A \in \mathcal{R}^{p \times m}$, $X \in \mathcal{R}^{m \times n}$, $C \in \mathcal{R}^{p \times n}$, reprezintă o colecție de ecuații vectoriale obișnuite, toate având

aceeași matrice A a coeficienților. Într-adevăr, partiționând matricele X și C pe coloane

$$X = [x_1 \ x_2 \ \dots \ x_n], \quad C = [c_1 \ c_2 \ \dots \ c_n] \quad (1.12)$$

cu $x_j = Xe_j$, $c_j = Ce_j$, $j = 1 : n$, ecuația matriceală $AX = C$ devine echivalentă cu setul de sisteme liniare

$$Ax_j = c_j, \quad j = 1 : n \quad (1.13)$$

care se rezolvă cu mijloacele clasice, cu precizarea că transformările necesare asupra matricei A a sistemelor (1.13) (cum ar fi, de exemplu, în cazul $p = m$ și A inversabilă, triangularizarea prin eliminare gaussiană cu pivotare parțială) se efectuează o singură dată. Având în vedere că rezolvarea sistemelor (1.13) nu ridică probleme deosebite, scrierea algoritmilor face obiectul exercițiilor 1.1 și 1.3.

Observația 1.1 Ecuațiile de tipul $XB = C$ se reduc la ecuații de tipul $AX = C$ prin transpunere. \diamond

Observația 1.2 Dacă matricea $A \in \mathcal{R}^{p \times m}$ este *monică* (i.e. are coloanele liniar independente) iar sistemele (1.13) sunt rezolvate în sensul celor mai mici pătrate (CMMP), atunci ansamblul *pseudosoluțiilor*

$$x_j^* = A^+c_j = (A^T A)^{-1} A^T c_j, \quad j = 1 : n \quad (1.14)$$

respectiv

$$X^* = A^+C, \quad A^+ = (A^T A)^{-1} A^T \quad (1.15)$$

are o semnificație asemănătoare și anume aceea că matricea $X^* \in \mathcal{R}^{m \times n}$ minimizează norma Frobenius a reziduuului matriceal

$$R = C - AX \quad (1.16)$$

adică

$$\|R^*\|_F = \|C - AX^*\|_F = \min_{X \in \mathcal{R}^{m \times n}} \|C - AX\|_F. \quad (1.17)$$

Într-adevăr,

$$\|R\|_F = \sqrt{\sum_{j=1}^n \sum_{i=1}^p r_{ij}^2} = \sqrt{\sum_{j=1}^n \|c_j - Ax_j\|^2}$$

este minimă dacă și numai dacă $\|c_j - Ax_j\|$ sunt minime pentru toți $j = 1 : n$. \diamond

Observația 1.3 Dacă matricea $A \in \mathcal{R}^{p \times m}$ este *epică* (adică are liniile liniar independente), atunci ansamblul soluțiilor *normale* ale sistemelor (1.13)

$$x_j^* = A^+ c_j = A^T (AA^T)^{-1} c_j, \quad j = 1 : n \quad (1.18)$$

respectiv

$$X^* = A^+ C, \quad A^+ = A^T (AA^T)^{-1} \quad (1.19)$$

reprezintă soluția de normă Frobenius minimă a sistemului $AX = C$:

$$\|X^*\|_F = \min_{\substack{AX=C \\ X \in \mathcal{R}^{m \times n}}} \|X\|_F. \quad (1.20)$$

Într-adevăr, este evident că X^* dat de (1.19) este o soluție a ecuației matriceale $AX = C$. În plus, dacă $A = LQ$ este o factorizare a matricei A , unde $Q \in \mathcal{R}^{m \times m}$ este ortogonală iar $L \in \mathcal{R}^{p \times m}$ are structura

$$L = [L' \ 0] \quad (1.21)$$

cu $L' \in \mathcal{R}^{p \times p}$ inferior triunghiulară nesingulară, atunci ecuația matriceală $AX = C$ se scrie

$$LQX = C. \quad (1.22)$$

Cu notația

$$QX = Y = \begin{bmatrix} Y' \\ Y'' \end{bmatrix}, \quad (1.23)$$

cu $Y' \in \mathcal{R}^{p \times n}$, $Y'' \in \mathcal{R}^{(m-p) \times n}$, având în vedere structura (1.21) a matricei L , rezultă

$$Y' = (L')^{-1} C.$$

În consecință, toate soluțiile sistemului $AX = C$ se scriu sub forma

$$X = Q^T \begin{bmatrix} (L')^{-1} C \\ Y'' \end{bmatrix} \quad (1.24)$$

cu $Y'' \in \mathcal{R}^{(m-p) \times n}$ arbitrar. Dar

$$\begin{aligned} \|X\|_F^2 &= \|Q^T \begin{bmatrix} (L')^{-1} C \\ Y'' \end{bmatrix}\|_F^2 = \left\| \begin{bmatrix} (L')^{-1} C \\ Y'' \end{bmatrix} \right\|_F^2 = \\ &= \|(L')^{-1} C\|_F^2 + \|Y''\|_F^2 \end{aligned}$$

minimul obținându-se pentru $Y'' = 0$. Prin urmare soluția de normă Frobenius minimă este

$$X^* = Q^T \begin{bmatrix} (L')^{-1}C \\ 0 \end{bmatrix} \quad (1.25)$$

verificându-se imediat, pe baza factorizării LQ , că aceasta coincide cu soluția (1.19). \diamond

Observația 1.4 În cazul în care matricea A este de rang nemaximal

$$\text{rang } A < \min(p, m),$$

notând cu x_j^* , $j = 1 : n$, *pseudosoluțiile normale*, în sensul celor mai mici pătrate, ale sistemelor (1.13), matricea $X^* = (x_j^*)_{j=1:n}$ poate fi interpretată drept pseudosoluția normală în sensul normei Frobenius a ecuației matriceale $AX = C$. Cu alte cuvinte $X^* \in \mathcal{R}^{m \times n}$ este matricea de normă Frobenius minimă dintre toate matricele $X \in \mathcal{R}^{m \times n}$ care minimizează norma Frobenius a reziduului matriceal $R = C - AX$:

$$\|X^*\|_F = \min_{X \in \mathcal{R}^{m \times n}} \|X\|_F. \quad (1.26)$$

$\|C - AX\|_F = \text{minim}$

Justificarea afirmațiilor de mai sus face obiectul exercițiului 1.2. \diamond

1.3 Rezolvarea ecuațiilor matriceale Sylvester

I. Vom începe cu prezentarea modalităților de rezolvare a *ecuației matriceale Sylvester continue* (1.8)

$$AX + XB = C$$

unde în general $A \in \mathcal{C}^{m \times m}$, $B \in \mathcal{C}^{n \times n}$, $C \in \mathcal{C}^{m \times n}$. Condițiile de existență și unicitate ale soluției acestei ecuații sunt date de

Propoziția 1.3 *Ecuația matriceală (1.8) admite o soluție $X \in \mathcal{C}^{m \times n}$ unică dacă și numai dacă*

$$\lambda_i + \mu_j \neq 0 \quad (1.27)$$

oricare ar fi $\lambda_i \in \lambda(A)$ și oricare ar fi $\mu_j \in \lambda(B)$ sau, altfel spus, dacă și numai dacă

$$\lambda(A) \cap \lambda(-B) = \emptyset. \quad (1.28)$$

Dacă matricele A , B , C sunt reale iar condițiile (1.27) sunt satisfăcute atunci soluția X rezultă și ea reală.

Demonstrație. Există matricele unitare $U \in \mathcal{C}^{m \times m}$ și $V \in \mathcal{C}^{n \times n}$ (i.e. $U^H U = U U^H = I_m$, $V^H V = V V^H = I_n$, indicele superior H semnificând dubla operație de transpunere și conjugare) astfel încât

$$U^H A U = S \in \mathcal{C}^{m \times m}, \quad (1.29)$$

$$V^H B V = T \in \mathcal{C}^{n \times n} \quad (1.30)$$

sunt *formele Schur complexe* ale matricelor A respectiv B și, deci, au o structură superior triunghiulară. Din (1.29), (1.30) avem relațiile

$$A = U S U^H \quad (1.31)$$

$$B = V T V^H \quad (1.32)$$

cu care ecuația (1.8) devine

$$U S U^H X + X V T V^H = C$$

de unde

$$S U^H X V + U^H X V T = U^H C V. \quad (1.33)$$

Notând cu

$$Y = U^H X V, \quad \tilde{C} = U^H C V \quad (1.34)$$

(1.33) se scrie

$$S Y + Y T = \tilde{C}. \quad (1.35)$$

Având în vedere natura nesingulară a transformărilor efectuate, ecuația matriceală (1.8) admite o soluție unică $X \in \mathcal{C}^{m \times n}$ dacă și numai dacă ecuația (1.35) admite o soluție unică $Y \in \mathcal{C}^{m \times n}$. Condițiile de existență și unicitate ale soluției ecuației (1.35) rezultă din aplicarea unei tehnici de rezolvare directă. Într-adevăr, ținând seama de structura superior triunghiulară a matricelor S și T rezultă că, scriind pe coloane, (1.35) devine

$$S y_j + Y t_j = \tilde{c}_j, \quad j = 1 : n \quad (1.36)$$

unde $y_j = Y e_j$, $t_j = T e_j$, $\tilde{c}_j = C e_j$. Dar

$$t_j = [t_{1j} \ t_{2j} \ \dots \ t_{jj} \ 0 \ \dots \ 0]^T$$

și, prin urmare, (1.36) devine

$$S y_j + \sum_{k=1}^j t_{kj} y_k = \tilde{c}_j, \quad j = 1 : n \quad (1.37)$$

sau

$$(S + t_{jj}I_m)y_j = \tilde{c}_j - \sum_{k=1}^{j-1} t_{kj}y_k, \quad j = 1 : n. \quad (1.38)$$

Se vede acum limpede că ecuațiile (1.38) sunt sisteme liniare triunghiulare având membrul drept calculabil dacă ordinea de rezolvare a acestor sisteme este $j = 1, 2, 3, \dots, n$. Mai mult, existența și unicitatea soluției este condiționată de nesingularitatea matricelor $(S + t_{jj}I_m)$, $j = 1 : n$. Matricea S fiind superior triunghiulară, această condiție este satisfăcută dacă și numai dacă

$$s_{ii} + t_{jj} \neq 0, \quad \forall i = 1 : m, \quad \forall j = 1 : n,$$

condiție identică cu (1.27) întrucât $\lambda(A) = \lambda(S)$ și $\lambda(B) = \lambda(T)$, observație care încheie demonstrația propoziției. \diamond

Tehnicile de rezolvare ale ecuației Sylvester continue fac apel la transformări care conduc la ecuații matriceale în care, în locul matricelor A și B , apar matrice cu o structură simplificată (Hessenberg, triunghiulară, diagonală) fapt care permite aplicarea metodelor consacrate de la sistemele clasice. O primă alternativă, numită ad-hoc *variante Schur-Schur*, are la bază procedura evidențiată de demonstrația propoziției 1.3. În cazul uzual, în care datele de intrare precum și soluția X sunt reale, redactarea algoritmului de calcul este următoarea.

Algoritmul 1.1 (Date matricele $A \in \mathcal{R}^{m \times m}$, $B \in \mathcal{R}^{n \times n}$, $C \in \mathcal{R}^{m \times n}$ cu $\lambda(A) \cap \lambda(-B) = \emptyset$, algoritmul calculează soluția $X \in \mathcal{R}^{m \times n}$ a ecuației Sylvester continue $AX + XB = C$ prin aducerea matricelor A și B la forma Schur complexă prin transformări unitare de asemănare. Algoritmul utilizează funcția **fsc** care calculează forma Schur complexă și matricea unitară de transformare corespunzătoare pentru o matrice de intrare dată.)

1. $[U, S] = \mathbf{fsc}(A)$
2. $[V, T] = \mathbf{fsc}(B)$
3. $C \leftarrow \tilde{C} = U^H C V$
4. Pentru $j = 1 : n$
 1. Dacă $j > 1$ atunci
 1. $c_j = c_j - \sum_{k=1}^{j-1} t_{kj}y_k$.
 2. Se rezolvă sistemul superior triunghiular $(S + t_{jj}I_m)y_j = c_j$.

$$5. X = \operatorname{Re}(UYV^H)$$

Comentarii. Referitor la algoritmul de mai sus se cuvin următoarele precizări:

1. Aducerea unei matrice reale la forma Schur complexă prin transformări de asemănare se face în două etape: într-o primă etapă cu algoritmul **QR** se obține forma Schur reală după care, într-o a doua etapă, utilizând rotații complexe sau reflectori complecși (exercițiile 1.4 : 1.7) blocurile 2×2 sunt reduse la forma superior triunghiulară complexă.

2. Efectuarea calculelor în format virgulă mobilă cu numere complexe conduce la un rezultat afectat de erori de rotunjire complexe deși rezultatul teoretic este o matrice reală. De aceea, în ultima instrucțiune, este inclusă eliminarea părților imaginare.

3. Evident, efortul de calcul cel mai important se consumă în execuția instrucțiunilor 1 și 2 de aducere la forma Schur complexă a matricelor A și B și de acumulare a transformărilor. \diamond

Din motive de eficiență se impune analiza alternativelor în care se renunță la aducerea la forma Schur a ambelor matrice A și B .

- Astfel, în așa numita *variantă Hessenberg-Schur* numai matricea B este adusă la forma Schur complexă apărând următoarele diferențe în raport cu algoritmul de mai sus:

- În instrucțiunea 1 matricea A este adusă printr-un algoritm de calcul direct (neiterativ) – algoritmul **HQ** – la forma superior Hessenberg. În acest fel se evită faza iterativă a algoritmului **QR** și trecerea de la forma Schur reală la forma Schur complexă.

- În compensație, la instrucțiunea 4.2, se rezolvă, în loc de n sisteme triunghiulare, n sisteme de tip Hessenberg incluzând eliminare gaussiană cu eventuală pivotare.

Scrierea explicită a algoritmului face obiectul exercițiului 1.10.

- În cazul matricelor A și B *simple* (adică al matricelor pentru care există un set complet de vectori proprii liniar independenți, în particular, al matricelor cu valori proprii distincte) metoda bazată pe calculul valorilor și vectorilor proprii poate conduce la o precizie superioară a rezultatelor simultan cu un timp de execuție rezonabil. Fie

$$\begin{aligned} \lambda(A) &= \{\lambda_1, \lambda_2, \dots, \lambda_m\} \subset \mathcal{C} \\ \lambda(B) &= \{\mu_1, \mu_2, \dots, \mu_n\} \subset \mathcal{C} \end{aligned} \tag{1.39}$$

spectrele de valori proprii ale matricelor A respectiv B și

$$\begin{aligned} U &= [u_1 \ u_2 \ \dots \ u_m], \\ V &= [v_1 \ v_2 \ \dots \ v_n] \end{aligned} \quad (1.40)$$

matricele corespunzătoare de vectori proprii (adică $u_j = Ue_j$ este un vector propriu al matricei A asociat valorii proprii $\lambda_j \in \lambda(A)$, etc.). Atunci

$$\begin{aligned} A &= ULU^{-1}, & L &= \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m) \\ B &= VMV^{-1}, & M &= \text{diag}(\mu_1, \mu_2, \dots, \mu_n) \end{aligned} \quad (1.41)$$

și ecuația Sylvester (1.8) se scrie sub forma

$$ULU^{-1}X + XVMV^{-1} = C \quad (1.42)$$

echivalentă cu

$$LU^{-1}XV + U^{-1}XVM = U^{-1}CV. \quad (1.43)$$

Notând

$$\begin{aligned} U^{-1}XV &= Y \\ U^{-1}CV &= Z \end{aligned} \quad (1.44)$$

din (1.43) rezultă ecuația

$$LY + YM = Z \quad (1.45)$$

cu necunoscuta matriceală Y ale cărei elemente, în condițiile propoziției 1.1, se calculează imediat cu relațiile

$$y_{ij} = z_{ij}/(\lambda_i + \mu_j), \quad i = 1 : m, \ j = 1 : n. \quad (1.46)$$

Matricea necunoscută X din ecuația inițială rezultă din (1.44)

$$X = UYV^{-1}. \quad (1.47)$$

Presupunând că dispunem de o funcție numită **vvp** care, pentru o matrice de intrare simplă T , furnizează o matrice V ale cărei coloane sunt vectorii proprii ai matricei T și un vector p al valorilor proprii corespunzătorii cu sintaxa

$$[V, p] = \mathbf{vvp}(T) \quad (1.48)$$

putem prezenta următorul algoritm.

Algoritm 1.2 (Date matricele simple $A \in \mathcal{R}^{m \times m}$, $B \in \mathcal{R}^{n \times n}$ cu $\lambda(A) \cap \lambda(-B) = \emptyset$ și matricea $C \in \mathcal{R}^{m \times n}$, algoritmul calculează soluția $X \in \mathcal{R}^{m \times n}$ a ecuației matriceale $AX + XB = C$ pe baza calculului valorilor și vectorilor proprii cu ajutorul funcției **vvp**.)

1. $[U, l] = \mathbf{vvp}(A)$
2. $[V, m] = \mathbf{vvp}(B)$
3. Se rezolvă sistemul matriceal $UZ = CV$ cu necunoscuta Z .
4. Pentru $i = 1 : m$
 1. Pentru $j = 1 : n$
 1. $y_{ij} = z_{ij}/(l_i + m_j)$.
5. Se rezolvă sistemul matriceal $XV = UY$ în raport cu X .
6. $X = \text{Re}(X)$.

Rezolvarea sistemelor matriceale de la instrucțiunile 3 și 5 ale algoritmului se realizează conform recomandărilor secțiunea 1.2.

II. Modalitățile de calcul ale soluției *ecuației matriceale Sylvester discrete* (1.9)

$$AXB + X = C$$

unde $A \in \mathcal{C}^{m \times m}$, $B \in \mathcal{C}^{n \times n}$, $C \in \mathcal{C}^{m \times n}$ sunt asemănătoare. Condițiile de existență și unicitate ale soluției sunt date de

Propoziția 1.4 *Ecuația matriceală (1.9) are o soluție $X \in \mathcal{C}^{m \times n}$ unică dacă și numai dacă*

$$\lambda_i \mu_j \neq -1 \quad (1.49)$$

oricare ar fi $\lambda_i \in \lambda(A)$ și oricare ar fi $\mu_j \in \lambda(B)$. Dacă matricele A, B, C sunt reale iar condițiile (1.49) sunt satisfăcute atunci soluția X rezultă și ea reală.

Demonstrație. Fie $U \in \mathcal{C}^{m \times m}$ și $V \in \mathcal{C}^{n \times n}$ matrice unitare care definesc reducerea matricelor A și B la *formele Schur complexe* S , respectiv T , ca în (1.29) și (1.30). Cu (1.31) și (1.32), ecuația (1.9) devine

$$USU^H XVTV^H + X = C \quad (1.50)$$

care, la rândul ei, este echivalentă cu

$$SYT + Y = \tilde{C} \quad (1.51)$$

unde Y și \tilde{C} sunt date de (1.34). Ecuația (1.9) are o soluție unică X dacă și numai dacă ecuația (1.51) are o soluție unică Y . Dacă $y_j = Ye_j$, $t_j =$

Te_j , $\tilde{c}_j = \tilde{C}e_j$, $j = 1 : n$ atunci, ținând seama de structura superior triunghiulară a matricelor T și S , din (1.51) obținem

$$\begin{aligned} (St_{11} + I_m)y_1 &= \tilde{c}_1, \\ (St_{jj} + I_m)y_j &= \tilde{c}_j - S \sum_{k=1}^{j-1} t_{kj}y_k, \quad j = 2 : n. \end{aligned} \tag{1.52}$$

Ecuatiile (1.52) sunt sisteme liniare superior triunghiulare care admit soluții unice dacă și numai dacă matricele $St_{jj} + I_m$, $j = 1 : n$ sunt nesingulare, respectiv

$$s_{ii}t_{jj} + 1 \neq 0, \quad \forall i \in 1 : m, \quad \forall j \in 1 : n. \tag{1.53}$$

Cum însă $\{s_{ii} \mid i = 1 : m\} = \lambda(A)$ și $\{t_{jj} \mid j = 1 : n\} = \lambda(B)$ și, în plus, rezolvarea sistemelor (1.52) în condițiile (1.53) este posibilă în ordinea $j = 1, 2, \dots, n$ în care termenul liber este calculabil, rezultă că propoziția este demonstrată. \diamond

Demonstrația de mai sus conduce imediat la algoritmul de rezolvare a ecuației matriceale Sylvester discrete în *variantea Schur-Schur*.

Algoritmul 1.3 (Fiind date matricele $A \in \mathcal{R}^{m \times m}$, $B \in \mathcal{R}^{n \times n}$ astfel încât $\lambda_i \mu_j \neq -1$ pentru toți $\lambda_i \in \lambda(A)$ și $\mu_j \in \lambda(B)$ și matricea $C \in \mathcal{R}^{m \times n}$ algoritmul calculează soluția X a ecuației Sylvester discrete $AXB + X = C$ prin reducerea matricelor A și B la forma Schur complexă. Algoritmul utilizează funcția **fsc** pentru calculul formei Schur complexe și a matricei de transformare corespunzătoare.)

1. $[U, S] = \mathbf{fsc}(A)$
2. $[V, T] = \mathbf{fsc}(B)$
3. $C \leftarrow \tilde{C} = U^H C V$
4. Pentru $j = 1 : n$
 1. Dacă $j > 1$ atunci
 1. $c_j \leftarrow c_j - S \left(\sum_{k=1}^{j-1} t_{kj} y_k \right)$
 2. Se rezolvă sistemul superior triunghiular $(St_{jj} + I_m)y_j = c_j$.
5. $X = \text{Re}(UYV^H)$.

Având în vedere similitudinea algoritmilor 1.1 și 1.2, similitudine care se extinde la celelalte variante, lăsăm în sarcina cititorului elaborarea *variantei Hessenberg-Schur* ca și justificarea faptului că următorul algoritm, cu notațiile din algoritmul 1.2, rezolvă problema în cazul matricelor de intrare A și B *simple* (diagonalizabile).

Algoritm 1.4 (Date matricele simple $A \in \mathcal{R}^{m \times m}$, $B \in \mathcal{R}^{n \times n}$ cu $\lambda_i \mu_j \neq -1$ pentru toți $\lambda_i \in \lambda(A)$, $\mu_j \in \lambda(B)$ și matricea $C \in \mathcal{R}^{m \times n}$, algoritmul calculează soluția X a ecuației $AXB + X = C$ pe baza determinării valorilor și vectorilor proprii ale matricelor A și B cu ajutorul funcției **vvp**.)

1. $[U, l] = \mathbf{vvp}(A)$
2. $[V, m] = \mathbf{vvp}(B)$
3. Se rezolvă sistemul matriceal $UZ = CV$ în raport cu Z .
4. Pentru $i = 1 : m$
 1. Pentru $j = 1 : n$
 1. $y_{ij} = z_{ij} / (l_i m_j + 1)$
5. Se rezolvă sistemul matriceal $XV = UY$ în raport cu X .
6. $X = \text{Re}X$.

1.4 Rezolvarea ecuațiilor matriceale Liapunov

Ecuațiile matriceale Liapunov continuă $A^T X + X A = C$ și discretă $A^T X A - X = C$ sunt cazuri particulare ale ecuațiilor Sylvester corespunzătoare și, prin urmare, rezultatele paragrafului precedent sunt aplicabile în întregime. Concret, având în vedere că $\lambda(A^T) = \lambda(A)$, obținem următoarele consecințe ale propozițiilor 1.3 și 1.4.

Corolar 1.1 *Ecuația matriceală Liapunov continuă (1.10) admite o soluție unică dacă și numai dacă matricea A nu are nici o pereche de valori proprii opuse, i.e.*

$$\lambda_i + \lambda_j \neq 0, \quad \forall \lambda_i, \lambda_j \in \lambda(A), \quad (1.54)$$

În particular, $0 \notin \lambda(A)$ i.e. matricea A este nesingulară.

Corolar 1.2 *Ecuația matriceală Liapunov discretă (1.11) admite o soluție unică dacă și numai dacă matricea A nu are nici o pereche de valori proprii inverse, i.e.*

$$\lambda_i \lambda_j \neq 1, \quad \forall \lambda_i, \lambda_j \in \lambda(A), \quad (1.55)$$

În particular, $\pm 1 \notin \lambda(A)$.

Algoritmii 1.1, 1.2, 1.3 și 1.4 se simplifică în sensul că este necesară o singură reducere la forma Schur, respectiv un singur calcul de valori și vectori proprii.

I. Într-adevăr, să considerăm mai întâi cazul *ecuației Liapunov continue* (1.10) și fie reducerea (1.31) a matricei $A \in \mathcal{R}^{n \times n}$ la *forma Schur complexă*. Rezultă

$$A^T = A^H = US^H U^H. \quad (1.56)$$

Cu (1.31), (1.56) ecuația Liapunov (1.10) devine

$$US^H U^H X + XUSU^H = C \quad (1.57)$$

respectiv

$$S^H Y + YS = \tilde{C} \quad (1.58)$$

unde

$$Y = U^H XU, \quad \tilde{C} = U^H CU. \quad (1.59)$$

Întrucât S^H are o structură inferior triunghiulară, scriind ecuația matriceală (1.58) pe coloane obținem sistemele inferior triunghiulare

$$(S^H + s_{jj}I_n)y_j = \tilde{c}_j - \sum_{k=1}^{j-1} s_{kj}y_k, \quad j = 1 : n. \quad (1.60)$$

În condițiile corolarului 1.1, aceste sisteme sunt nesingulare și pot fi rezolvate în ordinea $j = 1, 2, \dots, n$. Cu observația că soluția X a ecuației inițiale se calculează apoi din prima relație (1.59) rezultă următorul algoritm.

Algoritmul 1.5 (Date matricele $A \in \mathcal{R}^{n \times n}$, satisfăcând (1.54), și $C \in \mathcal{R}^{n \times n}$ oarecare, algoritmul calculează soluția X a ecuației Liapunov $A^T X + XA = C$ pe baza reducerii matricei A la forma Schur complexă cu ajutorul funcției **fsc**.)

1. $[U, S] = \mathbf{fsc}(A)$
2. $C \leftarrow U^H C U$
3. Pentru $j = 1 : n$
 1. Dacă $j > 1$ atunci
 1. $c_j = \tilde{c}_j - \sum_{k=1}^{j-1} s_{kj}y_k$
 2. Se rezolvă sistemul inferior triunghiular $(S^H + s_{jj}I_n)y_j = c_j$.
4. $X = \text{Re}(UYU^H)$.

În multe aplicații matricea C din (1.10) este *simetrică*. Într-un astfel de caz apar simplificări suplimentare datorită faptului că, în condițiile de existență și unicitate menționate, și matricea soluție este *simetrică*. Într-adevăr, dacă X este soluția, se verifică imediat că X^T satisface ecuația

matriceală (1.10) și, cum soluția este unică, rezultă $X = X^T$. Acest fapt poate fi exploatat observând că, datorită simetriei matricelor C și X , matricele Y și \tilde{C} din (1.58) sunt hermitice și în consecință este suficient să se calculeze, de exemplu, numai triunghiul superior al matricelor \tilde{C}, Y și, în final, X (exercițiul 1.11).

II. Algoritm de rezolvare a *ecuației Liapunov discrete* (1.11), bazat pe reducerea matricei A la *forma Schur complexă*, urmărește îndeaproape schema de mai sus.

Algoritm 1.6 (Date matricele $A \in \mathcal{R}^{n \times n}$, satisfăcând (1.55), și $C \in \mathcal{R}^{n \times n}$ oarecare, algoritmul calculează soluția X a ecuației Liapunov $A^T X A - X = C$ pe baza reducerii matricei A la forma Schur complexă cu ajutorul funcției **fsc**.)

1. $[U, S] = \mathbf{fsc}(A)$
2. $C \leftarrow U^H C U$
3. Pentru $j = 1 : n$
 1. Dacă $j > 1$ atunci
 1. $c_j = c_j - S^H (\sum_{k=1}^{j-1} s_{kj} y_k)$
 2. Se rezolvă sistemul inferior triunghiular $(S^H s_{jj} - I_n) y_j = c_j$.
4. $X = \text{Re}(U Y U^H)$.

Tratarea cazului în care matricele C și, drept consecință, X sunt simetrice face obiectul exercițiului 1.11.

Variantele bazate pe calculul valorilor și vectorilor proprii ale matricei A sunt, de asemenea, propuse cititorului (exercițiul 1.12).

Cazul important, în care matricea A este *stabilă* (vezi cap. 5) iar matricea C este *simetrică și negativ (semi)definită* este abordat în capitolul 7 (se pot consulta [X] și [6], [7]).

Încheiem prezentarea metodelor de rezolvare a ecuațiilor matriceale liniare prin evidențierea posibilității de a reduce rezolvarea unei ecuații Liapunov continue la rezolvarea unei ecuații corespondente discrete și reciproc. Pentru început să observăm că dacă ecuația Liapunov discretă (1.11) admite o soluție unică atunci, conform (1.55), $\pm 1 \notin \lambda(A)$ respectiv matricea $A - I_n$ este nesingulară și, deci, putem defini transformarea (omografică)

$$B = (A - I_n)^{-1}(A + I_n). \quad (1.61)$$

Apoi, din (1.61) rezultă că $1 \notin \lambda(B)$, i.e $B - I_n$ este nesingulară (exercițiul 1.13), deci

$$A = (B + I_n)(B - I_n)^{-1}, \quad (1.62)$$

respectiv,

$$A^T = (B^T - I_n)^{-1}(B^T + I_n) \quad (1.63)$$

expresii care, introduse în ecuația Liapunov discretă (1.11), conduc, după câteva prelucrări elementare care conservă soluția, la

$$B^T X + XB = D \quad (1.64)$$

unde

$$D = \frac{1}{2}(B^T - I_n)C(B - I_n). \quad (1.65)$$

În consecință, ecuația Liapunov continuă definită de (1.64), (1.61) și (1.65) are aceeași soluție cu ecuația Liapunov discretă (1.11) și poate servi ca suport pentru o modalitate de rezolvare a acesteia din urmă. Reciproc, dată ecuația Liapunov continuă (1.64), respectiv date matricele $n \times n$ B și D cu valorile proprii ale matricei B satisfăcând condiția $\mu_i + \mu_j \neq 0$ oricare ar fi $\mu_i, \mu_j \in \lambda(B)$ (în particular $0 \notin \lambda(B)$) putem găsi soluția X a acestei ecuații rezolvând ecuația Liapunov discretă (1.11) cu A dat de (1.62) și

$$C = 2(B^T - I_n)^{-1}D(B - I_n)^{-1} \quad (1.66)$$

obținut din (1.65).

Proprietăți numerice

Algoritmii prezentați utilizează în exclusivitate transformări ortogonale (unitare) ceea ce le conferă remarcabile proprietăți numerice. Toți algoritmii includ cel puțin o execuție a algoritmului **QR** de aducere iterativă a uneia din matricele de intrare la forma Schur reală sau complexă urmată de rezolvarea unor sisteme triunghiulare sau de tip Hessenberg. Algoritmii **QR** este admis a fi un algoritm numeric stabil ca și algoritmii de rezolvare prin substituție a sistemelor triunghiulare. Mai mult, se poate arăta că algoritmii din prezentul capitol sunt numeric stabili în întregul lor, i.e. soluția calculată reprezintă soluția exactă a problemei respective cu datele de intrare perturbate nesemnificativ. În consecință, preciziile efective ce se obțin depind de condiționările numerice ale matricelor inițiale și de nivelul practicat al toleranțelor.

Programe MATLAB disponibile

Pentru rezolvarea ecuațiilor matriceale Sylvester și Liapunov continue este disponibilă funcția **lyap**, care implementează algoritmi 1.1 și 1.5 (se apelează funcțiile **schur** pentru aducerea la forma Schur reală și **rsf2csf** pentru obținerea formei Schur complexe). În cazul discret se poate utiliza funcția **dlyap**, care impementează metoda transformării omografice (1.61).

Exerciții

E 1.1 Fie date matricele $A \in \mathcal{R}^{m \times m}$ nesingulară și $C \in \mathcal{C}^{m \times n}$. Scrieți algoritmi de rezolvare a sistemului matriceal $AX = C$ folosind eliminarea gaussiană cu pivotare parțială și completă. Evaluați numărul asimptotic de operații aritmetice. Cum procedați dacă matricea A este simetrică, eventual pozitiv definită ?

E 1.2 Se consideră date matricele $A \in \mathcal{R}^{p \times m}$, $C \in \mathcal{R}^{p \times n}$ și fie

$$A = U_1 \Sigma_1 V_1^T, \quad \Sigma_1 = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r), \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

descompunerea valorilor singulare a lui A . Arătați că $X^* = V_1 \Sigma_1^{-1} U_1^T C$ este matricea de normă Frobenius minimă dintre toate matricele $X \in \mathcal{R}^{m \times n}$ care minimizează norma Frobenius a rezidului matriceal $R = C - AX$.

E 1.3 Admițând că dispuneți de o funcție care calculează valorile singulare și matricele de transformare corespunzătoare pentru o matrice de intrare $A \in \mathcal{R}^{p \times m}$ dată, scrieți algoritmul de rezolvare, în sensul celor mai mici pătrate, a ecuației matriceale $AX = C$.

E 1.4 Dată o matrice $A \in \mathcal{R}^{2 \times 2}$ având $\lambda(A) \subset \mathcal{C} \setminus \mathcal{R}$, determinați o matrice de rotație bidimensională complexă $P \in \mathcal{C}^{2 \times 2}$, definită de

$$P = \begin{bmatrix} c & s \\ -\bar{s} & c \end{bmatrix}, \quad c^2 + |s|^2 = 1, \quad c \in \mathcal{R}$$

astfel încât matricea $S = P^H A P \in \mathcal{C}^{2 \times 2}$, $P^H \stackrel{\text{def}}{=} \bar{P}^T$ să fie superior triunghiulară.

E 1.5 Se dă o matrice $A \in \mathcal{R}^{n \times n}$ în formă Schur reală. Se cere un algoritm de calcul al unei matrice unitare $Q \in \mathcal{C}^{n \times n}$ și al matricei supe-

¹O matrice $P \in \mathcal{C}^{n \times n}$ se numește *unitară* dacă $P^H P = P P^H = I_n$ și joacă, în calculul matriceal complex, același rol cu cel jucat de matricele ortogonale în calculul matriceal real.

rior triunghiulare $S \in \mathcal{C}^{n \times n}$ astfel încât $Q^H A Q = S$. (Matricea S astfel obținută se numește *formă Schur complexă* a matricei A).

Indicație. Matricea Q poate fi un produs de matrice unitare de forma

$$Q_i = \begin{bmatrix} I & 0 & 0 \\ 0 & P_i & 0 \\ 0 & 0 & I \end{bmatrix} \in \mathcal{C}^{n \times n}$$

cu $P_i \in \mathcal{C}^{2 \times 2}$ o rotație complexă bidimensională.

E 1.6 Fie $A \in \mathcal{C}^{2 \times 2}$ superior triunghiulară. Să se calculeze o matrice unitară $U \in \mathcal{C}^{2 \times 2}$ astfel încât

$$U^H A U = \begin{bmatrix} a_{22} & * \\ 0 & a_{11} \end{bmatrix},$$

i.e. să se obțină permutarea elementelor diagonale.

E 1.7 Dacă $B \in \mathcal{C}^{n \times n}$ este o matrice superior triunghiulară dată, se cere un algoritm de calcul al unei matrice unitare $V \in \mathcal{C}^{n \times n}$ și a matricei superior triunghiulare $S \in \mathcal{C}^{n \times n}$ astfel încât $V^H B V = S$ și $|s_{11}| \leq |s_{22}| \leq \dots \leq |s_{nn}|$.

Observație. Dacă B este forma Schur complexă a unei matrice $A \in \mathcal{R}^{n \times n}$, atunci matricea S de mai sus poartă numele de *formă Schur complexă ordonată* a matricei A . Evident, se poate utiliza orice alt criteriu de ordonare a elementelor diagonale.

Indicație. Se folosește o secvență de transformări de forma

$$V_i = \begin{bmatrix} I & 0 & 0 \\ 0 & U_i & 0 \\ 0 & 0 & I \end{bmatrix}$$

unde U_i sunt matrice 2×2 determinate în exercițiul 1.6.

E 1.8 Admițând că numărul asimptotic de operații aritmetice în format virgulă mobilă necesar pentru execuția algoritmului **QR** de aducere a unei matrice reale $n \times n$ la forma Schur reală, (inclusiv calculul matricei de transformare) este, în medie statistică, conform [VI], $N_1 = 12n^3$, evaluați

a) numărul asimptotic de operații impus de execuția algoritmului de la exercițiul 1.7;

b) numărul asimptotic de operații necesar pentru rezolvarea, cu ajutorul algoritmului 1.2 (variante Schur–Schur), a ecuației Sylvester continue (1.8);

c) numărul asimptotic de operații necesar pentru rezolvarea ecuației Liapunov continue (1.10) cu ajutorul algoritmului 1.5, respectiv pentru rezolvarea ecuației Liapunov discrete (1.11) cu ajutorul algoritmului 1.6.

E 1.9 Admițând că sunt îndeplinite condițiile de existență și unicitate a soluției, scrieți un algoritm eficient de rezolvare a ecuației matriciale Sylvester continue (1.8) prin interpretarea ”desfășurată” a acesteia ca un sistem de mn ecuații cu mn necunoscute. Algoritmul nu trebuie să includă nici o procedură iterativă. Idem pentru ecuația Sylvester discretă (1.9).

Evaluati numărul de operații aritmetice necesar și comparați cu cel determinat la exercițiul 1.8. Ce concluzii se impun ?

Indicație. Dacă $\bar{x} \in \mathcal{R}^{mn}$ și $\bar{c} \in \mathcal{R}^{mn}$ sunt vectorii definiți, de exemplu, prin concatenarea coloanelor matricelor X respectiv C arătați că ecuația (1.8) se poate scrie în forma $(I_n \otimes A + B^T \otimes I_m)\bar{x} = \bar{c}$ iar (1.9) în forma $(B^T \otimes A + I_{mn})\bar{x} = \bar{c}$. În relațiile de mai sus \otimes semnifică produsul Kronecker definit în felul următor: dacă $U \in \mathcal{R}^{p \times q}$, $V \in \mathcal{R}^{r \times s}$ atunci $W \stackrel{\text{def}}{=} U \otimes V \in \mathcal{R}^{pr \times qs}$ este matricea având structura bloc $W = [W_{ij}]_{i=1:p, j=1:q}$ cu $W_{ij} = u_{ij}V$.

E 1.10 Fiind date matricele $A \in \mathcal{R}^{m \times m}$, $B \in \mathcal{R}^{n \times n}$ cu $\lambda(A) \cap \lambda(-B) = \emptyset$ și $C \in \mathcal{R}^{m \times n}$, să se scrie algoritmi Hessenberg – Schur de rezolvare a ecuației matriciale Sylvester continue $AX + XB = C$, respectiv a celei discrete $AXB + X = C$, bazate pe aducerea, prin transformări ortogonale, a matricei A la forma superior Hessenberg și a matricei B la forma Schur complexă. Presupuneți că se dispune de funcția **fsc** de calcul a formei Schur complexe, toate celelalte prelucrări urmând a fi scrise explicit.

E 1.11 Scrieți corespondenții algoritmilor 1.5 și 1.6 (variante Schur) pentru rezolvarea ecuațiilor matriciale Liapunov $A^T X + XA = C$ și $A^T XA - X = C$ în ipoteza că matricea C este simetrică, minimizând numărul de operații aritmetice și memoria ocupată.

E 1.12 Elaborați algoritmi de rezolvare a ecuațiilor Liapunov continuă și discretă (1.10) și (1.11) având matricele A simple, folosind o procedură de calcul a valorilor și vectorilor proprii. Exploatați cât mai eficient ipoteza simetriei matricei termenilor liberi. Admițând că efortul de calcul al valorilor și vectorilor proprii este $N_1 = kn^3$, determinați efortul total solicitat de rezolvarea ecuațiilor matriciale date.

E 1.13 Fie o matrice $A \in \mathcal{R}^{n \times n}$ cu proprietatea $\lambda_i \lambda_j \neq 1, \forall \lambda_i, \lambda_j \in \lambda(A)$. Se cere să se arate:

- a) Dacă $B = (A - I_n)^{-1}(A + I_n)$ atunci pentru orice $\lambda_i \in \lambda(A)$ avem $\mu_i = (\lambda_i + 1)/(\lambda_i - 1) \in \lambda(B)$;
- b) $\{0, 1\} \cap \lambda(B) = \emptyset$;
- c) $\lambda(B) \cap \lambda(-B) = \emptyset$.

E 1.14 Elaborați algoritmi de tip Schur-Schur pentru rezolvarea ecuațiilor Sylvester $AX + XB = C$ și $AXB + X = C$ care să utilizeze formele Schur reale ale matricelor A și B și exclusiv o aritmetică reală ².

Indicație. Utilizând reducerea matricelor A și B la forma Schur reală ecuațiile (1.35), respectiv (1.51), au matricele S și T cvasisuperior triunghiulare. Partiționând matricea Y în blocuri conform dimensiunilor blocurilor diagonale ale matricelor S și T , scriind (1.35), respectiv (1.51), pe bloc-coloane și utilizând o tehnică de bloc-substituție înapoi se obține că ecuațiile (1.35), respectiv (1.51), se reduc la un set de ecuații Sylvester continue, respectiv discrete, definite de blocurile diagonale ale matricelor S și T . Aceste ecuații se pot rezolva, în ordinea sugerată mai sus, într-o aritmetică reală, prin desfășurarea lor în sisteme uzuale de ordin cel mult patru (vezi exercițiul 1.9).

E 1.15 Elaborați algoritmi de tip Schur pentru rezolvarea ecuațiilor Liapunov $A^T X + XA = C$ și $A^T XA - X = C$ care să utilizeze exclusiv forma Schur reală a matricei A și o aritmetică reală ³.

Indicație. Adaptați indicațiile de la exercițiul precedent. Pentru detalii puteți consulta [4-6] și [X].

Bibliografie

- [1] **Davison E.J., Man F.T.** *The Numerical Solution of $A^T Q + QA = -C$* , IEEE Trans. Automat. Contr., vol. AC-13, pp. 448-449, 1968.
- [2] **Berger C.S.** *A Numerical Solution of the Matrix Equation $P = \Phi P \Phi^T + S$* , IEEE Trans. Automat. Contr., vol.AC-16, pp. 381-382, 1971.
- [3] **Bartell R.H., Stewart G.W.** *Solution of the Matrix Equation $AX + XB = C$* , Commun. Ass. Comput. Mach., vol 15, pp. 821-826, 1972.
- [4] **Barraud A.Y.** *A Numerical Algorithm to Solve $A^T XA - X = Q$* , IEEE Trans. Automat. Contr., vol.AC-22, pp. 883-885, 1977.

²Varianta Schur-Schur reală pentru rezolvarea ecuației Sylvester continue este referită în literatura de specialitate cu denumirea de *algoritmul Bartels-Stewart* [3].

³Varianta Schur reală pentru rezolvarea ecuației Liapunov discrete este referită în literatura de specialitate cu denumirea de *algoritmul Kitagawa-Barraud* [4, 5].

- [5] **Kitagawa G.** *An Algorithm for Solving the Matrix Equation $X = F^T X F + S$* , Int. J. Control, vol. 25, pp. 745–753, 1977.
- [6] **Sima V.** *Comparison of some Algorithms for Solving the Lyapunov-type Equations*, Rev. Roum. Sci. Techn. – Electrotechn. Energ., vol. 25, pp. 625–632, 1980.
- [7] **Hammarling S.J.** *Numerical Solution of the Stable, Non-negative Definite Lyapunov Equation*, IMA J. Numer. Anal., vol. 2, pp. 303–323, 1982.

Capitolul 2

Calculul funcțiilor de matrice

Exponențiala matriceală

Calculul funcțiilor de matrice este necesar în multe aplicații. În particular, analiza comportării dinamice a sistemelor liniare continue face apel la exponențiala matriceală, pentru calculul căreia a fost elaborată o întreagă serie de tehnici numerice mai mult sau mai puțin performante. În consecință, prezentăm în continuare unele metode generale de calcul pentru funcțiile de matrice, precum și metodele specifice, cele mai apreciate, de calcul pentru funcția exponențială de argument matriceal.

2.1 Funcții de matrice

Considerăm o funcție $f : D \subset \mathcal{C} \rightarrow \mathcal{C}$ definită pe o mulțime D din planul complex și fie $A \in \mathcal{C}^{n \times n}$ o matrice dată. Ne propunem mai întâi să definim noțiunea de *funcție de matrice* adică semnificația expresiei

$$F = f(A). \quad (2.1)$$

Pentru început observăm că dacă f este o funcție polinomială

$$f(z) = \sum_{i=0}^N c_i z^i, \quad c_i \in \mathcal{C}, \quad i = 0 : N \quad (2.2)$$

atunci matricea

$$F = \sum_{i=0}^N c_i A^i \quad (2.3)$$

este bine definită și poate fi numită *valoarea polinomului f în punctul (sau pe matricea) A* .

Fie acum $\mu_A(z)$ polinomul minimal al matricei A ¹ și $\lambda_i \in \mathcal{C}$, $i = 1 : l$, zerourile acestuia, zeroul λ_i având ordinul de multiplicitate m_i . Avem

$$\mu_A(z) = z^m + \sum_{k=0}^{m-1} \alpha_k z^k = \prod_{i=1}^l (z - \lambda_i)^{m_i} \quad (2.4)$$

unde

$$m \stackrel{\text{def}}{=} \text{grad}(\mu_A) = \sum_{i=1}^l m_i. \quad (2.5)$$

Conform teoremei împărțirii cu rest, există unic două polinoame q , r cu $\text{grad}(r) < m$ astfel încât

$$f = \mu_A q + r. \quad (2.6)$$

Polinomul minimal μ_A fiind un polinom anulator pentru matricea A , i.e. $\mu_A(A) = 0$, din (2.6) rezultă

$$F \stackrel{\text{def}}{=} f(A) = r(A). \quad (2.7)$$

În consecință, valoarea polinomului f în punctul A este aceeași cu cea a polinomului r și același lucru se poate spune despre orice alt polinom al cărui rest la împărțirea prin μ_A este r . Polinomul r poate fi determinat prin aplicarea algoritmului de împărțire cu rest a polinoamelor.

O altă modalitate de calcul a polinomului r se bazează pe faptul că din (2.4) rezultă

$$\mu_A^{(k)}(\lambda_i) = 0, \quad k = 0 : m_i - 1, \quad (2.8)$$

și, în consecință, cei (cel mult) m coeficienți ai polinomului r se pot determina, conform (2.6), prin rezolvarea sistemului de m ecuații liniare

$$r^{(k)}(\lambda_i) = f^{(k)}(\lambda_i), \quad i = 1 : l, \quad k = 0 : m_i - 1, \quad (2.9)$$

unde indicele superior indică ordinul derivatei.

¹Amintim că μ_A este polinomul monic de grad minim cu proprietatea $\mu_A(A) = 0$. Dacă $p(z) = \det(zI - A)$ este polinomul caracteristic al matricei A iar $d(z)$ este cel mai mare divizor comun (monic) al tuturor minorilor de ordinul $n-1$ ai matricei caracteristice $zI - A$ atunci $\mu_A(z) = p(z)/d(z)$.

Relațiile (2.7) și (2.9) servesc ca bază pentru definirea valorii în punctul A a oricărei funcții f care admite derivatele cerute în (2.9). Pentru o abordare formală introducem

Definiția 2.1 Fie $A \in \mathcal{C}^{n \times n}$ și

$$\Lambda = \{(\lambda_i, m_i) \mid i = 1 : l, \lambda_i \in \mathcal{C}, m_i \in \mathcal{N}\} \quad (2.10)$$

mulțimea zerourilor polinomului minimal μ_A al matricei A împreună cu multiplicitățile respective. Dacă funcția $f : D \subset \mathcal{C} \rightarrow \mathcal{C}$ este analitică pe o mulțime deschisă D ce conține punctele λ_i , $i = 1 : l$ atunci spunem că f este definită pe spectrul matricei A iar mulțimea valorilor funcției f pe spectrul matricei A este

$$f(\Lambda) = \left\{ f^{(k)}(\lambda_i) \mid i = 1 : l, k = 0 : m_i - 1 \right\}. \quad (2.11)$$

În particular, o funcție întreagă f (i.e. analitică pe $D = \mathcal{C}$) este definită pe spectrul oricărei matrice $A \in \mathcal{C}^{n \times n}$.

În condițiile definiției 2.1 introducem

Definiția 2.2 Fie date o matrice $A \in \mathcal{C}^{n \times n}$ și o funcție f definită pe spectrul lui A . Dacă polinomul minimal μ_A al matricei A are gradul m atunci polinomul r de grad cel mult $m-1$, unic determinat de sistemul liniar (2.9), se numește polinomul de interpolare Lagrange-Sylvester al funcției f pe spectrul matricei A .

Putem acum să introducem noțiunea de funcție de matrice extinzând relația (2.7) de la polinoame la orice funcție definită pe spectrul matricei argument.

Definiția 2.3 Dacă f este o funcție definită pe spectrul unei matrice A atunci valoarea funcției f în punctul A este

$$f(A) = r(A), \quad (2.12)$$

unde r este polinomul de interpolare Lagrange-Sylvester al funcției f pe spectrul lui A .

În scopul evidențierii cazurilor când o funcție de matrice este reală, în continuare vom utiliza următoarea terminologie

Definiția 2.4 O funcție $f : D \subset \mathcal{C} \rightarrow \mathcal{C}$ este reală pe spectrul matricei $A \in \mathcal{R}^{n \times n}$ dacă

$$f(\bar{\lambda}_i) = \overline{f(\lambda_i)}, \quad \forall \lambda_i, \bar{\lambda}_i \in \lambda(A),$$

unde \bar{z} este notația pentru conjugatul numărului complex z .

În particular $f(\lambda_i) \in \mathcal{R}$ pentru orice $\lambda_i \in \lambda(A) \cap \mathcal{R}$.

În acest context avem

Propoziția 2.1 *Dacă funcția f este reală pe spectrul matricei reale $A \in \mathcal{R}^{n \times n}$ atunci polinomul de interpolare Lagrange-Sylvester are coeficienții reali și, în consecință, matricea $F = f(A)$ este reală.*

Exemplul 2.1 Fie $f(z) = e^z$ și

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}.$$

Avem $\mu_A(z) = z^2 - 3z + 2$ și $\Lambda = \{(1, 1), (2, 1)\}$. Dacă $r(z) = r_1z + r_0$, atunci sistemul (2.9) se scrie

$$\begin{cases} r_1 + r_0 = e \\ 2r_1 + r_0 = e^2 \end{cases}$$

de unde rezultă $r_0 = 2e - e^2$, $r_1 = e^2 - e$ și, în consecință,

$$F = e^A = (e^2 - e)A + (2e - e^2)I_3 = \begin{bmatrix} e & e^2 - e & 0 \\ 0 & e^2 & 0 \\ 0 & 0 & e^2 \end{bmatrix}. \diamond$$

Definiția 2.3 a funcțiilor de matrice pe baza polinomului de interpolare Lagrange-Sylvester pune în evidență următoarele aspecte.

a) Evaluarea oricărei funcții de matrice se reduce la evaluarea unui polinom matriceal. Totuși, deși sunt stabilite expresii analitice pentru coeficienții polinomului de interpolare Lagrange-Sylvester, această cale nu este recomandată pentru calculul numeric al funcțiilor de matrice din considerente de eficiență și stabilitate numerică.

b) Valoarea funcției f în punctul A este determinată exclusiv de mulțimea valorilor funcției f pe spectrul matricei A .

c) Dacă matricea A nu are valori proprii multiple, atunci polinomul minimal coincide cu polinomul caracteristic iar sistemul liniar (2.9) a cărui soluție furnizează coeficienții polinomului de interpolare Lagrange-Sylvester devine

$$r(\lambda_i) = f(\lambda_i), \quad i = 1 : n. \quad (2.13)$$

Următoarele rezultate (ale căror demonstrații pot fi găsite, de exemplu, în [VII]), permit evidențierea unor proprietăți utile în elaborarea unor proceduri de calcul efectiv.

Teorema 2.1 Fie $\lambda(A)$ spectrul matricei $A \in \mathbb{C}^{n \times n}$ și $D \subset \mathbb{C}$ un domeniu cu frontiera Γ suficient de netedă astfel încât $\lambda(A) \subset D$. Dacă f este o funcție analitică pe $D \cup \Gamma$ atunci

$$f(A) = \frac{1}{2\pi i} \oint_{\Gamma} (zI - A)^{-1} f(z) dz. \quad (2.14)$$

Expresia (2.14) poate servi ca definiție pentru funcțiile analitice (pe un domeniu) iar calculul integralei Cauchy

$$f_{ij} = \frac{1}{2\pi i} \oint_{\Gamma} e_i^T (zI - A)^{-1} e_j f(z) dz \quad (2.15)$$

poate fi efectuat cu ajutorul teoremei reziduurilor.

Exemplul 2.2 Considerăm funcția f și matricea A din exemplul 2.1. Avem $\lambda(A) = \{1, 2, 2\}$ și

$$(zI - A)^{-1} = \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-2)} & 0 \\ 0 & \frac{1}{z-2} & 0 \\ 0 & 0 & \frac{1}{z-2} \end{bmatrix}.$$

Prin urmare, f fiind analitică în tot planul complex, putem alege un domeniu simplu conex oarecare ce satisface $\lambda(A) \subset D$. Obținem

$$f_{11} = \frac{1}{2\pi i} \oint_{\Gamma} \frac{e^z}{z-1} dz = \operatorname{Rez} \left[\frac{e^z}{z-1} \right] \Big|_{z=1} = e,$$

$$f_{22} = f_{33} = \frac{1}{2\pi i} \oint_{\Gamma} \frac{e^z}{z-2} dz = \operatorname{Rez} \left[\frac{e^z}{z-2} \right] \Big|_{z=2} = e^2,$$

$$\begin{aligned} f_{13} &= \frac{1}{2\pi i} \oint_{\Gamma} \frac{e^z}{(z-1)(z-2)} dz = \operatorname{Rez} \left[\frac{e^z}{(z-1)(z-2)} \right] \Big|_{z=1} + \\ &+ \operatorname{Rez} \left[\frac{e^z}{(z-1)(z-2)} \right] \Big|_{z=2} = -e + e^2, \end{aligned}$$

celelalte elemente ale matricei $F = e^A$ fiind, evident, nule. \diamond

În cele ce urmează vom aborda în exclusivitate cazul funcțiilor analitice și vom putea utiliza, în consecință, relația (2.14) ca o relație definitorie a funcțiilor de matrice.

Menționăm, de asemenea, posibilitatea exprimării funcțiilor de matrice prin serii matriceale de puteri. În acest sens avem următoarea

Teorema 2.2 *Dacă funcția $f(z)$ se poate dezvolta în serie de puteri în jurul punctului $z = z_0$*

$$f(z) = \sum_{k=0}^{\infty} \alpha_k (z - z_0)^k \quad (2.16)$$

și seria este convergentă în discul $|z - z_0| < r$ atunci această dezvoltare rămâne valabilă dacă argumentul scalar este înlocuit cu argumentul matriceal A

$$f(A) = \sum_{k=0}^{\infty} \alpha_k (A - z_0 I)^k \quad (2.17)$$

oricare ar fi matricea A al cărei spectru se află în interiorul discului de convergență.

Din această teoremă rezultă, printre altele, că dezvoltările

$$e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k,$$

$$\sin A = \sum_{k=0}^{\infty} \frac{(-1)^k}{(2k+1)!} A^{2k+1}, \quad \cos A = \sum_{k=0}^{\infty} \frac{(-1)^k}{(2k)!} A^{2k}$$

sunt valabile oricare ar fi matricea $A \in \mathcal{C}^{n \times n}$.

Exemplul 2.3 Reluăm matricea A și funcția f din exemplul 2.1. Prin inducție rezultă imediat

$$A^k = \begin{bmatrix} 1 & 2^k - 1 & 0 \\ 0 & 2^k & 0 \\ 0 & 0 & 2^k \end{bmatrix}$$

și deci

$$F = e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k = \begin{bmatrix} \sum_{k=0}^{\infty} \frac{1}{k!} & \sum_{k=0}^{\infty} \frac{2^k - 1}{k!} & 0 \\ 0 & \sum_{k=0}^{\infty} \frac{2^k}{k!} & 0 \\ 0 & 0 & \sum_{k=0}^{\infty} \frac{2^k}{k!} \end{bmatrix} =$$

$$= \begin{bmatrix} e & e^2 - e & 0 \\ 0 & e^2 & 0 \\ 0 & 0 & e^2 \end{bmatrix}. \diamond$$

În continuare, prezentăm câteva proprietăți ale funcțiilor de matrice, utile în dezvoltările procedurale care fac obiectul metodelor de calcul recomandate ca fiind cele mai bune în momentul actual.

Propoziția 2.2 *Dacă funcția f este definită pe spectrul matricei A și*

$$B = TAT^{-1} \quad (2.18)$$

unde T este o matrice nesingulară, atunci

$$f(B) = Tf(A)T^{-1}. \quad (2.19)$$

Demonstrația este imediată pentru funcțiile analitice pe baza relației definitorii (2.14) întrucât din (2.18) rezultă

$$(zI - B)^{-1} = (T(zI - A)T^{-1})^{-1} = T(zI - A)^{-1}T^{-1}. \diamond$$

Transferul transformărilor de asemănare (2.18) de la nivelul argumentelor matriceale la nivelul funcțiilor are o importanță decisivă în elaborarea unor tehnici de calcul adecvate care se pot focaliza asupra unor structuri matriceale remarcabile.

Propoziția 2.3 *Dacă matricea A este (bloc-) diagonală*

$$A = \text{diag}(A_{11}, A_{22}, \dots, A_{pp}) \quad (2.20)$$

atunci $f(A)$ este (bloc-) diagonală și

$$F \stackrel{\text{def}}{=} f(A) = \text{diag}(f(A_{11}), f(A_{22}), \dots, f(A_{pp})). \quad (2.21)$$

Demonstrație. Din faptul că

$$(zI - A)^{-1} = \text{diag}((zI - A_{11})^{-1}, (zI - A_{22})^{-1}, \dots, (zI - A_{pp})^{-1})$$

(2.21) rezultă direct din (2.14). \diamond

Observația 2.1 Conform propozițiilor 2.2 și 2.3, dacă

$$T^{-1}AT = \text{diag}(J_1, \dots, J_p)$$

este forma Jordan a lui A , atunci avem

$$f(A) = T \operatorname{diag}(f(J_1), \dots, f(J_p)) T^{-1},$$

și calculul lui $f(A)$ se reduce la calculul matricelor $f(J_k)$, unde J_k sunt blocurile Jordan. Considerând un bloc Jordan generic

$$J = \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ & & & \ddots & \\ & & & & \ddots & 1 \\ 0 & 0 & 0 & \cdots & \lambda \end{bmatrix}$$

și utilizând definiția 2.3 obținem

$$f(J) = \begin{bmatrix} f(\lambda) & f'(\lambda) & \cdots & \frac{f^{(m-1)}(\lambda)}{(m-1)!} \\ 0 & f(\lambda) & \cdots & \frac{f^{(m-2)}(\lambda)}{(m-2)!} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & f(\lambda) \end{bmatrix},$$

(exercițiul 2.4). \diamond

Propoziția 2.4 Dacă matricea $A \in \mathbb{C}^{n \times n}$ este superior (inferior) triunghiulară, iar f este o funcție definită pe spectrul lui A atunci

- a) $F = f(A)$ este superior (inferior) triunghiulară;
- b) $f_{ii} = f(a_{ii})$, $i = 1 : n$.

Demonstrație.

a) Se știe că inversa unei matrice superior (inferior) triunghiulare nesingulare este superior (inferior) triunghiulară. În consecință $(zI - A)^{-1}$ este superior (inferior) triunghiulară și, din (2.14), rezultă că aceeași structură o are și matricea $F = f(A)$.

b) Dacă A este o matrice triunghiulară atunci elementele diagonale ale matricei $(zI - A)^{-1}$ sunt

$$\frac{1}{z - a_{ii}}, \quad i = 1 : n.$$

Funcția f fiind definită pe spectrul matricei A , rezultă că $\lambda_i = a_{ii}$, $i = 1 : n$ nu sunt puncte singulare pentru f și, în consecință, utilizând formula

integrală a lui Cauchy pentru elementele diagonale din (2.14), rezultă

$$f_{ii} = \frac{1}{2\pi i} \oint_{\Gamma} \frac{f(z)}{z - a_{ii}} dz = f(a_{ii}), \quad i = 1 : n.$$

q.e.d. \diamond

Propoziția 2.5 Fie $\lambda(A) = \{\lambda_1, \dots, \lambda_n\}$ spectrul matricei $A \in \mathcal{C}^{n \times n}$. Atunci pentru orice funcție f definită pe spectrul lui A avem

$$\lambda(f(A)) = \{f(\lambda_1), \dots, f(\lambda_n)\}. \quad (2.22)$$

Demonstrație. Fie $B = Q^H A Q$ forma Schur complexă a matricei A . Atunci B este superior triunghiulară și $b_{ii} = \lambda_i$, $i = 1 : n$. Aplicând propozițiile 2.2 și 2.4 rezultă

$$f(A) = Q f(B) Q^H$$

și, deci, valorile proprii ale matricei $f(A)$ sunt $f(b_{ii}) = f(\lambda_i)$, $i = 1 : n$. \diamond

Propoziția 2.6 Matricele $A \in \mathcal{C}^{n \times n}$ și $F = f(A)$ comută i.e.

$$A f(A) = f(A) A. \quad (2.23)$$

Demonstrație. Existența matricei $F = f(A)$ presupune că funcția f este definită pe spectrul matricei A . În acest caz însă și funcția h definită de $h(z) = z f(z) = f(z) z$ este definită pe spectrul matricei A . Din faptul că $h(A)$ are semnificație rezultă (2.23). \diamond

Teorema 2.1 împreună cu proprietatea (2.23) stabilește un izomorfism $f(z) \rightsquigarrow f(A)$ între algebra comutativă a funcțiilor analitice $f(z)$ pe spectrul lui A și algebra (de asemenea comutativă) a funcțiilor de matrice $f(A)$.

Observația 2.2 Din propozițiile 2.2 și 2.4 rezultă că dacă $Q^H A Q = B$ este forma Schur a lui A atunci $f(A) = Q f(B) Q^H$, unde $G = f(B)$ este superior triunghiulară cu $g_{ii} = f(b_{ii})$. Ținând seama că forma Schur prezintă proprietăți numerice net superioare formei canonice Jordan și că propoziția 2.6 stă la baza unui algoritm fiabil de calcul al funcțiilor de matrice triunghiulare (vezi secțiunea următoare) deducem că rezultatele de mai sus fundamentează o categorie importantă de proceduri de calcul al funcțiilor de argument matriceal. \diamond

Proprietățile următoare, care ilustrează izomorfismul de algebre comutative amintit anterior, permit extinderea unor identități care pun în relație funcții de argument scalar cu funcții de matrice.

Propoziția 2.7 Fie $g(z_1, z_2, \dots, z_p)$ un polinom în nedeterminatele z_1, z_2, \dots, z_p și $f_i, i = 1 : p$ funcții definite pe spectrul matricei $A \in \mathbb{C}^{n \times n}$. Definim funcția

$$h(z) \equiv g(f_1(z), \dots, f_p(z)). \quad (2.24)$$

Atunci h este o funcție definită pe spectrul matricei A și dacă

$$h^{(k)}(\lambda_i) = 0, \quad i = 1 : l, k = 0 : m_i - 1 \quad (2.25)$$

unde $(\lambda_i, m_i), i = 1 : l$ sunt zerourile polinomului minimal μ_A cu multiplicitățile respective, atunci

$$H \stackrel{\text{def}}{=} h(A) = g(f_1(A), \dots, f_p(A)) = 0. \quad (2.26)$$

Demonstrație. Dacă $r_i, i = 1 : p$ sunt polinoamele de interpolare Lagrange-Sylvester ale funcțiilor f_i pe spectrul matricei A , atunci considerăm funcția u definită de

$$u(z) = g(r_1(z), \dots, r_p(z))$$

care este, evident, o funcție polinomială de argument scalar. În aceste condiții, din (2.25), (2.26) și modul de definire a polinoamelor de interpolare Lagrange-Sylvester (2.9) rezultă

$$u^{(k)}(\lambda_i) = 0, \quad i = 1 : l, k = 0 : m_i - 1,$$

respectiv polinomul de interpolare al funcției u este identic nul, de unde obținem

$$u(A) = g(r_1(A), \dots, r_p(A)) = g(f_1(A), \dots, f_p(A)) = 0$$

q.e.d. \diamond

Câteva consecințe imediate ale propoziției 2.7, care evidențiază transferul unor formule scalare în cazul matriceal sunt următoarele relații utile.

a) Fie polinomul $g(z_1, z_2) = z_1^2 + z_2^2 - 1$ și funcțiile $f_1(z) = \sin z, f_2(z) = \cos z$. Atunci funcția definită în (2.25)

$$h(z) \equiv g(\sin z, \cos z) = \sin^2 z + \cos^2 z - 1$$

este identic nulă și prin urmare (2.26) este satisfăcută pentru orice matrice A . De aici rezultă că, pentru toate matricele pătrate A , avem

$$\sin^2 A + \cos^2 A = I.$$

b) Similar, considerând polinomul $g(z_1, z_2) = z_1 z_2 - 1$ și funcțiile $f_1(z) = e^z$, $f_2(z) = e^{-z}$ obținem $e^A e^{-A} = I$, adică

$$e^{-A} = (e^A)^{-1}.$$

c) De asemenea, cu $g(z_1, z_2, z_3) = z_1 + z_2 + z_3$ și $f_1(z) = e^{iz}$, $f_2(z) = -\cos z$, $f_3(z) = -i \sin z$ rezultă că pentru orice matrice patrată A este adevărată formula

$$e^{iA} = \cos A + i \sin A.$$

Propoziția 2.8 *Dacă*

$$\rho(z) = \frac{g(z)}{h(z)} \quad (2.27)$$

este o funcție rațională ireductibilă, definită pe un domeniu care nu conține rădăcinile polinomului $h(z)$ iar $A \in \mathcal{C}^{n \times n}$ atunci funcția ρ este definită pe spectrul matricei A dacă și numai dacă

$$h(\lambda_i) \neq 0 \quad \forall \lambda_i \in \lambda(A). \quad (2.28)$$

Dacă (2.28) sunt satisfăcute atunci

$$\rho(A) = g(A)[h(A)]^{-1} = [h(A)]^{-1}g(A). \quad (2.29)$$

Demonstrație. Funcția ρ este indefinit derivabilă în toate punctele în care $h(z)$ nu se anulează. În consecința $\rho^{(k)}(\lambda_i)$ sunt definite dacă și numai dacă (2.28) sunt adevărate. Pe de altă parte din (2.28) și propoziția 2.5 rezultă că $0 \notin \lambda(h(A))$ i.e. matricea $h(A)$ este nesingulară. Aplicând acum propoziția 2.7, din identitatea

$$h(z)\rho(z) = \rho(z)h(z) = g(z)$$

rezultă

$$h(A)\rho(A) = \rho(A)h(A) = g(A)$$

care împreună cu nesingularitatea matricei $h(A)$ conduc la (2.29). \diamond

Propoziția 2.9 *Dacă funcția compusă $h = g \circ f$ este definită pe spectrul matricei $A \in \mathcal{C}^{n \times n}$ atunci*

$$h(A) = g(f(A)) \quad (2.30)$$

i.e. $h(A) = g(B)$ unde $B = f(A)$.

Demonstrația este propusă cititorului. \diamond

2.2 Calculul funcțiilor de matrice

Tehnicile numerice de evaluare a funcțiilor de matrice, recomandate de experiența numerică acumulată, evită determinarea efectivă a polinomului de interpolare Lagrange-Sylvester, în primul rând din motive de eficiență [1]. Metodele care s-au impus – și pentru care există și unele rezultate teoretice privind problemele legate de stabilitatea numerică – se pot împărți în două categorii:

- a) metode bazate pe calculul valorilor proprii;
- b) metode aproximative bazate pe trunchierea unor dezvoltări în serie.

Prezentăm, în continuare, procedurile bazate pe calculul valorilor proprii; metodele aproximative, bazate pe trunchierea unor dezvoltări în serie, vor fi expuse în paragraful următor, dedicat cazului particular dar important pentru teoria sistemelor liniare, al calculului exponențialei matriceale.

Pentru matricele simple (i.e. diagonalizabile) calculul funcțiilor de matrice prin evaluarea valorilor și vectorilor proprii se bazează pe propozițiile 2.2 și 2.3. Într-adevăr, în acest caz, dacă $V = [v_1 \ v_2 \ \dots \ v_n]$ este matricea vectorilor proprii ai matricei date $A \in \mathcal{C}^{n \times n}$ atunci $A = V\Lambda V^{-1}$ unde $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$.

Prin urmare, pentru orice funcție definită pe spectrul matricei A , aplicând (2.19) și (2.21), rezultă

$$F = f(A) = V f(\Lambda) V^{-1} = V \text{diag}(f(\lambda_1), f(\lambda_2), \dots, f(\lambda_n)) V^{-1}. \quad (2.31)$$

Introducând acum funcția **vvp** de calcul a valorilor și vectorilor proprii ai unei matrice (desigur, prin aplicarea algoritmului QR) cu sintaxa

$$[V, p] = \mathbf{vvp}(A)$$

unde p este vectorul valorilor proprii pentru matricea A , relația (2.31) stă la baza următorului algoritm de calcul al funcțiilor de matrice diagonalizabile.

Algoritm 2.1 (Date matricea simplă $A \in \mathcal{C}^{n \times n}$ și funcția $f : D \subset \mathcal{C} \rightarrow \mathcal{C}$ definită pe spectrul matricei A , algoritmul calculează $F = f(A)$ prin determinarea valorilor și vectorilor proprii.)

1. $[V, p] = \mathbf{vvp}(A)$
2. Pentru $i = 1 : n$
 1. $\beta_i = f(p_i)$
3. $D = V \text{diag}(\beta_1, \beta_2, \dots, \beta_n)$

4. Se rezolvă sistemul matriceal liniar nesingular $FV = D$ în raport cu F .

Evident, efortul de calcul principal este destinat calculului valorilor și vectorilor proprii și, într-o oarecare măsură, rezolvării sistemului matriceal liniar.

În cazul funcțiilor f reale pe spectrul matricei reale $A \in \mathcal{R}^{n \times n}$ (vezi definiția 2.4), conform propoziției 2.1 matricea rezultat F este reală. Întrucât matricele V , D și vectorul p din algoritmul 2.1 nu sunt în mod necesar reale, calculul (aproximativ) în aritmetică complexă conduce la apariția unor componente parazite imaginare în soluția calculată. De aceea, într-un astfel de caz, algoritmul 2.1 se completează cu o instrucțiune de eliminare a acestor componente:

5. $F = \text{Re}(F)$.

Dacă matricea A nu este diagonalizabilă, o extindere a algoritmului 2.1 apelează la calculul formei Jordan a matricei A , iar aplicarea propozițiilor 2.1 și 2.2 reduce calculul funcției de matrice la calculul funcțiilor având ca argumente blocurile Jordan. Această tehnică nu este însă recomandată în primul rând datorită unei sensibilități ridicate a structurii Jordan în raport cu perturbațiile ce apar la nivelul datelor de intrare.

Proprietăți numerice mult mai bune, în toate cazurile, se obțin dacă în locul formei canonice Jordan se utilizează forma Schur complexă (sau reală) a matricei A . Această alternativă este condiționată de existența unui algoritm eficient de calcul al funcțiilor de matrice triunghiulare (sau cvasitriunghiulare). Un astfel de algoritm, propus de B.N. Parlett [1], are la bază propoziția 2.3 și proprietatea de comutativitate formulată în propoziția 2.6. Vom deduce algoritmul pentru cazul generic al matricelor cu valori proprii distincte.

Fie $T \in \mathcal{C}^{n \times n}$ o matrice superior triunghiulară cu $t_{ii} \neq t_{jj}$ pentru orice $i \neq j$ și f o funcție definită pe spectrul lui T . Dacă $F = f(T)$ atunci, conform (2.23), avem

$$FT = TF. \quad (2.32)$$

Ținând cont, conform propoziției 2.4, că matricea F este, de asemenea, superior triunghiulară și scriind (2.32) pe elemente obținem

$$\sum_{k=i}^j t_{ik} f_{kj} = \sum_{k=i}^j f_{ik} t_{kj}, \quad j = 1 : n, \quad i = 1 : j, \quad (2.33)$$

de unde, în ipotezele menționate, avem

$$f_{ij} = \frac{1}{t_{jj} - t_{ii}} [t_{ij}(f_{jj} - f_{ii}) + \sum_{k=i+1}^{j-1} (t_{ik}f_{kj} - f_{ik}t_{kj})], \quad j = 1 : n, i = 1 : j. \quad (2.34)$$

Relația (2.34) este utilă numai în măsura în care se poate găsi o ordine de calcul a elementelor matricei F astfel încât, la fiecare moment al procesului de calcul, în membrul drept al expresiei (2.34) să apară numai elemente deja calculate. O astfel de ordine există și ea poate fi evidențiată observând că elementele diagonale sunt calculabile, conform (2.22), cu formula

$$f_{ii} = f(t_{ii}) \quad (2.35)$$

și că în membrul drept al relației (2.34) apar elementele $F(i, i : j - 1)$ din "stânga" elementului f_{ij} și $F(i + 1 : j, j)$ de "sub" elementul f_{ij} (vezi diagrama din figura 2.1).

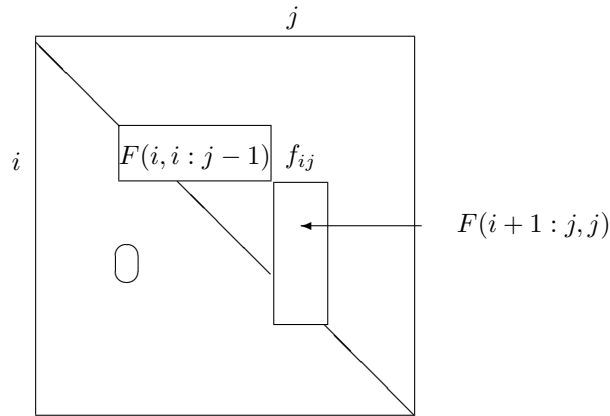


Figura 2.1: Matricea $F = f(T)$.

De exemplu, se poate adopta o ordine diagonală de calcul al elementelor triunghiului superior, după cum este indicat în diagrama din figura 2.2 pentru $n = 5$, numărul înscris în matrice marcând numărul de ordine pentru calculul elementului de pe poziția respectivă.

Aceasta nu este singura ordine posibilă. Într-adevăr, dacă efectuăm calculele pe coloane în ordinea $j = 1, 2, \dots, n$, pe fiecare coloană calculele

1	6	10	13	15
	2	7	11	14
		3	8	12
			4	9
				5

Figura 2.2: Ordinea diagonală de calcul a elementelor matricei $F = f(T)$.

efectuându-se de jos în sus (vezi diagrama a) din figura 2.3) sau pe linii în ordinea $i = n, n - 1, \dots, 1$, pe fiecare linie ordinea de calcul fiind de la stânga la dreapta (diagrama b) din figura 2.3), atunci în momentul calculului unui element curent oarecare, elementele din stânga și de sub elementul curent sunt deja calculate.

1	3	6	10	15
	2	5	9	14
		4	8	13
			7	12
				11

(a)

11	12	13	14	15
	7	8	9	10
		4	5	6
			2	3
				1

(b)

Figura 2.3: Ordinea "pe coloane" (a) și ordinea "pe linii" (b) de calcul a elementelor matricei $F = f(T)$.

Prezentăm în continuare algoritmul corespunzător ordinii diagonale de calcul, celelalte variante făcând obiectul unor exerciții. Pentru a urmări mai ușor indexarea, observăm că, atribuind indicele q direcțiilor paralele cu diagonala principală a matricei, rezultă următoarea schemă de calcul.

1. Se calculează f_{ii} , $i = 1 : n$ cu relația (2.35).
2. Pentru $q = 2 : n$
 1. Pentru $i = 1 : n - q + 1$
 1. Se calculează $f_{i,i+q-1}$ cu relația (2.34).

Pentru a nu modifica indexarea folosită în expresia (2.34) a elementului f_{ij} facem următoarele schimbări de indici:

$$p = q - 1, \quad j = i + p,$$

cu care schema de mai sus ne conduce la următorul algoritm de calcul al funcțiilor de matrice triunghiulare cu elementele diagonale distincte.

Algoritmul 2.2 (Parlett)(Date o matrice superior triunghiulară $T \in \mathcal{C}^{n \times n}$ cu $t_{ii} \neq t_{jj}$ pentru orice $i \neq j$ și o funcție $f : D \subset \mathcal{C} \rightarrow \mathcal{C}$ definită pe spectrul matricei T , algoritmul calculează, în ordine diagonală, elementele matricei superior triunghiulare $F = f(T)$.)

1. Pentru $i = 1 : n$
 1. $f_{ii} = f(t_{ii})$
2. Pentru $p = 1 : n - 1$
 1. Pentru $i = 1 : n - p$
 1. $j = i + p$
 2. $s = t_{ij}(f_{jj} - f_{ii})$
 3. Dacă $p > 1$ atunci
 1. Pentru $k = i + 1 : j - 1$
 1. $s = s + t_{ik}f_{kj} - f_{ik}t_{kj}$
 4. $f_{ij} = \frac{s}{t_{jj} - t_{ii}}$

Schema recurentă care stă la baza algoritmului 2.2 este deosebit de eficientă, în afara celor n evaluări de funcții scalare de la instrucțiunea 1, algoritmul necesitând un număr asimptotic de $N_{op} \approx \frac{2n^3}{3}$ operații în virgulă mobilă.

Cu acest algoritm, completat cu procedura de aducere a unei matrice date la forma Schur complexă (superior triunghiulară) prin transformări ortogonale de asemănare, se obține o procedură cu bune calități numerice pentru calculul funcțiilor de matrice. Introducem funcția **fsc** de calcul a formei Schur complexe având sintaxa

$$[S, U] = \mathbf{fsc}(A)$$

unde datele de ieșire sunt matricea unitară de transformare U și forma Schur complexă S a matricei A legate prin relațiile

$$S = U^H A U, \quad A = U S U^H, \quad (2.36)$$

indicele superior H semnificând dubla operație de transpunere și conjugare. Notăm cu **funmt** ("funcție de matrice triunghiulară") funcția de calcul realizată de algoritmul 2.2 de mai sus, cu sintaxa (neformală)

$$F = \mathbf{funmt}(f, T)$$

având ca parametri de intrare funcția f și matricea superior triunghiulară T , iar ca parametru de ieșire matricea $F = f(T)$. Cu aceste notații prezentăm

Algoritmul 2.3 (Date o matrice A cu valori proprii distincte și o funcție f definită pe spectrul matricei A , algoritmul calculează matricea $F = f(A)$ prin metoda aducerii matricei A la forma Schur complexă.)

1. $[S, U] = \mathbf{fsc}(A)$
2. $T = \mathbf{funmt}(f, S)$
3. $F = UTU^H$

În cazul în care matricea A are valori proprii multiple, algoritmul 2.3 nu mai este funcțional pentru că formula (2.34) nu mai este aplicabilă. De asemenea, în situația existenței unor valori proprii foarte apropiate, aplicarea formulei (2.34) conduce la fenomene de instabilitate numerică. Soluția problemei într-un astfel de caz poate consta în utilizarea unei variante "bloc" a algoritmului Parlett după o prealabilă grupare a valorilor proprii apropiate în cadrul așa numitei forme Schur complexe ordonate. Pentru detalii recomandăm [1], [VII].

2.3 Calculul exponențialei matriceale

Deși metodele de calcul prezentate în secțiunea precedentă sunt aplicabile pentru toate funcțiile definite pe spectrul matricei argument, pentru anumite funcții, de un interes aplicativ deosebit, au fost dezvoltate proceduri alternative, cu calități numerice superioare.

În acest paragraf prezentăm principalele metode pentru calculul exponențialei matriceale

$$\Phi(t) = e^{tA}, \quad (2.37)$$

unde $t > 0$ este un parametru scalar. După cum se știe, funcția $\Phi(t)$ este matricea de tranziție a stărilor sistemului liniar $\dot{x} = Ax$, adică satisface ecuația diferențială matriceală liniară $\dot{\Phi} = A\Phi$ cu condiția inițială $\Phi(0) = I$.

Având în vedere importanța unui calcul fiabil al exponențialei matriciale pentru studiul sistemelor liniare, o discuție prealabilă a condiționării numerice a problemei calculului lui $\Phi(t) = e^{tA}$, adică a sensibilității lui $\Phi(t)$ în raport cu variațiile lui A , este necesară, întrucât erorile de calcul ale lui $\Phi(t)$ se traduc, prin mecanismul de propagare inversă, în perturbații la nivelul datelor de intrare, adică al matricii A .

Notând cu $\hat{\Phi}(t)$ exponențiala asociată matricii perturbate $A + E$ avem $\dot{\hat{\Phi}} = (A + E)\hat{\Phi}$ cu $\hat{\Phi}(0) = I$. Punând $\hat{\Phi} = \Phi + \Delta\Phi$, obținem

$$\dot{\hat{\Phi}} + \Delta\dot{\Phi} = A\Phi + A\Delta\Phi + E\Phi + E\Delta\Phi,$$

unde $\Phi(0) = I$, deci $\Delta\Phi(0) = 0$. Pentru simplitate, neglijăm termenul $E\Delta\Phi$ și, prin integrare, deducem

$$\Delta\Phi(t) = \int_0^t e^{(t-\tau)A} E e^{\tau A} d\tau = \Phi(t) \int_0^t e^{-\tau A} E e^{\tau A} d\tau.$$

Utilizând norma matriceală $\|\cdot\| = \|\cdot\|_2$ putem scrie

$$\|\Delta\Phi(t)\| \leq \|\Phi(t)\| \cdot \left\| \int_0^t e^{-\tau A} E e^{\tau A} d\tau \right\|$$

relație din care rezultă

$$\frac{\|\Delta\Phi(t)\|}{\|\Phi(t)\|} \leq \text{cond}(\Phi(t)) \frac{\|E\|}{\|A\|}, \quad (2.38)$$

unde

$$\text{cond}(\Phi(t)) \stackrel{\text{def}}{=} \max_{\|E\|=1} \left\| \int_0^t e^{-\tau A} E e^{\tau A} d\tau \right\| \cdot \|A\| \quad (2.39)$$

este *numărul de condiție* al lui A relativ la calculul lui $\Phi(t) = e^{tA}$.

Obținem o evaluare grosieră a lui $\text{cond}(\Phi(t))$ observând că

$$\left\| \int_0^t e^{-\tau A} E e^{\tau A} d\tau \right\| \leq \int_0^t \|e^{-\tau A}\| \cdot \|E\| \cdot \|e^{\tau A}\| d\tau$$

unde

$$\|e^{\pm\tau A}\| \leq e^{\tau\|A\|}, \quad \tau > 0.$$

De asemenea, avem

$$\int_0^t e^{2\tau\|A\|} d\tau = \frac{1}{2\|A\|} (e^{2t\|A\|} - 1),$$

deci

$$\text{cond}\Phi(t) \leq \frac{1}{2}(e^{2t\|A\|} - 1). \quad (2.40)$$

În concluzie, pentru a limita sensibilitatea lui $\Phi(t)$ în raport cu variațiile lui A e necesar să limităm $t\|A\|$, de exemplu astfel încât

$$t\|A\| \leq 1, \quad (2.41)$$

fie alegând t suficient de mic, fie (în cazul în care t este impus) utilizând proprietatea

$$e^{tA} = (e^{\frac{t}{2^m}A})^{2^m}, \quad m \geq 1. \quad (2.42)$$

Într-adevăr, oricare ar fi $t\|A\|$ există m astfel încât $\frac{t}{2^m}\|A\| \leq 1$ iar dacă $e^{\frac{t}{2^m}A}$ e cunoscut atunci, conform (2.42), e^{tA} poate fi calculat printr-un proces de ridicare succesivă la pătrat. O evaluare a lui m rezultă observând că dacă $t\|A\| \geq 1$ atunci trebuie să avem $\log_2(t\|A\|) < m$, deci

$$m = 1 + [\log_2(t\|A\|)], \quad (2.43)$$

unde $[\cdot]$ este partea întreagă a argumentului.

În practică m se poate alege printr-un procedeu de înjumătățire succesivă. Schema de calcul, numită ”de înjumătățire și ridicare la pătrat” (*scaling and squaring*), pe scurt $(1/2)^2$, este următoarea:

Schema $(1/2)^2$

1. Se calculează $\|A\|$.
2. Se inițializează $m = 0$.
3. Cât timp $t\|A\| \geq 1$
 1. $t \leftarrow t/2$
 2. $m \leftarrow m + 1$
4. Se calculează $F = e^{tA}$.
5. Cât timp $m \geq 1$
 1. $F \leftarrow F^2$
 2. $m \leftarrow m - 1$

Observația 2.3 Având în vedere necesitatea limitării normei lui tA la calculul exponențialei, schema $(1/2)^2$ se aplică întotdeauna indiferent de metoda utilizată pentru calculul efectiv (Taylor, Padé, etc.). Trebuie însă

menționat că ridicarea la pătrat (adică înmulțirea a două matrice) introduce erori de rotunjire suplimentare, mai precis calculul lui F^2 se face cu erori mici față de $\|F\|^2$, dar nu neapărat față de $\|F^2\|$. Din această cauză utilizarea schemei $(1/2)^2$ ridică probleme de precizie dacă norma $\|e^{tA}\|$ e mică ². Pe de altă parte, dacă exponențiala e utilizată pentru simularea evoluției unui sistem liniar, o condiție de tipul $t\|A\| < 1$ orientează alegerea pasului de discretizare și, în consecință, poate fi asigurată de la început. \diamond

În ipoteza că $t\|A\|$ e limitată (de exemplu cu schema $(1/2)^2$), metodele uzuale de calcul al exponențialei matriceale constau în construcția unei aproximări *locale* în jurul lui $z = 0$ a funcției e^z ³. De regulă aproximările utilizate sunt de tip polinomial (Taylor) sau rațional (Padé) și permit calculul lui e^{tA} prin evaluarea aproximării considerate pentru $z = tA$.

I. Aproximația polinomială (Taylor) de ordin p este de forma

$$T_p(z) = \sum_{k=0}^p c_k z^k \quad (2.44)$$

²Acest lucru e posibil chiar dacă $t\|A\| \gg 1$ deoarece marginea $\|e^{tA}\| \leq e^{t\|A\|}$ utilizată mai sus este în general supraestimată (de exemplu $A = -5$, $t = 1$). De aceea este util să dispunem pentru evaluarea lui $\text{cond}(\Phi(t))$ de estimări mai precise pentru $\|e^{tA}\|$, eventual capabile să țină seama de repartiția spectrului lui A și în special de "dispersia" acestuia ("stiffness"). Se știe că dacă A are valorile proprii λ_i , atunci e^{tA} are valorile proprii $e^{\lambda_i t}$ (propoziția 2.5); pe de altă parte pentru orice matrice A și orice normă matriceală $\|\cdot\|$ avem $\|A\| \geq \rho(A)$, unde $\rho(A) \stackrel{\text{def}}{=} \max |\lambda_i|$, deci $\|e^{tA}\| \geq \max |e^{\lambda_i t}|$. Se poate arăta că dacă $\text{Re } \lambda_i < \alpha$ atunci există M astfel încât $\|e^{tA}\| \leq M e^{\alpha t}$, $t \geq 0$. În mod analog, dacă $\beta < \text{Re } \lambda_i$ atunci există N astfel încât $\|e^{-tA}\| \leq N e^{-\beta t}$, $t \geq 0$. Prin urmare dacă spectrul lui A e cuprins în banda $\beta < \text{Re } \lambda_i < \alpha$ atunci

$$\left\| \int_0^t e^{-tA} E e^{tA} dt \right\| \leq MN \|E\| \int_0^t e^{(\alpha-\beta)\tau} d\tau = MN \|E\| \frac{e^{(\alpha-\beta)t} - 1}{\alpha - \beta}$$

deci

$$\text{cond } \Phi(t) \leq MN \frac{e^{\delta t} - 1}{\delta} \|A\|$$

unde

$$\delta \stackrel{\text{def}}{=} \alpha - \beta > 0$$

De aici se vede influența nefavorabilă a dispersiei δ a spectrului (cu cât δ e mai mare cu atât $\text{cond}\Phi(t)$ e mai mare) precum și faptul că numai pentru δ mic (adică $\frac{e^{\delta t} - 1}{\delta} \approx t$) $\text{cond}\Phi(t)$ e realmente limitat de $t\|A\|$.

³În cazuri speciale (matrice simetrice cu valori proprii reale și negative, cum sunt cele rezultate prin discretizarea spațială a problemelor la limite eliptice) se pot utiliza aproximări *globale* (uniforme în sens Cebâșev sau în medie pătratică) valabile pe domeniul ce conține spectrul lui A .

unde coeficienții c_k se determină din condițiile

$$\left. \frac{d^k}{dz^k} (e^z - T_p(z)) \right|_{z=0} = 0, \quad k = 0 : p. \quad (2.45)$$

Obținem

$$c_k = \frac{1}{k!}, \quad (2.46)$$

adică cea mai bună aproximare polinomială locală de ordin p coincide cu polinomul Taylor corespunzător (fapt binecunoscut, de altfel).

Exemplul 2.4 Pentru $p = 1$ din (2.44) obținem $T_1(z) = 1 + z$ iar aproximația $T_1(tA) = I + tA$ a lui e^{tA} corespunde integrării numerice a ecuației diferențiale $\dot{x}(t) = Ax(t)$ prin metoda Euler explicită de ordin 1. În mod similar se poate arăta că metoda Runge-Kutta de ordin 4 corespunde alegerii $p = 4$ în (2.44). \diamond

Putem evalua eroarea de trunchiere comisă prin utilizarea aproximației polinomiale (2.44) observând că

$$e^z - T_p(z) \stackrel{\text{def}}{=} \sum_{k \geq p+1} \frac{z^k}{k!} = \frac{z^{p+1}}{(p+1)!} \left[1 + \frac{z}{p+2} + \frac{z^2}{(p+2)(p+3)} + \dots \right]$$

unde paranteza dreaptă e majorată prin

$$1 + \frac{|z|}{1} + \frac{|z|^2}{1 \cdot 2} + \dots = e^{|z|}$$

sau (dacă $|z| < p+2$) prin

$$1 + \frac{|z|}{p+2} + \frac{|z|^2}{(p+2)^2} + \dots = \frac{1}{1 - \frac{|z|}{p+2}}.$$

Punând $z = tA$ și utilizând prima evaluare (evident, acoperitoare) obținem eroarea relativă de trunchiere sub forma

$$\frac{\|e^{tA} - T_p(tA)\|}{\|e^{tA}\|} \leq \frac{(t\|A\|)^{p+1}}{(p+1)!}. \quad (2.47)$$

Altfel spus, avem

$$\frac{\|e^{tA} - T_p(tA)\|}{\|e^{tA}\|} \leq tol, \quad (2.48)$$

unde tol este precizia de calcul impusă, dacă alegem p astfel încât

$$\frac{(t\|A\|)^{p+1}}{(p+1)!} < tol. \quad (2.49)$$

Pentru orientare, menționăm că dacă $t\|A\| < \frac{1}{2}$ și tol este de ordinul 10^{-5} , atunci condiția (2.49) impune $p = 5 \div 6$.

Evaluarea polinomului matriceal $T_p(A)$ are loc după schema

$$T_{k+1}(tA) = T_k(tA) + X_k, \quad X_k = \frac{1}{k!} t^k A^k, \quad (2.50)$$

unde între doi termeni succesivi X_k și X_{k-1} are loc relația

$$X_k = \frac{1}{k} t A X_{k-1}, \quad (2.51)$$

iar inițializările sunt, evident, $T_0(tA) = I$ și $X_0 = I$.

Ținând seama de cele spuse mai sus, algoritmul de calcul al exponențialei matriceale poate fi formulat în felul următor.

Algoritmul 2.4 (Date $A \in \mathcal{R}^{n \times n}$ și $t \in \mathcal{R}$ algoritmul calculează $F = e^{tA}$, în ipoteza $t\|A\| < 1$, utilizând aproximația Taylor).

1. Se determină p din condiția (2.49).
2. $X = I$
3. $F = I$
4. Pentru $k = 1 : p$
 1. $X \leftarrow \frac{1}{k} t A X$
 2. $F \leftarrow F + X$

Numărul de operații necesare este de ordinul $N_{op} \sim pn^3$ (pentru efectuarea produsului matriceal de la pasul 4.1).

Observația 2.4 În absența limitării lui $t\|A\|$, la pasul 4.2 pot avea loc fenomene de anulare prin scădere cu consecințe catastofale asupra preciziei. Prin urmare, algoritmul de mai sus se asociază *în mod obligatoriu* cu **schema** $(1/2)^2$. În orice caz, pentru siguranța calculului, se recomandă monitorizarea normei termenilor X . \diamond

II. Aproximația rațională (Padé) de grade (p, q) și de ordin $p + q$ este de forma

$$R_{pq}(z) = \frac{N_p(z)}{D_q(z)} = \frac{c_0 + c_1 z + \cdots + c_p z^p}{1 + d_1 z + \cdots + d_q z^q}, \quad (d_0 \stackrel{\text{def}}{=} 1), \quad (2.52)$$

unde coeficienții c_k, d_k se determină din condițiile

$$\frac{d^k}{dz^k} (e^z - R_{pq}(z)) = 0, \quad k = 0 : p + q. \quad (2.53)$$

Echivalent, se poate utiliza metoda coeficienților nedeterminați, scriind formal

$$D_q(z) \left(\sum_{k \geq 0} \frac{z^k}{k!} \right) = N_p(z)$$

și egalând coeficienții gradelor egale din cei doi membri până la gradul $p+q$ necesar. Rezultă

$$c_k = \frac{(p+q-k)!p!}{(p+q)!k!(p-k)!}, \quad d_k = \frac{(p+q-k)!q!}{(p+q)!k!(q-k)!}(-1)^k. \quad (2.54)$$

Exemplul 2.5 Pentru $q = 0$ obținem aproximația polinomială (Taylor). Primele aproximații raționale "diagonale", i.e. cu $p = q$, sunt

$$R_{11}(z) = \frac{1 + \frac{z}{2}}{1 - \frac{z}{2}}$$

cu polul $z = 2$, respectiv

$$R_{22}(z) = \frac{1 + \frac{z}{2} + \frac{z^2}{12}}{1 - \frac{z}{2} + \frac{z^2}{12}}$$

cu polii $z_{1,2} = 3 \pm i\sqrt{3}$ având $\operatorname{Re} z_{1,2} = 3$, plasați pe cercul de rază $2\sqrt{3}$. \diamond

Prin urmare, trecând la cazul matriceal, constatăm că în general existența aproximației raționale $R_{pq}(tA)$ nu este asigurată necondiționat (vezi propoziția 2.8), întrucât matricea $D_q(tA)$ poate fi singulară; de exemplu,

$$D_1(tA) = I - \frac{1}{2}tA$$

este singulară dacă $t = 1$ și A are o valoare proprie $\lambda = 2$. (În legătură cu aceasta, observăm că dacă este îndeplinită condiția $t\|A\| < 1$, adică $\|A\| < 1$, atunci cu atât mai mult $|\lambda_i| < 1$, deci această situație este imposibilă). În general, se poate arăta că toate zerourile lui $D_q(z)$ sunt situate în $\operatorname{Re} z > 0$ și se acumulează la ∞ pentru $q \rightarrow \infty$ astfel încât $D_q(tA)$ rezultă sigur nesingulară dacă este satisfăcută cel puțin una din următoarele condiții

- i) $\operatorname{Re} \lambda_i < 0$, i.e. matricea A este stabilă;
- ii) q este suficient de mare;
- iii) $t\|A\| < 2$.

În consecință, condiția de limitare a normei lui tA este utilă pentru a asigura nu numai buna condiționare a problemei calculului lui e^{tA} ci și

existența aproximației raționale considerate mai sus, independent de caracteristicile de stabilitate ale matricei date A sau de ordinul q adoptat în calcule.

Pe de altă parte, pentru a conecta exemplele 2.4 și 2.5, este instructiv să observăm că aproximația rațională $R_{11}(tA)$ corespunde integrării numerice a ecuației diferențiale $\dot{x}(t) = Ax(t)$ prin metoda Euler implicită de ordin 1.

Eroarea de trunchiere poate fi limitată alegînd convenabil gradele p, q ale aproximării (2.51). Se poate arăta că (vezi [4]) dacă $t\|A\| < \frac{1}{2}$ atunci

$$R_{pq}(tA) = e^{t(A+E)} \quad (2.55)$$

unde $AE = EA$ iar

$$\|tE\| \leq 8 \frac{p!q!}{(p+q)!} \cdot \frac{(t\|A\|)^{p+q+1}}{(p+q+1)!}. \quad (2.56)$$

Mai mult, deoarece $AE = EA$ implică $e^{t(A+E)} = e^{tA} \cdot e^{tE}$ (exercițiul 2.12), eroarea relativă de trunchiere comisă prin utilizarea aproximației Padé $R_{pq}(tA)$ poate fi evaluată acoperitor prin

$$\frac{\|e^{t(A+E)} - e^{tA}\|}{\|e^{tA}\|} \leq \|e^{tE} - I\| = \|tE + \frac{t^2 E^2}{2!} + \dots\| \leq t\|E\| e^{t\|E\|}. \quad (2.57)$$

De exemplu, în cazul aproximației polinomiale ($q = 0$) avem

$$\|tE\| \leq 8 \frac{(t\|A\|)^{p+1}}{(p+1)!} \quad (2.58)$$

și această inegalitate poate fi utilizată în locul lui (2.47). De menționat că, în orice caz, evaluarea (2.56), valabilă în cazul general $p \neq q$, este simetrică în raport cu p și q , deci nu permite alegerea lor univocă.

Avem nevoie de un criteriu suplimentar care în cazul de față este efortul de calcul necesar pentru realizarea unei aproximații de ordin $p+q$ dat. Într-adevăr, dacă $p+q$ e fixat, atunci trebuie să calculăm polinoamele $N_p(tA)$ și $D_q(tA)$ (printr-o schemă similară cu cea utilizată mai sus pentru $T_p(tA)$), iar în acest scop sunt necesare pn^3 și, respectiv, qn^3 operații. Având în vedere expresiile coeficienților c_k, d_k , cazul diagonal $p = q$ este evident optimal deoarece corespunde la $d_k = (-1)^k c_k$ adică $N_p(tA)$ și $D_p(tA)$ au (până la semn) aceiași termeni care se calculează – o singură dată pentru ambele polinoame – în pn^3 operații. Odată fixată relația $p = q$ valoarea concretă se alege din condiția $\|tE\| < eps$, mai precis

$$8 \frac{(p!)^2}{(2p)!} \cdot \frac{(t\|A\|)^{2p+1}}{(2p+1)!} \leq eps, \quad (2.59)$$

adică, deoarece exponentul din (2.59) este $2p+1$ în loc de $p+1$ din (2.49), p corespunzător aproximației raționale e aproximativ jumătate din p corespunzător aproximației polinomiale. Altfel spus, cu același p aproximația rațională asigură un ordin de precizie dublu.

Polinoamele $N_p(tA)$ și $D_p(tA)$ se evaluează după schema cunoscută, adică

$$N_{k+1}(tA) = N_k(tA) + X_k, \quad X_k = c_k t^k A^k \quad (2.60)$$

unde

$$X_k = \frac{p-k+1}{k(2p-k+1)} t A X_{k-1}. \quad (2.61)$$

Rezulta următorul algoritm.

Algoritmul 2.5 (Padé) (Date $A \in \mathcal{R}^{n \times n}$ și $t \in \mathcal{R}$ astfel încât $t\|A\| < \frac{1}{2}$, algoritmul calculează $F = e^{tA}$ utilizând aproximația Padé).

1. Se determină p din condiția (2.59).
2. $X = I$
3. $N = I$
4. $D = I$
5. Pentru $k = 1 : p$
 1. $X \leftarrow \frac{p-k+1}{k(2p-k+1)} t A X$
 2. $N \leftarrow N + X$
 3. $D \leftarrow D + (-1)^k X$
6. Se rezolvă ecuația matriceală $DF = N$ în raport cu F .

Numărul de operații necesar este de ordinul $N_1 = pn^3$ (pentru efectuarea produsului matriceal de la pasul 5.1.), $N_2 = \frac{n^3}{3}$ (pentru triangularizarea lui D) și $N_3 = n^3$ (pentru calculul celor n coloane ale lui F), adică, în total, $N_{op} \approx (p + \frac{4}{3})n^3$, comparabil cu cel obținut anterior.

Observația 2.5 Pentru siguranța calculului, în afară de monitorizarea normei termenilor X , se recomandă testarea lui $\text{cond}(D)$ la pasul 6. \diamond

Programe MATLAB disponibile

Pentru calculul unei funcții de matrice cu valori proprii distincte se poate folosi funcția **funm**, care implementează algoritmul 2.3. Pentru calculul exponențialei matriceale este disponibilă funcția **expm**, care implementează algoritmul 2.5, bazat pe aproximarea Padé. (Opțional, pentru

compararea rezultatelor, se poate recurge la utilizarea unor versiuni ale funcției **expm** ce implementează algoritmi 2.1 și 2.3).

Alte proceduri de discretizare aproximativă a unor sisteme de ecuații diferențiale liniare, utilizate tradițional în analiză și simulare (aproximațiile de ordin zero și unu, aproximația Tustin, etc.) fac obiectul funcției MATLAB **c2dm**.

Exerciții

E 2.1 Se consideră dată matricea

$$A = \begin{bmatrix} 3 & -3 & 2 \\ -1 & 5 & -2 \\ -1 & 3 & 0 \end{bmatrix}.$$

Să se calculeze următoarele funcții de matrice folosind definiția 2.3 și verificând în prealabil că funcțiile respective sunt definite pe spectrul matricei argument: a) $F = A^{-1}$, b) $G = \sqrt{A}$, c) $H = e^{A \ln 2}$, d) $S_1 = \sin(\frac{\pi}{4}A)$, e) $C_1 = \cos(\frac{\pi}{4}A)$, f) $S_2 = \sin(\frac{\pi}{8}A)$, g) $C_2 = \cos(\frac{\pi}{8}A)$.

E 2.2 Verificați că matricele calculate în exercițiul precedent satisfac relațiile: a) $G^2 = A$ (există matrice $X \neq G$ care satisfac $X^2 = A$?), b) $S_1 = 2S_2C_2 = 2C_2S_2$, c) $C_1^2 = C_2^2 - S_2^2$, etc. Puteți generaliza aceste rezultate arătând, e.g., că $(\sqrt{A})^2 = A$, $\sin 2A = 2 \sin A \cos A$, $\cos 2A = 2(\cos A)^2 - I$ etc, oricare ar fi matricea patrată A pe spectrul căreia funcțiile respective sunt definite ?

E 2.3 Reluați matricea A din exercițiul 2.1 și calculați funcțiile de matrice cerute pe baza teoremei 2.1 (relațiile (2.15)).

E 2.4 Să se deducă expresia analitică a matricei $F = f(J)$ unde $J \in \mathcal{C}^{m \times m}$ este un bloc Jordan iar f este o funcție definită pe spectrul lui J (vezi observația 2.1).

E 2.5 Verificați că matricea A din exercițiul 2.1 este simplă și determinați un set complet de vectori proprii al acesteia. Pe această bază calculați funcțiile de matrice cerute utilizând algoritmul 2.1.

E 2.6 Fie $A \in \mathcal{C}^{n \times n}$ o matrice simplă și v_j , $j = 1 : n$ un set complet de vectori proprii al acestei matrice. Dacă $u_j = V^{-1}(:, j)$, unde V este matricea vectorilor proprii menționați, iar f este o funcție definită pe spectrul

$\lambda(A) = \{ \lambda_1, \lambda_2, \dots, \lambda_n \}$ al matricei date să se arate că

$$F \stackrel{\text{def}}{=} f(A) = \sum_{j=1}^n f(\lambda_j) v_j u_j^H.$$

Folosind setul de vectori proprii al matricei A din exercițiul 2.1, calculat în exercițiul precedent, aplicați formula de mai sus pentru calculul funcțiilor de matrice cerute.

E 2.7 Rescrieți algoritmul Parlett de calcul al funcțiilor de matrice superior triunghiulare cu valori proprii distincte calculând elementele matricei $F = f(T)$ în ordinea "pe coloane", respectiv "pe linii" (vezi figura 2.3). Adaptați acest algoritm pentru calculul funcțiilor de matrice inferior triunghiulare cu elementele diagonale distincte.

E 2.8 Scrieți un algoritm eficient pentru calculul funcției $F = f(A)$, unde A este superior bidiagonală cu elemente diagonale distincte.

E 2.9 Scrieți un algoritm eficient pentru calculul unei funcții $F = f(A)$, unde $A \in \mathcal{R}^{n \times n}$ este simetrică. Considerați cazurile $f(z) = e^z$, $f(z) = \ln z$, $f(z) = \sqrt{z}$.

Indicație. Utilizând algoritmul **QR** simetric se obține $A = Q^T \Lambda Q$, unde Q este ortogonală iar $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ conține valorile proprii (reale) ale lui A . Apoi se aplică (2.19).

E 2.10 Scrieți un algoritm eficient pentru calculul matricei de tranziție $\Phi(k) = A^k$ a sistemului linear discret $x(k+1) = Ax(k)$.

Indicație. Executând $A \leftarrow A * A$ de zece ori se obține $A \leftarrow A^{1024}$. Cum procedați dacă, de exemplu, $k = 1453$, $k = 1789$ sau $k = 1917$?

E 2.11 Discutați metodele de calcul ale exponențialei matriceale e^{tA} , unde A este o matrice companion, e.g. are forma

$$A = \begin{bmatrix} 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ -a_1 & -a_2 & \cdots & -a_n \end{bmatrix},$$

și scrieți algoritmul de calcul pe care îl recomandați.

Indicație. Polinomul caracteristic al lui A este $p(s) \stackrel{\text{def}}{=} \det(sI_n - A) = s^n + a_n s^{n-1} + \cdots + a_1$. Utilizând teorema Cayley - Hamilton, conform

căreia $p(A) = 0$, precum și dezvoltarea în serie a lui e^z , găsim

$$e^{tA} = \sum_{i=1}^n \alpha_i(t) A^{i-1},$$

unde coeficienții funcțiilor analitice $\alpha_i(t)$ pot fi calculați recurent. Alternativ, utilizând transformata Laplace a exponențialei, putem scrie

$$(sI - A)^{-1} = \sum_{k=1}^n \frac{s^{k-1}}{p(s)} A_k,$$

unde atât matricele A_k cât și transformatele inverse ale funcțiilor $s^{k-1}/p(s)$ pot fi, în principiu, evaluate numeric. Experiența numerică arată că, în ambele situații, efectul erorilor de rotunjire și trunchiere poate fi catastrofal, ceea ce confirmă încă o dată faptul că o metodă aplicabilă în regim de lucru "cu hârtia și creionul" este aproape sigur contraindicată în calitate de algoritm de uz curent (pe calculator). În cazul de față se recomandă utilizarea algoritmului 2.5, în care structura rară a lui A reduce considerabil numărul de operații la pasul 5.1.

E 2.12 Arătați că $e^{A+B} = e^A e^B$ este adevărată dacă și numai dacă $AB = BA$. Puteți da o expresie (aproximativă) pentru e^{A+B} dacă $AB \neq BA$?

E 2.13 Scrieți *schema* $(1/2)^2$ pentru calculul funcțiilor $\cos(tA)$ și $\sin(tA)$ și argumentați utilitatea ei. Idem pentru funcțiile $\cosh(tA)$ și $\sinh(tA)$.

E 2.14 Scrieți analogul algoritmilor 2.4 și 2.5 pentru calculul funcțiilor din exercițiul 2.13.

Observație. Soluția ecuației diferențiale $\ddot{x}(t) + Ax(t) = 0$ cu condițiile inițiale $x(0) = x$, $\dot{x}(0) = y$ se scrie sub forma

$$x(t) = (\cos A^{\frac{1}{2}}t) x + (A^{-\frac{1}{2}} \sin A^{\frac{1}{2}}t) y.$$

Examinând seriile corespunzătoare, explicați cum se aplică rezultatele de mai sus în această situație. Cum procedați dacă A este simetrică și pozitiv definită ?

E 2.15 Stabiliți legătura dintre funcțiile considerate mai sus și exponențiala matriceală e^{tA^0} , unde

$$A^0 = \begin{bmatrix} 0 & I_n \\ -A & 0 \end{bmatrix}.$$

Ce implicații calculatorii are constatarea făcută?

E 2.16 Fie $A \in \mathcal{R}^{n \times n}$ și $B \in \mathcal{R}^{n \times m}$ două matrice date. Utilizând definiția matricei de tranziție, precizați structura exponențialei matriceale e^{tA^0} , unde

$$A^0 = \begin{bmatrix} A & B \\ 0 & 0 \end{bmatrix}.$$

Indicație. Pe baza rezultatelor din § 2.1, se constată imediat că e^{tA^0} are structura

$$e^{tA^0} = \begin{bmatrix} F(t) & G(t) \\ 0 & I_m \end{bmatrix}.$$

Cine sunt $F(t)$ și $G(t)$?

Observație. Aceste rezultate sunt utilizate în secțiunea 4.3 în legătură cu problema discretizării sistemelor liniare.

E 2.17 Generalizați constatarea făcută în exercițiul 2.16 pentru a obține o procedură eficientă de calcul a integralelor

$$G_k(t) = \int_0^t e^{\eta A} \frac{(t-\eta)^{k-1}}{(k-1)!} d\eta, \quad k = 1 : p.$$

Indicație. Consultați [6].

E 2.18 Utilizând seria lui $e^{\eta A}$ și efectuând integrarea, obțineți dezvoltările în serie ale funcțiilor $G_k(t)$ și scrieți un algoritm de tip 2.4 pentru evaluarea lor. Estimați eroarea de trunchiere.

Bibliografie

- [1] **Parlett B. N.** *A Recurrence among the Elements of Functions of Triangular Matrices*, Lin. Alg. and Its Applic., vol. 14, pp. 117–121, 1976.
- [2] **Ward R. C.** *Numerical Computation of the Matrix Exponential with Accuracy Estimate*, SIAM J. Numer. Anal., vol. 14, pp. 600–610, 1977.
- [3] **Van Loan C. F.** *The Sensitivity of the Matrix Exponential*, SIAM J. Numer. Anal., vol. 6, pp. 971–981, 1977.

- [4] **Moler C. B., Van Loan C. F.** *Nineteen Dubious Ways to Compute the Exponential of a Matrix*, SIAM Review, vol. 20, pp. 801–836, 1978.
- [5] **Van Loan C. F.** *A Note on the Evaluation of Matrix Polynomials*, IEEE Trans. Automat. Contr., vol. AC-24, pp. 320–321, 1978.
- [6] **Van Loan C. F.** *Computing Integrals Involving the Matrix Exponential*, IEEE Trans. Automat. Contr., vol. AC-24, pp. 395–404, 1978.

Capitolul 3

Tehnici de procesare a modelelor sistemice liniare

În acest capitol ne vom ocupa de reprezentările numerice ale modelelor matematice ale sistemelor liniare și vom prezenta principalele proceduri de calcul privitoare la utilizarea și manipularea acestora în rezolvarea problemelor de analiză și sinteză sistemică.

3.1 Modele sistemice liniare

În general, un sistem liniar este definit printr-un model de stare de forma

$$(\mathbf{S}) \quad \begin{cases} \dot{x}(t) = Ax(t) + Bu(t) \\ y(t) = Cx(t) + Du(t) \end{cases} \quad (3.1)$$

unde vectorii $x(t) \in \mathcal{R}^n$, $u(t) \in \mathcal{R}^m$ și $y(t) \in \mathcal{R}^l$ sunt *starea*, *intrarea* și, respectiv, *ieșirea* sistemului la momentul curent $t \in \mathcal{R}$ iar A, B, C, D sunt matrice constante de dimensiuni corespunzătoare. Dimensiunea n a spațiului stărilor $\mathcal{X} = \mathcal{R}^n$ se numește *ordinul* (sau *dimensiunea*) lui (\mathbf{S}) .

Pe scurt, un model de stare al unui sistem liniar este definit de un cvartet de matrice și în consecință se notează $S = (A, B, C, D)$.

Semnificația dinamică a modelului de stare (\mathbf{S}) poate fi rezumată astfel. Dacă starea inițială $x(0) = x$ și funcția de intrare $u(t)$, $t \geq 0$ sunt precizate, atunci traiectoria de stare $x(t)$, $t \geq 0$ este soluția ecuației diferențiale din (3.1) iar funcția de ieșire $y(t)$, $t \geq 0$ rezultă din a doua relație (3.1). (Metodele de calcul corespunzătoare sunt expuse în capitolul 4). Din acest

punct de vedere, matricea $A \in \mathcal{R}^{n \times n}$ caracterizează dinamica internă a sistemului (adică evoluția acestuia la nivelul spațiului stărilor), coloanele matricei $B \in \mathcal{R}^{n \times m}$ definesc canalele de intrare, liniile matricei $C \in \mathcal{R}^{l \times n}$ definesc canalele de ieșire iar elementele matricei $D \in \mathcal{R}^{l \times m}$ reprezintă căile de transfer direct intrare-ieșire. Dacă $D = 0$ atunci sistemul S , notat pe scurt $S = (A, B, C)$, este pur dinamic, în sensul că orice transfer intrare-ieșire este posibil numai prin intermediul modificării (dinamice) a stării $x(t)$, $t \geq 0$. Aceste idei sunt evidențiate de schema bloc asociată lui (S) (figura 3.1).

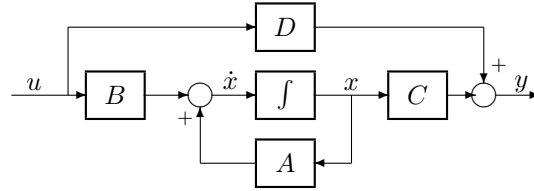


Figura 3.1: Structura modelului de stare al unui sistem liniar.

Modelul de stare (S) este invariant la o transformare liniară de stare

$$\tilde{x} = Tx, \quad (3.2)$$

unde $T \in \mathcal{R}^{n \times n}$ este o matrice *nesingulară* arbitrară, în sensul că noul vector de stare \tilde{x} satisface ecuații (\tilde{S}) de același tip cu (S) , în care matricele corespunzătoare sunt

$$\begin{aligned} \tilde{A} &= TAT^{-1}, & \tilde{B} &= TB, \\ \tilde{C} &= CT^{-1}, & \tilde{D} &= D. \end{aligned} \quad (3.3)$$

Modelele (S) și (\tilde{S}) legate prin relațiile (3.3) se numesc *echivalente* (sau *asemenea*) și sunt indiscernabile prin experimente intrare-ieșire, deci reprezintă un același sistem liniar considerat modulo relația de echivalență (3.3). Pe scurt, au sens sistemic numai acele proprietăți ale sistemului (S) care sunt *invarianti* ai lui (S) în raport cu transformările (3.3).

Din motive de eficiență și siguranță a calculului, deseori vom restrânge clasa transformărilor utilizând în (3.3) numai matrice $T = U$ *ortogonale*. În acest caz, vom spune că modelele (S) și (\tilde{S}) , legate prin relațiile

$$\begin{aligned} \tilde{A} &= UAU^T, & \tilde{B} &= UB, \\ \tilde{C} &= CU^T, & \tilde{D} &= D, \end{aligned} \quad (3.4)$$

sunt *ortogonal echivalente* și, în mod corespunzător, vom acorda o atenție specială *invariantilor ortogonali* ai lui (\mathbf{S}) . În legătură cu aceste noțiuni, observăm că orice echivalență ortogonală este o echivalență, deci clasa invariantilor ortogonali este mai largă decât clasa invariantilor (orice invariant este invariant ortogonal dar reciproca este falsă!). De aceea, invariantii ortogonali au deseori proprietăți de robustețe și permit evaluări cantitative, ceea ce le conferă avantaje de calcul și interpretare sistemică importante. Pentru concretizarea acestor idei în legătură cu noțiunile fundamentale de controlabilitate și observabilitate, cititorul poate consulta capitolul 5.

În practică se consideră deseori modele de stare generalizate (de tip „descriptor”) de forma

$$(\mathbf{S}') \quad \begin{cases} E\dot{x}(t) = Ax(t) + Bu(t) \\ y(t) = Cx(t) + Du(t) \end{cases}, \quad (3.5)$$

unde în plus față de (\mathbf{S}) , aici apare și matricea *nesingulară* $E \in \mathcal{R}^{n \times n}$.

Interesant este cazul în care E este aproape singulară astfel încât (la fel ca în cazul fasciculelor) formarea produselor $E^{-1}A$, $E^{-1}B$ (respectiv AE^{-1} , CE^{-1}) este nedorită. În acest caz transformările (3.3), (3.4) se scriu

$$\begin{aligned} \tilde{E} &= TES, & \tilde{A} &= TAS, & \tilde{B} &= TB, \\ \tilde{C} &= CS, & \tilde{D} &= D, \end{aligned} \quad (3.6)$$

în care T, S sunt matrice nesingulare, respectiv

$$\begin{aligned} \tilde{E} &= UEV, & \tilde{A} &= UAV, & \tilde{B} &= UB, \\ \tilde{C} &= CV, & \tilde{D} &= D, \end{aligned} \quad (3.7)$$

în care U, V sunt matrice ortogonale.

Observăm că dacă $E = I$, adică (\mathbf{S}') este de forma (\mathbf{S}) , atunci obținem $\tilde{E} = I$ dacă și numai dacă $S = T^{-1}$ și în acest caz (3.6), (3.7) se reduc respectiv la (3.3), (3.4). Pe de altă parte, dacă $E \neq I$ atunci obținem $\tilde{E} = I$ dacă și numai dacă $TES = I$, de exemplu $S = I$ și $T = E^{-1}$ sau $S = E^{-1}$ și $T = I$, în ambele cazuri transformarea (care implică calculul lui E^{-1}) fiind rău condiționată.

Din punct de vedere sistemic, modelele de tip (\mathbf{S}') cu E aproape singulară corespund sistemelor liniare cu constante de timp mici (sau dinamică internă parazită) și tehnicile de aproximare a comportării acestui tip de sisteme pe baza unor modele de ordin redus se înscriu în clasa așa numitelor metode de perturbații singulare.

În ultimul timp se studiază și modele de tip (\mathbf{S}') în care matricea E este

singulară. În cazul în care fascicolul (E, A) este regulat ¹ aceste modele corespund unor sisteme „singulare” care, de exemplu, conțin elemente de derivare pură.

În practica inginerescă proprietățile de transfer intrare-ieșire ale unui sistem liniar în starea inițială nulă, adică pentru $x(0) = 0$, sunt caracterizate prin intermediul matricei de transfer $T(s)$.

De exemplu, dacă sistemul este definit printr-un model de tip **(S)** atunci avem

$$T(s) = C(sI - A)^{-1}B + D \quad (3.8)$$

și se constată ușor (pe baza expresiei lui $(sI - A)^{-1}$) că $T(s)$ este o matrice cu l linii și m coloane ale cărei elemente sunt funcții raționale de variabilă complexă s , pe scurt $T(s) \in \mathcal{R}^{l \times m}(s)$ ². De asemenea avem

$$\lim_{s \rightarrow \infty} T(s) = D, \quad (3.9)$$

deci $T(s)$ este proprie, adică fiecare element are gradul numitorului mai mare decât sau cel mult egal cu gradul numărătorului; în particular $T(s)$ este strict proprie dacă și numai dacă $D = 0$. În sfârșit, $T(s)$ este evident invariantă în raport cu transformările de echivalență (3.3) astfel încât ea are într-adevăr caracter sistemic (caracterizează efectiv sistemul, și nu doar modelul **(S)** considerat).

În teoria realizării sistemelor liniare se arată că, reciproc, orice matrice de transfer $T(s)$ admite o realizare de stare, adică există un model de tip **(S)** ale cărui matrice A, B, C, D satisfac (3.8). Mai mult, acest model este în esență unic determinat (modulo echivalența (3.3)) și de ordin minim dacă și numai dacă el este controlabil și observabil. (Altfel spus, matricea de transfer este un invariant complet pentru sistemele controlabile și observabile). Prin urmare, în toate chestiunile care privesc exclusiv proprietățile de transfer intrare-ieșire ale sistemelor liniare, modelele de stare pot fi suplinite de matricele de transfer corespunzătoare ³.

¹Un fascicol matriceal (E, A) se numește *regulat* dacă polinomul $\det(\lambda E - A)$ nu este identic nul.

²Notăm cu $\mathcal{R}[s]$ inelul polinoamelor și cu $\mathcal{R}(s)$ corpul său de fracții, i.e. corpul funcțiilor raționale, ambele cu coeficienți reali. Indicăm, ca de obicei, prin indici superiori dimensiunea vectorilor și matricelor formate cu aceste elemente.

³Afirmația nu e valabilă pentru modelele de stare necontrolabile sau/și neobservabile (utilizate în multe aplicații), pentru care nu există (și nu poate exista) o descriere adecvată în termenii matricelor de transfer. Altfel spus, în general o matrice de transfer identifică nu un sistem ci o întreagă clasă de sisteme liniare *echivalente intrare-ieșire* cu un sistem dat.

În acest sens în continuare vom considera și modele sistemice „de transfer”, definite prin matrice de forma

$$(\mathbf{T}) \quad T(s) = [T_{ij}(s)], \quad T_{ij}(s) = \frac{N_{ij}(s)}{p_{ij}(s)} \quad (3.10)$$

în care $p_{ij}, N_{ij} \in \mathcal{R}[s]$ și $\partial p_{ij} \geq \partial N_{ij}$, $i = 1 : l$, $j = 1 : m$ iar ∂ denotă gradul.

Alte forme sub care se poate prezenta modelul (\mathbf{T}) precum și procedurile de conversie corespunzătoare vor fi discutate în continuare.

Similar cu (3.1), un sistem liniar discret este definit tot printr-un cvartet de matrice $S = (A, B, C, D)$, iar matricea de transfer corespunzătoare $T(z)$ are aceeași expresie (3.8). Deoarece toate rezultatele procedurale (pur algebrice) din acest capitol se aplică identic și în cazul discret, orice referire explicită la acesta este lipsită de interes.

3.2 Conexiuni

În acest paragraf arătăm cum se construiesc modelele sistemice de tip (\mathbf{S}) și (\mathbf{T}) pentru conexiunile uzuale a două sisteme S_i , $i = 1, 2$ definite prin modele de stare

$$(\mathbf{S}_i) \quad \begin{cases} \dot{x}_i(t) = A_i x_i(t) + B_i u_i(t) \\ y_i(t) = C_i x_i(t) + D_i u_i(t), \end{cases} \quad (3.11)$$

respectiv prin matrice de transfer

$$(\mathbf{T}_i) \quad y_i(s) = T_i(s) u_i(s). \quad (3.12)$$

Construcția modelelor sistemice pentru structuri oricât de complexe se reduce, în ultimă instanță, la aplicarea repetată a procedurilor prezentate mai jos.

În toate cazurile considerate vectorul de stare al sistemului realizat prin conexiune este

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathcal{R}^n, \quad n = n_1 + n_2, \quad (3.13)$$

adică întotdeauna vectorul de stare rezultă prin agregarea vectorilor de stare ai sistemelor componente iar ordinul conexiunii este egal cu suma ordinilor sistemelor conectate.

Conexiunea paralel

Pentru a putea fi conectate în paralel, sistemele S_1 și S_2 trebuie să satisfacă următoarele condiții structurale:

$$m_1 = m_2, \quad l_1 = l_2, \quad (3.14)$$

adică să aibă același număr de intrări și același număr de ieșiri. Relațiile de interconexiune (vezi figura 3.2 (b)) sunt

$$u_1 = u_2 = u, \quad y = y_1 + y_2. \quad (3.15)$$

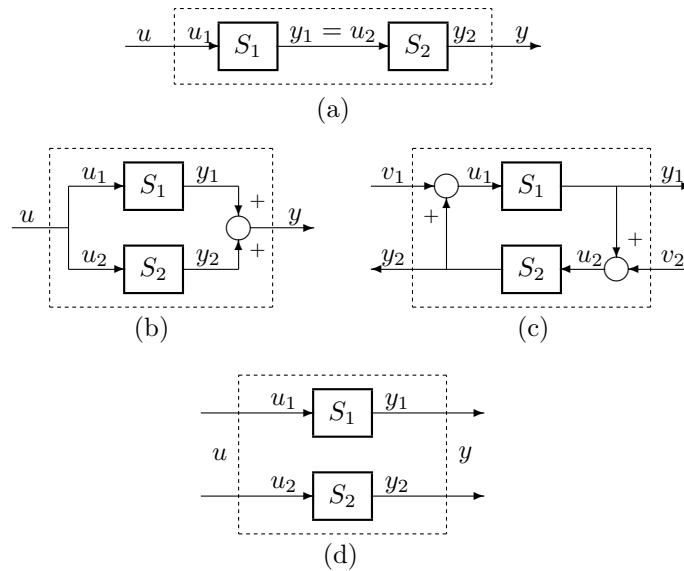


Figura 3.2: Conexiunile sistemice fundamentale: serie (a), paralel (b), reacție (c) și produs direct (d).

În consecință, ecuațiile de stare se scriu

$$\begin{cases} \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} A_1 & 0 \\ 0 & A_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} u \\ y = [C_1 \quad C_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + [D_1 + D_2] u, \end{cases} \quad (3.16)$$

unde matricele lui $S = (A, B, C, D)$ au expresii evidente, în particular A este bloc - diagonală. Matricea de transfer este

$$T(s) = T_1(s) + T_2(s), \quad (3.17)$$

deci, în cazul SISO, cu notații evidente, avem

$$T(s) \stackrel{\text{def}}{=} \frac{N(s)}{p(s)} = \frac{N_1(s)p_2(s) + N_2(s)p_1(s)}{p_1(s)p_2(s)}. \quad (3.18)$$

Observăm că $T(s)$ nu rezultă neapărat ireductibilă chiar dacă $T_i(s)$, $i = 1 : 2$, au această proprietate.

Conexiunea serie (cascadă sau tandem)

Pentru a putea fi conectate în serie, sistemele trebuie să satisfacă următoarea condiție structurală:

$$l_1 = m_2. \quad (3.19)$$

Relațiile de interconexiune (vezi figura 3.2 (a)) sunt

$$u_1 = u, \quad u_2 = y_1, \quad y = y_2. \quad (3.20)$$

În consecință, ecuațiile de stare se scriu

$$\begin{cases} \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} A_1 & 0 \\ B_2C_1 & A_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} B_1 \\ B_2D_1 \end{bmatrix} u \\ y = \begin{bmatrix} D_2C_1 & C_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} D_2D_1 \end{bmatrix} u, \end{cases} \quad (3.21)$$

unde, din nou, matricele lui S au expresii evidente, în particular A este bloc - inferior - triunghiulară. Matricea de transfer este

$$T(s) = T_2(s)T_1(s), \quad (3.22)$$

(în această ordine a matricelor factor!).

Conexiunea în circuit închis (în buclă sau cu reacție)

Pentru a putea fi conectate în circuit închis, sistemele trebuie să satisfacă următoarele condiții structurale:

$$l_2 = m_1, \quad l_1 = m_2, \quad (3.23)$$

precum și condiția „de bună formulare” a conexiunii,

$$\det(I - D_1 D_2) \neq 0. \quad (3.24)$$

Această condiție, necesară pentru a putea exprima y în mod unic funcție de x_1, x_2 precum și de intrările externe v_1, v_2 (vezi figura 3.2 (c)), este generic satisfăcută în raport cu D_1, D_2 și este satisfăcută în mod sigur dacă fie $D_1 = 0$, fie $D_2 = 0$, adică cel puțin unul dintre sistemele conectate în circuit închis este strict propriu (pur dinamic).

Relațiile de interconexiune sunt

$$u_1 = v_1 + y_2, \quad u_2 = v_2 + y_1. \quad (3.25)$$

În consecință, ecuațiile de stare ale conexiunii

$$\begin{cases} \dot{x}_1 = A_1 x_1 + B_1(v_1 + y_2) \\ \dot{x}_2 = A_2 x_2 + B_2(v_2 + y_1) \end{cases}$$

$$\begin{cases} y_1 = C_1 x_1 + D_1(v_1 + y_2) \\ y_2 = C_2 x_2 + D_2(v_2 + y_1) \end{cases}$$

se scriu sub forma evidentă

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} A_1 & 0 & | & B_1 & 0 \\ 0 & A_2 & | & 0 & B_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ v_1 \\ v_2 \end{bmatrix} + \begin{bmatrix} 0 & B_1 \\ B_2 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad (3.26)$$

$$\begin{bmatrix} I & -D_1 \\ -D_2 & I \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} C_1 & 0 & | & D_1 & 0 \\ 0 & C_2 & | & 0 & D_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ v_1 \\ v_2 \end{bmatrix}. \quad (3.27)$$

Introducând matricele produsului direct $S_d = (A_d, B_d, C_d, D_d)$ al celor două sisteme, vezi fig. 3.2 (d), unde

$$A_d = \begin{bmatrix} A_1 & 0 \\ 0 & A_2 \end{bmatrix}, \quad B_d = \begin{bmatrix} B_1 & 0 \\ 0 & B_2 \end{bmatrix},$$

$$C_d = \begin{bmatrix} C_1 & 0 \\ 0 & C_2 \end{bmatrix}, \quad D_d = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix},$$

precum și matricele de interconexiune

$$B^* = \begin{bmatrix} 0 & B_1 \\ B_2 & 0 \end{bmatrix}, \quad D^* = \begin{bmatrix} 0 & D_1 \\ D_2 & 0 \end{bmatrix}$$

ecuațiile (3.26) și (3.27) devin

$$\dot{x} = \begin{bmatrix} A_d & B_d \end{bmatrix} \begin{bmatrix} x \\ v \end{bmatrix} + B^* y$$

$$(I - D^*)y = \begin{bmatrix} C_d & D_d \end{bmatrix} \begin{bmatrix} x \\ v \end{bmatrix},$$

unde, în virtutea condiției de bună formulare, matricea $I - D^*$ este inversabilă. Prin urmare avem

$$\begin{cases} \dot{x} = \left(\begin{bmatrix} A_d & B_d \end{bmatrix} + B^*(I - D^*)^{-1} \begin{bmatrix} C_d & D_d \end{bmatrix} \right) \begin{bmatrix} x \\ v \end{bmatrix}, \\ y = \begin{matrix} (I - D^*)^{-1} \begin{bmatrix} C_d & D_d \end{bmatrix} \end{matrix} \begin{bmatrix} x \\ v \end{bmatrix}. \end{cases} \quad (3.28)$$

Matricele modelului de stare agregat $S = (A, B, C, D)$ al conexiunii în circuit închis au expresiile

$$\begin{aligned} A &= A_d + B^*(I - D^*)^{-1}C_d, & B &= B_d + B^*(I - D^*)^{-1}D_d, \\ C &= (I - D^*)^{-1}C_d, & D &= (I - D^*)^{-1}D_d. \end{aligned} \quad (3.29)$$

iar calculul se face în ordinea

$$\begin{bmatrix} C & D \end{bmatrix} = (I - D^*)^{-1} \begin{bmatrix} C_d & D_d \end{bmatrix},$$

$$\begin{bmatrix} A & B \end{bmatrix} = \begin{bmatrix} A_d & B_d \end{bmatrix} + B^* \begin{bmatrix} C & D \end{bmatrix}.$$

Matricea de transfer este

$$T(s) = (I - T^*(s))^{-1}T_d(s), \quad (3.30)$$

unde

$$T_d(s) = \begin{bmatrix} T_1(s) & 0 \\ 0 & T_2(s) \end{bmatrix}, \quad T^*(s) = \begin{bmatrix} 0 & T_1(s) \\ T_2(s) & 0 \end{bmatrix}. \quad (3.31)$$

Relațiile de mai sus pentru construcția modelelor de stare ale sistemelor agregat evidențiază faptul că în toate cazurile se face apel exclusiv la calcule matriceale elementare precum și la manipularea unor structuri matriceale organizate pe blocuri (ceea ce, de exemplu în MATLAB, este, de asemenea, elementar). În consecință, scrierea algoritmilor corespunzători este lăsată în sarcina cititorului (vezi exercițiile 3.1 și 3.2).

3.3 Realizări

În acest paragraf prezentăm modalitățile de construcție a realizărilor de stare $S = (A, B, C, D)$, nu nepărat de ordin minim, pentru sisteme definite de matrice de transfer $T(s)$ raționale *proprii*. Aceste proceduri sunt importante în aplicații, deoarece deseori sistemele liniare se descriu utilizând limbajul intrare - ieșire al matricelor de transfer dar se tratează numeric utilizând, aproape exclusiv, realizările de stare corespunzătoare.

În toate cazurile considerate, matricea D a realizării de stare se obține "extrăgând partea întreagă" a lui $T(s)$, adică scriind

$$T(s) = T'(s) + D, \quad (3.32)$$

unde $T'(s)$ este strict proprie. Această operație se poate efectua separat pentru fiecare element al lui T , mai precis, dacă

$$T_{ij}(s) \stackrel{\text{def}}{=} \frac{N_{ij}(s)}{p_{ij}(s)}, \quad (3.33)$$

unde polinomul p_{ij} este monic ⁴ și de grad n , atunci d_{ij} este coeficientul termenului de grad n al lui N_{ij} astfel încât

$$\frac{N_{ij}(s)}{p_{ij}(s)} = d_{ij} + \frac{N'_{ij}(s)}{p_{ij}(s)}, \quad N'_{ij}(s) = N_{ij}(s) - d_{ij}p_{ij}(s). \quad (3.34)$$

Altfel spus, dacă notăm $d_{ij} \stackrel{\text{def}}{=} \nu_0$, atunci dispunând coeficienții $a_0 = 1$, a_k , $k = 1 : n$ ai polinomului p_{ij} și ν_k , $k = 0 : n$ ai polinomului N_{ij} în ordinea descrescătoare a puterilor

$$\begin{array}{cccc} 1 & a_1 & \dots & a_n \\ \nu_0 & \nu_1 & \dots & \nu_n \end{array}$$

și făcând eliminare gaussiană pentru a anula ν_0 , obținem coeficienții polinomului N'_{ij} , adică

$$\nu_k \leftarrow \nu'_k = \nu_k - \nu_0 a_k, \quad k = 1 : n. \quad (3.35)$$

Operația de mai sus se repetă pentru $i = 1 : l$, $j = 1 : m$.

În consecință, rămâne să construim realizarea de stare a matricei $T'(s)$ *strict proprie* cu elementele

$$T'_{ij}(s) = \frac{N'_{ij}(s)}{p_{ij}(s)}, \quad (3.36)$$

⁴Un polinom se numește *monic* dacă are coeficientul termenului de grad maxim egal cu 1.

adică să construim $S' = (A, B, C)$ astfel încât

$$T'(s) = C(sI - A)^{-1}B. \quad (3.37)$$

Pentru simplificarea notațiilor, mai departe eliminăm indicele superior ' și considerăm succesiv câteva cazuri reprezentative de complexitate crescândă. (Subliniem din nou că, deși procedurile de realizare obținute au un caracter elementar și, ca atare, ar justifica o tratare mai expeditivă, ele sunt esențiale pentru aplicarea metodelor de calcul numeric în analiza și sinteza sistemică).

A. Sisteme simple, cu o singură intrare și o singură ieșire (SISO)

În cazul sistemelor cu o singură intrare și o singură ieșire sau, pe scurt, SISO (Single-Input, Single-Output) se dă funcția de transfer (scalară) strict proprie

$$T(s) = \frac{N(s)}{p(s)}, \quad (3.38)$$

cu

$$\begin{aligned} p(s) &= s^n + a_1 s^{n-1} + \dots + a_n, \\ N(s) &= \nu_1 s^{n-1} + \dots + \nu_n. \end{aligned} \quad (3.39)$$

O realizare de stare a lui $T(s)$ se scrie, prin inspecție, într-una din cele două forme duale următoare

1⁰. *Forma standard controlabilă*

$$(C) \quad \begin{cases} A_c = \begin{bmatrix} -a_1 & \dots & -a_{n-1} & -a_n \\ 1 & \dots & 0 & 0 \\ & \ddots & & \\ & & 1 & 0 \end{bmatrix}, & B_c = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \\ C_c = [\nu_1 \quad \dots \quad \nu_{n-1} \quad \nu_n], \end{cases} \quad (3.40)$$

sau, echivalent

$$(\mathbf{C}') \quad \begin{cases} A'_c = \begin{bmatrix} 0 & 1 & & \\ & & \ddots & \\ -a_n & -a_{n-1} & \cdots & -a_1 \end{bmatrix}, & B'_c = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \\ C'_c = [\nu_n \quad \nu_{n-1} \quad \cdots \quad \nu_1]. \end{cases} \quad (3.41)$$

2^o. Forma standard observabilă

$$(\mathbf{O}) \quad \begin{cases} A_o = \begin{bmatrix} -a_1 & 1 & & \\ \vdots & & \ddots & \\ -a_{n-1} & & & 1 \\ -a_n & 0 & \cdots & 0 \end{bmatrix}, & B_o = \begin{bmatrix} \nu_1 \\ \vdots \\ \nu_{n-1} \\ \nu_n \end{bmatrix}, \\ C_o = [1 \quad 0 \quad \cdots \quad 0], \end{cases} \quad (3.42)$$

sau, echivalent

$$(\mathbf{O}') \quad \begin{cases} A'_o = \begin{bmatrix} 0 & 0 & -a_n \\ 1 & & -a_{n-1} \\ & \ddots & \vdots \\ & & 1 & -a_1 \end{bmatrix}, & B'_o = \begin{bmatrix} \nu_n \\ \nu_{n-1} \\ \vdots \\ \nu_1 \end{bmatrix}, \\ C'_o = [0 \quad \cdots \quad 0 \quad 1]. \end{cases} \quad (3.43)$$

Formele (\mathbf{C}) și (\mathbf{C}') (respectiv (\mathbf{O}) și (\mathbf{O}')) diferă neesențial între ele printr-o permutare evidentă a variabilelor de stare. Formele (\mathbf{C}) și (\mathbf{O}) (respectiv (\mathbf{C}') și (\mathbf{O}')) se numesc *duale* deoarece au loc relațiile

$$A_o = A_c^T, \quad B_o = C_c^T, \quad C_o = B_c^T. \quad (3.44)$$

(Implicațiile calculatorii ale acestui tip de relații vor fi exploatate intensiv în capitolele următoare).

Matricele A_c și A'_c (respectiv A_o și A'_o) sunt într-o formă superior (inferior) Hessenberg particulară, numită *forma superior (inferior) Frobenius*. Sugestiv, aceste matrice se mai numesc *matrice companion* ale polinomului $p(s)$, cu ai cărui coeficienți sunt formate.

Prin calcul direct se poate arăta că forma standard controlabilă (\mathbf{C}) se bucură de următoarele proprietăți:

I. Are loc egalitatea

$$(sI - A_c)^{-1}B_c = \frac{L(s)}{p(s)}, \quad (3.45)$$

unde

$$L(s) = \begin{bmatrix} s^{n-1} \\ \vdots \\ s \\ 1 \end{bmatrix}.$$

În consecință

- a) $p(s)$ este polinomul caracteristic al lui A_c .
- b) (A_c, B_c, C_c) e o realizare a lui $T(s)$.

II. Matricea de controlabilitate

$$R_c \stackrel{\text{def}}{=} [B_c \quad A_c B_c \quad \cdots \quad A_c^{n-1} B_c]$$

este superior triunghiulară nesingulară, deci perechea (A_c, B_c) este controlabilă.

III. Perechea (C_c, A_c) este observabilă (rang $Q_c = n$), deci matricea de observabilitate

$$Q_c = \begin{bmatrix} C_c \\ C_c A_c \\ \vdots \\ C_c A_c^{n-1} \end{bmatrix}$$

este nesingulară, dacă și numai dacă $T(s)$ este ireductibilă, adică polinoamele $N(s)$ și $p(s)$ sunt coprime.

Proprietățile I – II de mai sus justifică denumirea de *formă standard controlabilă* atribuită realizărilor **(C)** (sau **(C')**), iar proprietatea III arată că ordinul minim al realizării coincide cu gradul numitorului funcției de transfer $T(s)$, presupusă ireductibilă.

Proprietățile formei standard observabile **(O)** se formulează prin dualitate. În particular avem

$$R_o = Q_c^T, \quad Q_o = R_c^T, \quad (3.46)$$

deci, la fel ca mai sus, ordinul minim coincide cu gradul numitorului funcției de transfer, presupusă ireductibilă.

B. Sisteme cu o singură intrare și mai multe ieșiri (SIMO)

În cazul sistemelor SIMO (Single-Input, Multi-Output) matricea de transfer, considerată dată, are o structură tip coloană

$$T(s) = \begin{bmatrix} T_1(s) \\ \vdots \\ T_l(s) \end{bmatrix} \in \mathcal{R}^l(s), \quad (3.47)$$

unde fiecare $T_i(s)$ este o funcție rațională strict proprie de forma

$$T_i(s) = \frac{N_i(s)}{p_i(s)}, \quad i = 1 : l \quad (3.48)$$

cu

$$\begin{aligned} p_i(s) &= s^{n(i)} + a_1(i)s^{n(i)-1} + \dots + a_{n(i)}(i), \\ N_i(s) &= \nu_1(i)s^{n(i)-1} + \dots + \nu_{n(i)}(i). \end{aligned} \quad (3.49)$$

Există două posibilități de scriere a unei realizări, una "naivă", care conduce la o realizare de ordin

$$n_I = \sum_{i=1}^l n(i), \quad (3.50)$$

alta "eficientă", care conduce la o realizare de ordin n_{II} egal cu gradul celui mai mic multiplu comun (cmmmc) al numitorilor $p_i(s)$, $i = 1 : l$ (evident, $n_{II} \leq n_I$).

B1. Realizarea "naivă"

Scriem relația $y(s) = T(s)u(s)$ pe componente sub forma

$$y_i(s) = T_i(s)u(s), \quad i = 1 : l \quad (3.51)$$

și realizăm fiecare funcție de transfer $T_i(s)$ printr-un sistem $S_i = (A_i, B_i, C_i)$ utilizând metodele expuse la punctul A, de exemplu putem considera $S_i =$

(A_i, B_i, C_i) în forma standard controlabilă (\mathbf{C}) , i.e.

$$A_i \stackrel{\text{def}}{=} \begin{bmatrix} -a_1(i) & \cdots & -a_{n(i)-1}(i) & -a_{n(i)}(i) \\ 1 & \cdots & 0 & 0 \\ & \ddots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 \end{bmatrix}, \quad B_i \stackrel{\text{def}}{=} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

$$C_i \stackrel{\text{def}}{=} [\quad \nu_1(i) \quad \cdots \quad \nu_{n-1}(i) \quad \nu_n(i) \quad].$$
(3.52)

Apoi realizăm conexiunea "paralel la intrare" (vezi figura 3.3 (a)), echivalentă cu (3.51) prin

$$\begin{cases} A = \begin{bmatrix} A_1 & & \\ & \ddots & \\ & & A_l \end{bmatrix}, & B = \begin{bmatrix} B_1 \\ \vdots \\ B_l \end{bmatrix}, \\ C = \begin{bmatrix} C_1 & & \\ & \ddots & \\ & & C_l \end{bmatrix}, \end{cases}$$
(3.53)

blocurile nedigonale ale matricelor A și C fiind nule.

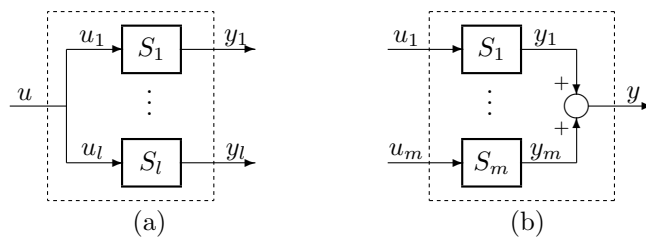


Figura 3.3: Conexiunile "paralel la intrare" (a) și "paralel la ieșire" (b).

Realizarea (3.53) se numește *cu ieșiri decuplate* (sau *decuplată la ieșire*) datorită formei bloc diagonale a perechii (C, A) .

B2. Realizarea ”eficientă”

Aducând la același numitor toate elementele lui $T(s)$, scriem

$$T(s) = \begin{bmatrix} T_1(s) \\ \vdots \\ T_l(s) \end{bmatrix} = \begin{bmatrix} N_1^*(s) \\ \vdots \\ N_l^*(s) \end{bmatrix} p^{-1}(s) \stackrel{\text{def}}{=} N(s)p^{-1}(s), \quad (3.54)$$

unde

$$p(s) = s^n + a_1 s^{n-1} + \dots + a_n \stackrel{\text{def}}{=} \text{cmmmc}(p_1(s), \dots, p_l(s)), \quad (3.55)$$

$$N_i^*(s) = \nu_1^*(i) s^{n-1} + \dots + \nu_n^*(i) \quad (3.56)$$

deci numărătorul $N \in \mathcal{R}^l[s]$ se scrie sub forma:

$$N(s) = \nu_1 s^{n-1} + \dots + \nu_n \quad (3.57)$$

unde vectorii $\nu_k \in \mathcal{R}^l$ au componentele $\nu_k^*(i)$, $i = 1 : l$, $k = 1 : n$.

Matricea de transfer factorizată $T(s) = N(s)p^{-1}(s)$ se realizează direct în forma *standard controlabilă* (**C**), unde acum polinomul caracteristic al matricei A_c coincide cu $p(s)$ iar C_c conține, pe fiecare coloană, coeficienții ν_k ai lui N , adică $C_c \in \mathcal{R}^{l \times n}$.

Observația 3.1 Realizarea obținută este controlabilă. Dacă factorizarea Np^{-1} este coprimă la dreapta, adică nu există un polinom D (divizor) care divide simultan p și fiecare N_i^* , atunci realizarea astfel obținută este și observabilă, deci de ordin minim. \diamond

C. Sisteme cu mai multe intrări și o singură ieșire (MISO)

În acest caz matricea de transfer are o structură linie

$$T(s) = [T_1(s) \quad \dots \quad T_m(s)] \in \mathcal{R}^{1 \times m}(s). \quad (3.58)$$

unde fiecare $T_i(s)$ este o funcție rațională strict proprie de forma

$$T_i(s) = \frac{N_i(s)}{p_i(s)}, \quad i = 1 : m.$$

Ca și în cazul precedent, prezentăm pentru problema realizării o soluție ”naivă” și una ”eficientă”.

C1. Realizarea ”naivă”

Scriem relația $y(s) = T(s)u(s)$ pe componente sub forma

$$y(s) = \sum_{i=1}^m T_i(s)u_i(s) \quad (3.59)$$

și realizăm fiecare funcție de transfer $T_i(s)$ printr-un model de stare $S_i = (A_i, B_i, C_i)$, utilizând metodele expuse la punctul A. Apoi realizăm conexiunea ”paralel la ieșire” (vezi figura 3.3 (b)), echivalentă cu (3.59) prin

$$\left\{ \begin{array}{l} A = \begin{bmatrix} A_1 & & \\ & \ddots & \\ & & A_m \end{bmatrix}, \quad B = \begin{bmatrix} B_1 & & \\ & \ddots & \\ & & B_m \end{bmatrix}, \\ C = [C_1 \quad \cdots \quad C_m], \end{array} \right. \quad (3.60)$$

blocurile nediagonale ale matricelor A și B fiind nule.

Realizarea (3.60) se numește *cu intrări decuplate* (sau *decuplată la intrare*).

C2. Realizarea ”eficientă”

Procedăm „prin dualitate” față de cazul B2. Aceasta înseamnă că, aducând la același numitor toate elementele lui $T(s)$, scriem

$$\begin{aligned} T(s) &= [T_1(s) \quad \cdots \quad T_m(s)] = \\ &= p^{-1}(s) [N_1^*(s) \quad \cdots \quad N_m^*(s)] = p^{-1}(s)N(s), \end{aligned} \quad (3.61)$$

unde $p(s) \stackrel{\text{def}}{=} \text{cmmmc}(p_1(s), \dots, p_m(s))$, $N_i^*(s)$ au expresii similare cu (3.55) și (3.56), deci numărătorul $N \in \mathcal{R}^{1 \times m}[s]$ se scrie sub forma (3.57), unde $\nu_k \in \mathcal{R}^{1 \times m}$ au componentele $\nu_k^*(i)$, $i = 1 : m$, $k = 1 : n$.

Matricea de transfer factorizată $T(s) = p^{-1}(s)N(s)$ se realizează direct în forma *standard observabilă* (**O**), unde acum polinomul caracteristic al matricei A_o coincide cu $p(s)$, iar B_o conține, pe fiecare linie, coeficienții lui N , adică $B_o \in \mathcal{R}^{n \times m}$.

Observația 3.2 Realizarea obținută este observabilă. Dacă factorizarea $p^{-1}N$ este coprimă la stânga (i.e. polinoamele $p(s)$ și fiecare $N_i^*(s)$ nu au un polinom divizor comun) atunci realizarea astfel obținută este și controlabilă, deci de ordin minim. \diamond

D. Sisteme cu mai multe intrări și mai multe ieșiri (MIMO)

Se dă matricea de transfer cu l linii și m coloane $T(s) = [T_{ij}(s)]$.

D1. Realizarea "naivă"

Această realizare se bazează pe scrierea

$$y_i(s) = \sum_{j=1}^m T_{ij}(s)u_j(s), \quad i = 1 : l \quad (3.62)$$

și este de ordin

$$n = \sum_{i=1}^l \sum_{j=1}^m n(i, j) \quad (3.63)$$

unde $n(i, j)$ e gradul numitorului lui $T_{ij}(s)$.

D2. Realizarea "eficientă"

În acest caz, o realizare eficientă se poate construi în două moduri.

1⁰. Partiționăm $T(s)$ pe coloane

$$T(s) = [T_1(s) \quad \dots \quad T_m(s)] \quad (3.64)$$

și aducem fiecare coloană $T_j(s) \in \mathcal{R}^l(s)$ la același numitor ca la punctul **B**, adică scriem $T_j(s) = N_j(s)p_j^{-1}(s)$, unde $N_j(s) \in \mathcal{R}^l[s]$, $p_j(s) \in \mathcal{R}[s]$. Deci avem

$$\begin{aligned} T(s) &= [N_1(s)p_1^{-1}(s) \quad \dots \quad N_m(s)p_m^{-1}(s)] = \\ &= [N_1(s) \quad \dots \quad N_m(s)] \begin{bmatrix} p_1(s) & & \\ & \ddots & \\ & & p_m(s) \end{bmatrix}^{-1} = \\ &\stackrel{\text{def}}{=} N(s)p^{-1}(s), \end{aligned} \quad (3.65)$$

unde

$$\begin{aligned} p_j(s) &= s^{n(j)} + a_1(j)s^{n(j)-1} + \dots + a_{n(j)}(j), \quad a_k(j) \in \mathcal{R}, \\ N_j(s) &= \nu_1(j)s^{(j)-1} + \dots + \nu_{n(j)}(j), \quad \nu_k(j) \in \mathcal{R}^l. \end{aligned} \quad (3.66)$$

Întrucât relația $y(s) = T(s)u(s)$ se scrie

$$y(s) = \sum_{j=1}^m N_j(s)p_j^{-1}(s)u_j(s), \quad (3.67)$$

iar fiecare coloană $N_j(s)p_j^{-1}(s)$ se realizează în forma standard controlabilă $(A_c(j), B_c(j), C_c(j))$, ca în cazul B2, obținem realizarea *standard controlabilă* (decuplată la intrare)

$$\left\{ \begin{array}{l} A_c = \begin{bmatrix} A_c(1) & & \\ & \ddots & \\ & & A_c(m) \end{bmatrix}, \quad B_c = \begin{bmatrix} B_c(1) & & \\ & \ddots & \\ & & B_c(m) \end{bmatrix}, \\ C_c = [C_c(1) \quad \cdots \quad C_c(m)]. \end{array} \right. \quad (3.68)$$

Ordinul $n = n_c$ este suma gradelor $n(j)$ ale numitorilor comuni pe coloană. (Se poate arăta că aceste grade coincid cu indicii de controlabilitate ai perechii (A_c, B_c)).

2⁰. Partiționând $T(s)$ pe linii și procedând prin dualitate față de punctul precedent, obținem realizarea *standard observabilă* (decuplată la ieșire)

$$\left\{ \begin{array}{l} A_o = \begin{bmatrix} A_o(1) & & \\ & \ddots & \\ & & A_o(l) \end{bmatrix}, \quad B_o = \begin{bmatrix} B_o(1) \\ \vdots \\ B_o(l) \end{bmatrix} \\ C_o = \begin{bmatrix} C_o(1) & & \\ & \ddots & \\ & & C_o(l) \end{bmatrix}. \end{array} \right. \quad (3.69)$$

Ordinul $n = n_o$ al realizării (3.69) este suma gradelor numitorilor comuni pe linie. (Se poate arăta că aceste grade coincid cu indicii de observabilitate ai perechii (C_o, A_o)).

Observația 3.3 În general, niciuna dintre cele două realizări (3.68) și (3.69) nu este minimală căci, în general, perechea (C_c, A_c) nu rezultă observabilă iar perechea (A_o, B_o) nu rezultă controlabilă. \diamond

În finalul acestei secțiuni precizăm că obținerea efectivă a realizărilor de mai sus necesită, în primul rând, un mecanism de manipulare a matricelor bloc cum este cel oferit de MATLAB (sau cel care poate fi creat

cu relativă ușurință de utilizator în orice limbaj de programare). Cu un astfel de mecanism la dispoziție, scrierea algoritmilor de construcție a realizărilor prezentate devine extrem de simplă și, din acest motiv, face obiectul exercițiilor.

3.4 Conversii de modele

Conversia de la un model de tip matrice de transfer la un model de stare, pe scurt $(\mathbf{T}) \rightarrow (\mathbf{S})$, constă din aplicarea uneia dintre procedurile de realizare, care au fost discutate în paragraful anterior. Conversia inversă $(\mathbf{S}) \rightarrow (\mathbf{T})$, respectiv calculul matricei de transfer asociate unui model de stare dat, se face separat pentru fiecare pereche (u_j, y_i) de intrări și ieșiri, adică în esență se referă la un sistem simplu (SISO). Deoarece formele alternative de reprezentare ale unei matrice de transfer se referă tot la sisteme simple, adică la funcții de transfer, în acest paragraf ne vom referi numai la aceste obiecte.

Calculul funcției de transfer (Conversia $(\mathbf{S}) \rightarrow (\mathbf{T})$)

Se dă un sistem simplu definit prin

$$\begin{cases} \dot{x} = Ax + Bu \\ y = Cx \end{cases}, \quad (3.70)$$

unde matricele $A \in \mathcal{R}^{n \times n}$, $B \in \mathcal{R}^n$, $C \in \mathcal{R}^{1 \times n}$ sunt cunoscute.

Pentru a calcula funcția de transfer, adică două polinoame $p(s)$, $N(s)$ astfel încât

$$T(s) \stackrel{\text{def}}{=} C(sI - A)^{-1}B = \frac{N(s)}{p(s)}, \quad (3.71)$$

constatăm că funcția de transfer a unui sistem cu reacție unitară (vezi figura 3.4) este

$$T_0(s) = \frac{T(s)}{1 + T(s)} = \frac{N(s)}{p(s) + N(s)} = \frac{N(s)}{p_0(s)} \quad (3.72)$$

unde $p_0(s) = p(s) + N(s)$. Deci

$$N(s) = p_0(s) - p(s). \quad (3.73)$$

(Menționăm că rădăcinile polinomului $N(s)$ sunt zerourile sistemului $S = (A, B, C)$).

Problema s-a redus la a calcula $p(s)$ și $p_0(s)$.

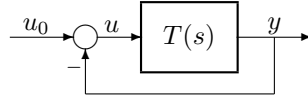


Figura 3.4: Conexiune cu reacție unitară.

Dar $p(s)$ este polinomul caracteristic al lui A adică

$$p(s) = \prod_{i=1}^n (s - \lambda_i), \quad (3.74)$$

unde λ_i sunt polii lui $T(s)$, adică valorile proprii ale matricei A . Analog, $p_0(s)$ este polinomul caracteristic al matricei A_0 a sistemului în circuit închis obținut punând $u = u_0 - y$, adică $u = u_0 - Cx$. Deci $A_0 = A - BC$.

Prin urmare, pentru a calcula funcția de transfer $T(s)$ se procedează astfel:

1. Se calculează valorile proprii λ_i ale matricei A utilizând algoritmul QR.
2. Se formează $A_0 = A - BC$.
3. Se calculează valorile proprii λ_{0i} ale matricei A_0 utilizând algoritmul QR.
4. Se calculează (coeficienții lui) $p(s) = \prod_{i=1}^n (s - \lambda_i)$.
5. Se calculează $p_0(s) = \prod_{i=1}^n (s - \lambda_{0i})$.
6. Se calculează $N(s) = p_0(s) - p(s)$.

Dacă sistemul $S = (A, B, C)$ nu este controlabil și observabil (ceea ce este probabil să se întâmple dacă provine dintr-un sistem multiplu) atunci p și N , adică p și p_0 au un cel mai mare divizor comun (cmMdc) $\neq 1$, care coincide cu produsul $\prod (s - \lambda_{fi})$, unde $\lambda_{fi} \in \sigma(A) \cap \sigma(A_0)$ sunt polii ficși sistemului (A, B, C) . Acest divizor poate fi eliminat prin inspecție înainte de pasul 4, dar decizia poate fi afectată de erorile de rotunjire inerente.

Pentru o funcție de transfer strict proprie definim, în afară de reprezentarea "rațională"

$$T(s) = \frac{N(s)}{p(s)}, \quad (3.75)$$

ca raport de două polinoame, încă două reprezentări:

$$(\mathbf{T}_1) \quad T(s) = \delta \frac{\prod_{j=1}^{n-k} (s - \mu_j)}{\prod_{i=1}^n (s - \lambda_i)} \quad (3.76)$$

și

$$(\mathbf{T}_2) \quad T(s) = \sum_{i=1}^n \frac{r_i}{s - \lambda_i}, \quad (3.77)$$

unde λ_i sunt polii, μ_j sunt zerourile, δ este amplificarea la infinit, numărul întreg $k \geq 1$ este excesul poli - zerouri sau „tipul” lui $T(s)$, iar r_i sunt reziduurile lui $T(s)$ în polii λ_i (a doua reprezentare există ca atare numai dacă toți polii lui $T(s)$ sunt simpli.)

Relativ la reprezentarea (\mathbf{T}_1) discutăm în primul rând:

Conversia $(\mathbf{S}) \rightarrow (\mathbf{T}_1)$

Se dă modelul de stare $S = (A, B, C)$. Pentru a calcula $T(s)$ sub forma (\mathbf{T}_1) se procedează astfel:

1. Se calculează valorile proprii λ_i ale matricei A .
2. Se calculează zerourile μ_j ale sistemului (A, B, C) .
3. Se alege $\sigma_0 \in \mathcal{R}$, astfel încât $\sigma_0 \neq \lambda_i, \mu_j, \forall (i, j)$ și se calculează

$$\delta = \frac{C(\sigma_0 I - A)^{-1} B}{\frac{\prod_{j=1}^{n-k} (\sigma_0 - \mu_j)}{\prod_{i=1}^n (\sigma_0 - \lambda_i)}}.$$

Observația 3.4 La pasul 2 se pot utiliza încă cel puțin două metode.

a) *Metoda generală* ține seama de faptul că zerourile sistemului (A, B, C) (considerat controlabil și observabil) coincid cu valorile proprii generalizate finite ale fascicolului

$$\lambda \left[\begin{array}{cc} I_n & 0 \\ 0 & 0 \end{array} \right] - \left[\begin{array}{cc} A & B \\ C & 0 \end{array} \right], \quad (3.78)$$

deci se calculează utilizând algoritmul **QZ**, vezi [3].

b) *Metoda specifică* în cazul $l = 1, m = 1$ considerat ”simulează”, în limbajul descrierii de stare $S = (A, B, C)$, procedura de inversare a sistemului propriu având funcția de transfer $s^k T(s)$, unde (A, B, C) este

o realizare a lui $T(s)$. Fie k ordinul relativ al lui (A, B, C) , adică primul întreg ≥ 1 astfel încât $\delta \stackrel{\text{def}}{=} CA^{k-1}B \neq 0$. Din ecuațiile de stare

$$\begin{cases} \dot{x} = Ax + Bu \\ y = Cx, \end{cases} \quad (3.79)$$

prin derivarea succesivă a ieșirii y ținând seama de faptul că $CA^i B = 0$, $i = 0 : k - 2$, obținem

$$\begin{aligned} \dot{y} &= CAx \\ &\vdots \\ y^{(k-1)} &= CA^{k-1}x \\ y^{(k)} &= CA^k x + CA^{k-1}Bu \stackrel{\text{def}}{=} CA^k x + \delta u, \end{aligned} \quad (3.80)$$

unde, prin ipoteză, $\delta \neq 0$. În consecință, sistemul

$$\begin{cases} \dot{x} &= Ax + Bu \\ y^{(k)} &= CA^k x + \delta u \end{cases} \quad (3.81)$$

este inversabil iar inversul său are ecuațiile

$$\begin{cases} \dot{x} = (A - \delta^{-1}BCA^k)x + \delta^{-1}By^{(k)} \\ u = -\delta^{-1}CA^k x + \delta^{-1}y^{(k)}. \end{cases} \quad (3.82)$$

Mai mult, matricea $A - \delta^{-1}BCA^k$ are exact k valori proprii nule (corespunzătoare celor k valori proprii generalizate infinite ale fascicolului (3.78)) precum și $n - k$ valori proprii egale cu zerourile lui (A, B, C) . Prin urmare ordinul relativ k coincide cu excesul poli - zerouri iar $\delta = CA^{k-1}B$ este amplificarea la infinit, adică $s^k T(s) \rightarrow \delta$, când $s \rightarrow \infty$. În definitiv, zerourile (finite) μ_j , $j = 1 : n - k$ ale sistemului (A, B, C) coincid cu cele $n - k$ valori proprii nenule ale matricei $A - \delta^{-1}BCA^k$. \diamond

În ceea ce privește celelalte conversii, dăm numai câteva sugestii. (Menționăm că procedurile de calcul polinomial care apar în schemele de calcul de mai jos sunt discutate în secțiunea finală a acestui capitol.)

Conversia $(\mathbf{T}_1) \rightarrow (\mathbf{T})$

1. Se calculează $p(s) = \prod_{i=1}^n (s - \lambda_i)$.
2. Se calculează $N(s) = \prod (s - \mu_i)$.
2. $N(s) \leftarrow \delta N(s)$

Observația 3.5 Uneori în loc de algoritmul $(\mathbf{S}) \rightarrow (\mathbf{T})$ prezentat mai sus se utilizează $(\mathbf{S}) \rightarrow (\mathbf{T}_1)$ urmat de $(\mathbf{T}_1) \rightarrow (\mathbf{T})$. \diamond

Conversia $(\mathbf{T}) \rightarrow (\mathbf{T}_1)$

1. Se construiește o realizare de stare a lui $T(s)$.
2. Se aplică procedura de conversie $(\mathbf{S}) \rightarrow (\mathbf{T}_1)$.

Conversia $(\mathbf{T}_1) \rightarrow (\mathbf{S})$

1. Se construiește reprezentarea rațională a lui $T(s)$.
2. Se aplică algoritmul de realizare $(\mathbf{T}) \rightarrow (\mathbf{S})$.

Relativ la reprezentarea (\mathbf{T}_2) , menționăm numai că ea se obține din (\mathbf{T}) calculând polii λ_i cu algoritmul QR aplicat matricei companion A_c (respectiv A_o) și apoi calculând reziduurile r_i cu formula cunoscută

$$r_i = \frac{N(\lambda_i)}{p'(\lambda_i)}, \quad (3.83)$$

evaluările polinoamelor realizându-se cu algoritmul (schema) lui Horner (vezi secțiunea următoare).

Alternativ, se construiește o realizare de stare (A, B, C) a lui $T(s)$ și se calculează valorile proprii λ_i (presupuse distincte) precum și vectorii proprii x_i ai lui A . Considerând matricea $X = [x_1 \ \dots \ x_n]$ și efectuând transformarea

$$X^{-1}AX = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}, \quad X^{-1}B = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}, \quad (3.84)$$

$$CX = [\gamma_1 \ \dots \ \gamma_n],$$

avem

$$r_i = \gamma_i \beta_i, \quad i = 1 : n. \quad (3.85)$$

3.5 Algoritmi de calcul polinomial

În secțiunile precedente ale capitolului de față s-a evidențiat faptul că procesarea modelelor sistemice de tip matrice de transfer se reduce, în ultimă instanță, la efectuarea unui set de operații cu polinoame. De aceea reunim aici procedurile numerice de bază pentru rezolvarea problemelor de calcul polinomial.

Notăm un polinom $a(s)$ de grad n prin

$$a(s) = a_0 s^n + a_1 s^{n-1} + \cdots + a_n, \quad a_0 \neq 0 \quad (3.86)$$

și îi asociem vectorul

$$a = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} \quad (3.87)$$

de lungime $n + 1$.⁵

Se știe că polinoamele de grad $\leq m$, cu operațiile de adunare și înmulțire cu scalari definesc un *spațiu vectorial* de dimensiune $m + 1$. În acest spațiu, cu m fix⁶ și $n \leq m$, polinomul $a(s)$ se reprezintă prin vectorul

$$T^{m-n} a = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ a_0 \\ \vdots \\ a_n \end{bmatrix} \in \mathcal{R}^{m+1}, \quad (3.88)$$

unde T este operatorul de „deplasare (*shift*) cu un pas în jos urmat de adăugarea unui zero pe prima poziție” definit prin

$$(Ta)_i = a_{i-1}, \quad i \geq 1, \quad (Ta)_0 = 0. \quad (3.89)$$

Reținând acest fapt simplu, în continuare trecem în revistă modalitățile de efectuare a operațiilor principale cu polinoame.

1⁰. **Suma** $c(s) \stackrel{\text{def}}{=} a(s) + b(s)$ a două polinoame $a(s)$, $b(s)$ de grade n , m cu $n \leq m$ ⁷ este de grad m , deci

$$c = T^{m-n} a + b. \quad (3.90)$$

Rezultă următorul algoritm.

⁵Numerotarea elementelor unui vector (tablou unidimensional) se face în mod diferit în diverse limbaje formale de programare, respectiv începând cu 0, ca în limbajul **C**, sau cu 1 ca în **MATLAB**. Trecerea la cea de a doua variantă reprezintă un exercițiu util pentru cititor.

⁶De fiecare dată când efectuăm o operație alegem m minim necesar.

⁷Altfel facem suma $b(s) + a(s)$.

Algoritmul 3.1 (Se dau vectorii $a \in \mathcal{R}^{n+1}$ și $b \in \mathcal{R}^{m+1}$, cu $m \geq n$, ai coeficienților polinoamelor $a(s)$ și $b(s)$. Algoritmul calculează vectorul $c \in \mathcal{R}^{m+1}$ al coeficienților polinomului sumă $c(s) = a(s) + b(s)$.)

1. Pentru $i = 0 : m - n - 1$

1. $c_i = b_i$

2. Pentru $i = m - n : m$

1. $c_i = a_{i-(m-n)} + b_i$

2⁰. **Produsul** $p(s) \stackrel{\text{def}}{=} a(s)b(s)$ a două polinoame $a(s)$, $b(s)$ de grade $n \geq m$ ⁸ este de grad $m + n$. Avem

$$a(s)b(s) = \left(\sum_{i=0}^n a_i s^{n-i} \right) \left(\sum_{j=0}^m b_j s^{m-j} \right) = \sum_{k=0}^{m+n} p_k s^{n+m-k}, \quad (3.91)$$

unde

$$p_k = \sum_{i+j=k} a_i b_j = \begin{cases} \sum_{j=0}^k a_{k-j} b_j \\ \sum_{i=0}^k a_i b_{k-i} \end{cases}. \quad (3.92)$$

(Pentru comoditatea scrierii am presupus că $a_i = 0$ pentru $i > n$ și $b_j = 0$ pentru $j > m$. În general, operația (3.92) reprezintă *convoluția discretă* a șirurilor $(a_i)_{i \geq 0}$, $(b_j)_{j \geq 0}$).

Din (3.92) se constată ușor ca vectorul $p \in \mathcal{R}^{m+n+1}$, al coeficienților polinomului produs, se poate exprima într-o formă matriceală în modul următor

$$\begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_m \\ \vdots \\ p_n \\ p_{n+1} \\ \vdots \\ p_{n+m} \end{bmatrix} = \begin{bmatrix} a_0 & 0 & \cdots & 0 \\ a_1 & a_0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_m & a_{m-1} & \cdots & a_0 \\ \vdots & \vdots & & \vdots \\ a_n & a_{n-1} & \cdots & a_{n-m} \\ 0 & a_n & \cdots & a_{n-m+1} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_n \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_m \end{bmatrix}, \quad (3.93)$$

⁸ Altfel facem produsul $b(s)a(s)$.

adică, pentru a putea descrie produsul unui polinom $a(s)$ de grad n cu un polinom $b(s)$ de grad m , îi asociem lui $a(s)$ matricea cu $m + 1$ coloane

$$A = \begin{bmatrix} a_0 & & & & & \\ \vdots & & \ddots & & & \\ a_n & & & & a_0 & \\ & & & \ddots & \vdots & \\ & & & & & a_n \end{bmatrix} \in \mathcal{R}^{(m+n+1) \times (m+1)}, \quad (3.94)$$

și scriem

$$p = Ab. \quad (3.95)$$

În mod analog, datorită comutativității produsului de polinoame, avem și $p = Ba$ unde matricea cu $n + 1$ coloane $B \in \mathcal{R}^{(m+n+1) \times (n+1)}$ se formează, în aceeași manieră, cu coeficienții lui $b(s)$.

Obținem următorul algoritm.

Algoritmul 3.2 (Se dau vectorii $a \in \mathcal{R}^{n+1}$ și $b \in \mathcal{R}^{m+1}$, cu $m \geq n$, ai coeficienților polinoamelor $a(s)$ și $b(s)$. Algoritmul calculează vectorul $c \in \mathcal{R}^{m+n+1}$ al coeficienților polinomului produs $c(s) = a(s)b(s)$.)

1. Pentru $k = 0 : m + n$

1. $p_k = 0$

2. Pentru $j = 0 : m$

1. Pentru $i = 0 : n$

1. $p_{i+j} = p_{i+j} + a_i b_j$

Observația 3.6 Operațiile 1⁰ și 2⁰ definesc *inelul comutativ* $\mathcal{R}[s]$ al polinoamelor cu coeficienți reali. \diamond

3⁰. **Teorema împărțirii întregi** pentru două polinoame $a(s), b(s)$ de grade $n \leq m$ afirmă existența unică a două polinoame $q(s)$ de grad $m - n$ și $r(s)$ de grad cel mult $n - 1$ astfel încât

$$b(s) = a(s)q(s) + r(s). \quad (3.96)$$

Dacă $n > m - n$, ținând seama de exprimarea matriceală a produsului de

polinoame, (3.96) se scrie

$$\begin{bmatrix} b_0 \\ \vdots \\ b_{m-n} \\ \hline b_{m-n+1} \\ \vdots \\ b_n \\ \vdots \\ b_m \end{bmatrix} = \begin{bmatrix} a_0 & & & \\ \vdots & \ddots & & \\ a_{m-n} & \cdots & a_0 & \\ \hline a_{m-n+1} & \cdots & a_1 & \\ \vdots & \cdots & \vdots & \\ a_n & \cdots & a_{2n-m} & \\ & \ddots & \vdots & \\ & & a_n & \end{bmatrix} \begin{bmatrix} q_0 \\ \vdots \\ q_{m-n} \end{bmatrix} + \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \hline r_0 \\ \vdots \\ r_{2n-m-1} \\ \vdots \\ r_{n-1} \end{bmatrix} \quad (3.97)$$

deci, cu partiții evidente, avem

$$b_1 = A_1 q \quad (3.98)$$

$$b_2 = A_2 q + r \quad (3.99)$$

adică

$$r = b_2 - A_2 q. \quad (3.100)$$

Cu alte cuvinte, vectorul q al coeficienților polinomului cât se calculează rezolvând sistemul inferior triunghiular nesingular format din primele $m - n + 1$ ecuații

$$A_1 q = b_1, \quad (a_0 \neq 0), \quad (3.101)$$

după care r se calculează cu (3.100) ca reziduu al sistemului $A_2 q = b_2$, format cu celelalte n ecuații.

4⁰. Reprezentarea prin fracție continuă finită a unei funcții raționale este o aplicație directă a algoritmului de împărțire întreagă. Fie

$$T(s) = \frac{N(s)}{p(s)} \quad (3.102)$$

strict proprie ⁹ și fie k_1 primul coeficient nenul al lui N astfel încât $N(s) = k_1 p_1(s)$ unde $p_1(s)$ e monic și de grad strict mai mic decât $p(s)$. Avem

$$T(s) = \frac{k_1 p_1(s)}{p(s)} = \frac{k_1}{\frac{p(s)}{p_1(s)}} = \frac{k_1}{q_1(s) + \frac{N_1(s)}{p_1(s)}}, \quad (3.103)$$

⁹ Altfel extragem partea întreagă ν_0 și punem $q_0(s) = \nu_0$.

unde q_1 e câtul iar N_1 e restul împărțirii lui p la p_1 , deci

$$T_1(s) \stackrel{\text{def}}{=} \frac{N_1(s)}{p_1(s)}$$

e strict proprie. Putem scrie

$$T(s) = \frac{k_1}{q_1(s) + T_1(s)} = \frac{k_1 \frac{1}{q_1(s)}}{1 + T_1(s) \frac{1}{q_1(s)}}. \quad (3.104)$$

Dacă privim funcția rațională $T(s)$ ca funcție de transfer a unui sistem simplu atunci putem asocia relației (3.104) schema bloc din figura 3.5. Procedura de mai sus se repetă pentru $T_1(s)$. În general,

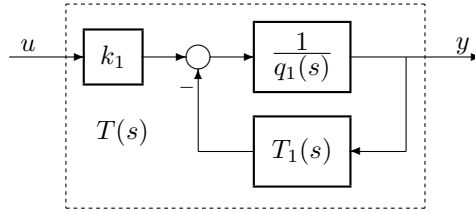


Figura 3.5: Corespondentul sistemic al primei etape de descompunere în fracție continuă finită a funcției raționale $T(s)$.

$$T_{i-1}(s) \stackrel{\text{def}}{=} \frac{k_i p_i(s)}{p_{i-1}(s)} = \frac{k_i}{q_i + \frac{k_{i+1} p_{i+1}(s)}{p_i(s)}}, \quad i = 1, 2, \dots \quad (3.105)$$

unde $T_0(s) = T(s)$. Deoarece gradele lui p_i scad strict cu i și $\text{grad } p_0(s) = n < \infty$, procedura are un număr finit de pași. (În mod analog, fracția continuă asociată unui număr rațional e finită și reciproc.)

Calculul reprezentării prin fracție continuă finită a unei funcții raționale date constă în determinarea tripletelor $(k_i, q_i(s), p_i(s))$, $i = 1, 2, \dots$ prin aplicarea repetată a algoritmului de împărțire întreagă. Scrierea efectivă a algoritmului este propusă ca exercițiu cititorului.

5⁰. **Calculul celui mai mare divizor comun (cmMdc)** a două polinoame $a(s)$ și $b(s)$ se face cel mai bine cu algoritmul lui Euclid, deci consti-

tuie o altă aplicație directă a algoritmului de împărțire întreagă. Reamintim schema de calcul care descrie algoritmul lui Euclid, unde presupunem $n = \text{grad } a(s) \leq \text{grad } b(s) = m$ ¹⁰.

1. $r(s) = 1$
2. Cât timp $r(s) \neq 0$
 1. Se calculează polinoamele $q(s)$ și $r(s)$ astfel încât $b(s) = a(s)q(s) + r(s)$ utilizând algoritmul de împărțire întreagă.
 2. $b(s) \leftarrow a(s)$
 3. $a(s) \leftarrow r(s)$
3. $\text{cmMdc} = b(s)$

De notat faptul că la implementarea schemei de mai sus condiția din instrucțiunea 2 trebuie relaxată introducându-se o toleranță corespunzătoare.

6⁰. **Calculul celui mai mic multiplu comun (cmmmc)** a două polinoame $a(s)$ și $b(s)$ se face pe baza relației

$$\text{cmmmc}(a(s), b(s)) = \frac{a(s)b(s)}{\text{cmMdc}(a(s), b(s))}. \quad (3.106)$$

7⁰. **Calculul rădăcinilor unui polinom** dat prin coeficienți, i.e. rezolvarea ecuațiilor polinomiale $p(s) = 0$, este o problema a căreia i s-a acordat, de foarte mult timp, o atenție cu totul deosebită și pentru care există numeroși algoritmi „candidat”, e.g. metoda bisecției, metode de interpolare (liniară sau pătratică), metodele Newton, Laguerre, Bairstow, etc. Subliniem, încă odată, că metoda cu cele mai bune performanțe numerice este aplicarea algoritmului **QR** matricelor companion A_c sau A_o din (3.40), respectiv (3.42) (evident, după împărțirea prin coeficientul termenului de grad maxim).

8⁰. **Calculul coeficienților unui polinom** atunci când se cunosc rădăcinile λ_i , $i = 1 : n$ se face pe baza relației

$$p(s) = \prod_{i=1}^n (s - \lambda_i), \quad (3.107)$$

i.e. prin ”acumularea” produselor după schema:

¹⁰ Altfel facem $a(s) \leftrightarrow b(s)$.

1. $p(s) = 1$
2. Pentru $i = 1 : n$
 1. $p(s) \leftarrow p(s)(s - \lambda_i)$

unde la începutul pasului curent 2.1 $p(s) \stackrel{\text{def}}{=} p_{i-1}(s)$ este un polinom monic de grad $i - 1$, iar după execuția lui un polinom monic $p(s) \stackrel{\text{def}}{=} p_i(s)$ de grad i . Deci

$$\begin{aligned} p_i(s) &\stackrel{\text{def}}{=} s^i + a_1^* s^{i-1} + \dots + a_i^* = \\ &= (s^{i-1} + a_1 s^{i-2} + \dots + a_{i-1})(s - \lambda_i) = \\ &= s^i + (a_1 - \lambda_i) s^{i-1} + (a_2 - \lambda_i a_1) s^{i-2} + \dots + a_i - \lambda_i a_{i-1}. \end{aligned}$$

Obținem următorul algoritm.

Algoritmul 3.3 (Se dau rădăcinile $\lambda_i \in \mathcal{C}$, $i = 1 : n$ ale unui polinom. Algoritmul calculează coeficienții $a_0 = 1$, a_i , $i = 1 : n$ ai acestuia).

1. $a_0 = 1$
2. Pentru $k = 1 : n$
 1. $a_k = 0$
3. Pentru $i = 1 : n$
 1. Pentru $k = i : -1 : 1$
 1. $a_k \leftarrow a_k - \lambda_i a_{k-1}$

9⁰. **Evaluarea valorii unui polinom** $p(s)$ sau/și a derivatei sale pentru o valoare $s = \sigma \in \mathcal{C}$ precizată se face cu *algoritmul lui Horner*. Reamintim acest algoritm, cu precizarea că, în contextul scalar discutat, el este optimal din punctul de vedere al eficienței.

Algoritmul 3.4 (Horner) (Se dau coeficienții a_i , $i = 0 : n$ ai unui polinom $p(s)$, în ordinea descrescătoare a puterilor nedeterminate, și numărul $\sigma \in \mathcal{C}$. Algoritmul calculează valoarea $\alpha = p(\sigma)$ a polinomului și valoarea $\delta = p'(\sigma)$ a derivatei acestuia).

1. $\alpha = a_0$
2. $\delta = n a_0$
3. Pentru $i = 1 : n - 1$

1. $\alpha \leftarrow \alpha\sigma + a_i$
2. $\delta \leftarrow \delta\sigma + (n - i)a_i$
4. $\alpha \leftarrow \alpha\sigma + a_n$

Algoritmul Horner generalizat care rezolvă problema calculului coeficienților polinomului $p(s + \sigma)$ face obiectul exercițiului 3.15.

Observația 3.7 Așa cum s-a mai precizat în capitolul 2, algoritmul Horner este util și pentru evaluarea polinoamelor matriceale. Într-un astfel de context algoritmul nu mai este însă optimal din punctul de vedere al numărului de operații aritmetice scalare (vezi [VI] și exercițiul 3.16). \diamond

Programe MATLAB disponibile

Pentru realizarea unor conexiuni specifice sunt disponibile funcțiile **append** (care efectuează produsul direct), **parallel**, **series**, **feedback** și **cloop** (sistem cu reacție unitară). Conexiunile de tip general se realizează apelând **bklbuild** și **conect**. Pentru construcția realizării de stare în forma standard controlabilă a unei matrice de transfer cu o singură intrare este disponibilă funcția **tf2ss** ("transfer function **t(w)o(!)** state space"). Conversia inversă se efectuează cu funcția **ss2tf**. Numele celorlalte funcții de conversie se formează similar, ținându-se seama că pentru codificarea formei (\mathbf{T}_1) este folosită sigla **zp** ("zerouri – poli"). Calculul descompunerii în fracții simple (forma (\mathbf{T}_2)) se poate face cu funcția **residue**. Transformările nesingulare de stare (3.3) se pot calcula cu ajutorul funcției **ss2ss**.

Pentru implementarea algoritmilor de calcul polinomial din § 3.5, menționăm în primul rând convenția MATLAB de reprezentare a polinoamelor prin vectorii *linie* ai coeficienților în ordinea descrescătoare a puterilor nedeterminatei (primul element are indicele 1). Zerourile (rădăcinile) unui polinom se reprezintă prin vectori coloană. Adunarea polinoamelor se face prin adunarea vectorilor coeficienților după egalarea dimensiunilor prin completarea lor corespunzătoare cu zerouri. Pentru înmulțire se utilizează funcția **conv** iar pentru împărțire întreagă **deconv**. Funcția **roots** calculează rădăcinile aplicând algoritmul **QR** matricei companion iar funcția **poly** calculează coeficienții din rădăcini (pentru argument matriceal calculează polinomul caracteristic). În sfârșit, calculul valorii unui polinom, inclusiv pentru argument matriceal, se efectuează cu funcția **polyval**.

Exerciții

E 3.1 Se dau două sisteme $S_i = (A_i, B_i, C_i, D_i)$, $i = 1, 2$. Să se scrie algoritmi de construcție a matricelor sistemelor $S = (A, B, C, D)$, rezultate prin interconectarea paralel, serie și în circuit închis a sistemelor date.

Indicație. Se utilizează relațiile (3.16), (3.21) și (3.28).

E 3.2 Se dau trei sisteme $S_i = (A_i, B_i, C_i, D_i)$, $i = 1 : 3$. Să se scrie algoritmi de construcție a matricelor sistemelor $S = (A, B, C, D)$, rezultate prin interconectarea sistemelor date conform schemelor din figura 3.6.

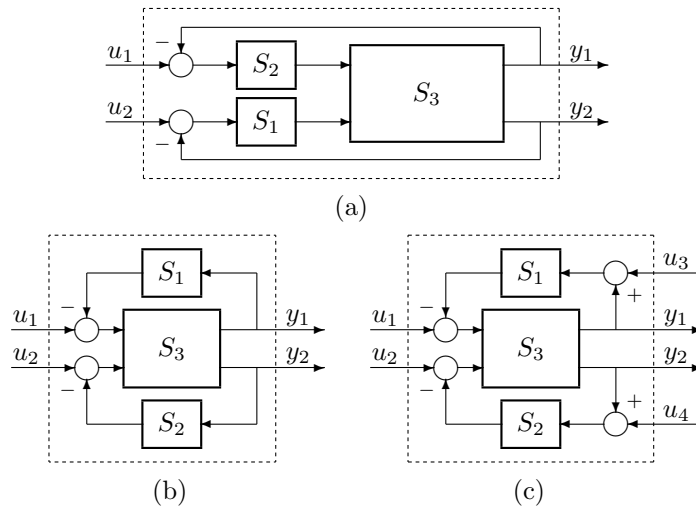


Figura 3.6: Conexiuni pentru exercițiul 3.2.

Cum procedați dacă sistemele S_i , $i = 1 : 3$ sunt date prin reprezentările lor de transfer ?

E 3.3 Se dă sistemul $S = (A, B, C, D)$ ale cărui matrice au structura

$$A = \begin{bmatrix} A_1 & 0 \\ A_{21} & A_2 \end{bmatrix}, \quad B = \begin{bmatrix} B_1 \\ G_2 \end{bmatrix},$$

$$C = [H_1 \quad C_2], \quad D = [D].$$

Să se scrie o procedură de reprezentare a lui S sub forma unei conexiuni
a) paralel

b) serie

a două sisteme S_1, S_2 , *convenabil* definite.

Indicație. a) Cu notația (3.13), considerați transformarea $\tilde{x}_2 = x_2 + Vx_1$ și determinați matricea V astfel încât $\tilde{A}_{21} = 0$. În ce condiții acest lucru este posibil?

b) Cu referire la (3.21), considerați factorizarea

$$\begin{bmatrix} A_{21} & G_2 \\ H_1 & D \end{bmatrix} = \begin{bmatrix} B_2 \\ D_2 \end{bmatrix} [C_1 \quad D_1] .$$

Pentru calculul factorizării utilizați **DVS**. Prin urmare numărul minim al terminalelor comune $l_1 = m_2$ coincide cu rangul matricei din membrul stâng.

E 3.4 Se dă o matrice de transfer $T(s)$, nu neapărat strict proprie, ale cărei coloane (linii) sunt aduse la același numitor comun. Să se scrie un program MATLAB de construcție a unei realizări de stare $S = (A, B, C, D)$ a lui $T(s)$ utilizând structura (3.68) (respectiv (3.69)).

E 3.5 Se dă o matrice de transfer $T(s)$, nu neapărat strict proprie, sub forma $T(s) = N(s)/p(s)$, unde $p(s)$ este un polinom de grad n iar $N(s) \in \mathcal{R}^{l \times m}[s]$. Scrieți prin inspecție două realizări ale lui $T(s)$. Precizați ce relație există între ordinul n al acestor realizări și ordinele n_I, n_{II} obținute în secțiunea 3.3D.

Indicație. Realizările cerute se obțin generalizând convenabil formele standard **(C)** și **(O)**.

E 3.6 Se dă un sistem $S = (A, B, C, D)$ precum și doi întregi $i \in 1 : l$, $j \in 1 : m$. Să se scrie algoritmi care calculează funcția de transfer T_{ij} în formele **(T)**, **(T₁)** și **(T₂)**.

Indicație. Vezi procedurile de conversie **(S)** \rightarrow **(T)**, **(S)** \rightarrow **(T₁)** și **(S)** \rightarrow **(T₂)** din §3.4.

E 3.7 Se dă o funcție de transfer $T(s)$, nu neapărat strict proprie. Să se scrie și să se discute comparativ câțiva algoritmi de conversie **(T)** \rightarrow **(T₁)** și **(T)** \rightarrow **(T₂)**.

E 3.8 Analizați posibilitățile de construcție a unei realizări de stare a unei funcții de transfer $T(s)$ date în forma **(T₁)**, *fără a construi*, în prealabil, reprezentarea rațională a lui $T(s)$.

E 3.9 Să se scrie, la nivel de operații scalare, algoritmul de împărțire întreagă a două polinoame.

E 3.10 Utilizând algoritmul lui Euclid de calcul al cmMdc a două polinoame, scrieți o procedură de calcul a cmMdc a l polinoame.

Indicație. Operația $p_1 \circ p_2 = \text{cmMdc}(p_1, p_2)$ este asociativă.

E 3.11 Utilizând algoritmi stabiliți în exercițiile precedente scrieți o procedură de calcul a cmmmc a l polinoame $p_i(s)$, $l \geq 2$.

E 3.12 Se dau l funcții raționale, fiecare definită ca raport de două polinoame. Să se scrie algoritmul de aducere a celor l funcții la același numitor comun ¹¹.

E 3.13 Să se scrie algoritmul de adunare a două sau mai multe funcții (matrice) raționale ¹².

E 3.14 Să se scrie algoritmul de înmulțire a două sau mai multe funcții (matrice) raționale ¹³.

E 3.15 Să se scrie algoritmul (lui Horner generalizat) de calcul al coeficienților polinomului $p(s + \sigma)$, unde $p(s)$ este dat prin vectorul coeficienților iar σ este un scalar dat, utilizând un volum minim de memorie. Elaborati o procedură de aplicare a algoritmului "cu hârtia și creionul".

Indicație. Toate calculele se pot desfășura în locațiile de memorie asociate vectorului coeficienților polinomului inițial.

E 3.16 Se dau un polinom $p(s)$ de grad m și o matrice $A \in \mathcal{R}^{n \times n}$. Să se scrie un algoritm de calcul al matricei $B = p(A)$. Căutați un exemplu care arată că în cazul matriceal schema lui Horner nu este optimală din punctul de vedere al numărului de operații.

E 3.17 Să se studieze problema inversării unei matrice de transfer (pătrate) $T(s)$, precizând condițiile în care inversa este tot o matrice de transfer. Formulați aceeași problemă relativ la o realizare $S = (A, B, C, D)$ a lui $T(s)$ și comentați.

E 3.18 Ce puteți afirma relativ la problema anterioară în cazul $l \neq m$?

¹¹Vezi procedurile de realizare "eficiente".

¹²Vezi conexiunea paralel.

¹³Vezi conexiunea serie.

Bibliografie

- [1] **Ionescu V.** TEORIA SISTEMELOR. SISTEME LINIARE. EDP, București 1985.
- [2] **Varga A., Sima V.** *Numerically Stable Algorithm for Transfer Function Matrix Evaluation*, Int J. Control, Vol. 33, pp. 1123 – 1133, 1981.
- [3] **Emami A., Naeini A., Van Dooren P.** *Computation of Zeros of Linear Multivariable Systems*, Automatica, Vol. 4, pp. 415 – 430, 1982.
- [4] **Barnett S.** *Matrices, Polynomials and Linear Time - Invariant Systems*, IEEE Trans. AC-18, No. 1, pp. 1 – 10, 1973.

Capitolul 4

Calculul răspunsului în timp al sistemelor liniare

Analiza comportării sistemelor dinamice include, în mod necesar, determinarea evoluției temporale a ieșirilor (și, eventual, a stării) la diverse tipuri de intrări externe. În faza de concepție a unui sistem, astfel de determinări nu pot fi realizate decât într-un mediu de simulare în care testele se fac sau pe modele la scară sau prin rezolvarea ecuațiilor ce descriu sistemul respectiv. Un mediu de simulare pentru sisteme cu timp continuu ce folosește tehnica numerică de calcul presupune, în mod obligatoriu, o discretizare a timpului și calculul răspunsului în momentele de timp discret corespunzătoare. În acest scop se utilizează o procedură de discretizare adecvată a cărei alegere poate influența în mod determinant precizia rezultatelor.

4.1 Răspunsul liber al sistemelor liniare

Problema determinării răspunsului liber (la intrare nulă) al unui sistem liniar $S = (A, B, C, D)$ constă în calculul funcției de ieșire $y(t)$ definită prin

$$\begin{cases} \dot{x}(t) = Ax(t), & x(0) = x \\ y(t) = Cx(t), \end{cases} \quad (4.1)$$

pe un interval fixat $[0, t_f]$.

Soluția ecuației diferențiale omogene (4.1) este

$$x(t) = e^{tA}x(0), \quad (4.2)$$

unde, prin definiție, $\Phi(t) = e^{tA}$ este matricea de tranziție a sistemului considerat. Prin urmare, ținând seama de ecuația ieșirii, răspunsul liber are expresia simplă

$$y(t) = Ce^{tA}x(0).$$

Pentru a evalua numeric în mod eficient valorile $y(t)$, divizăm intervalul $[0, t_f]$ într-un număr N de subintervale egale, i.e. considerăm $t_f = Nh$, unde h este pasul de simulare, utilizat pentru calculul valorilor dorite $y(kh)$, $k = 0 : N - 1$, ale răspunsului, iar N este numărul de pași necesar pentru acoperirea intervalului prescris. Scriind relația (4.2) la două momente de timp discret $t = kh$ și $t = (k + 1)h$, obținem evident

$$\begin{cases} x((k + 1)h) = e^{hA}x(kh) \\ y(kh) = Cx(kh), \end{cases} \quad (4.3)$$

adică $y(kh)$ rezultă direct în funcție de $x(kh)$ care se determină recurent pentru $k = 0 : N - 1$, ținând seama că stare inițială $x(0) = x$ este cunoscută. Prin urmare, totul se reduce la calculul lui $F = e^{hA}$, utilizând, de exemplu, aproximația Padé. De fapt, sistemul (4.3) constituie reprezentarea discretizată exactă – cu element de extrapolare de ordinul 0 – a sistemului (4.1) (vezi § 4.3).

Observația 4.1 Din acest punct de vedere, dacă pasul de simulare h nu este impus, atunci el se poate alege egal cu pasul de discretizare, utilizând condiția cunoscută $h\|A\| < \frac{1}{2}$. Pe de altă parte, dacă t_f nu este nici el precizat dar se urmărește obținerea regimului tranzitoriu ”esențial”, atunci – în ipoteza că matricea A este stabilă cu valorile proprii satisfăcând condiția $\text{Re } \lambda(A) < -\alpha$, unde $\alpha > 0$ – putem evalua constanta de timp dominantă prin

$$T = \frac{1}{\alpha}$$

iar timpul de răspuns (pentru banda de 5% din jurul regimului staționar constant) prin

$$t_f = (4 \div 5) T = \frac{4 \div 5}{\alpha}.$$

În definitiv, un N ales astfel încât

$$Nh \approx \frac{4 \div 5}{\alpha}.$$

asigură dezideratul propus. (În legătură cu aceasta, vezi și observația 5.1).

◇

În concluzie, convenind să memorăm valorile $y_k \stackrel{\text{def}}{=} y(kh)$ într-o matrice Y de forma $l \times N$ astfel încât $Y(:, k) = y_{k-1}$, $k = 1 : N$, calculul răspunsului liber al sistemelor liniare se poate face utilizând următoarea procedură.

Algoritmul 4.1 (Se dau sistemul liber (A, C) , starea inițială $x(0) = x$, intervalul de simulare $[0, t_f]$ precum și parametrii h , N astfel încât $Nh = t_f$. Se calculează răspunsul liber Y în momentele $t_k = kh$, $k = 0 : N - 1$ și se returnează $x = x(t_f)$.)

1. Se calculează $F = e^{hA}$.

2. Pentru $k = 1 : N$

1. $Y(:, k) = Cx$

2. $x \leftarrow Fx$

Comentarii. Starea finală $x(t_f) = x$, furnizată de algoritmul 4.1, este necesară în calitatea de nouă stare inițială pentru o eventuală continuare a simulării pe un interval de timp $[t_f, t'_f]$, ulterior lui t_f . \diamond

4.2 Răspunsul la intrări liniar generate

Problema determinării răspunsului unui sistem liniar $S = (A, B, C, D)$ constă în calculul funcției de ieșire $y(t)$, definită prin

$$(S) \quad \begin{cases} \dot{x}(t) = Ax(t) + Bu(t), & x(0) = x \\ y(t) = Cx(t) + Du(t), \end{cases} \quad (4.4)$$

pe un interval fixat $[0, t_f]$, pe care funcția de intrare $u(t)$ se consideră cunoscută.

Soluția ecuației diferențiale (4.4) este

$$x(t) = e^{tA}x(0) + \int_0^t e^{(t-\tau)A}Bu(\tau) d\tau. \quad (4.5)$$

În cazul unei intrări de tip impuls $u(t) = u_0\delta(t)$, $u_0 \in \mathcal{R}^m$ (v. fig. 4.1 a)) se obține

$$x(t) = e^{tA}x(0) + e^{tA}Bu_0, \quad (4.6)$$

respectiv

$$y(t) = Ce^{tA}x(0) + [Ce^{tA}B + D\delta(t)]u_0, \quad (4.7)$$

deci răspunsul poate fi calculat (mai puțin termenul impulsiv $Du_0\delta(t)$) utilizând algoritmul 4.1 cu datele A, C și $x = x(0) \leftarrow x(0) + Bu_0$, unde Bu_0

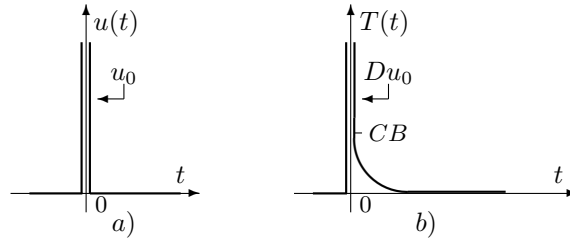


Figura 4.1: Funcție de intrare tip impuls (a) și funcția pondere a unui sistem liniar (b).

este saltul traiectoriei de stare (4.6) datorat impulsului aplicat la intrare în momentul $t = 0$.

Totodată, din (4.7) rezultă că matricea pondere (v. fig. 4.1 b))

$$T(t) \stackrel{\text{def}}{=} Ce^{tA}B1(t) + D\delta(t), \quad (4.8)$$

poate fi calculată pe coloane determinând succesiv (prin metoda indicată mai sus) răspunsurile sistemului în stare inițială nulă $x(0) = 0$ la impulsurile unitate $u_0 = e_i \in \mathcal{R}^m$, $i = 1 : m$. Odată $T(t)$ cunoscut, răspunsul la o intrare $u(t)$ oarecare, e.g. continuă pe porțiuni, poate fi, în principiu, obținut utilizând relația

$$y(t) = Ce^{tA}x(0) + \int_0^t T(t-\tau)u(\tau) d\tau. \quad (4.9)$$

Deoarece nu există metode generale de reprezentare a unei funcții de intrare $u(t)$ „arbitrară”, iar evaluarea numerică a integralei din (4.9) este dificilă ¹, în cele ce urmează vom considera situația (particulară, dar frecvent întâlnită în aplicații) în care $u(t)$ este *liniar generată*, adică poate fi modelată ca ieșire a unui sistem liniar liber (S_F) de forma

$$(S_F) \quad \begin{cases} \dot{x}_F(t) = A_F x_F(t), & x_F(0) = x_F \\ u(t) = C_F x_F(t), \end{cases} \quad (4.10)$$

unde A_F , C_F sunt matrice date, iar starea inițială $x_F \in \mathcal{R}^{n_F}$ este un vector arbitrar, dar fixat. (Altfel spus, precizând x_F fixăm implicit intrarea în

¹Câteva metode de reprezentare aproximativă a unor funcții de intrare de forma generală, precum și metodele corespunzătoare de calcul al răspunsului, sunt prezentate într-o altă secțiune a acestui capitol.

clasa tuturor intrărilor $u(t)$ ce pot fi obținute ca ieșiri ale lui S_F). Sistemul S_F se numește *generator* sau *filtru de formare* al funcției de intrare u . Conectând sistemele (S) și (S_F) în serie putem scrie

$$(S^0) \left\{ \begin{array}{l} \begin{bmatrix} \dot{x}(t) \\ \dot{x}_F(t) \end{bmatrix} = \begin{bmatrix} A & BC_F \\ 0 & A_F \end{bmatrix} \begin{bmatrix} x(t) \\ x_F(t) \end{bmatrix}, \quad \begin{bmatrix} x(0) \\ x_F(0) \end{bmatrix} = \begin{bmatrix} x \\ x_F \end{bmatrix} \\ y(t) = [C \quad DC_F] \begin{bmatrix} x(t) \\ x_F(t) \end{bmatrix}, \end{array} \right. \quad (4.11)$$

i.e. calculul răspunsului sistemului (S) la intrarea liniar generată $u(t)$ se reduce la calculul răspunsului liber al sistemului (S^0) , prin metoda (exactă) expusă în paragraful anterior.

Având în vedere că aplicând transformarea Laplace ecuațiilor (4.10) ale sistemului (S_F) obținem

$$u(s) = C_F(sI - A_F)^{-1}x_F, \quad (4.12)$$

rețeta generală de construcție a sistemului generator (S_F) constă în realizarea transformatei $u(s)$, interpretată ca funcție de transfer a unui sistem SIMO, sub forma unui triplet (A_F, b_F, C_F) cu o singură intrare, unde vectorul $b_F = x_F$ este stare inițială a lui (S_F) . De aici rezultă că funcția de intrare $u(t)$ este liniar generată dacă și numai dacă transformata sa Laplace este o matrice rațională strict proprie.

În concluzie, calculul răspunsului sistemelor liniare la intrări liniar generate se poate face utilizând următoarea procedură.

Algoritm 4.2 (Se dau sistemul (A, B, C, D) , starea inițială $x(0) = x$, intervalul de simulare $[0, t_f]$, parametrii h, N astfel încât $Nh = t_f$, precum și transformata $u(s)$ sub forma matricei de transfer a unui sistem SIMO. Se calculează răspunsul Y în momentele $t_k = kh, k = 0 : N-1$ și se returnează $x = x(t_f)$.)

1. Se realizează $u(s)$ sub forma (A_F, b_F, C_F) .
2. Se construiesc matricele

$$A \leftarrow A^0 = \begin{bmatrix} A & BC_F \\ 0 & A_F \end{bmatrix}, \quad x \leftarrow x^0 = \begin{bmatrix} x \\ b_F \end{bmatrix}$$

$$C \leftarrow C^0 = [C \quad DC_F].$$

3. Se aplică algoritmul 4.1 de calcul al răspunsului liber cu datele de intrare A, C și $x(0) = x$.

4. $x = x(1 : n)$.

Pentru exemplificare, vom considera câteva cazuri particulare importante.

A. Intrări polinomiale

Intrările polinomiale sunt de forma

$$u(t) = \sum_{k=1}^p \alpha_k t^{k-1}, \quad t \geq 0, \quad (4.13)$$

unde coeficienții $\alpha_k \in \mathcal{R}^m$ sunt vectori dați. Notând $u_k \stackrel{\text{def}}{=} \alpha_k(k-1)!$ avem

$$u(t) = \sum_{k=1}^p u_k \frac{t^{k-1}}{(k-1)!}, \quad t \geq 0, \quad (4.14)$$

iar transformata Laplace corespunzătoare este

$$u(s) = \sum_{k=1}^p \frac{u_k}{s^k} = \frac{u_1 s^{p-1} + \dots + u_p}{s^p}. \quad (4.15)$$

Deci, o realizare de stare a lui (S_F) , în forma standard controlabilă, (vezi capitolul 3) este

$$A_F = \begin{bmatrix} 0 & \dots & 0 & 0 \\ 1 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \end{bmatrix}, \quad b_F = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (4.16)$$

$$C_F = [u_1 \quad \dots \quad u_{p-1} \quad u_p].$$

În particular, dacă $p = 1$ atunci $u(t) = u_1 1(t)$ sau

$$u(t) = u_1, \quad t \geq 0 \quad (4.17)$$

este o intrare *treaptă* de amplitudine $u_1 \in \mathcal{R}^m$ dată (v. fig. 4.2 a)). În acest caz filtrul de formare este

$$\begin{aligned} A_F &= [0], & b_F &= [1] \\ C_F &= [u_1] \end{aligned} \quad (4.18)$$

iar matricele sistemului (S^0) sunt

$$A^0 = \begin{bmatrix} A & Bu_1 \\ 0 & 0 \end{bmatrix}, \quad x^0 = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad (4.19)$$

$$C^0 = [C \quad Du_1].$$

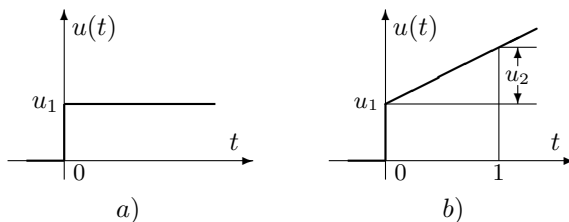


Figura 4.2: Funcțiile de intrare tip treaptă (a) și rampă (b).

În mod analog, dacă $p = 2$ atunci

$$u(t) = u_1 + u_2 t, \quad t \geq 0 \quad (4.20)$$

este o intrare *rampă* de amplitudine inițială $u(0) = u_1$ și pantă $u_2 \in \mathcal{R}^m$ date (v. fig. 4.2 b)). În acest caz filtrul de formare este

$$A_F = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad b_F = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (4.21)$$

$$C_F = [u_1 \quad u_2].$$

iar matricele sistemului (S^0) sunt

$$A^0 = \begin{bmatrix} A & Bu_1 & Bu_2 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad x^0 = \begin{bmatrix} x \\ 1 \\ 0 \end{bmatrix} \quad (4.22)$$

$$C^0 = [C \quad Du_1 \quad Du_2].$$

Clasa tuturor intrărilor polinomiale de forma (4.14), în care vectorii u_k , $k = 1 : p$ sunt arbitrari, se obține considerând în loc de (4.16) forma

standard observabilă

$$A_F = \begin{bmatrix} 0 & I_m & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & I_m \\ 0 & 0 & \cdots & 0 \end{bmatrix}, \quad b_F = \begin{bmatrix} u_1 \\ \cdots \\ u_{p-1} \\ u_p \end{bmatrix} \quad (4.23)$$

$$C_F = [I_m \quad 0 \quad \cdots \quad 0],$$

unde vectorul $b_F = x_F$ se alege pentru a preciza intrarea în clasa considerată. (Observăm că, acum, S_F este de ordin $n_F = pm$, unde p este ordinul transformatei Laplace (4.15), iar m este numărul de intrări ale sistemului (4.4)).

În particular, dacă $u(t) = u_1 1(t)$ este o intrare treaptă de amplitudine arbitrară, atunci filtrul de formare este

$$\begin{aligned} A_F &= [0], & b_F &= [u_1] \\ C_F &= [I_m], \end{aligned} \quad (4.24)$$

iar matricele sistemului (S^o) sunt

$$\begin{aligned} A^o &= \begin{bmatrix} A & B \\ 0 & 0 \end{bmatrix}, & x^o &= \begin{bmatrix} x \\ u_1 \end{bmatrix} \\ C^o &= [C \quad D]. \end{aligned} \quad (4.25)$$

B. Intrări armonice

Intrările armonice sunt de forma

$$u(t) = \gamma_1 \cos \omega t + \gamma_2 \sin \omega t, \quad t \geq 0, \quad (4.26)$$

unde pulsația ω este fixată iar coeficienții $\gamma_1, \gamma_2 \in \mathcal{R}^m$ sunt vectori dați. Avem

$$u(s) = \frac{\gamma_1 s + \gamma_2 \omega}{s^2 + \omega^2}, \quad (4.27)$$

deci o realizare a lui (S_F) este

$$\begin{aligned} A_F &= \begin{bmatrix} 0 & -\omega^2 \\ 1 & 0 \end{bmatrix}, & b_F &= \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ C_F &= [\gamma_1 \quad \omega \gamma_2]. \end{aligned} \quad (4.28)$$

Cazurile relativ mai complicate, de tipul

$$u = \sum_{k=1}^p \gamma_k \cos k\omega t + \delta_k \sin k\omega t, \quad (4.29)$$

$$u = e^{\alpha t}(\gamma_1 \cos \omega t + \gamma_2 \sin \omega t), \quad (4.30)$$

în care coeficienții sunt eventual arbitrari, se tratează analog.

4.3 Răspunsul la intrări etajate

Domeniul de aplicație al metodei bazate pe utilizarea sistemului generator, expusă în paragraful anterior, poate fi sensibil lărgit, considerând că funcția de intrare $u(t)$ este *liniar generată pe porțiuni*, e.g. are o expresie de tip (4.14) pe fiecare subinterval $[kh, (k+1)h)$, $k = 0 : N-1$ al unei diviziuni uniforme cu pasul h al intervalului $[0, t_f]$.

În cel mai simplu caz (care în același timp este și cel mai frecvent întâlnit în aplicații) funcția de intrare $u(t)$ este constantă pe porțiuni (sau *etajată*), i.e. avem

$$u(t) = u_k, \quad t \in [kh, (k+1)h) \quad (4.31)$$

unde amplitudinile $u_k \in \mathcal{R}^m$ sunt vectori dați, în general variabili de la pas la pas (v. fig. 4.3).

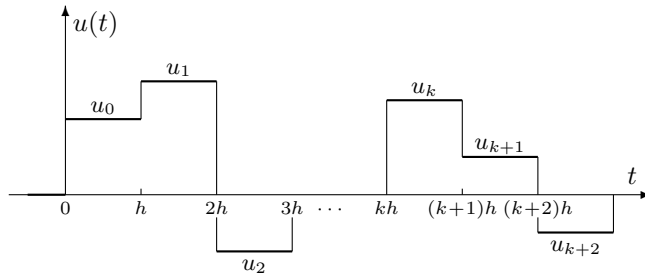


Figura 4.3: Funcție de intrare etajată.

Utilizând relația (4.6) scrisă sub forma

$$x(t) = e^{(t-kh)A}x(kh) + \int_{kh}^t e^{(t-\tau)A}Bu(\tau) d\tau,$$

în care punem $t = (k + 1)h$, în virtutea lui (4.31) obținem ușor

$$\begin{cases} x((k + 1)h) = Fx(kh) + Gu_k, \\ y(kh) = Cx(kh) + Du_k, \end{cases} \quad (4.32)$$

unde

$$F = e^{hA}, \quad G = \int_0^h e^{\eta A} d\eta B. \quad (4.33)$$

Prin definiție, sistemul liniar discret $S_d = (F, G, C, D)$ constituie *reprezentarea discretizată* (pe scurt discretizatul) cu pasul h al sistemului liniar $S = (A, B, C, D)$. Matricele F și G , necesare pentru construcția lui S_d , se determină calculând exponențiala matriceală $F^o = e^{hA^o}$, unde matricea A^o are expresia (4.25), corespunzătoare cazului în speță. (Amintim că, pe fiecare subinterval, funcția de intrare (4.31) este de tip treaptă). Procedând în acest mod, se obține

$$e^{hA^o} = \begin{bmatrix} F & G \\ 0 & I_m \end{bmatrix}. \quad (4.34)$$

În concluzie, convenind să dispunem valorile date u_k într-o matrice U de forma $m \times N$, astfel încât $U(:, k) = u_{k-1}$, $k = 1 : N$, calculul răspunsului sistemelor liniare la intrări etajate se poate face utilizând următoarea procedură

Algoritmul 4.3 (Se dau sistemul (A, B, C, D) , starea inițială $x(0) = x$, intervalul de simulare $[0, t_f]$, parametrul h , N astfel încât $Nh = t_f$ precum și tabloul U ce conține palierele funcției de intrare etajate $u(t)$. Se calculează răspunsul Y în momentele $t_k = kh$, $k = 0 : N - 1$ și se returnează $x = x(t_f)$).

1. Se calculează $F = e^{hA}$ și G în acord cu (4.34), unde A^o are expresia (4.25).
2. Pentru $k = 1 : N$
 1. $Y(:, k) = Cx + DU(:, k)$
 2. $x \leftarrow Fx + GU(:, k)$.

Idei asemănătoare se aplică în cazul utilizării unor reprezentări ale funcției de intrare $u(t)$, mai complicate decât (4.31). În cazul general, corespunzător lui (4.14), putem scrie

$$u(t) = \sum_{i=1}^p u_k^{(i)} \frac{(t - kh)^{i-1}}{(i-1)!}, \quad t \in [kh, (k+1)h) \quad (4.35)$$

unde coeficienții $u_k^{(i)}$ au o semnificație evidentă, în particular $u_k^{(1)} = u(kh)$ coincide cu u_k din (4.31). În locul lui (4.32), acum se obține o realizare discretizată de forma

$$\begin{cases} x((k+1)h) = Fx(k) + \sum_{i=1}^p G_i u_k^{(i)} \\ y(kh) = Cx(kh) + Du_k^{(1)}, \end{cases} \quad (4.36)$$

în care matricele $F = e^{hA}$ și

$$G_i \stackrel{\text{def}}{=} \int_0^h e^{(h-\tau)A} \frac{\tau^{i-1}}{(i-1)!} d\tau, \quad i = 1 : p$$

se determină în mod corespunzător, calculând exponențiala matriceală $F^o = e^{hA^o}$, unde A^o se construiește cu datele din (4.23). (Stabilirea formei lui F^o , în cazul de față, este propusă cititorului ca exercițiu. Evident, avem $G_1 = G$ și în cazul $p = 1$ ecuațiile (4.36) se reduc la (4.32)).

Pentru determinarea coeficienților $u_k^{(i)}$, în special dacă asigurarea „raccordării” aproximațiilor (4.35) pe intervale adiacente este importantă, se recomandă utilizarea metodelor de interpolare tip spline, [IX]. Subliniem că în acest caz, sistemul (4.36) nu este causal în sensul că, în general, coeficienții $u_k^{(i)}$, deci starea $x((k+1)h)$ rezultată din prima ecuație (4.36), depind de valorile viitoare ale intrării $u(t)$ ². În schimb, procedurile de calcul astfel obținute (net superioare tehnicilor de extrapolare clasice) combină precizia aproximărilor spline cu eficiența și siguranța metodelor de calcul ale exponențialei matriceale, concurând deseori cu succes procedurile alternative de tip general, bazate pe integrarea numerică a ecuației diferențiale (4.4).

Fără a intra în detalii, menționăm numai că utilizarea metodelor de integrare numerică (tip Runge - Kutta, predictor - corector, etc., în general cu pas variabil) devine obligatorie în cazul general, în care matricea A a sistemului liniar $S = (A, B, C, D)$ este variabilă în timp. În locul exponențialei matriceale $F = e^{hA}$, pe fiecare subinterval $[t_k, t_{k+1})$, acum se utilizează matricea de tranziție $F_k = \Phi(t_{k+1}, t_k)$, calculată pe coloane prin integrarea ecuației diferențiale omogene (4.1) cu condițiile inițiale $x(t_k) = e_j \in \mathcal{R}^n$, $j = 1 : n$.

² Acest fenomen poate fi evitat numai în cazul $p = 2$, când în (4.35) se obține $u_k^{(1)} = u_k$ și $u_k^{(2)} = (u_{k+1} - u_k)/h$. Efectuând în (4.36) schimbarea de variabilă de stare $x((k+1)h) \leftarrow x((k+1)h) - G_2 u_{k+1}/h$, ecuațiile (4.36) pot fi scrise sub forma uzuală (4.32). În cazul $p = 4$, corespunzător funcțiilor spline cubice, frecvent utilizate în practică, acest artificiu nu mai este eficient.

4.4 Răspunsul staționar

Considerăm contextul din § 4.2, unde presupunem, în plus, că matricea A este stabilă, adică $\operatorname{Re} \lambda(A) < 0$, iar funcția de intrare u este *persistentă*, adică $\operatorname{Re} \lambda(A_F) \geq 0$.

Procedura generală de separare a componentelor *tranzitorie* și „*staționară*” (*asimptotică* sau *permanentă*) ale răspunsului $y(t)$ constă în a scrie

$$x(t) = \xi(t) + Vx_F(t), \quad (4.37)$$

unde $\xi(t)$ este componenta *tranzitorie* a traiectoriei de stare, adică $\xi(t) \rightarrow 0$ când $t \rightarrow \infty$, iar matricea constantă $V \in \mathcal{R}^{n \times n_F}$ trebuie determinată. Derivând relația (4.37) și ținând seama de ecuațiile sistemelor (S) și (S_F) , obținem

$$\begin{aligned} \dot{x} &= \dot{\xi} + V\dot{x}_F = \dot{\xi} + VA_Fx_F = \\ &= Ax + Bu = \\ &= A(\xi + Vx_F) + BC_Fx_F = \\ &= A\xi + (AV + BC_F)x_F. \end{aligned} \quad (4.38)$$

Prin urmare, dacă matricea V este aleasă astfel încât

$$VA_F = AV + BC_F \quad (4.39)$$

atunci din (4.38) rezultă $\dot{\xi}(t) = A\xi(t)$, deci $\xi(t) \rightarrow 0$ când $t \rightarrow \infty$, $\forall \xi(0)$ în virtutea ipotezei de stabilitate $\operatorname{Re} \lambda(A) < 0$. (În fapt, dacă $x(0)$ este precizat, atunci din (4.37) avem $\xi(0) = x(0) - Vx_F(0)$).

Cu alte cuvinte, considerând un moment inițial flotant $t_0 \in \mathcal{R}$, presupunând mai sus $t \geq t_0$ și făcând $t_0 \rightarrow -\infty$, obținem din nou $\xi(t) \rightarrow 0$, deci $x(t) \rightarrow Vx_F(t)$, $\forall x(t_0)$. În virtutea acestui fapt, spunem că $Vx_F(t)$ este componenta „*staționară*” (*asimptotică* sau *permanentă*) a traiectoriei de stare $x(t)$.

Observația 4.2 După cum se știe, ecuația matriceală Sylvester (4.39) are soluție unică V oricare ar fi membrul drept BC_F dacă și numai dacă $\lambda(A) \cap \lambda(A_F) = \emptyset$, (vezi capitolul 1) adică nici o valoare proprie a lui A nu coincide cu nici o valoare proprie a lui A_F . În cazul de față această condiție este evident satisfăcută în virtutea ipotezelor de stabilitate $\operatorname{Re} \lambda(A) < 0$ și persistență $\operatorname{Re} \lambda(A_F) \geq 0$. \diamond

În continuare obținem

$$y(t) = C(\xi(t) + Vx_F(t)) + DC_Fx_F(t), \quad (4.40)$$

deci avem

$$y(t) = \eta(t) + Wx_F(t), \quad (4.41)$$

unde

$$\eta(t) = C\xi(t) \quad (4.42)$$

este (prin definiție) componenta *tranzitorie* a răspunsului, iar matricea W este definită prin

$$W = CV + DC_F. \quad (4.43)$$

Având în vedere că, în esență, $x_F(t) = e^{tA_F}x_F(0)$ este cunoscut, componenta *staționară* $Wx_F(t)$ a răspunsului este complet determinată prin (4.43) în funcție de soluția V a ecuației matriceale algebrice liniare (4.39).

În concluzie, calculul răspunsului staționar al sistemelor liniare la intrări persistente liniar - generate se poate face utilizând următoarea procedură.

Algoritm 4.4 (Se dau sistemul (A, B, C, D) precum și transformata $u(s)$ sub forma matricei de transfer a unui sistem SIMO. Se calculează răspunsul staționar $Wx_F(t)$ pe un interval $[0, t_f]$ convenabil precizat).

1. Se realizează $u(s)$ sub forma (A_F, b_F, C_F) .
2. Se calculează matricea V rezolvând ecuația matriceală Sylvester (4.39).
3. Se calculează matricea W în acord cu (4.43).
4. Se determină răspunsul staționar aplicând algoritmul 4.1 cu datele de intrare A_F, W și $x_F = b_F$.

Pentru exemplificare, vom considera câteva cazuri particulare importante.

Răspunsul staționar la intrări polinomiale

Intrările polinomiale au un filtru de formare S_F definit prin matricele A_F, C_F date în (4.16). Prin urmare, partiționând V pe coloane sub forma $V = [v_1 \dots v_p]$ unde $v_i \in \mathcal{R}^n$, ecuația (4.39) se scrie

$$\begin{aligned} [v_1 \quad v_2 \quad \dots \quad v_p] \begin{bmatrix} 0 & \dots & 0 & 0 \\ 1 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \end{bmatrix} &= \\ &= A [v_1 \quad v_2 \quad \dots \quad v_p] + B [u_1 \quad \dots \quad u_{p-1} \quad u_p]. \end{aligned} \quad (4.44)$$

Obținem

$$\begin{cases} v_2 = Av_1 + Bu_1 \\ \vdots \\ v_p = Av_{p-1} + Bu_{p-1} \\ 0 = Av_p + Bu_p \end{cases} \quad (4.45)$$

deci vectorii $v_p \dots v_1$ rezultă succesiv rezolvând sistemele liniare

$$\begin{aligned} Av_p &= -Bu_p \\ Av_k &= -Bu_k + v_{k+1}, \quad k = p-1 : 1. \end{aligned} \quad (4.46)$$

(Evident, matricea A trebuie triangularizată o singură dată.)

În sfârșit, din (4.43), unde $W = [w_1 \dots w_p]$ rezultă ușor

$$w_k = Cv_k + Du_k, \quad k = 1 : p. \quad (4.47)$$

Observația 4.3 Unicitatea soluției V este în acest caz garantată dacă A e inversabilă (adică $0 \notin \lambda(A)$), ceea ce evident rezultă din ipoteza A stabilă, adoptată mai sus. Altfel spus, S poate avea soluția polinomială $Vx_F(t)$, dar ea nu este soluție asimptotică (pentru $t \rightarrow \infty$), decât dacă A este stabilă. \diamond

În particular, răspunsul staționar la o intrare treaptă de amplitudine $u_1 \in \mathcal{R}^m$ corespunde la $p = 1$ și $x_F(t) = 1$, deci are expresia

$$W = [-CA^{-1}B + D]u_1 \stackrel{\text{def}}{=} T(0)u_1,$$

unde

$$T(0) = [C(sI - A)^{-1}B + D] \Big|_{s=0}.$$

Răspunsul staționar la intrări armonice

Intrările armonice au un filtru de formare definit prin

$$\begin{aligned} A_F &= \begin{bmatrix} 0 & -\omega^2 \\ 1 & 0 \end{bmatrix}, \\ C_F &= [\gamma_1 \quad \omega\gamma_2]. \end{aligned} \quad (4.48)$$

Partiționând V pe coloane sub forma $V = [v_1 \quad \omega v_2]$ ecuația (4.39) se scrie

$$[v_1 \quad \omega v_2] \begin{bmatrix} 0 & -\omega^2 \\ 1 & 0 \end{bmatrix} = A[v_1 \quad \omega v_2] + B[\gamma_1 \quad \omega\gamma_2] \quad (4.49)$$

Obținem

$$\begin{cases} \omega v_2 = Av_1 + B\gamma_1 \\ -\omega v_1 = Av_2 + B\gamma_2 \end{cases} \quad (4.50)$$

deci vectorii v_i , $i = 1 : 2$ rezultă rezolvând sistemul liniar

$$\begin{bmatrix} -A & \omega I \\ -\omega I & -A \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} B\gamma_1 \\ B\gamma_2 \end{bmatrix}. \quad (4.51)$$

În sfârșit, din (4.43), unde $W = [w_1 \ \omega w_2]$, rezultă ușor

$$w_i = Cv_i + D\gamma_i, \quad i = 1 : 2 \quad (4.52)$$

Observația 4.4 Unicitatea soluției V este în acest caz garantată dacă $i\omega \notin \lambda(A)$, adică în absența rezonanței pe frecvență ω a intrării, ceea ce rezultă evident din ipoteza A stabilă, adoptată mai sus. \diamond

4.5 Calculul caracteristicilor de frecvență

Deseori, problema de calcul abordată în finalul paragrafului precedent intervine sub o formă modificată, adică se cere calculul răspunsului armonic pentru diverse valori ale frecvenței ω renunțând la ipoteza A stabilă în favoarea condiției mai slabe $i\omega \notin \lambda(A)$ de absență a rezonanțelor pe frecvențele ω prescrise. Evident, aceasta înseamnă calculul caracteristicilor de frecvență ale sistemului considerat, ceea ce, tradițional, se face recurgând la "complexificarea" acestuia.

Concret, aceasta înseamnă că în locul funcției de intrare reale (4.26), i.e. $u(t) = \gamma_1 \cos \omega t + \gamma_2 \sin \omega t$ se consideră intrarea complexă

$$u_c(t) = \gamma_c e^{i\omega t}, \quad (4.53)$$

unde $\gamma_c \in \mathcal{C}^m$ este un vector dat cu componente complexe. Avem

$$u_c(s) = \frac{\gamma_c}{s - i\omega} \quad (4.54)$$

deci o realizare (complexă) a lui (S_F) este

$$\begin{aligned} A_F &= i\omega, & b_F &= 1, \\ C_F &= \gamma_c. \end{aligned} \quad (4.55)$$

Dacă scriem acum (4.39) sub forma

$$i\omega V = AV + B\gamma_c \quad (4.56)$$

obținem imediat

$$V = (i\omega I - A)^{-1} B\gamma_c, \quad (4.57)$$

iar din (4.43) rezultă

$$W = CV + D\gamma_c = [C(i\omega I - A)^{-1} B + D]\gamma_c \stackrel{\text{def}}{=} T(i\omega)\gamma_c \quad (4.58)$$

unde

$$T(i\omega) = C(i\omega I - A)^{-1} B + D, \quad \omega \in \mathcal{R} \quad (4.59)$$

reprezintă caracteristica (complexă) de frecvență a sistemului (S).

Problema calculului lui $T(i\omega)$ pentru diverse valori $\omega_k \in \mathcal{R}$, $k = 1 : N$ ale lui ω presupune deci

1. Pentru $k = 1 : N$

1. Rezolvarea sistemelor matriceale liniare $(i\omega_k I - A)V_k = B$.
2. Calculul matricei $T(i\omega_k) = CV_k + D$.

Deosebirea importantă față de cazul considerat anterior al intrărilor polinomiale constă aici, (în afara necesității evidente de a utiliza aritmetica complexă), în faptul că, aparent, la pasul 1.1 matricele $(i\omega_k I - A)$ trebuie să fie triangularizate pentru fiecare k în parte, ceea ce în cazul în care A este de forma generală presupune, în total, $N\frac{n^3}{3}$ operații (în afara celor Nmn^2 necesare pentru determinarea coloanelor lui V). Pentru a evita aceasta, amintim că rezolvarea unui sistem $Ax = b$ cu A superior (sau inferior) Hessenberg cere numai n^2 operații. Prin urmare, în cazul de față se recomandă determinarea unei transformări prelabile M astfel încât

$$\tilde{A} = MAM^{-1}$$

să fie, de exemplu, superior Hessenberg. În acest scop, se utilizează eliminarea gaussiană cu pivotare parțială în aritmetica complexă și rezultă $M = M_{n-1}P_{n-1} \dots M_2P_2$, unde M_k sunt matrice inferior triunghiulare elementare, iar P_k sunt matrice de permutare adecvat alese. (Alternativ, dar mai costisitor, se poate utiliza procedura ortogonală în care $M = U = U_{n-1} \dots U_2$ este o secvență de reflectori). Aplicând transformarea M matricelor B și C , adică $\tilde{B} = MB$, $\tilde{C} = CM^{-1}$ și definind $\tilde{V} = MV$, problema calculului lui $T(i\omega_k)$ se reduce la

1. Pentru $k = 1 : N$

1. Rezolvarea sistemelor matriceale liniare superior Hessenberg $(i\omega_k I - \tilde{A})\tilde{V}_k = \tilde{B}$.
2. Calculul matricei $T(i\omega_k) = \tilde{C}\tilde{V}_k + D$.

Aplicarea acestei scheme de calcul necesită, la pasul 1.1, în total, nu mai Nmn^2 operații, mai precis, $Nm\frac{n^2}{2}$ operații pentru triangularizare plus $Nm\frac{n^2}{2}$ pentru rezolvarea efectivă.

Procedura de calcul poate fi rezumată astfel.

Algoritmul 4.5 (Se dau A, B, C, D și $\omega_k \in \mathcal{R}$, $k = 1 : N$, de regulă în scară logaritmică. Se calculează $T(i\omega_k)$, $k = 1 : N$).

1. Se determină M astfel încât $A \leftarrow \tilde{A} = MAM^{-1}$ este superior Hessenberg.
2. $B \leftarrow MB$
3. $C \leftarrow CM^{-1}$
4. Pentru $k = 1 : N$
 1. Se formează matricea $A_0 = i\omega_k I - A$.
 2. Se face triangularizarea $A_0 \leftarrow R = NA_0$.
 3. $B_0 = NB$
 4. Se rezolvă $A_0\tilde{V} = B_0$ în raport cu \tilde{V} .
 5. $T(i\omega_k) = C\tilde{V} + D$

Comentarii. Matricele A_0 și B_0 sunt tablouri de lucru. La pasul 4.2. se efectuează eliminarea gaussiană cu pivotare parțială în aritmetică complexă astfel încât R rezultă superior triunghiulară. La pasul 4.3 transformarea N se aplică lui B iar la pasul 4.4 se procedează ca în cap. 1, §2. \diamond

Observația 4.5 Aducerea prealabilă a lui A la forma Schur (pentru a evita pașii 4.2, 4.3) este mai puțin eficientă. \diamond

4.6 Răspunsul sistemelor liniare discrete

Datorită specificului recurent al ecuațiilor de stare în cazul discret, matricea de tranziție este $\Phi(k) = A^k$, $k \in \mathcal{N}$. În particular, matricea de tranziție

într-un pas $\Phi(1) = A$ este dată, deci procedurile de calcul prezentate mai sus se simplifică. În algoritmi 4.1 și 4.3 se elimină pașii 1, înlocuind peste tot perechea discretizată (F, G) cu perechea (A, B) a sistemului discret considerat. În algoritmul 4.2, în locul lui $u(s)$ se consideră transformata \mathcal{Z} corespunzătoare $u(z)$. Considerațiile din § 4.4 rămân valabile înlocuind condițiile de stabilitate și persistență prin versiunile lor discrete $|\lambda(A)| < 1$ și, respectiv, $|\lambda(A_F)| \geq 1$. (Vezi § 5.1). În sfârșit, în versiunea discretă a algoritmului 4.5, numerele $i\omega_k$ situate pe axa imaginară se înlocuiesc cu numere $e^{i\theta_k}$, $k = 1 : N$, convenabil plasate pe cercul unitate T . Detaliile de redactare sunt propuse cititorului ca exerciții.

Programe MATLAB disponibile

Pentru calculul și reprezentarea grafică a răspunsului liber al unui sistem liniar este disponibilă funcția **initial**, care folosește algoritmul 4.1. Răspunsul în stare inițială nulă la o intrare (scalară) de tip impuls sau treaptă unitate se poate calcula utilizând funcțiile **impulse** și **step**, care implementează algoritmi corespunzători 4.1 și 4.2 considerând $x(0) = b$ și respectiv $x(0) = 0$. Funcția **lsim** implementează algoritmul 4.3 utilizând o aproximare liniară pe porțiuni a funcției de intrare $u(t)$, vezi formula (4.35) în cazul $p = 2$. În cazul discret, în aceleași scopuri, se utilizează funcțiile **dinitial**, **dimpulse**, **dstep** și **dlsim**.

Pentru construcția reprezentării discretizate cu pasul h a unei perechi (A, B) este disponibilă funcția **c2d**, bazată pe formula de calcul exactă (4.33). Diverse metode de discretizare aproximativă pot fi utilizate apelând funcția **c2dm**.

Pentru calculul și reprezentarea grafică a caracteristicilor de frecvență ale unui sistem liniar sunt disponibile funcțiile **nyquist**, **bode** și **nichols** care apelează procedura **freqresp**, bazată pe algoritmul 4.5. În cazul discret, în același scop se utilizează funcțiile **dnyquist**, **dbode** și **dnichols**. Pentru funcții de transfer sunt disponibile versiunile simple **freqs** și respectiv **freqz**.

Exerciții

E 4.1 Prin analogie cu algoritmi din text, scrieți algoritmi de calcul al răspunsului în timp al unui sistem liniar, discret $S = (A, B, C, D)$ definit

prin ecuațiile de stare

$$\begin{cases} x(k+1) = Ax(k) + Bu(k), & x(0) = x \\ y(k) = Cx(k) + Du(k) \end{cases} \quad (4.60)$$

E 4.2 Reluați exercițiul 4.1 în ipoteza că matricele sistemului discret S au formele particulare obținute prin aplicarea procedurilor de realizare expuse în capitolul 1, §3. Considerați succesiv cazurile SISO, MISO, SIMO și în fiecare caz scrieți algoritmi propuși cât mai amănunțit posibil, evaluând memoria necesară precum și numărul de operații. Formulați concluziile rezultate.

E 4.3 Scrieți algoritmul de calcul al răspunsului unui sistem liniar discret S descris printr-o relație intrare - ieșire de forma

$$y(k) = \sum_{i=1}^n a_i y(k-i) + \sum_{i=0}^d b_i u(k-i), \quad k = 0 : N-1 \quad (4.61)$$

precizând cu grijă datele inițiale necesare.

E 4.4 Reluați exercițiul anterior aplicând transformata \mathcal{Z} ecuației (4.61) pentru a obține $T(z)$ și utilizați în acest context procedurile stabilite la exercițiul 4.2. Formulați concluziile rezultate.

E 4.5 Cum se calculează răspunsul în timp al unui sistem liniar discret având funcția de transfer dată sub forma

$$T(z) = k \frac{\prod(z - \mu_i)}{\prod(z - \lambda_i)}, \quad \text{respectiv} \quad T(z) = \sum \frac{r_i}{z - \lambda_i} ? \quad (4.62)$$

E 4.6 Reluați exercițiul 4.1 pentru sistemul liniar discret definit prin

$$\begin{cases} x(k+1) = Ax(k) + \sum_{i=0}^d B_i u(k-i), & x(0) = x \\ y(k) = Cx(k) + \sum_{i=0}^d D_i u(k-i) \end{cases} \quad (4.63)$$

în care $d \geq 1$ reprezintă întârzierea (pură) maximă pe canalul de intrare, iar A , C și B_i , D_i , $i = 1 : d$ sunt matrice date.

E 4.7 Calculați matricele de transfer $T(z)$ pentru sistemele din exercițiile 4.1 și 4.6, indicând o procedură eficientă de calcul al răspunsului staționar la intrare treaptă.

E 4.8 Elaborați proceduri eficiente de calcul pentru caracteristicile de frecvență $T(e^{i\theta})$, $\theta \in [0, 2\pi)$, unde $T(z)$ sunt matricele de transfer calculate în exercițiul precedent.

E 4.9 Considerați sistemul liniar descris prin

$$\begin{cases} \dot{x}(t) = Ax(t) + \sum_{i=0}^d B_i u(t - ih), & x(0) = x_0 \\ y(t) = Cx(t) + \sum_{i=0}^d D_i u(t - ih) \end{cases} \quad (4.64)$$

în care $d \geq 1$, iar funcția de intrare $u(t)$ este etajată cu pasul h . Scrieți algoritmul de calcul al răspunsului $y(kh)$, $k = 0 : N - 1$.

E 4.10 Formulați și rezolvați versiunile continue ale exercițiilor 4.2 – 4.5, 4.7 și 4.8.

E 4.11 Scrieți algoritmul de calcul al răspunsului sistemului liniar (4.4) la o intrare impuls $u(t) = u_0 \delta(t)$, unde vectorul $u_0 \in \mathcal{R}^m$ este dat. Cum procedați în cazul $u(t) = \sum_{k=0}^{N-1} u_k \delta(t - kh)$, $t_f = Nh$?

E 4.12 Formulați și rezolvați versiunea discretă a exercițiului 4.11.

E 4.13 Scrieți algoritmul de calcul al răspunsului sistemului liniar (4.4) la un semnal de intrare periodic de formă dreptunghiulară, presupunând că perioada acestuia este ph , $p \geq 1$.

E 4.14 Elaborați algoritmul de calcul al răspunsului în timp al sistemului liniar (4.4) pentru un semnal de intrare periodic având formă de dinți de fierăstrău cu perioada ph , $p \geq 1$ și panta u_2 i.e. $u(t) = u_2(t - kph)$, $t \in [kph, (k+1)ph)$, $k \in \mathcal{N}$.

E 4.15 Analizați posibilitățile de calcul al regimului staționar rezultat în condițiile exercițiilor 4.13 și 4.14.

E 4.16 Propuneți un analog al formulei (4.3) pentru sistemul liniar cu întârziere descris de ecuația diferențială $\dot{x}(t) = Ax(t) + Bx(t - \tau)$ cu condiția inițială $x(t) = u(t)$, considerată cunoscută pentru $t \in [-\tau, 0]$, unde $\tau > 0$ reprezintă întârzierea (internă), iar A și B sunt matrice date. Indicați o procedură de calcul aproximativ al răspunsului liber $y(kh) = Cx(kh)$, $k = 0 : N - 1$.

Indicație. Considerați mai întâi ecuația corespunzătoare cazului $n = 1$ și alegeți h astfel încât $\tau = dh$, unde, de exemplu, $d = 10$. Cum se pot calcula valorile $x(kh)$, $h = 1 : d$?

E 4.17 Scrieți algoritmul de calcul al răspunsului $y(kh)$, $k = 1 : N$, al sistemului (4.1) dacă în locul stării inițiale $x(0) = x_0$ se dă starea finală $x(t_f) = \xi$.

E 4.18 Scrieți algoritmul de calcul al răspunsului $y(kh)$, $k = 0 : N$, al sistemului (4.1) dacă între $x(0)$ și $x(t_f)$ se dă o relație de forma $Lx(0) + Mx(t_f) = d$, unde L, M sunt două matrice de ordin n iar $d \in \mathcal{R}^n$ este un vector dat. (Cazul $L = I_n, M = 0$ corespunde algoritmului 4.1, iar cazul $L = 0, M = I_n$ corespunde exercițiului 4.17).

Indicație. Spre deosebire de algoritmul 4.1 și exercițiul 4.17, care în esență efectuează integrarea sistemului $\dot{x} = Ax$ în timp direct și, respectiv, în timp invers, problema bilocală considerată acum nu are întotdeauna soluție unică. Considerând $t = t_f$ în expresia (4.2) a răspunsului liber, obținem $x(t_f) = e^{t_f A} x(0)$, deci trebuie să avem $(L + Me^{t_f A})x(0) = d$. Acest sistem de n ecuații liniare cu n necunoscute $x(0) \in \mathcal{R}^n$ are soluție unică dacă și numai dacă matricea $\bar{L} = L + Me^{t_f A}$ rezultă inversabilă. În acest caz, se calculează $x(0) = \bar{L}^{-1}d$ și se aplică algoritmul 4.1.

E 4.19 Formulați și rezolvați exercițiile 4.17 și 4.18 în cazul $u(\cdot) \neq 0$.

E 4.20 Formulați și rezolvați versiunile discrete ale exercițiilor 4.17 – 4.19, precizând ipotezele necesare pentru buna formulare a problemelor de calcul astfel obținute.

E 4.21 Scrieți un algoritm pentru rezolvarea *ecuației diferențiale matriciale Sylvester*

$$\dot{X}(t) = AX(t) + X(t)B, \quad (4.65)$$

cu condiția inițială $X(0) = X$, unde A, B, X sunt matrice date de dimensiuni corespunzătoare.

Indicație. Analogul formulei (4.2) în cazul de față este

$$X(t) = e^{tA} X(0) e^{tB}.$$

În cazul $B = A^T$ se obține *ecuația diferențială matriceală Liapunov*.

E 4.22 Reluați exercițiul precedent pentru ecuația

$$\dot{X}(t) = AX(t) + X(t)B + C \quad (4.66)$$

cu $X(0) = X$, unde, în plus, C este o matrice dată, eventual variabilă în timp.

Indicație. Analogul formulei (4.5) este

$$X(t) = e^{tA} X(0) e^{tB} + \int_0^t e^{(t-\tau)A} C e^{(t-\tau)B} d\tau.$$

Bibliografie

- [1] **Lupaș L., Bogdan M.** *Forced Response Computation of Linear Time - Invariant Multivariable Systems*, Rev. Roum. Sci. Tehn. – Electrotehn. et Energ., Vol. 19, No. 3, pp. 475 – 490, 1974.
- [2] **Lupaș L.** *A Numerical Method of Evaluating the Transient Response of Linear Time - Invariant Multivariable Systems*, Rev. Roum. Sci. Tehn. – Electrotehn. et Energ., Vol. 19, No. 4, pp. 639 - 648, 1974.
- [3] **Enright W.** *On the Efficient and Reliable Numerical Solution of Large Linear Systems of ODE's*, IEEE Trans. AC - 24, No. 6, pp. 905 – 908, 1979.
- [4] **Laub A.J.** *Efficient Multivariable Frequency Response Computations*, IEEE Trans. AC - 26, No. 2, pp. 407 - 408, 1981.

Capitolul 5

Proceduri de analiză sistemică

În acest capitol vom discuta procedurile de analiză numerică a proprietăților sistemice fundamentale (stabilitate, controlabilitate și observabilitate, etc.) precum și procedurile conexe de minimizare dimensională și echilibrare a modelelor sistemice liniare.

5.1 Stabilitatea sistemelor liniare

Stabilitatea unui sistem liniar $S = (A, B, C)$ este o proprietate ce depinde numai de matricea A a sistemului și, în ultimă instanță, numai de matricea de tranziție $\Phi(t)$ asociată lui A . Mai precis, sistemul S este *stabil* dacă $\Phi(t) \rightarrow 0$ când $t \rightarrow \infty$. Deoarece matricea de tranziție are expresii diferite în cazurile continuu și discret, adică $\Phi(t) = e^{tA}$, $t \in \mathcal{R}$, respectiv $\Phi(t) = A^t$, $t \in \mathcal{N}$, proprietatea de stabilitate are și ea exprimări diferite în cele două cazuri.

Definiția 5.1 *Matricea A a unui sistem continuu este stabilă dacă toate valorile proprii ale lui A au partea reală strict negativă, deci sunt plasate în semiplanul stâng deschis \mathcal{C}^- al planului complex \mathcal{C} .*

Pe scurt, criteriul (condiția necesară și suficientă) de stabilitate în cazul continuu este $\lambda(A) \subset \mathcal{C}^-$ sau, mai sugestiv, $\operatorname{Re}\lambda(A) < 0$.

Definiția 5.2 Matricea A a unui sistem discret este stabilă dacă toate valorile proprii ale lui A au modul strict subunitar, deci sunt plasate în discul unitate deschis $\mathcal{D}_1(0)$ din planul complex \mathcal{C} .

Pe scurt, criteriul de stabilitate în cazul discret este $|\lambda(A)| < 1$.

În general, o valoare proprie λ a lui A satisfăcând condiția de stabilitate $\operatorname{Re}\lambda < 0$, respectiv $|\lambda| < 1$, se numește valoare proprie stabilă. Prin urmare, matricea A este stabilă dacă și numai dacă toate valorile sale proprii sunt stabile.

Deoarece în ambele cazuri, continuu și discret, proprietatea de stabilitate vizează o anumită plasare a valorilor proprii, iar acestea se calculează eficient utilizând algoritmul **QR**, obținem următorul test de stabilitate pentru sisteme liniare.

Algoritmul 5.1 (Se testează stabilitatea matricei A).

1. Se calculează valorile proprii $\lambda(A)$ utilizând algoritmul QR (fără acumularea transformărilor).
2. Se determină $\alpha = \max \operatorname{Re}\lambda(A)$ în cazul continuu, respectiv $\rho = \max |\lambda(A)|$ în cazul discret. Matricea A este stabilă dacă și numai dacă $\alpha < 0$, respectiv $\rho < 1$.

Observația 5.1 Proprietatea de stabilitate este o proprietate binară, i.e. poate fi numai adevărată sau falsă. În aplicații interesează deseori evaluări cantitative ale unor "margini" sau "rezerve" de stabilitate. Asemenea evaluări sunt furnizate de parametri α și ρ calculați mai sus. De exemplu, dacă un sistem continuu este stabil, adică $\alpha < 0$, atunci $|\alpha|$ reprezintă inversul constantei de timp dominante, deci relația

$$t_r \approx \frac{4 \div 5}{|\alpha|}$$

oferă o evaluare orientativă a timpului de răspuns (definit pentru banda de 5% din jurul regimului staționar constant) al sistemului considerat. Pe de altă parte, dacă matricea discretă A rezultă prin discretizarea cu pasul h a unei matrice A_c , aparținând unui sistem continuu atunci avem $\lambda(A) = e^{h\lambda A_c}$, în particular $\rho = e^{h\alpha_c}$, deci numărul N_r de pași necesar pentru stabilizarea răspunsului discret poate fi evaluat cu formula

$$N_r \approx \frac{4 \div 5}{|\ln \rho|}.$$

◇

În general, numărul $\rho = \rho(A)$ se numește *rază spectrală* a matricei A .

5.2 Descompunerea spectrală

În anumite probleme de analiză și sinteză structurală în care matricea A nu este stabilă, interesează evidențierea subspațiilor invariante ale matricei A generate de vectorii proprii asociați valorilor proprii stabile și respectiv nestabile.

Pentru precizare, fie sistemul linear $S = (A, B, C)$ definit prin ecuațiile de stare

$$\begin{cases} \dot{x} = Ax + Bu \\ y = Cx \end{cases} \quad (5.1)$$

unde $x \in \mathcal{R}^n$ este un vector cu n componente.

Considerăm transformarea de coordonate $\tilde{x} = Ux$, în care $U \in \mathcal{R}^{n \times n}$ este matricea ortogonală ce aduce matricea A la forma Schur reală $S \stackrel{\text{def}}{=} \tilde{A} = UAU^T$. Punem în evidență partițiile conforme

$$U = \left[\begin{array}{c} U_a \\ U_b \end{array} \right] \begin{matrix} \} n_a \\ \} n_b \end{matrix}, \quad S \stackrel{\text{def}}{=} \tilde{A} = \left[\begin{array}{cc} \overbrace{A_a}^{n_a} & \overbrace{A_{ab}}^{n_b} \\ 0 & A_b \end{array} \right] \begin{matrix} \} n_a \\ \} n_b \end{matrix} \quad (5.2)$$

unde dimensiunile blocurilor satisfac relația $n_a + n_b = n$. Se constată ușor că un vector \tilde{x} de forma

$$\tilde{x} = \left[\begin{array}{c} x_a \\ 0 \end{array} \right] \begin{matrix} \} n_a \\ \} n_b \end{matrix}, \quad (5.3)$$

unde x_a este arbitrar, se transformă prin \tilde{A} într-un vector de aceeași formă, mai precis avem

$$\tilde{A}\tilde{x} = \left[\begin{array}{c} A_a x_a \\ 0 \end{array} \right] \begin{matrix} \} n_a \\ \} n_b \end{matrix},$$

Acest fapt arată că

(i) subspațiul $\mathcal{X}_a \subset \mathcal{R}^n$ al vectorilor de forma (5.3) este invariant în raport cu \tilde{A} și restricția lui \tilde{A} la acest subspațiu coincide cu blocul stânga - sus A_a al lui \tilde{A} .

În coordonatele inițiale $x = U^T \tilde{x}$ obținem

$$x = \left[\begin{array}{cc} U_a^T & U_b^T \end{array} \right] \left[\begin{array}{c} x_a \\ 0 \end{array} \right] = U_a^T x_a,$$

deci

(j) coloanele matricei U_a^T formează o bază ortogonală a subspațiului \mathcal{X}_a .

Proprietățile (i) și (j) de mai sus sunt esențiale pentru orice procedură de calcul a subspațiilor invariante \mathcal{X}_a , definite prin condiții spectrale impuse restricției A_a corespunzătoare.

De exemplu, dacă se cere ca evoluția sistemului S pe subspațiul \mathcal{X}_a să fie stabilă atunci matricea A_a trebuie să fie stabilă, deci $Re\lambda(A_a) < 0$. *Cel mai mare* subspațiu A -invariant \mathcal{X}_a cu această proprietate se notează cu $\mathcal{X}^-(A)$ și corespunde situației în care matricea A_a din (5.2) conține *toate* valorile proprii ale lui A cu parte reală negativă.

Această situație se realizează ușor dacă presupunem că matricea S din (5.2) reprezintă *forma Schur reală ordonată* a lui A , în care ordonarea valorilor proprii pe diagonala lui S se face *în sensul crescător al părților reale*, deci valorile proprii λ_i ale lui A apar pe diagonala lui S în ordinea

$$Re\lambda_1 \leq Re\lambda_2 \leq \dots \leq Re\lambda_n. \quad (5.4)$$

În acest caz submatricea A_a satisface condiția de stabilitate și are dimensiunea maximă dacă n_a se consideră egal cu numărul n^- al valorilor proprii stabile ale lui A .

Procedura de determinare a subspațiului A -invariant $\mathcal{X}^-(A)$ asociat valorilor proprii stabile ale lui A poate fi rezumată astfel.

Algoritmul 5.2 (Se determină o bază ortogonală a subspațiului A -invariant stabil $\mathcal{X}^-(A)$).

1. Se calculează forma Schur reală $S = UAU^T$ a lui A , utilizând algoritmul QR cu acumularea transformărilor U .
2. Se ordonează S , în acord cu criteriul (5.4), utilizând transformări de asemănare ortogonală $S \leftarrow PSP^T$ și efectuând acumularea acestora conform schemei $U \leftarrow PU$.
3. Se determină numărul $n_a \stackrel{\text{def}}{=} n_-$ al valorilor proprii stabile inspectând diagonala lui S . O bază ortogonală a subspațiului $\mathcal{X}^-(A)$ este formată din primele n_- coloane ale matricei U^T .

Comentarii. Algoritmul 5.2 utilizează exclusiv transformări ortogonale și, ca atare, este numeric stabil. În cazul discret se obțin concluzii similare utilizând *forma Schur reală ordonată* a lui A , în care ordonarea valorilor proprii pe diagonala lui S se face *în sensul crescător al modulelor*, deci astfel încât

$$|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_n|. \quad (5.5)$$

◇

Revenim acum la sistemul liniar $S = (A, B, C)$ adus la forma (5.2) cu transformarea ortogonală $\tilde{x} = Ux$. Considerăm transformarea suplimentară $z = V_0^{-1}\tilde{x}$ sau, echivalent, $\tilde{x} = V_0z$, în care matricea $V_0 \in \mathcal{R}^{n \times n}$ are structura

$$V_0 = \begin{bmatrix} \overbrace{I}^{n_a} & \overbrace{Y}^{n_b} \\ 0 & I \end{bmatrix} \begin{matrix} \} n_a \\ \} n_b \end{matrix} \quad (5.6)$$

Se constată imediat că avem

$$\tilde{A}V_0 = V_0A_0, \quad \text{unde } A_0 \stackrel{\text{def}}{=} \begin{bmatrix} A_a & 0 \\ 0 & A_b \end{bmatrix}, \quad (5.7)$$

dacă și numai dacă matricea Y satisface ecuația matriceală Sylvester

$$A_aY - YA_b = -A_{ab}. \quad (5.8)$$

Pe de altă parte, se știe că ecuația (5.8) are o soluție unică Y dacă și numai dacă matricele A_a și A_b nu au valori proprii comune, pe scurt

$$\lambda(A_a) \cap \lambda(A_b) = \emptyset. \quad (5.9)$$

Deoarece condiția (5.9) este evident satisfăcută în situația creată aplicând algoritmul 5.2, în care A_a și A_b colectează toate valorile proprii stabile și, respectiv, nestabile ale lui A , concludem că prin schimbarea

$$x \stackrel{\text{def}}{=} Vz = U^T\tilde{x} = U^TV_0z,$$

sistemul $S = (A, B, C)$ poate fi întotdeauna adus la forma

$$\begin{cases} \dot{z} = A_0z + B_0u \\ y = C_0z \end{cases} \quad (5.10)$$

unde

$$A_0 \stackrel{\text{def}}{=} V^{-1}AV = \begin{bmatrix} A_a & 0 \\ 0 & A_b \end{bmatrix}, \quad B_0 \stackrel{\text{def}}{=} V^{-1}B = \begin{bmatrix} B_a \\ B_b \end{bmatrix}, \quad (5.11)$$

$$C_0 \stackrel{\text{def}}{=} CV = [C_a \ C_b],$$

în care matricea A_a colectează toate valorile proprii ale matricei A cu $\text{Re}\lambda(A) < 0$, deci, este *stabilă*, iar matricea A_b colectează toate valorile proprii ale matricei A cu $\text{Re}\lambda(A) \geq 0$, deci este *total nestabilă*, pe scurt

$$\lambda(A_a) \subset \mathcal{C}^-, \quad \lambda(A_b) \subset \bar{\mathcal{C}}^+, \quad (5.12)$$

unde $\bar{\mathcal{C}}^+$ este semiplanul drept închis al planului complex \mathcal{C} .

Definiția 5.3 Forma bloc-diagonală (5.10), (5.11) se numește descompunerea spectrală a sistemului $S = (A, B, C)$ și corespunde reprezentării lui S prin conexiunea paralel a două sisteme, $S_a = (A_a, B_a, C_a)$ și $S_b = (A_b, B_b, C_b)$.

În acord cu (5.12), S_a este stabil de ordin $n_a = n^-$, unde n^- este numărul valorilor proprii stabile ale lui A iar S_b este total nestabil de ordin $n_b = \bar{n}^+$, unde \bar{n}^+ este numărul valorilor proprii nestabile (cu parte reală strict pozitivă sau nulă) ale lui A . În consecință, S_a se numește *partea stabilă* iar S_b se numește *partea total nestabilă* a sistemului S .

În particular, dacă matricea A este *dihotomică*, adică nu are valori proprii pe axa imaginară (cu $Re\lambda = 0$), atunci matricea A_b din (5.11) este *antistabilă*, adică $\lambda(-A_b) \subset \mathcal{C}^-$, iar S_b se numește *partea antistabilă* a sistemului dihotomic S .

Conform celor de mai sus, matricea V din (5.11) este

$$V = U^T V_0 = \begin{bmatrix} U_a^T & U_b^T \end{bmatrix} \begin{bmatrix} I & Y \\ 0 & I \end{bmatrix},$$

deci avem

$$V \stackrel{\text{def}}{=} \begin{bmatrix} V_a & V_b \end{bmatrix} = \begin{bmatrix} U_a^T & U_b^T + U_a^T Y \end{bmatrix}. \quad (5.13)$$

Punând

$$z = \begin{bmatrix} z_a \\ z_b \end{bmatrix}$$

obținem

$$x = Vz = V_a z_a + V_b z_b,$$

unde coloanele matricei $V_a = U_a^T$ constituie o bază ortogonală a *subspațiului stabil* $\mathcal{X}^-(A) = \mathcal{X}_a$, corespunzător vectorilor z cu $z_b = 0$ (vezi (5.3)).

În mod analog, coloanele matricei V_b constituie o bază (în general ne-ortogonală, deoarece $Y \neq 0$) a *subspațiului total nestabil* $\bar{\mathcal{X}}^+(A) = \mathcal{X}_b$, corespunzător vectorilor z cu $z_a = 0$. (Spre deosebire de (5.2), în (5.11) avem $A_{ab} = 0$, astfel încât acum și \mathcal{X}_b este \tilde{A} -invariant cu restricția A_b).

Procedura de determinare a descompunerii spectrale (5.10), (5.11) poate fi rezumată astfel.

Algoritmul 5.3 (Se construiește descompunerea spectrală a sistemului $S = (A, B, C)$ și se determină o bază ortogonală a subspațiului A -invariant stabil $\mathcal{X}^-(A)$).

1. Se aplică algoritmul 5.2.

2. Se face

$$B \leftarrow UB = \begin{bmatrix} B_a \\ B_b \end{bmatrix}, \quad C \leftarrow CU^T = [C_a \quad C_b].$$

3. Se rezolvă ecuația matriceală Sylvester (5.8), ținând seama că matricele A_a, A_b sunt deja în forma Schur reală (vezi (5.2)).

4. Se face

$$B_a \leftarrow B_a - YB_b, \quad C_b \leftarrow C_a Y + C_b.$$

Comentarii. Spre deosebire de algoritmul 5.2, algoritmul 5.3 utilizează transformări de asemănare neortogonală (vezi forma (5.13) a lui V_0), care în general ridică probleme de condiționare. În practică, numărul de condiție $\text{cond}(V_0)$ se poate evalua eficient utilizând estimatorul de condiție pentru matrice triunghiulare propus în LINPACK [III]. În ceea ce privește însăși condiționarea problemei rezolvării ecuației matriceale Sylvester (5.8), aceasta poate fi apreciată (destul de nesigur) evaluând apropierea dintre spectrele matricelor A_a și A_b , e.g. prin

$$\delta \stackrel{\text{def}}{=} \min \text{Re}\lambda(A_b) - \max \text{Re}\lambda(A_a).$$

(Relativ la această chestiune, vezi cap. 1). În cazul discret se obțin concluzii similare pe baza aceleiași ecuații (5.8). Dacă este necesar, aceasta poate fi scrisă sub forma specifică

$$A_a Y A_b^{-1} - Y = -A_{ab} A_b^{-1}, \quad (5.14)$$

căreia i se pot aplica rezultatele corespunzătoare din capitolul 1. \diamond

Opțional, algoritmul 5.3 poate fi completat pentru a furniza baza $V_b = U_b^T + U_a^T Y$ a subspațiului $\mathcal{X}^+(A)$ sau/și matricea de transformare (neortogonală) $T \in \mathcal{R}^{n \times n}$ a tripletului $S = (A, B, C)$ la forma (5.11), (5.12). Prin definiție avem

$$z = Tx, \quad (5.15)$$

unde, conform notațiilor de mai sus

$$T = V^{-1} = V_0^{-1} U = \begin{bmatrix} I & -Y \\ 0 & I \end{bmatrix} \begin{bmatrix} U_a \\ U_b \end{bmatrix} = [U_a - YU_b \quad U_b]. \quad (5.16)$$

Prin urmare, matricea T poate fi calculată pe loc în tabloul U furnizat de algoritmul 5.2, efectuând actualizarea

$$U_a \leftarrow U_a - YU_b.$$

Alternativ, pentru a evita distrugerea bazei ortogonale U_a a subspațiului $\mathcal{X}^-(A)$, se poate recurge la calculul matricii $V = T^{-1}$ din (5.13), corespunzătoare schimbării de coordonate $x = Vz$.

5.3 Controlabilitatea și observabilitatea sistemelor liniare

Controlabilitatea unui sistem liniar $S = (A, B, C)$ este o proprietate ce depinde numai de perechea (A, B) . Pe scurt, S este *controlabil* dacă orice tranziție de stare dorită poate fi realizată prin alegerea adecvată a funcției de intrare.

Pentru precizare, considerăm sistemul liniar discret

$$\begin{aligned} x(t+1) &= Ax(t) + Bu(t) \\ y(t) &= Cx(t) \end{aligned} \quad (5.17)$$

în care, ca de obicei, $A \in \mathcal{R}^{n \times n}$, $B \in \mathcal{R}^{n \times m}$ și $C \in \mathcal{R}^{l \times n}$ sunt matrice constante.

Presupunând $x(0) = 0$, din (5.17) rezultă succesiv

$$\begin{aligned} x(1) &= Bu(0), \\ x(2) &= ABu(0) + Bu(1), \end{aligned}$$

și în general

$$x(k) = A^{k-1}Bu(0) + \dots + ABu(k-2) + Bu(k-1). \quad (5.18)$$

Introducând notațiile

$$R_k = [B \quad AB \quad \dots \quad A^{k-1}B], \quad u_k = \begin{bmatrix} u(k-1) \\ u(k-2) \\ \vdots \\ u(0) \end{bmatrix}, \quad (5.19)$$

în care vectorul u_k conține k blocuri, fiecare cu m componente, iar R_k este o matrice bloc de forma $n \times mk$, relațiile (5.18) se scriu pe scurt

$$R_k u_k = x(k). \quad (5.20)$$

Considerând că în (5.20) pasul $k \geq 1$ este *fixat* iar $x(k) \in \mathcal{R}^n$ este un vector arbitrar care reprezintă starea dorită la pasul k , vom spune că perechea

discretă (A, B) este *controlabilă în k pași* dacă sistemul de ecuații liniare (5.20) are o soluție u_k , adică *există* un șir de intrări $u(0), u(1), \dots, u(k-1)$ care realizează tranziția de stare $0 \rightarrow x(k)$. După cum se știe, sistemul (5.20) este *compatibil oricare ar fi membrul drept* dacă și numai dacă matricea R_k este *epică*, adică

$$\text{rang} R_k = n. \quad (5.21)$$

Mai mult, în acest caz *soluția normală* u_k^* a sistemului (5.20) este unic determinată și are expresia

$$u_k^* = R_k^T (R_k R_k^T)^{-1} x(k), \quad (5.22)$$

deci avem

$$\|u_k^*\|^2 = \min_{R_k u = x(k)} \|u_k\|^2, \quad (5.23)$$

unde norma (euclidiană) a unui vector u_k de forma (5.19) este evident

$$\|u_k\|^2 = \sum_{t=0}^{k-1} \|u(t)\|^2. \quad (5.24)$$

Interpretând expresia (5.24) ca "efort" de comandă necesar pentru realizarea tranziției $0 \rightarrow x(k)$, soluția normală (5.22) se numește sugestiv comandă de normă sau efort minim.

În rezumat, perechea (A, B) este controlabilă în k pași dacă și numai dacă este satisfăcută condiția de rang (5.21), și în acest caz comanda de normă minimă are expresia (5.22).

În general, considerând mai sus un număr de pași k arbitrar, obținem noțiunea de controlabilitate. Prin urmare, perechea discretă (A, B) este *controlabilă* dacă există $k \geq 1$ astfel încât oricare ar fi $x(k) \in \mathcal{R}^n$ există o soluție u_k a sistemului (5.20) sau, echivalent, este îndeplinită condiția de rang (5.21).

Mai mult, în acest caz avem întotdeauna $k \leq n$, deoarece, în virtutea teoremei Cayley-Hamilton, eventualele blocuri $A^{k-1}B$ cu $k > n$ din (5.19) sunt în mod necesar liniar dependente de precedentele, deci nu contribuie la satisfacerea condiției de rang (5.21). Pe scurt, un sistem liniar discret controlabil este întotdeauna controlabil în cel mult n pași.

Notând pentru simplitate $R_n = R$, deci

$$R = [B \quad AB \quad \dots \quad A^{n-1}B], \quad (5.25)$$

precum și $u_n = \mathbf{u}$, relațiile (5.20) pentru $k = n$ se scriu

$$R\mathbf{u} = x(n) \quad (5.26)$$

iar condiția de rang (5.21), echivalentă cu controlabilitatea perechii (A, B) , devine

$$\text{rang} R = n. \quad (5.27)$$

În cazul sistemelor cu timp continuu, considerații similare conduc în ultimă instanță tot la condiția (5.27). De aceea, în continuare vom adopta următoarea definiție unificatoare, pur algebrică.

Definiția 5.4 *Perechea (A, B) este controlabilă dacă și numai dacă este satisfăcută condiția de rang (5.27).*

Matricele R și R_k definite prin relațiile (5.25) și (5.19) se numesc matrice de controlabilitate și respectiv matrice de controlabilitate în k pași ale perechii (A, B) .

În cazul în care perechea (A, B) este controlabilă, cel mai mic întreg $k \geq 1$ pentru care condiția de rang (5.21) este satisfăcută se numește indice de controlabilitate al perechii (A, B) și se notează cu ν .

Observăm că, în cazul discret, ν reprezintă cel mai mic număr k de pași (deci timpul minim) în care este posibilă tranziția de stare $0 \rightarrow x(k)$, oricare ar fi starea dorită $x(k) \in \mathcal{R}^n$. De asemenea, dacă perechea (A, B) are o singură intrare, deci $m = 1$, atunci R este pătrată și nesingulară iar $\nu = n$. Din contră, dacă $m > 1$ atunci în general avem $\nu < n$.

Observabilitatea unui sistem liniar $S = (A, B, C)$ este proprietatea duală controlabilității și, în consecință, depinde numai de perechea (C, A) . Pe scurt, S este *observabil* dacă starea inițială este unic determinată cunoscând funcțiile de intrare și ieșire, obținute e.g. prin măsurători efectuate la terminalele lui S .

Pentru precizare, considerăm din nou sistemul liniar discret descris prin ecuațiile de stare (5.17). Presupunând că $u(t) = 0, t \geq 0$, din (5.17) rezultă $x(t) = A^t x(0)$, astfel încât șirul de ieșiri corespunzător stării inițiale $x(0)$ este

$$y(t) = CA^t x(0), \quad t = 0 : k - 1. \quad (5.28)$$

Introducând notațiile

$$Q_k = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{k-1} \end{bmatrix}, \quad y_k = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(k-1) \end{bmatrix}, \quad (5.29)$$

în care vectorul y_k conține k blocuri, fiecare cu l componente, iar Q_k este o matrice bloc de forma $kl \times n$, relațiile (5.28) se scriu pe scurt

$$Q_k x(0) = y_k. \quad (5.30)$$

Considerând că în (5.30) pasul $k \geq 1$ este *fixat*, vom spune că perechea discretă (C, A) este *observabilă în k pași* dacă sistemul de ecuații liniare (5.30) are cel mult o soluție $x(0)$, adică starea inițială $x(0)$ este *unic* determinată de șirul de ieșiri $y(0), y(1), \dots, y(k-1)$. După cum se știe, sistemul (5.30) este *unic determinat* dacă și numai dacă matricea Q_k este *monică*, adică

$$\text{rang } Q_k = n. \quad (5.31)$$

Mai mult, în acest caz *pseudosoluția* $x^*(0)$ în sensul CMMP a sistemului (5.30) este unic determinată și are expresia

$$x^*(0) = (Q_k^T Q_k)^{-1} Q_k^T y_k, \quad (5.32)$$

deci avem

$$\|y_k - Q_k x^*(0)\|^2 = \min_{x(0) \in \mathcal{R}^n} \|y_k - Q_k x(0)\|^2, \quad (5.33)$$

unde evident

$$\|y_k - Q_k x(0)\|^2 = \sum_{t=0}^{k-1} \|y_k - Cx(k)\|^2. \quad (5.34)$$

Interpretând reziduurile $r(k) = y(k) - Cx(k)$ ca erori sau "zgomote" de măsură care afectează ieșirile măsurate $y(k)$, pseudosoluția $x^*(0)$ se numește sugestiv *estimare optimală* (în sensul CMMP) a stării inițiale $x(0)$.

În rezumat, perechea (C, A) este observabilă în k pași dacă și numai dacă este îndeplinită condiția de rang (5.31) și în acest caz estimarea optimală a stării inițiale $x(0)$ are expresia (5.32).

În general, considerând mai sus un număr de pași k arbitrar, obținem noțiunea de observabilitate. Prin urmare, perechea discretă (C, A) este *observabilă* dacă există $k \geq 1$ astfel încât este îndeplinită condiția de rang (5.31). Mai mult, în acest caz, la fel ca în cazul controlabilității, putem considera $k \leq n$.

Notând pentru simplitate $Q_n = Q$, deci

$$Q = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix}, \quad (5.35)$$

precum și $y_n = \mathbf{y}$, relațiile (5.30) pentru $k = n$ se scriu

$$Qx(0) = \mathbf{y} \quad (5.36)$$

iar condiția de rang (5.31), echivalentă cu observabilitatea perechii (C, A) , devine

$$\text{rang } Q = n. \quad (5.37)$$

În cazul sistemelor cu timp continuu, considerații similare conduc în ultimă instanță tot la condiția (5.37). În consecință, mai departe vom adopta următoarea definiție unificatoare.

Definiția 5.5 *Perechea (C, A) este observabilă dacă și numai dacă este satisfăcută condiția de rang (5.37).*

Matriceile Q și Q_k definite prin relațiile (5.35) și (5.29) se numesc matrice de observabilitate și respectiv matrice de observabilitate în k pași ale perechii (C, A) .

În cazul în care perechea (C, A) este observabilă, cel mai mic întreg $k \geq 1$ pentru care condiția de rang (5.31) este satisfăcută se numește indice de observabilitate al perechii (C, A) .

Examinând în paralel relațiile (5.18)-(5.27) și (5.28)-(5.37), concludem că noțiunile de controlabilitate și observabilitate sunt *duale*, în sensul că ele se corespund prin corespondența

$$A \leftarrow A^T, \quad C \leftarrow B^T,$$

în particular avem

$$Q(C, A) = R^T(A^T, C^T).$$

În consecință, mai departe vom prezenta procedurile de calcul asociate controlabilității, rămânând ca dualele lor să fie formulate "prin dualitate" și discutate în detaliu de către cititorul interesat.

5.4 Teste elementare de controlabilitate

În general procedurile de testare a controlabilității unei perechii date (A, B) sunt de două tipuri.

Primul tip de proceduri, numite ad-hoc *elementare*, testează controlabilitatea perechii (A, B) utilizând direct definiția 5.3 sau proprietăți echivalente cu aceasta. Aceste proceduri sunt în general neeficiente și pot fi aplicate cu succes numai pentru perechi de dimensiuni mici.

Al doilea tip de proceduri utilizează transformări de asemănare pentru a aduce perechea inițială (A, B) la o formă "simplă", pe care proprietatea de controlabilitate să poată fi testată direct ("prin inspecție"). Aceste proceduri, numite în continuare *de transformare*, sunt cele mai utilizate în aplicații, cu atât mai mult cu cât ele facilitează și rezolvarea altor probleme de calcul conexe vizând, de exemplu, descompunerea controlabilă, testarea stabilizabilității, alocarea polilor etc. Cea mai eficientă procedură de transformare, care constă în aducerea perechii (A, B) la forma (superior) Hessenberg prin transformări de asemănare ortogonale, va fi prezentată în paragraful următor.

Observația 5.2 Proprietatea de controlabilitate este o proprietate binară, exprimată prin condiția de rang (5.27). De aceea, pe de o parte nu pot exista proceduri numerice absolut sigure de testare a controlabilității iar, pe de altă parte, pentru reducerea riscurilor implicate de deciziile de rang, se recomandă utilizarea transformărilor ortogonale și, în special, a descompunerii valorilor singulare (**DVS**). \diamond

Din clasa procedurilor elementare face parte în primul rând procedura "directă", bazată pe construcția matricei de controlabilitate R și testarea condiției de rang (5.27).

Construcția matricei R se face recurent, ținând seama că prin definiție avem $R = R_n$, unde matricele șirului R_k , $k \geq 1$ au expresia (5.19) iar $R_1 = B$. Pentru a obține relația de recurență care leagă doi termeni succesivi R_k și R_{k+1} , observăm că avem

$$R_{k+1} \stackrel{\text{def}}{=} \begin{bmatrix} B & : & AB & A^2B & \dots & A^k B \end{bmatrix} = \begin{bmatrix} B & : & AR_k \end{bmatrix}, \quad (5.38)$$

respectiv

$$R_{k+1} \stackrel{\text{def}}{=} \begin{bmatrix} B & AB & \dots & A^{k-1}B & : & A^k B \end{bmatrix} = \begin{bmatrix} R_k & : & B_{k+1} \end{bmatrix}, \quad (5.39)$$

unde submatricea $B_{k+1} = A^k B$, formată din ultimele m coloane ale lui R_{k+1} , rezultă în funcție de submatricea corespunzătoare B_k a lui R_k prin relația evidentă

$$B_{k+1} = AB_k, \quad k \geq 1 \quad (5.40)$$

iar $B_1 = B$.

Se constată imediat că procedura de construcție a matricei R bazată pe relația (5.38) este relativ neeconomică, întrucât necesită efectuarea produselor AR_k , unde matricea R_k are mk coloane, $k \geq 1$, în timp ce procedura bazată pe relația (5.39) necesită efectuarea produselor AB_k , unde la fiecare

pas $k \geq 1$ matricea B_k are numai m coloane. De asemenea, este clar că la fiecare pas matricea B_k este disponibilă în R_k , mai precis avem

$$B_k = R_k(\cdot, (k-1)m+1 : km) \quad (5.41)$$

astfel încât nu este necesară nici alocarea unui tablou suplimentar de lucru pentru memorarea matricei curente B_k , nici calculul pe loc în B , conform schemei $B \leftarrow AB$, al termenilor șirului B_k , ceea ce ar implica distrugerea matricei inițiale B .

În rezumat, procedura de testare a controlabilității perechii (A, B) bazată pe utilizarea relațiilor (5.39) – (5.41) poate fi formulată astfel.

Algoritmul 5.4 (Se testează controlabilitatea perechii (A, B) construind matricea de controlabilitate R).

1. $R(\cdot, 1 : m) \leftarrow B_1 = B$
2. Pentru $k = 2 : n$
 1. $R(\cdot, (k-1)m+1 : km) \leftarrow B_k = AB_{k-1}$
3. Se calculează $r = \text{rang} R$ utilizând algoritmul DVS.
4. Dacă $r = n$ atunci
 1. **Tipărește** "Perechea (A, B) este controlabilă."

În general, în ciuda simplității sale, algoritmul 5.4 *nu este recomandabil* întrucât matricea R este de mari dimensiuni iar calculul coloanelor sale la pasul 2.1 implică riscuri la efectuarea testului de rang de la pasul 3.

De exemplu, dacă $m = 1$, deci matricea $B = b$ se reduce la o singură coloană, atunci, conform relației (5.40), coloanele $b, Ab, \dots, A^k b, \dots$ ale matricei (patrate) R coincid cu vectorii obținuți aplicând *metoda puterii* vectorului inițial $b_1 = b$. Prin urmare, dacă $n \gg 1$ iar matricea A are o valoare proprie dominantă, deci

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|,$$

atunci ultimele coloane ale lui R sunt practic aliniată cu vectorul propriu x_1 al lui A asociat cu λ_1 . Aceasta înseamnă că, chiar dacă perechea (A, b) este controlabilă, matricea de controlabilitate R este numeric aproape singulară, adică

$$\text{cond}(R) = \frac{\sigma_1}{\sigma_n} \gg 1,$$

unde $\text{cond}(R)$ este numărul de condiționare la inversare al matricei R (în raport cu norma spectrală) iar σ_i , $i = 1 : n$ sunt valorile singulare ale lui R , ordonate astfel încât $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$. Fenomenul poartă numele de

necontrolabilitate numerică a perechii (A, b) și în general nu poate fi evitat, indiferent de testul de rang utilizat la pasul 3 al algoritmului de mai sus.

Din acest punct de vedere, în cazul $m > 1$ se recomandă efectuarea testului de rang în bucla 2, conform schemei următoare, care, în plus, furnizează și indicele de controlabilitate ν al perechii (A, B) .

Algoritmul 5.5 (Se testează controlabilitatea perechii (A, B) construind matricea de controlabilitate R . Se presupune $m > 1$.)

1. $R(:, 1 : m) \leftarrow B_1 = B$
2. Pentru $k = 2 : n$
 1. Se calculează $r = \text{rang} R$ utilizând algoritmul DVS.
 2. Dacă $r = n$ atunci
 1. $\nu = k$
 2. **Tipărește** "Perechea (A, B) este controlabilă."
 3. **stop**
3. $R(:, (k-1)m + 1 : km) \leftarrow B_k = AB_{k-1}$

Subliniem însă că schema de mai sus este în orice caz neeficientă din punctul de vedere al efortului de calcul datorită apelării repetate a algoritmului **DVS** într-o formă brută, care nu exploatează rezultatele calculelor anterioare. Cel puțin aceeași siguranță a testului este asigurată de algoritmul 5.7 prezentat mai departe, cu un efort de calcul considerabil mai redus.

În afară de procedura "directă", bazată pe testarea condiției de rang (5.27), din grupul metodelor elementare fac parte o serie de proceduri derivate pe baza unor criterii echivalente cu aceasta. Cele mai importante sunt rezumate mai jos pentru a simplifica referirile ulterioare.

Propoziția 5.1 *Controlabilitatea perechii (A, B) este echivalentă cu oricare dintre condițiile*

1. Nu există nici un vector $h \neq 0$ astfel încât $A^T h = \lambda h, \lambda \in \mathcal{C}$ și $B^T h = 0$.
2. $\text{rang} \begin{bmatrix} A - \lambda I & B \end{bmatrix} = n, \quad \forall \lambda \in \lambda(A)$.
3. Oricare ar fi vectorul $g \in \mathcal{R}^m$ cu proprietatea $Bg \neq 0$, există o matrice $F \in \mathcal{R}^{m \times n}$ astfel încât perechea cu o singură intrare $(A - BF, Bg)$ este controlabilă.

În general, testele de controlabilitate derivate prin aplicarea directă a criteriilor de mai sus nu sunt recomandabile din punct de vedere numeric. Criteriul 1 necesită calculul valorilor și vectorilor proprii $x_i \in \mathcal{C}^n$ ai matricei A (ceea ce implică aplicarea algoritmului **QR** cu acumularea transformărilor de asemănare) precum și testarea condițiilor $B^T x_i \neq 0$, $i = 1 : n$. Criteriul 2 necesită numai calculul valorilor proprii, în schimb presupune aplicarea repetată a algoritmului **DVS** pentru testarea condițiilor de rang. În sfârșit, criteriul 3 este interesant deoarece reduce cazul general $m > 1$ la cazul aparent mai simplu $m = 1$. Deoarece controlabilitatea este o proprietate generică, matricea F precum și vectorul g se pot alege aleator. Această idee poate fi deseori exploatată cu succes în legătură cu algoritmul 5.6, prezentat mai departe, dar în legătură cu aceasta trebuie să facem următoare observație de principiu.

Observația 5.3 De obicei în aplicații testele de controlabilitate nu se aplică *per se*, ci fac parte din anumite proceduri de sinteză mai generale, în care constituie condiții de calcul sau/și de validare a rezultatelor. (Un exemplu în acest sens îl constituie procedurile de alocare a polilor, prezentate în capitolul 6). De aceea este important ca secvențele de calcul să se înlănțuie în mod natural, transformările datelor necesare la un moment dat fiind utilizate ca atare în fazele ulterioare ale procedurii considerate. \diamond

Din perspectiva observației 5.3, criteriul 3 devine mai puțin semnificativ, întrucât presupune distrugerea perechii date (A, B) și înlocuirea ei cu o pereche $(A - BF, Bg)$ transformată aleator precum și renunțarea nemotivată la libertățile suplimentare de sinteză existente în cazul $m > 1$ al perechilor cu mai multe intrări față de cazul $m = 1$.

În concluzie, testele "elementare" de controlabilitate sunt în general nesatisfăcătoare din punctul de vedere al eficienței și siguranței numerice. Aplicarea lor practică este limitată la sisteme de mici dimensiuni, pentru care, eventual, pot constitui mijloace de calcul "cu hârtia și creionul".

Teste elementare de observabilitate

În general, dacă o procedură **proc** calculează un rezultat $R = \mathbf{proc}(A, B)$ privind controlabilitatea perechii (A, B) atunci rezultatul dual Q , privind observabilitatea perechii (C, A) , se obține cu secvența

1. $R_d = \mathbf{proc}(A^T, C^T)$
2. $Q = R_d^T$

De exemplu, dacă **proc** (A, B) este algoritmul 5.4, care produce matricea de controlabilitate R a perechii (A, B) , atunci secvența de mai sus produce

matricea de observabilitate Q a perechii (C, A) .

5.5 Forma Hessenberg controlabilă

Considerăm sistemul liniar $S = (A, B, C)$ și fie $\tilde{S} = (\tilde{A}, \tilde{B}, \tilde{C})$ un sistem asemenea cu S , având matricele

$$\tilde{A} = TAT^{-1}, \quad \tilde{B} = TB, \quad \tilde{C} = CT^{-1}, \quad (5.42)$$

unde $T \in \mathcal{R}^{n \times n}$ este o matrice nesingulară arbitrară. Avem

$$\tilde{A}\tilde{B} = TAB, \quad \tilde{A}^2\tilde{B} = \tilde{A}(\tilde{A}\tilde{B}) = TA^2B, \quad \dots$$

astfel încât între matricele de controlabilitate și observabilitate atașate sistemelor S și \tilde{S} există relațiile

$$\tilde{R} = TR, \quad \tilde{Q} = QT^{-1}. \quad (5.43)$$

Deoarece multiplicarea cu o matrice nesingulară nu modifică rangul, din relațiile (5.43) rezultă că matricele R și \tilde{R} au același rang, deci perechile (A, B) și (\tilde{A}, \tilde{B}) sunt simultan controlabile sau necontrolabile. În mod similar (sau prin dualitate), perechile (C, A) și (\tilde{C}, \tilde{A}) sunt simultan observabile sau neobservabile. Pe scurt, proprietățile de controlabilitate și observabilitate, precum și rangurile matricelor corespunzătoare de controlabilitate și observabilitate, sunt invariante în raport cu transformările de asemănare (5.42).

Dacă în (5.42) presupunem că matricea $T = U$ este ortogonală, deci sistemele S și \tilde{S} sunt *ortogonal asemenea*

$$\tilde{A} = UAU^T, \quad \tilde{B} = UB, \quad \tilde{C} = CU^T, \quad (5.44)$$

atunci relațiile (5.43) devin

$$\tilde{R} = UR, \quad \tilde{Q} = QU^T. \quad (5.45)$$

Prin urmare, nu numai proprietățile (binare!) de controlabilitate și observabilitate ci și valorile singulare $\sigma_i(R)$ și $\sigma_i(Q)$, $i \in 1 : n$ ale matricelor de controlabilitate și observabilitate sunt invariante ai sistemului (A, B, C) în raport cu transformările de asemănare ortogonală (5.44).

Observația 5.4 Invariantii ortogonali $\sigma_i(R)$ și $\sigma_i(Q)$, $i \in 1 : n$ au o deosebită importanță în analiza sistemelor liniare, întrucât reprezintă o

măsură cantitativă a "gradelor" de controlabilitate și observabilitate ale sistemului $S = (A, B, C)$. În special, cea mai mică valoare singulară $\sigma_n(R) \geq 0$ a lui R reprezintă distanța (măsurată în norma spectrală) dintre perechea dată (A, B) și mulțimea perechilor necontrolabile de aceeași formă. \diamond

În cele ce urmează vom descrie procedurile de aducere a perechii (A, B) la o formă simplă utilizând transformările de asemănare ortogonală (5.44), cu scopul identificării invariantilor (ortogonali ai) acestei perechi și, în particular, al testării controlabilității acesteia.

În consecință, vom considera o pereche (A, B) cu m intrări, nu neapărat controlabilă, și pentru claritate vom discuta succesiv cazurile $m = 1$ și $m > 1$. Peste tot vom presupune $B \neq 0$ și vom nota

$$r \stackrel{\text{def}}{=} \text{rang} R, \quad \bar{r} \stackrel{\text{def}}{=} n - r, \quad (5.46)$$

unde evident r satisface inegalitatea $1 \leq r \leq n$.

Cazul $m = 1$

În cazul unei perechi (A, b) de ordin n cu o singură intrare, forma "simplă", rezultată prin aplicarea transformării (5.44), se numește *forma superior Hessenberg* și este definită prin

$$\tilde{A} = UAU^T = \left[\begin{array}{cc} \overbrace{A_R}^r & \overbrace{A_{R\bar{R}}}^{\bar{r}} \\ 0 & A_{\bar{R}} \end{array} \right] \left. \vphantom{\begin{array}{c} r \\ \bar{r} \end{array}} \right\} \begin{array}{l} r \\ \bar{r} \end{array}, \quad \tilde{b} = Ub = \left[\begin{array}{c} b_R \\ 0 \end{array} \right] \left. \vphantom{b_R} \right\} \bar{r} \quad (5.47)$$

în care perechea (A_R, b_R) de ordin r se află în forma superior Hessenberg ireductibilă

$$A_R = \left[\begin{array}{cccccc} x & x & \dots & x & x \\ h_2 & x & \dots & x & x \\ & \ddots & \ddots & \vdots & \vdots \\ & & \ddots & x & \vdots \\ & & & h_r & x \end{array} \right], \quad b_R = \left[\begin{array}{c} h_1 \\ 0 \\ \vdots \\ 0 \end{array} \right], \quad (5.48)$$

unde

$$h_i \neq 0, \quad i = 1 : r \quad (5.49)$$

iar x denotă elemente a căror valoare numerică nu are importanță în context.

Dacă perechea inițială (A, b) este *controlabilă* atunci în (5.47) avem $r = n$, astfel încât blocurile $A_{\bar{R}}$ și $A_{R\bar{R}}$ dispar iar perechea $(\tilde{A}, \tilde{b}) = (A_R, b_R)$ are structura ireductibilă (5.48) cu $r = n$, adică

$$\tilde{A} = UAU^T = \begin{bmatrix} x & x & \dots & x & x \\ h_2 & x & \dots & x & x \\ & \ddots & \ddots & \vdots & \vdots \\ & & \ddots & x & \vdots \\ & & & h_n & x \end{bmatrix}, \quad \tilde{b} = Ub = \begin{bmatrix} h_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad (5.50)$$

unde

$$h_i \neq 0, \quad i = 1 : n. \quad (5.51)$$

De aceea, structura ireductibilă din (5.50) se numește *forma superior Hessenberg controlabilă* a perechii (A, b) .

Pe de altă parte, dacă perechea (\tilde{A}, \tilde{b}) are forma (5.50), nu neapărat ireductibilă, iar $k \geq 2$ este primul întreg pentru care $h_k = 0$, atunci această pereche admite evident o descompunere (deflație) de tip (5.47), în care $r = k - 1$. În acest sens, mai departe vom spune că structura (5.50), în general reductibilă, constituie *forma superior Hessenberg* (completă) a perechii (A, b) .

Afirmația următoare relativă la (5.50) se demonstrează ușor prin calcul direct.

Propoziția 5.2 *Matricea de controlabilitate \tilde{R} a unei perechi (\tilde{A}, \tilde{b}) de forma (5.50) este superior triunghiulară, mai precis avem*

$$\tilde{R} \stackrel{\text{def}}{=} [\tilde{b} \quad \tilde{A}\tilde{b} \quad \dots \quad \tilde{A}^{n-1}\tilde{b}] = \begin{bmatrix} \tilde{r}_{11} & x & \dots & x \\ & \tilde{r}_{22} & \dots & x \\ & & \ddots & \vdots \\ & & & \tilde{r}_{nn} \end{bmatrix}, \quad (5.52)$$

în care elementele diagonale ale lui \tilde{R} sunt

$$\tilde{r}_{kk} = \prod_{i=1}^k h_i, \quad k = 1 : n. \quad (5.53)$$

Prin urmare, perechea inițială (A, b) este controlabilă, deci matricea $R = U\tilde{R}$ este nesingulară, dacă și numai dacă sunt satisfăcute condițiile (5.51).

Pe o cale încă mai simplă, aceeași concluzie rezultă utilizând criteriul de controlabilitate furnizat de punctul 2 al propoziției 5.1. Într-adevăr, avem

$$[\tilde{b} \quad \tilde{A} - \lambda I] = \begin{bmatrix} h_1 & x & \dots & x & x \\ & h_2 & \dots & x & x \\ & & \ddots & \vdots & \vdots \\ & & & h_n & x \end{bmatrix}, \quad (5.54)$$

deci (5.51) implică evident $\text{rang} [\tilde{b} \quad \tilde{A} - \lambda I] = n$, $\forall \lambda \in \mathcal{C}$. Demonstrația implicației inverse este propusă cititorului ca exercițiu (de atenție!).

La fel de simplu se demonstrează afirmația următoare relativă la (5.47).

Propoziția 5.3 *Matricea de controlabilitate \tilde{R} a unei perechi (\tilde{A}, \tilde{b}) de forma (5.47) are structură superior trapezoidală*

$$\tilde{R} = \begin{bmatrix} R_R & x \\ 0 & 0 \end{bmatrix}, \quad (5.55)$$

în care $R_R \in \mathcal{R}^{r \times r}$ este matricea de controlabilitate a perechii (A_R, b_R) .

Prin urmare, ordinul r al perechii (A_R, b_R) coincide cu rangul matricei de controlabilitate $R = U\tilde{R}$ al perechii (A, b) dacă și numai dacă sunt satisfăcute condițiile (5.49), deci perechea (A_R, b_R) din (5.47) este controlabilă.

În rezumat, sunt posibile două cazuri:

1. Perechea (A, b) este controlabilă, deci $r \stackrel{\text{def}}{=} \text{rang} R = n$.

În acest caz perechea (A, b) este ortogonal asemenea cu o pereche (\tilde{A}, \tilde{b}) în forma superior Hessenberg controlabilă (5.50), (5.51).

2. Perechea (A, b) nu este controlabilă, deci $r \stackrel{\text{def}}{=} \text{rang} R < n$.

În acest caz perechea (A, b) este ortogonal asemenea cu o pereche (\tilde{A}, \tilde{b}) în forma bloc superior triunghiulară (5.47), în care perechea (A_R, b_R) este controlabilă de ordin r , deci are structura (5.48), (5.49). (În particular, (5.47) poate coincide cu forma superior Hessenberg completă (5.50), în care r este cel mai mic întreg ≥ 1 astfel încât $h_{r+1} = 0$).

Aducerea unei perechi (A, b) cu o singură intrare la forma superior Hessenberg completă (5.50) se face printr-o procedură directă, a cărei schemă de calcul este următoarea.

1. Se determină un reflector $U_1 \in \mathcal{R}^{n \times n}$ astfel încât ultimele $n - 1$ componente ale vectorului transformat $b \leftarrow U_1 b$ să fie nule.
2. Se aplică U_1 lui A , deci se calculează $A \leftarrow U_1 A U_1$.

3. Se aduce matricea A la forma superior Hessenberg $A \leftarrow \tilde{A} = \tilde{U}A\tilde{U}^T$, utilizând algoritmul **HQ**, [VI,IX].

Pasul 1 este elementar. Reflectorul U_1 este definit prin

$$U_1 = I_n - \frac{u_1 u_1^T}{\beta_1}$$

în care componentele u_{i1} , $i = 1 : n$ ale vectorului $u_1 \in \mathcal{R}^n$ precum și scalarul β_1 se determină cu formulele cunoscute

Procedura $U = \mathbf{H0}(b)$

1. $\sigma = \text{sgn}(b_1)(\sum_{i=1}^n b_i^2)^{1/2}$
2. $u_{11} = b_1 + \sigma$; $u_{i1} = b_i$, $i = 2 : n$
3. $\beta_1 = \sigma u_{11}$
4. $b_1 \leftarrow h_1 = -\sigma$; $b_i = 0$, $i = 2 : n$

iar ultima instrucțiune efectuează transformarea $b \leftarrow U_1 b$.

Pasul 2 se desfășoară în două faze. În prima fază se calculează $A \leftarrow U_1 A$, aplicând U_1 matricei A partiționate pe coloane. Obținem

Procedura $A = \mathbf{H1}(U, A)$

1. Pentru $j = 1 : n$
 1. $\tau = (\sum_{i=1}^n u_{i1} a_{ij}) / \beta_1$
 2. $a_{ij} \leftarrow a_{ij} - \tau u_{i1}$, $i = 1 : n$.

În a doua fază se calculează $A \leftarrow AU_1$ aplicând U_1 matricei A partiționate pe linii, unde U_1 e simetric deci $AU_1 = (U_1 A^T)^T$. Obținem

Procedura $A = \mathbf{H2}(A, U)$

1. Pentru $i = 1 : n$
 1. $\tau = (\sum_{j=1}^n u_{i1} a_{ji}) / \beta_1$
 2. $a_{ij} \leftarrow a_{ij} - \tau u_{j1}$, $j = 1 : n$.

Pasul 3 constă din $n - 2$ etape. La fiecare etapă se efectuează calcule similare cu cele descrise la pașii 1 și 2 de mai sus, operând efectiv asupra unor perechi "restante" (\bar{A}, \bar{b}) de ordin din ce în ce mai mic. Pentru a simplifica referirile ulterioare și, în special, redactarea algoritmului general, corespunzător cazului $m > 1$, mai jos vom trece în revistă principalele momente ale calculului.

Presupunem că înainte de etapa k , $k = 2 : n-1$, perechea (A, b) (parțial transformată în etapele anterioare) are structura

$$A = \left[\begin{array}{ccc|ccc} x & \dots & x & x & x & \dots & x \\ h_2 & \dots & x & x & x & \dots & x \\ & \ddots & \vdots & \vdots & \vdots & & \vdots \\ & & h_{k-1} & x & x & \dots & x \\ \hline & & & a_{k,k-1} & a_{kk} & \dots & a_{kn} \\ & & & \vdots & \vdots & & \vdots \\ & & & a_{n,k-1} & a_{nk} & \dots & a_{nn} \end{array} \right] = \left[\begin{array}{c|c} \tilde{A}_k & X_k \\ \hline 0 & \tilde{b}_k \mid \tilde{A}_k \end{array} \right], \quad (5.56)$$

$$b = \left[\begin{array}{c} h_1 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \end{array} \right] = \left[\begin{array}{c} \tilde{b}_k \\ 0 \end{array} \right]$$

în care perechea $(\tilde{A}_k, \tilde{b}_k)$ este deja în formă superior Hessenberg completă de ordin $k-1$.

La etapa k se efectuează următoarele calcule:

- 3.1. Se determină un reflector $U_k \in \mathcal{R}^{n \times n}$ astfel încât ultimele $n-k$ elemente din coloana $k-1$ a matricei A transformate să fie nule.
- 3.2. Se aplică U_k lui A , deci se calculează $A \leftarrow U_k A U_k$.

La pasul 3.1 reflectorul U_k este definit prin

$$U_k = I_n - \frac{u_k u_k^T}{\beta_k},$$

în care vectorul $u_k \in \mathcal{R}^n$ are primele $k-1$ componente nule. În virtutea acestui fapt, U_k are structura bloc

$$U_k = \left[\begin{array}{cc} I_{k-1} & 0 \\ 0 & \tilde{U}_k \end{array} \right], \quad (5.57)$$

iar datorită acesteia premultiplicarea cu U_k nu modifică primele $k-1$ linii iar postmultiplicarea cu U_k nu modifică primele $k-1$ coloane ale tablourilor

transformate. Prin urmare, la pasul 3.2 obținem

$$U_k A = \left[\begin{array}{c|c} \tilde{A}_k & X_k \\ \hline 0 & \bar{U}_k \bar{b}_k \end{array} \right] \stackrel{\text{def}}{=} \left[\begin{array}{c|c} \tilde{A}_k & \\ \hline 0 & \bar{U}_k \bar{b}_k \end{array} \right] Y_k \quad (5.58)$$

precum și

$$(U_k A) U_k = \left[\begin{array}{c|c} \tilde{A}_k & \\ \hline 0 & \bar{U}_k \bar{b}_k \end{array} \right] Y_k \bar{U}_k \stackrel{\text{def}}{=} \left[\begin{array}{c|c} \tilde{A}_{k+1} & X_{k+1} \\ \hline 0 & \bar{b}_{k+1} \end{array} \right] \bar{A}_{k+1} \quad (5.59)$$

$$U_k b = b = \left[\begin{array}{c} \tilde{b}_{k+1} \\ \hline 0 \end{array} \right]$$

unde, datorită alegerii lui U_k de la pasul 3.1, perechea $(\tilde{A}_{k+1}, \tilde{b}_{k+1})$ rezultă în forma superior Hessenberg completă de ordin k .

Astfel procedura a progresat, adică dimensiunea părții transformate (marcate cu tildă) a crescut, iar cea a părții restante (barate) a scăzut cu o unitate.

După a $(n - 1)$ -a etapă vom obține (A, b) în forma dorită.

Pentru organizarea procedurală a calculelor, este esențial să observăm că, în acord cu (5.57), reflectorul U_k este complet determinat de blocul său \bar{U}_k , care este el însuși un reflector de ordin $n - k + 1$. Ținând seama de această observație, precum și de modul de acțiune al lui \bar{U}_k în (5.58), rezultă că pasul 3.1 se reduce la aplicarea procedurii **H0** vectorului

$$\text{(H0)} \quad \bar{b}_k = A(k : n, k - 1)$$

(În acest fel se determină componentele nenule u_{ik} , $i = k : n$ ale vectorului u_k , scalarul β_k precum și vectorul transformat $\bar{b}_k \leftarrow \bar{U}_k \bar{b}_k$ din (5.59)).

Similar, pasul 3.2 se reduce la efectuarea transformărilor $\bar{A}_k \leftarrow \bar{U}_k \bar{A}_k$ și $Y_k \leftarrow Y_k \bar{U}_k$ din (5.58) și (5.59), ceea ce necesită aplicarea procedurilor **H1** și **H2** submatricelor lui A definite prin

$$\text{(H1)} \quad \bar{A}_k = A(k : n, k : n)$$

și respectiv

$$\text{(H2)} \quad Y_k = A(1 : n, k : n).$$

În sfârșit, ultima relație (5.59) arată că reflectorul U_k lasă b invariant, deci fiecare etapă k , $k = 2 : n - 1$ conservă automat asemănarea ortogonală a perechilor transformate (A, b) .

În rezumat, procedura de aducere a perechii (A, b) cu o singură intrare la forma superior Hessenberg completă poate fi formulată astfel.

1. $U_1 = \mathbf{H0}(b)$.
2. $A = \mathbf{H1}(U_1, A)$; $A = \mathbf{H2}(A, U_1)$.
3. Pentru $k = 2 : n - 1$
 1. $U_k = \mathbf{H0}(A(k : n, k - 1))$.
 2. $A(k : n, k : n) = \mathbf{H1}(U_k, A(k : n, k : n))$;
 $A(:, k : n) = \mathbf{H2}(A(:, k : n), U_k)$

Construcția matricei ortogonale $U = U_{n-1} \cdots U_2 U_1$ din (5.50), i.e. acumularea transformărilor, se face după schema

1. $U = U_1$
2. Pentru $k = 2 : n - 1$
 1. $U(k : n, 1 : n) = \mathbf{H1}(U_k, U(k : n, 1 : n))$.

În practică, se preferă deseori obținerea formei superior Hessenberg (5.47), (5.48), (5.49), ceea ce evident presupune testarea condiției $\|\bar{b}_k\| \neq 0$ înainte de parcurgerea etapei k , $k = 2 : n - 1$. Se obține următoarea procedură de calcul.

Algoritmul 5.6 (Construiește forma superior Hessenberg (5.47), (5.48), (5.49) a unei perechi (A, b) de ordin n cu o singură intrare și acumulează transformările. Se presupune $b \neq 0$).

1. $k = 1$
2. $U_1 = \mathbf{H0}(b)$
3. $A = \mathbf{H1}(U_1, A)$, $A = \mathbf{H2}(A, U_1)$
4. $U = U_1$
5. CONTINUĂ='da'
6. Cât timp CONTINUĂ='da'
 1. $r = k$
 2. Dacă $k = n$
atunci
 1. CONTINUĂ='nu'
 2. tipărește "Perechea (A, b) este controlabilă"

```

altfel
1.  $g = A(k + 1 : n, k)$ 
2. Dacă  $\|g\| = 0$ 
   atunci
   1. CONTINUĂ = 'nu'
   altfel
   1.  $k \leftarrow k + 1$       % Etapa  $k, k \geq 2$ 
   2.  $U_k = \mathbf{H0}(g)$ 
   3.  $A(k : n, k : n) = \mathbf{H1}(U_k, A(k : n, k : n))$ ,
       $A(:, k : n) = \mathbf{H2}(A(:, k : n), U_k)$ 
   4.  $U(k : n, 1 : n) = \mathbf{H1}(U_k, U(k : n, 1 : n))$ 

```

Comentarii. CONTINUĂ este o variabilă logică ce determină avansarea contorului de etapă k . În final, variabila r furnizează rangul matricei de controlabilitate R a perechii date (A, b) . În cazul $r = n$, deci dacă perechea (A, b) este controlabilă (vezi pasul 6.2) numărul de operații necesar este de ordinul $5\frac{n^3}{3}$. \diamond

Cazul $m > 1$

În cazul unei perechi (A, B) de ordin n cu m intrări, $m \geq 1$, procedura de transformare are la bază următorul rezultat general.

Propoziția 5.4 (Lema de deflație controlabilă). *Oricare ar fi perechea (A, B) cu $\text{rang} B = r_1 > 0$, există o matrice ortogonală $U_1 \in \mathcal{R}^{n \times n}$ astfel încât*

$$A \leftarrow \tilde{A} = U_1 A U_1^T = \begin{bmatrix} A_1 & X \\ G & F \end{bmatrix}, \quad B \leftarrow \tilde{B} = U_1 B = \begin{bmatrix} H_1 \\ 0 \end{bmatrix} \quad (5.60)$$

unde matricea $H_1 \in \mathcal{R}^{r_1 \times m}$ este epică, i.e. $\text{rang} H_1 = r_1$. Mai mult, perechea (A, B) este controlabilă dacă și numai dacă perechea redusă (F, G) , de ordin $n - r_1 < n$, este controlabilă.

Demonstrație. Matricea U_1 se determină aplicând lui B procedura de descompunere a valorilor singulare (**DVS**) sau procedura de triangularizare ortogonală cu pivotarea coloanelor. În primul caz avem

$$U_1 B V = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix}, \quad (5.61)$$

unde U_1 și $V \in \mathcal{R}^{m \times m}$ sunt ortogonale iar Σ_1 este diagonală de ordin $r_1 > 0$ cu elementele diagonale pozitive. Partiționând $V = [V_1 \ V_2]$, unde V_1 are r_1 coloane, obținem

$$\tilde{B} = U_1 B = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} = \begin{bmatrix} \Sigma_1 V_1^T \\ 0 \end{bmatrix},$$

unde $H_1 = \Sigma_1 V_1^T$ este evident epică de rang r_1 . În al doilea caz obținem

$$U_1 B P = \begin{bmatrix} R_1 & R_2 \\ 0 & 0 \end{bmatrix}, \quad (5.62)$$

unde U_1 este ortogonală iar $P \in \mathcal{R}^{m \times m}$ este o matrice de permutare a coloanelor lui B , aleasă pe parcursul procedurii de triangularizare ortogonală cu reflectori astfel încât matricea superior triunghiulară R_1 de ordin r_1 să rezulte nesingulară. În acest caz matricea $\tilde{B} = U_1 B$ rezultă din nou în forma (5.60) cu $H_1 = [R_1 \ R_2] P^T$ ¹.

După determinarea lui U_1 printr-una din procedurile (5.61) sau (5.62) se aplică U_1 lui A , i.e. se calculează $\tilde{A} = U_1 A U_1^T$, și în acord cu (5.60) se evidențiază perechea redusă (F, G) de ordin $n - r_1$ cu r_1 intrări.

A doua afirmație din leamnă se demonstrează ca în cazul $m = 1$, utilizând punctul 2 al propoziției 5.1. \diamond

În ipoteza $G \neq 0$, perechii (F, G) cu $\text{rang } G = r_2 > 0$ i se poate aplica din nou o transformare de tip (5.60). Se obține

$$F \leftarrow \tilde{F} = \tilde{U}_2 F \tilde{U}_2^T = \begin{bmatrix} A_2 & X \\ G_{nou} & F_{nou} \end{bmatrix}, \quad G \leftarrow \tilde{G} = \tilde{U}_2 G = \begin{bmatrix} H_2 \\ 0 \end{bmatrix}$$

sau echivalent

$$A \leftarrow \tilde{A} = U_2 A U_2^T = \left[\begin{array}{c|cc} A_1 & X & X \\ H_2 & A_2 & X \\ 0 & G_{nou} & F_{nou} \end{array} \right], \quad (5.63)$$

$$B \leftarrow \tilde{B} = U_2 B = \begin{bmatrix} H_1 \\ 0 \\ 0 \end{bmatrix},$$

¹În cazul $m = 1$, considerat anterior, avem $r_1 = 1$, U_1 este reflectorul care anulează ultimele $n - 1$ elemente ale matricei $B = b \in \mathcal{R}^n$ iar $P = 1$. Această constatare evidențiază clar complicațiile de calcul ce apar în cazul $m > 1$ față de cazul simplu $m = 1$.

unde matricele \bar{U}_2 precum și

$$U_2 = \begin{bmatrix} I_{r_1} & 0 \\ 0 & \bar{U}_2 \end{bmatrix} \quad (5.64)$$

sunt ortogonale iar matricea H_2 este epică de rang r_2 .

Forma finală a perechii (A, B) , obținută prin aplicarea repetată a procedurii de deflație descrisă în propoziția 5.4, se numește *forma bloc-superior Hessenberg* și este definită prin

$$A \leftarrow \tilde{A} = UAU^T = \begin{bmatrix} A_R & A_{R\bar{R}} \\ 0 & A_{\bar{R}} \end{bmatrix}, \quad B \leftarrow \tilde{B} = UB = \begin{bmatrix} B_R \\ 0 \end{bmatrix}, \quad (5.65)$$

în care perechea (A_R, B_R) , de ordin r , se găsește în forma *bloc-superior Hessenberg controlabilă*

$$A_R = \begin{bmatrix} A_1 & X & \cdots & X & X \\ H_2 & A_2 & \cdots & X & X \\ & \ddots & \ddots & \vdots & \vdots \\ & & \ddots & \ddots & \vdots \\ & & & H_k & A_k \end{bmatrix}, \quad B_R = \begin{bmatrix} H_1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}, \quad (5.66)$$

unde toate blocurile $H_i \in \mathcal{R}^{r_i \times r_{i-1}}$ sunt epice, i.e.

$$\text{rang} H_i = r_i > 0, \quad i = 1 : k, \quad r_0 \stackrel{\text{def}}{=} m. \quad (5.67)$$

Evident avem

$$r = \sum_{i=1}^k r_i \leq n \quad (5.68)$$

și în virtutea condițiilor (5.67) perechea (A_R, B_R) din (5.65) este controlabilă. Mai mult, *perechea inițială* (A, B) este controlabilă dacă și numai dacă $r = n$, i.e. $(\tilde{A}, \tilde{B}) = (A_R, B_R)$, și în acest caz numărul k de blocuri din (5.66) coincide cu *indicele de controlabilitate* ν al lui (A, B) .

Matricea ortogonală $U \in \mathcal{R}^{n \times n}$ din (5.65) rezultă prin acumularea transformărilor parțiale, i.e.

$$U = U_k \cdots U_2 U_1, \quad (5.69)$$

în care U_1 se determină ca în demonstrația lemei 5.4, U_2 are forma (5.64) etc. Subliniem că în virtutea structurii lui U_2 , premultiplicarea cu U_2 nu modifică primele r_1 linii, iar postmultiplicarea cu U_2^T nu modifică primele

r_1 coloane ale matricei asupra căreia acționează transformarea corespunzătoare. De asemenea, proprietăți asemănătoare au toate matricele U_i , $i = 2 : k$ din (5.69).

În rezumat, aducerea pe loc a perechii inițiale (A, B) la forma bloc-superior Hessenberg (5.65), (5.66) se face printr-o succesiune de transformări ortogonale de asemănare

$$A \leftarrow U_i A U_i^T, \quad B \leftarrow U_i B, \quad i = 1 : k. \quad (5.70)$$

La etapa $i = 1$ se operează asupra perechii inițiale (A, B) ca în demonstrația propoziției 5.4, iar la etapele $i \geq 2$ se operează analog asupra perechii reduse corespunzătoare (F, G) , conținute în tabloul A . În particular, transformările U_i , $i \geq 2$ nu modifică forma lui B obținută după prima etapă, deci pentru $i \geq 2$ relațiile (5.70) se reduc la

$$A \leftarrow U_i A, \quad A \leftarrow A U_i^T. \quad (5.71)$$

Pentru a descrie precis algoritmul astfel obținut, notăm cu

$$[r_1, U_1, H_1] = \mathbf{red}(B) \quad (5.72)$$

procedura de reducere a lui B decrișă în demonstrația propoziției 5.4, vezi relațiile (5.60) și (5.61) sau (5.62). (Amintim că la etapele $i \geq 2$ procedura **red** se aplică unui bloc G , localizat în tabloul A). De asemenea, notăm transformările (5.71) respectiv prin

$$A = \mathbf{mul1}(U_i, A), \quad A = \mathbf{mul2}(A, U_i),$$

convenind să indicăm explicit de fiecare dată elementele lui A care sunt efectiv modificate.

Acumularea transformărilor din (5.69) se face conform schemei

$$U \leftarrow U_i U, \quad i = 2 : k \quad (5.73)$$

cu inițializarea $U = U_1$.

Algoritmul 5.7 (Dată perechea (A, B) de ordin n cu m intrări și $B \neq 0$ algoritmul construiește forma bloc-superior Hessenberg (5.65), (5.66), (5.67) și acumulează transformările. Perechea rezultată suprascrie pe cea inițială).

1. $k = 1$
2. $[r_1, U_1, H_1] = \mathbf{red}(B)$

3. $A = \mathbf{mul1}(U_1, A), \quad A = \mathbf{mul2}(A, U_1)$
4. $U = U_1$
5. $r_v = 0, \quad r = r_1$
6. CONTINUĂ = 'da'
7. Cât timp CONTINUĂ = 'da'
 1. $\nu = k$
 2. Dacă $r = n$
atunci
 1. CONTINUĂ = 'nu'
 2. Tipărește "Perechea (A, B) este controlabilă"
 - altfel
 1. $G = A(r+1 : n, r_v+1 : r)$
 2. Dacă $\|G\| = 0$
atunci
 1. CONTINUĂ = 'nu'
 - altfel
 1. $k \leftarrow k+1$ % (etapa $k, k \geq 2$)
 2. $[r_k, U_k, H_k] = \mathbf{red}(G)$
 3. $A(r+1 : n, r+1 : n) =$
 $\mathbf{mul1}(U_k, A(r+1 : n, r+1 : n)),$
 $A(1 : n, r+1 : n) = \mathbf{mul2}(A(1 : n, r+1 : n), U_k)$
 4. $U(r+1 : n, 1 : n) = \mathbf{mul1}(U_k, U(r+1 : n, 1 : n))$
 5. $r_v = r, \quad r \leftarrow r + r_k$

Comentarii. Având în vedere că numărul k al etapelor de reducere nu este apriori cunoscut, utilizăm și aici variabila binară CONTINUĂ cu semnificația cunoscută. În final, variabila r furnizează rangul matricei de controlabilitate R a perechii date (A, B) . În cazul $r = n$, i.e. dacă perechea (A, B) este controlabilă, algoritmul se oprește la instrucțiunea 7.2. iar ν este indicele de controlabilitate. Din contră, dacă algoritmul se oprește datorită faptului că blocul curent G este nul, atunci perechea (A, B) nu este controlabilă, iar perechea transformată (\tilde{A}, \tilde{B}) , furnizată de algoritmul 5.7, are forma (5.65), în care $r < n$. \diamond

Forma Hessenberg observabilă

Pentru analiza proprietăților de observabilitate a unei perechi (C, A) de ordin n cu l ieșiri, $l \geq 1$, se procedează prin dualitate. Presupunem $C \neq 0$

și fie $q_1 = \text{rang } C > 0$. Lema de deflație observabilă (vezi propoziția 5.4) afirmă existența unei matrice ortogonale $U_1 \in \mathcal{R}^{n \times n}$, astfel încât

$$A \leftarrow \tilde{A} = U_1 A U_1^T = \begin{bmatrix} A_1 & H \\ X & F \end{bmatrix}, \quad (5.74)$$

$$C \leftarrow \tilde{C} = C U_1^T = \begin{bmatrix} G_1 & 0 \end{bmatrix},$$

unde matricea $G_1 \in \mathcal{R}^{l \times q_1}$ este *monică*, i.e. $\text{rang } G_1 = q_1$. În plus, perechea (C, A) este observabilă dacă și numai dacă perechea redusă (H, F) este observabilă.

Prin aplicarea repetată a acestei proceduri de deflație, sau echivalent, prin aplicarea schemei de dualizare relativ la algoritmul 5.7, se obține *forma bloc-inferior Hessenberg* a perechii (C, A) , definită prin

$$A \leftarrow \tilde{A} = U A U^T = \begin{bmatrix} A_Q & 0 \\ A_{\bar{Q}Q} & A_{\bar{Q}} \end{bmatrix}, \quad (5.75)$$

$$C \leftarrow \tilde{C} = C U^T = \begin{bmatrix} C_Q & 0 \end{bmatrix},$$

în care perechea (C_Q, A_Q) de ordin q se află în forma *bloc-inferior Hessenberg observabilă*

$$A_Q = \begin{bmatrix} A_1 & G_2 & & & \\ X & A_2 & \ddots & & \\ \vdots & \vdots & \ddots & \ddots & \\ X & X & \cdots & \ddots & G_k \\ X & X & \cdots & \cdots & A_k \end{bmatrix}, \quad C_Q = \begin{bmatrix} G_1 & 0 & \cdots & 0 \end{bmatrix}, \quad (5.76)$$

unde toate blocurile $G_i \in \mathcal{R}^{q_{i-1} \times q_i}$ sunt monice, i.e.

$$\text{rang } G_i = q_i > 0, \quad i = 1 : k, \quad q_0 \stackrel{\text{def}}{=} l. \quad (5.77)$$

Evident avem

$$q = \sum_{i=1}^k q_i \leq n \quad (5.78)$$

și în virtutea condițiilor (5.77) perechea (C_Q, A_Q) din (5.75) este observabilă. În plus, *perechea inițială* (C, A) este *observabilă* dacă și numai dacă $q = n$, i.e. $(\tilde{C}, \tilde{A}) = (C_Q, A_Q)$, și în acest caz numărul k de blocuri din (5.76) coincide cu *indicele de observabilitate* al lui (C, A) .

În sfârșit, la fel ca în cazul controlabilității, matricea ortogonală $U \in \mathcal{R}^{n \times n}$ din (5.75) rezultă din acumularea transformărilor parțiale, i.e.

$$U = U_k \cdots U_2 U_1,$$

în care U_1 se determină ca în (5.74).

5.6 Descompunerea controlabilă

Fie (A, B) o pereche, în general necontrolabilă, de ordin n cu m intrări, unde $m \geq 1$. Perechea bloc-superior Hessenberg $(\tilde{A}, \tilde{B}) = (UAU^T, UB)$, obținută aplicând algoritmul 5.7 expus în paragraful anterior, are forma (5.65) cu $r \leq n$ și constituie prin definiție, *descompunerea controlabilă* a perechii date (A, B) . (În cazul $m = 1$ se utilizează (5.47)). Perechea controlabilă (A_R, B_R) de ordin r cu m intrări reprezintă *partea controlabilă*, iar matricea $A_{\bar{R}}$ de ordin $\bar{r} = n - r$ reprezintă *partea total necontrolabilă* a lui (A, B) .

Semnificația sistemică a descompunerii controlabile este evidențiată de graful de semnal asociat, vezi figura 5.1(a), din care se vede că partea necontrolabilă nu este influențată de intrarea u nici direct, nici indirect (i.e. prin intermediul părții controlabile), iar matricea $A_{R\bar{R}}$ caracterizează acțiunea părții total necontrolabile asupra părții controlabile. (De aceea, în aplicații, partea necontrolabilă modelează perturbațiile exterioare ce acționează asupra sistemului dat, considerat controlabil). Prin definiție, valorile proprii ale matricei $A_{\bar{R}}$ constituie valorile proprii necontrolabile sau *polii ficși ai perechii* (A, B) .

Principalele elemente de analiză a proprietăților de controlabilitate asociate unei perechi oarecare (A, B) sunt furnizate direct de algoritmul 5.7 și pot fi rezumate astfel.

1. *Testarea controlabilității* se face pe baza egalității $r = n$. Altfel spus, perechea (A, B) este *controlabilă* dacă și numai dacă partea sa total necontrolabilă este vidă.

2. *Spațiul controlabil* $\mathcal{R} \subset R^n$ al perechii (A, B) este definit prin $\mathcal{R} = \text{Im}R$, unde R este matricea de controlabilitate corespunzătoare. (Amintim că prin $\text{Im}R$, unde matricea R are n linii, notăm subspațiul liniar al lui \mathcal{R}^n generat de coloanele lui R , vezi (5.26)). Avem evident $r = \text{rang}R = \dim\mathcal{R}$, unde r este ordinul perechii (A_R, B_R) din (5.65). Mai mult, partiționând conform matricea de transformare U , i.e.

$$U = \left[\begin{array}{c} U_R \\ U_{\bar{R}} \end{array} \right] \left. \begin{array}{l} \} r \\ \} \bar{r} = n - r \end{array} \right\} \quad (5.79)$$

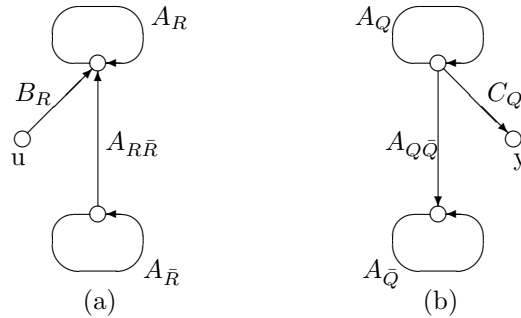


Figura 5.1: Semnificația sistemică a *descompunerii controlabile* (a) și a *descompunerii observabile* (b).

și, utilizând aceleași raționamente ca în § 2, constatăm că – în coordonatele inițiale – coloanele matricei U_R^T formează o *bază ortogonală* a subspațiului controlabil \mathcal{R} .

În coordonatele lui \mathcal{R}^n corespunzătoare reprezentării (5.65) subspațiul controlabil \mathcal{R} este format din totalitatea vectorilor de forma

$$\tilde{x} = \left[\begin{array}{c} x_R \\ 0 \end{array} \right] \left. \begin{array}{l} \} r \\ \} \bar{r} \end{array} \right\} , \quad (5.80)$$

iar structura perechii transformate (\tilde{A}, \tilde{B}) din (5.65) arată că

- (i) subspațiul \mathcal{R} este A -invariant și
- (ii) conține $\text{Im}B$.

Mai mult, în virtutea faptului că perechea (A_R, B_R) este controlabilă (deci nu poate fi descompusă la rândul ei ca în (5.65)), subspațiul \mathcal{R} este *cel mai mic* subspațiu al lui \mathcal{R}^n cu proprietățile de mai sus.

În legătură cu aceasta, subliniem că noțiunea de subspațiu necontrolabil al perechii (A, B) nu are semnificație geometrică invariantă deoarece orice subspațiu $\tilde{\mathcal{R}}$ (vezi figura 5.2) cu proprietatea $\mathcal{R}^n = \tilde{\mathcal{R}} \oplus \mathcal{R}$ poate fi numit astfel. În aplicații se consideră de obicei complementul ortogonal $\tilde{\mathcal{R}} = \mathcal{R}^\perp$, pentru care matricea U_R^T definește evident o bază ortogonală. (Simplitatea acestor construcții se datorează faptului că matricea de transformare (5.80) este ortogonală, i.e. coloanele lui U^T constituie o bază ortogonală a întregului spațiu al stărilor \mathcal{R}^n).

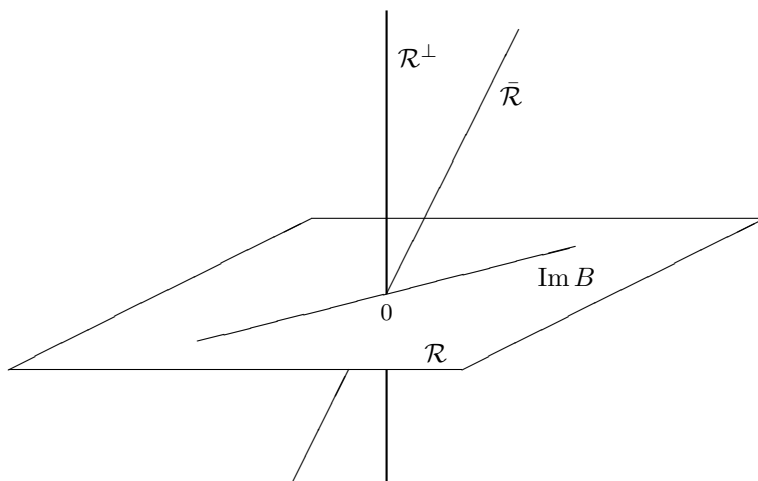


Figura 5.2: Reprezentarea geometrică a spațiului controlabil.

Descompunerea observabilă

Fie (C, A) o pereche în general neobservabilă de ordin n cu l ieșiri, unde $l \geq 1$. Perechea bloc-inferior Hessenberg $(\tilde{C}, \tilde{A}) = (CU^T, UAU^T)$, obținută aplicând schema de dualizare relativ la algoritmul 5.7, expus în paragraful anterior, are forma (5.75) cu $q \leq n$, și constituie, prin definiție, *descompunerea observabilă* a perechii date (C, A) .

Semnificația sistemică a descompunerii observabile este evidențiată de graficul de semnal asociat, vezi figura 5.1(b). Prin definiție, valorile proprii ale matricei $A_{\tilde{O}}$ constituie valorile proprii neobservabile sau *polii ficși ai perechii* (C, A) .

Principalele elemente de analiză a proprietăților de observabilitate asociate unei perechi oarecare (C, A) sunt propuse ca exerciții.

5.7 Stabilizabilitatea și detectabilitatea sistemelor liniare

În problemele de sinteză vizând stabilizarea sistemelor liniare prin reacție după stare sau prin compensare dinamică, interesează o proprietate mai

slabă decât controlabilitatea, numită *stabilizabilitate* și definită cu referire la descompunerea controlabilă (5.65).

Definiția 5.6 *Perechea (A, B) este stabilizabilă dacă partea sa total necontrolabilă $A_{\bar{R}}$ este stabilă sau, echivalent, polii fiși ai perechii (A, B) sunt stabili.*

În consecință, principalele elemente de analiză a proprietăților de stabilizabilitate asociate unei perechi (A, B) rezultă combinând procedurile de analiză a proprietăților fundamentale de stabilitate și controlabilitate, stabilite în §§ 1, 2 și §§ 5, 6.

1. *Testarea stabilizabilității* se face pe baza definiției 5.6, i.e. construind descompunerea controlabilă (5.65) a perechii (A, B) și apoi aplicând algoritmul 5.1 matricei $A_{\bar{R}}$. Alternativ, se poate construi descompunerea spectrală a perechii (A, B) ca în §2, după care se aplică algoritmul 5.7 perechii total nestabile (A_b, B_b) .

2. *Spațiul stabilizabil $\mathcal{S} \subset R^n$* al perechii (A, B) este definit prin $\mathcal{S} = \mathcal{R} + \mathcal{X}^-(A)$, astfel încât evidențierea sa se face construind descompunerea controlabilă (5.65) a perechii (A, B) și, apoi, aplicând algoritmul 5.2 matricei $A_{\bar{R}}$.

În urma acestor două transformări se obține, evident, *descompunerea stabilizabilă* a perechii date (A, B) , care evidențiază partea stabilizabilă (controlabilă sau necontrolabilă dar stabilă), respectiv partea total nestabilizabilă (necontrolabilă și nestabilă) a perechii considerate. (În aplicații, partea nestabilizabilă modelează perturbațiile exterioare persistente ce acționează asupra sistemului dat, considerat stabilizabil).

În problemele de proiectare a estimatoarelor de stare stabile interesează o proprietate mai slabă decât observabilitatea, numită *detectabilitate* și definită prin dualitate față de stabilizabilitate cu referire la descompunerea observabilă (5.75).

Definiția 5.7 *Perechea (C, A) este detectabilă dacă partea total neobservabilă $A_{\bar{Q}}$ este stabilă sau, echivalent, polii fiși ai perechii (C, A) sunt stabili.*

Principalele elemente de analiză a proprietăților de detectabilitate asociate unei perechi (C, A) sunt propuse ca exerciții.

5.8 Realizări minimale

Principalul rezultat al teoriei realizării sistemelor liniare afirmă că orice sistem $S = (A, B, C)$ este echivalent intrare-ieșire, i.e. are aceeași matrice de transfer, cu un sistem de ordin minim $S_m = (A_m, B_m, C_m)$ care este simultan controlabil și observabil. Mai mult, sistemul S_m este determinat până la o transformare de asemănare și, în esență, coincide cu partea simultan controlabilă și observabilă a sistemului dat S .

Definiția 5.8 *Un sistem $S_m = (A_m, B_m, C_m)$ controlabil și observabil, echivalent intrare-ieșire cu S se numește realizare minimală a lui S .*

Mai general, orice sistem $S_m = (A_m, B_m, C_m)$ controlabil și observabil se numește minimal.

În consecință, o realizare minimală (i.e. de ordin minim) a unei matrice de transfer date $T(s)$ poate fi construită aplicând următoarea procedură.

Algoritmul 5.8 (Construiește o realizare minimală a matricei de transfer $T(s)$).

1. Se construiește o realizare observabilă $S = (A, B, C)$ a lui $T(s)$, utilizând formele standard observabile, expuse în capitolul 3.
2. Se construiește descompunerea controlabilă

$$\tilde{A} = UAU^T = \begin{bmatrix} A_R & A_{R\bar{R}} \\ 0 & A_{\bar{R}} \end{bmatrix}, \quad \tilde{B} = UB = \begin{bmatrix} B_R \\ 0 \end{bmatrix},$$

utilizând algoritmul 5.7 și se aplică transformarea U matricei C , i.e.

$$\tilde{C} = CU^T = [C_R \ C_{\bar{R}}].$$

3. Se reține partea controlabilă $S_m = (A_R, B_R, C_R)$ a triplețului transformat, unde

$$A_R = \tilde{A}(1:r, 1:r), \quad B_R = \tilde{B}(1:r, :)$$

este partea controlabilă a perechii (A, B) iar

$$C_R = \tilde{C}(:, 1:r).$$

Procedurile de identificare sistemică, bazate pe teoria realizării sistemelor liniare, sunt prezentate în ANEXA A.

Procedurile de reducere dimensională aproximativă a unor sisteme controlabile și observabile (deci minimale) au la bază conceptul fundamental de realizare echilibrată și vor fi expuse în §10.

5.9 Gramieni de controlabilitate și observabilitate

Pentru a introduce cât mai simplu noțiunea de gramian în contextul sistemic al discuției noastre, considerăm din nou sistemul liniar discret $S = (A, B, C)$ definit prin relațiile (5.17).

Utilizând matricea de controlabilitate în k pași R_k a perechii (A, B) și ținând seama de relațiile (5.19) și (5.22) definim *gramianul de controlabilitate în k pași* prin

$$Y_k \stackrel{\text{def}}{=} R_k R_k^T = \sum_{i=0}^{k-1} A^i B B^T (A^T)^i. \quad (5.81)$$

Evident, Y_k este o matrice $n \times n$ simetrică pozitiv semidefinită, pe scurt $Y_k = Y_k^T \geq 0$, $k \geq 1$. Ținând seama de relația (5.38) putem scrie

$$Y_{k+1} = R_{k+1} R_{k+1}^T = \begin{bmatrix} B & A R_k \end{bmatrix} \begin{bmatrix} B^T \\ R_k^T A^T \end{bmatrix},$$

prin urmare șirul Y_k satisface ecuația matriceală recurentă Liapunov discretă (DEML)

$$Y_{k+1} = B B^T + A Y_k A^T, \quad (5.82)$$

cu condiția inițială $Y_1 = B B^T$ (sau $Y_0 = 0$).

Prin dualitate, utilizând matricea de observabilitate în k pași Q_k a perechii (C, A) , definim *gramianul de observabilitate în k pași* prin

$$X_k \stackrel{\text{def}}{=} Q_k^T Q_k = \sum_{i=0}^{k-1} (A^T)^i C^T C A^i, \quad (5.83)$$

unde $X_k = X_k^T \geq 0$, $k \geq 1$ satisface DEML

$$X_{k+1} = C^T C + A^T X_k A \quad (5.84)$$

cu condiția inițială $X_1 = C^T C$ (sau $X_0 = 0$).

Dacă matricea discretă A este *stabilă*, i.e. $|\lambda(A)| < 1$, atunci se poate arăta ușor că pentru $k \rightarrow \infty$ șirurile (5.81) și (5.83) converg către limitele $Y = Y^T \geq 0$ și respectiv $X = X^T \geq 0$ care satisfac evident ecuațiile matriceale algebrice Liapunov discrete (DEMAL)

$$Y = B B^T + A Y A^T \quad (5.85)$$

și respectiv

$$X = C^T C + A^T X A. \quad (5.86)$$

În cazul sistemelor cu timp continuu *stabile*, i.e. $\operatorname{Re}\lambda(A) < 0$, considerații similare conduc la ecuațiile matriceale algebrice Liapunov (EMAL)

$$0 = BB^T + AY + YA^T, \quad (5.87)$$

respectiv

$$0 = C^T C + A^T X + XA. \quad (5.88)$$

Având în vedere că, în condițiile de *stabilitate* formulate mai sus, ecuațiile (5.85) – (5.88) au soluții unice, în continuare vom adopta următoarea definiție unificatoare.

Definiția 5.9 *Matricele simetrice $Y, X \in \mathcal{R}^{n \times n}$, unic definite ca soluții ale EMAL (5.87), (5.88) (respectiv DEMAL (5.85), (5.86) în cazul discret), se numesc gramian de controlabilitate al perechii (A, B) , respectiv gramian de observabilitate al perechii (C, A) .*

Rezultatul următor se demonstrează ușor utilizând propoziția 5.1 și are o importanță deosebită pentru tot restul capitolului.

Propoziția 5.5 *Sistemul liniar stabil $S = (A, B, C)$ este controlabil, respectiv observabil, dacă și numai dacă gramianul corespunzător Y , respectiv X , este o matrice pozitiv definită.*

Observația 5.5 Utilizând rezultatele obținute în capitolul 1 precum și în §2 al acestui capitol se poate vedea ușor că ecuațiile (5.85) – (5.88) au soluții unice chiar dacă, în locul condiției de stabilitate, este îndeplinită numai condiția (mai slabă) de *dihotomie*, i.e. $\operatorname{Re}\lambda(A) \neq 0$ (respectiv $|\lambda(A)| \neq 1$ în cazul discret). Mai mult, în acest caz $S = (A, B, C)$ este controlabil, respectiv observabil, dacă și numai dacă gramianul corespunzător este o matrice simetrică inversabilă a cărei inerție este dictată de numerele valorilor proprii stabile și nestabile ale lui A . \diamond

În cele ce urmează nu vom exploata implicațiile importante ale observației 5.5, ci vom presupune mereu că sistemul considerat $S = (A, B, C)$ este stabil.

Adoptând această ipoteză, profităm de ocazie pentru a efectua o scurtă trecere în revistă a unor proceduri iterative de rezolvare a DEMAL. (Procedurile ”de transformare”, bazate pe utilizarea formei Schur, au fost discutate în capitolul 1).

Pentru precizare, în continuare ne referim la ecuația (5.86).

Schema de iterare "directă", care utilizează ecuația recurentă corespunzătoare (5.84), este neeficientă deoarece, de regulă, convergența $X_k \rightarrow X$ este lentă (cu atât mai lentă cu cât raza spectrală $\rho(A) < 1$ a matricei discrete A este mai apropiată de 1, vezi observația 5.1). Pe de altă parte, matricele X_k sunt pozitiv semidefinite, ceea ce permite utilizarea avantajă pe parcursul calculelor a unor factorizări de tip Cholesky. Având în vedere aceste observații de principiu, în practică se recomandă două tipuri de scheme iterative.

Algoritmii "de dublare" urmăresc accelerarea convergenței. Ideea constă în a efectua iterarea cu pas mereu dublat, i.e. conform schemei

$$X_1 \rightarrow X_2 \rightarrow X_4 \rightarrow X_8 \rightarrow \dots,$$

fără a mai calcula iterațiile intermediare. Notăm

$$X^{(s)} = X_k, \quad k = 2^s, \quad s \geq 0. \quad (5.89)$$

Prin calcul direct, din (5.83) sau (5.84) rezultă

$$\begin{aligned} X^{(0)} &= X_1 = C^T C, \\ X^{(1)} &= X_2 = A^T X^{(0)} A + X^{(0)}, \\ X^{(2)} &= X_4 = (A^2)^T X^{(1)} A^2 + X^{(1)} \end{aligned}$$

și în general

$$X^{(s+1)} = (A^k)^T X^{(s)} A^k + X^{(s)}. \quad (5.90)$$

Prin urmare, schema de iterare accelerată, bazată pe aplicarea formulei de dublare (5.90), poate fi rezumată astfel.

Algoritm 5.9 (Rezolvă DEMAL (5.86) în cazul A stabilă utilizând un algoritm de dublare).

1. CONTINUĂ = 'da'
 2. $X = C^T C$
 3. Cât timp CONTINUĂ = 'da'
 1. $X(+1) = AXA^T + X$
 2. Dacă $\|X(+1) - X\| < tol\|X\|$
atunci
 1. CONTINUĂ = 'nu'
- altfel

1. $X = X(+1)$
2. $A \leftarrow A^2$.

La fiecare iterație sunt necesare $2n^3$ înmulțiri (se exploatează simetria lui X), deci schema se justifică dacă $\rho(A)$ nu este prea aproape de 1 i.e. convergența este suficient de rapidă.

Algoritmii "de rădăcină pătrată" urmăresc exploatarea caracterului pozitiv (semi)definit al matricelor X_k din (5.83) prin determinarea factorilor Cholesky

$$X_k = L_k^T L_k \geq 0, \quad (5.91)$$

unde, evident, se poate lua $L_1 = C$. Dacă pentru $k \geq 1$ matricea L_k este cunoscută atunci L_{k+1} se determină din (5.84), observând că trebuie să avem

$$X_{k+1} = L_{k+1}^T L_{k+1} = C^T C + A^T L_k^T L_k A = \begin{bmatrix} C^T & (L_k A)^T \end{bmatrix} \begin{bmatrix} C \\ L_k A \end{bmatrix}.$$

Este clar că putem lua

$$\begin{bmatrix} L_{k+1} \\ 0 \end{bmatrix} = U_k \begin{bmatrix} C \\ L_k A \end{bmatrix}, \quad (5.92)$$

unde U_k este o matrice *ortogonală* (e.g. o secvență de reflectori) care aduce matricea bloc din dreapta la forma *superior triunghiulară*. (În consecință, toate matricele șirului L_k , $k \geq 2$ rezultă superior triunghiulare).

Prin urmare, schema de iterare "compactă", bazată pe aplicarea relației (5.92) poate fi rezumată astfel.

Algoritmul 5.10 (Rezolvă DEMAL (5.86) în cazul A stabilă utilizând un algoritm de rădăcină pătrată, bazat pe triangularizarea ortogonală).

1. CONTINUĂ = 'da'
 2. $L = C$
 3. Cât timp CONTINUĂ = 'da'
 1. $\begin{bmatrix} L(+1) \\ 0 \end{bmatrix} = U \begin{bmatrix} C \\ AL \end{bmatrix}$,
 2. Dacă $\|L(+1) - L\| < tol \|L\|$
atunci
 1. CONTINUĂ = 'nu'
- altfel

1. $L = L(+1)$.

Datorită convergenței în general lente, schema de rădăcină pătrată, a cărei idee a fost expusă mai sus relativ la DEMAL (5.86), se utilizează avantajos mai ales pentru rezolvarea DEML (5.84) și în special a unor ecuații matriceale recurente de tip Riccati, la care păstrarea caracterului pozitiv (semi-)definit al soluției prin iterare directă este problematică. În legătură cu această mențiune, observăm că schemele de rădăcină pătrată rămân valabile ca atare chiar dacă sistemul discret $S = (A, B, C)$ este variabil în timp, i.e. matricele $A = A_k$, $B = B_k$, $C = C_k$ variază de la pas la pas.

5.10 Echilibrarea sistemelor liniare

Considerăm sistemul liniar stabil, controlabil și observabil $S = (A, B, C)$ și fie $Y = Y^T > 0$ și $X = X^T > 0$ gramienii de controlabilitate și observabilitate, definiți ca soluții unice ale EMAL (5.87) și respectiv (5.88).

Utilizând transformarea de asemănare

$$\tilde{A} = TAT^{-1}, \quad \tilde{B} = TB, \quad \tilde{C} = CT^{-1}, \quad (5.93)$$

obținem ușor

$$\begin{aligned} 0 &= \tilde{B}\tilde{B}^T + \tilde{A}\tilde{Y} + \tilde{Y}\tilde{A}^T, \\ 0 &= \tilde{C}^T\tilde{C} + \tilde{A}^T\tilde{X} + \tilde{X}\tilde{A}, \end{aligned}$$

unde gramienii transformați au expresiile

$$\tilde{Y} = TYT^T, \quad \tilde{X} = T^{-T}XT^{-1}. \quad (5.94)$$

În particular avem

$$\tilde{Y}\tilde{X} = T(YX)T^{-1}, \quad (5.95)$$

deci valorile proprii $\lambda(YX)$ ale produsului YX sunt *invarianti* ai sistemului $S = (A, B, C)$ în raport cu transformarea (5.93). (Se poate arăta că aceste valori proprii sunt pozitive și coincid cu pătratul valorilor singulare Hankel ale lui S , $[X]$). Mai mult, aceste valori proprii pot fi interpretate – într-un sens ce va fi precizat mai departe – ca expresii cantitative ale ”gradelor” de transfer intrare-ieșire atașate (celor n componente ale stării $x \in \mathcal{R}^n$) sistemului $S = (A, B, C)$.

Pe de altă parte, utilizând transformarea de asemănare (5.93), în care matricea $T = U$ este ortogonală, relațiile (5.94) devin

$$\tilde{Y} = UYU^T, \quad \tilde{X} = UXU^T. \quad (5.96)$$

Prin urmare, în acest caz se conservă *separat* valorile proprii (evident pozitive) ale matricelor simetrice pozitiv definite Y și X . Interpretând aceste valori proprii ca expresii cantitative ale "gradelor" de controlabilitate și respectiv observabilitate atașate lui S , putem spune sugestiv că aceste grade sunt *invarianti ortogonali* ai sistemului S . De aceea, în problemele de analiză *separată* a controlabilității și observabilității, abordate în §§3–8, am utilizat cu succes transformări ortogonale.

Din contră, în acest paragraf, dedicat analizei proprietăților de transfer ale sistemului S , rezultate (într-un sens precizat) prin "compunerea" proprietăților de controlabilitate și observabilitate, vom utiliza transformări (5.93) de formă generală, în particular, cu scopul "echilibrării" (egalării) gradelor corespunzătoare de controlabilitate și observabilitate.

Enunțul precis al problemei echilibrării este următorul.

Se dă sistemul $S = (A, B, C)$, presupus stabil, controlabil și observabil. Să se determine o matrice de transformare $T \in \mathcal{R}^{n \times n}$ inversabilă astfel încât gramienii transformați (5.94) să fie diagonali și egali, i.e.

$$\tilde{X} = \tilde{Y} = \Sigma, \quad (5.97)$$

unde Σ este o matrice diagonală cu elemente pozitive. (Conform relației (5.95) aceste elemente reprezintă valorile singulare Hankel ale lui S).

Definiția 5.10 *Dacă matricea inversabilă $T \in \mathcal{R}^{n \times n}$ este o soluție a problemei echilibrării formulate mai sus relativ la sistemul S , atunci sistemul transformat $\tilde{S} = (\tilde{A}, \tilde{B}, \tilde{C})$, definit prin relațiile (5.93), se numește reprezentare echilibrată a lui S .*

Mai general, orice sistem $\tilde{S} = (\tilde{A}, \tilde{B}, \tilde{C})$ cu proprietatea (5.97) se numește echilibrat.

Soluția problemei echilibrării în formularea de mai sus este simplă din punct de vedere principal și, în esență, se reduce la diagonalizarea simultană a două matrice simetrice, ambele pozitiv definite, prin transformări de congruență de forma (5.94).

Vom determina matricea căutată T sub forma

$$T = T_3 T_2 T_1, \quad (5.98)$$

unde transformările T_k , $k = 1 : 3$ se determină succesiv.

În prima etapă considerăm factorizarea Cholesky a matricei $Y > 0$, i.e.

$$Y = MM^T, \quad (5.99)$$

unde M este *inferior triunghiulară inversabilă* și luăm $T_1 = M^{-1}$. Aplicând (5.94) obținem

$$\begin{aligned} Y_1 &\stackrel{\text{def}}{=} T_1 Y T_1^T = M^{-1} (M M^T) M^{-T} = I, \\ X_1 &\stackrel{\text{def}}{=} M^T X M. \end{aligned} \quad (5.100)$$

În a doua etapă considerăm descompunerea spectrală a matricei $X_1 > 0$, obținută aplicând algoritmul QR simetric cu acumularea transformărilor, i.e.

$$X_1 = U^T \Sigma^2 U, \quad (5.101)$$

unde coloanele matricei ortogonale U^T sunt vectorii proprii ai matricei (simetrice) X_1 iar notația $\Sigma^2 \stackrel{\text{def}}{=} \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ se justifică ținând seama că valorile proprii $\lambda_i \stackrel{\text{def}}{=} \sigma_i^2$ ale matricei X_1 sunt strict pozitive. Convenim să alegem $\sigma_i > 0$, $i = 1 : n$ și luăm $T_2 = U$. Aplicând încă o dată (5.94) obținem

$$\begin{aligned} Y_2 &\stackrel{\text{def}}{=} T_2 Y_1 T_2^T = U U^T = I, \\ X_2 &\stackrel{\text{def}}{=} U (U^T \Sigma^2 U) U^T = \Sigma^2. \end{aligned} \quad (5.102)$$

Acum ambele matrice X_2 și Y_2 sunt diagonale. Pentru a le face egale, luăm $T_3 = \Sigma^{1/2}$ și aplicând ultima dată (5.94) obținem într-adevăr

$$\begin{aligned} \tilde{Y} &= T_3 Y_2 T_3^T = \Sigma^{1/2} \Sigma^{1/2} = \Sigma, \\ \tilde{X} &= \Sigma^{-1/2} (\Sigma^2) \Sigma^{-1/2} = \Sigma. \end{aligned} \quad (5.103)$$

În definitiv transformarea căutată (5.98) este

$$T = \Sigma^{1/2} U M^{-1}, \quad (5.104)$$

iar elementele pozitive ale matricei diagonale Σ coincid cu invariantii $\lambda^{1/2}(YX)$.

Este clar că procedura de mai sus admite o versiune duală, care începe cu factorizarea Cholesky $X = L^T L$ a matricei $X > 0$, unde L este *superior triunghiulară inversabilă*, etc.

Pentru a evidenția "simetria" de principiu a procedurii de echilibrare, considerăm ambele factorizări Cholesky

$$Y = M M^T, \quad X = L^T L \quad (5.105)$$

și formăm produsul LM ². Aplicând algoritmul de descompunere a valorilor singulare (**DVS**) obținem

$$LM = V^T \Sigma U, \quad (5.106)$$

unde U, V sunt matrice ortogonale, $\Sigma \stackrel{\text{def}}{=} (\sigma_1, \dots, \sigma_n)$ este diagonală, iar

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0 \quad (5.107)$$

sunt valorile singulare corespunzătoare. Transformarea căutată T poate fi exprimată prin formulele echivalente

$$T = \Sigma^{1/2} U M^{-1}, \quad \text{respectiv} \quad T = \Sigma^{-1/2} V L \quad (5.108)$$

și conduce la realizarea egalității dorite (5.97), unde în plus Σ satisface condiția de ordonare (5.107).

Prin urmare, procedura de echilibrare "simetrică", bazată pe relațiile (5.105) și (5.106), poate fi rezumată astfel.

Algoritmul 5.11 (Se echilibrează sistemul $S = (A, B, C)$, presupus stabil, controlabil și observabil).

1. Se aduce matricea A la forma Schur, i.e. $A \leftarrow U_0 A U_0^T$, utilizând algoritmul **QR** cu acumularea transformărilor.
2. Se face $B \leftarrow U_0 B$ și $C \leftarrow C U_0^T$.
3. Se rezolvă EMAL (5.87) și (5.88), ținând seama că matricele A și A^T sunt triunghiulare.
4. Se calculează factorii Cholesky L și M conform relației (5.105).
5. Se efectuează produsul LM și se calculează DVS (5.106).
6. Se calculează matricea de transformare T conform (5.108).
7. Se face $A \leftarrow T A T^{-1}$, $B \leftarrow T B$ și $C \leftarrow C T^{-1}$.

Comentarii. Pașii 1 și 2 au drept scop reducerea efortului de calcul implicat de rezolvarea ambelor EMAL (5.87) și (5.88) prin metodele expuse în capitolul 1, întrucât matricea A este adusă la forma Schur o singură dată. Utilizând o procedură specială de tip rădăcină pătrată, propusă de Hammarling [17], este posibilă contragerea pașilor 3 și 4, i.e. calculul direct al factorilor Cholesky L și M fără a determina explicit soluțiile X și Y

²Având în vedere definițiile (5.81) și (5.83), matricele M și L din (5.105) constituie versiuni "prescurtate" ale matricelor de controlabilitate și observabilitate, deci matricea $H = LM$ are o semnificație similară în raport cu matricea Hankel $H_{\alpha\beta} = Q_\alpha R_\beta$.

ale EMAL. În general, procedura de echilibrare rezumată prin algoritmul 5.11, care în esență reprezintă analogul sistemic al procedurii **DVS**, este susceptibilă de ameliorări vizând eliminarea a cât mai multe cantități intermediare în favoarea operării directe asupra datelor primare $S = (A, B, C)$ ale sistemului considerat. Căile prin care acest obiectiv de principiu poate fi cel mai bine realizat urmează încă a fi descoperite. \diamond

Dintre aplicațiile procedurii de echilibrare, expunem mai jos numai tehnicile de reducere dimensională a unui sistem $S = (A, B, C, D)$ în scopul eliminării stărilor neesențiale din punctul de vedere al transferului intrare-ieșire.

Aprecierea caracterului esențial, respectiv neesențial, al variabilelor de stare $x_i \in \mathcal{R}$, $i = 1 : n$ se face pe forma echilibrată $\tilde{S} = (\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ satisfăcând condiția (5.97), în care elementele matricei diagonale Σ sunt ordonate în acord cu (5.107). Subliniem că matricele $\tilde{A}, \tilde{B}, \tilde{C}$ ale lui \tilde{S} rezultă aplicând algoritmul 5.11 tripletului (A, B, C) iar $\tilde{D} = D$.

Fie $\varepsilon > 0$ un număr dat suficient de mic, înțeles ca prag de semnificație, și fie t cel mai mare număr întreg astfel încât

$$0 \leq t \leq n \quad \text{cu} \quad \sigma_i \geq \varepsilon, \quad i = 1 : t. \quad (5.109)$$

Sistemul de ordin redus S_t va fi de ordin t și, în esență, se obține eliminând stările x_j , $j = t + 1 : n$ corespunzătoare valorilor singulare neesențiale $\sigma_j < \varepsilon$, $j = t + 1 : n$.

Pentru a descrie succint cele două modalități principale de calcul ale matricelor sistemului redus S_t , considerăm partiția matricelor sistemului echilibrat \tilde{S} definită prin

$$\begin{aligned} \tilde{A} &= \begin{bmatrix} A_1 & A_{12} \\ A_{21} & A_2 \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, \\ \tilde{C} &= [C_1 \quad C_2], \quad \tilde{D} = D, \end{aligned} \quad (5.110)$$

unde A_1 și B_1 au t linii iar A_1 și C_1 au t coloane.

I. Prima modalitate de reducere corespunde *trunchierii* directe a lui \tilde{S} , conform relației

$$S_t = (A_1, B_1, C_1, D_1), \quad (5.111)$$

prin urmare, matricele sistemului S_t , dispuse sugestiv sub forma

$$S_t = \begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix}, \quad (5.112)$$

rezultă prin inspecție din (5.110).

II. A doua modalitate de reducere, preferențial utilizată în practică, corespunde unei *trunchieri compensate* în scopul conservării transferului în regim staționar constant al sistemelor S și \tilde{S} , i.e. al matricei de amplificare

$$K \stackrel{\text{def}}{=} T(0) = D + (sI - A)^{-1} B \Big|_{s=0}. \quad (5.113)$$

Utilizând idei specifice metodei perturbațiilor singulare, scriem ecuațiile lui \tilde{S} sub forma

$$\begin{aligned} \dot{x}_1 &= A_1 x_1 + A_{12} x_2 + B_1 u \\ \dot{x}_2 &= A_{21} x_1 + A_2 x_2 + B_2 u \\ y &= C_1 x_1 + C_2 x_2 + D u \end{aligned} \quad (5.114)$$

și presupunând $x_2 = \text{const}$, i.e. $\dot{x}_2 = 0$, din a doua ecuație găsim

$$x_2 = -A_2^{-1}(A_{21}x_1 + B_2u).$$

Introducând această expresie în celelalte două ecuații, obținem

$$\begin{aligned} \dot{x}_1 &= (A_1 - A_{12}A_2^{-1}A_{21})x_1 + (B_1 - A_{12}A_2^{-1}B_2)u \\ y &= (C_1 - C_2A_2^{-1}A_{21})x_1 + (D - C_2A_2^{-1}B_2)u \end{aligned}$$

deci trunchierea compensată S_t a lui S este

$$S_t = (A_1 - A_{12}A_2^{-1}A_{21}, B_1 - A_{12}A_2^{-1}B_2, C_1 - C_2A_2^{-1}A_{21}, D - C_2A_2^{-1}B_2). \quad (5.115)$$

Prin urmare, matricele lui S_t se calculează acum utilizând (5.112) și efectuând actualizarea

$$S_t \leftarrow S_t - \begin{bmatrix} A_{12} \\ C_2 \end{bmatrix} A_2^{-1} \begin{bmatrix} A_{21} & B_2 \end{bmatrix}. \quad (5.116)$$

În cazul sistemelor discrete (presupuse de asemenea stabile, controlabile și observabile) raționamentele și rezultatele de mai sus rămân valabile în cea mai mare parte. În particular, algoritmul 5.11 se păstrează ca atare, înlocuind peste tot EMAL cu DEMAL corespunzătoare. De asemenea, pentru a construi trunchierea compensată discretă, care conservă matricea de amplificare

$$K \stackrel{\text{def}}{=} T(1) = D + (zI - A)^{-1} B \Big|_{z=1}, \quad (5.117)$$

scriem ecuațiile lui \tilde{S} sub forma

$$\begin{aligned} x_1(+1) &= A_1 x_1 + A_{12} x_2 + B_1 u \\ x_2(+1) &= A_{21} x_1 + A_2 x_2 + B_2 u \\ y &= C_1 x_1 + C_2 x_2 + D u \end{aligned} \quad (5.118)$$

și presupunând $x_2 = \text{const}$, i.e. $x_2(+1) = x_2$, din a doua ecuație găsim

$$x_2 = (I - A_2)^{-1}(A_{21}x_1 + B_1u)$$

Prin urmare, trunchierea compensată a unui sistem discret este definită tot de (5.115) sau (5.116), în care matricea $(sI - A_2)|_{s=0} = -A_2$ se înlocuiește cu $(zI - A_2)|_{z=1} = I - A_2$. Redactarea algoritmilor este propusă cititorului sub formă de exerciții.

Programe MATLAB disponibile

Testarea stabilității se face conform algoritmului 5.1, i.e. calculând valorile proprii ale matricei A a sistemului cu funcția **eig**. Pentru obținerea formei Schur ordonate, necesară în algoritmul 5.2, se utilizează funcția **schord**. Construcția matricei de controlabilitate R a perechii (A, B) se face cu funcția **ctrb**, care implementează algoritmul 5.4, mai precis versiunea acestuia bazată pe formula (5.38). Pentru construcția formei bloc-superior Hessenberg a unei perechi (A, B) este disponibilă funcția **ctrbf**, care implementează algoritmul 5.7 cu blocurile permutate în ordine inversă. Construcția unei realizări minimale (vezi algoritmul 5.8) se poate face utilizând funcția **minreal**. Pentru calculul gramienilor de controlabilitate și observabilitate sunt disponibile funcțiile **gram** și **dgram**, care apelează, evident, **lyap** și **dlyap**. Echilibrarea unei realizări stabile și minimale se efectuează cu **balreal**, care implementează o versiune neoptimizată a algoritmului 5.11. (În particular, valorile singulare nu rezultă ordonate ca în (5.107)). Trunchierea compensată a realizării echilibrate furnizate de **balreal** se face interactiv, pe bază opțiunii utilizatorului ce apelează funcția corespunzătoare **modred**.

Exerciții

E 5.1 Conform teoremei lui Liapunov, matricea A a unui sistem continuu (discret) este stabilă dacă și numai dacă oricare ar fi matricea $Q = Q^T > 0$ ecuația matriceală algebrică Liapunov (discretă)

$$A^T X + X A = -Q \quad (A^T X A - X = -Q)$$

are o soluție $X = X^T > 0$. Explicați de ce enunțul de mai sus (de mare valoare teoretică) nu conduce la un test de stabilitate numeric eficient în raport cu algoritmul 5.1.

E 5.2 Un polinom $p(s) = a_0s^n + a_1s^{n-1} + \dots + a_n$ este stabil (sau "hurwitzian") dacă toate rădăcinile sale au partea reală strict negativă. Scrieți algoritmul de testare a stabilității polinomului $p(s)$ utilizând criteriul Routh-Hurwitz.

Indicație. Se utilizează algoritmul de eliminare gaussiană *fără pivotare* relativ la matricea Hurwitz asociată lui $p(s)$, ceea ce corespunde exact formulării propuse de Routh pentru criteriul în discuție.

Observație. Întrucât în absența unei strategii de pivotare adecvate, algoritmul de eliminare gaussiană este *numeric nestabil* (aparitia pivotilor nuli corespunde așa-numitelor cazuri "critice" ale algoritmului Routh), în practică testarea stabilității polinomului $p(s)$ se face construind matricea companion și aplicând algoritmul 5.1.

E 5.3 Formulați și rezolvați versiunea discretă a exercițiului 5.2.

E 5.4 Matricea A a unui sistem continuu se numește *admisibilă* dacă $\operatorname{Re}\lambda(A) < -\eta$, unde $\eta > 0$ este un număr real dat. Formulați testul corespunzător de admisibilitate. (Algoritmul 5.1 corespunde evident cazului $\eta = 0$).

E 5.5 Formulați și rezolvați versiunea discretă a exercițiului 5.4.

E 5.6 Utilizând algoritmul **QZ** în locul algoritmului **QR**, scrieți algoritmul de testare a stabilității sistemului generalizat (de tip descriptor), descris prin

$$\begin{aligned} E\dot{x} &= Ax + Bu, \\ y &= Cx, \end{aligned}$$

unde matricea $E \in \mathcal{R}^{n \times n}$ este nesingulară dar poate fi rău condiționată la inversare. (Prin urmare, trebuie evitată formarea matricelor $A_0 = E^{-1}A$ etc).

Indicație. Valorile proprii ale matricei A_0 coincid cu valorile proprii generalizate ale fascicolului $\lambda E - A$.

E 5.7 Scrieți algoritmul de calcul a unei baze ortogonale a subspațiului A -invariant total nestabil $\bar{\mathcal{X}}^+(A)$.

E 5.8 Utilizând rezultatele furnizate de algoritmul 3 (vezi observația 5.4), indicați o procedură de calcul a unei baze ortogonale a lui $\bar{\mathcal{X}}^+(A)$.

E 5.9 Scrieți versiunile discrete ale algoritmilor 5.2 și 5.3 (vezi observația 5.2).

E 5.10 Presupunând dată o funcție de transfer $T(s)$ într-una din formele

$$T(s) = \delta \frac{\prod (s - \mu_j)}{\prod (s - \lambda_i)}, \quad \text{respectiv} \quad T(s) = \sum \frac{r_i}{s - \lambda_i},$$

scrieți o realizare de stare $S = (A, b, c^T)$ în forma (5.11) și indicați baze (ortogonale) pentru subspațiile $\mathcal{X}^-(A)$ și $\mathcal{X}^+(A)$. Cum trebuie procedat dacă $T(s) = N(s)/p(s)$ este dată ca raport a două polinoame ?

E 5.11 Formulați și rezolvați versiunea discretă a exercițiului 5.10.

E 5.12 Presupunând că matricea A este simplă, scrieți algoritmul de descompunere spectrală bazat pe calculul valorilor și vectorilor proprii ai matricei A și comparați-l cu algoritmul 5.3. Care este preferabil și de ce ?

E 5.13 Utilizând algoritmul **QZ** în locul algoritmului **QR**, analizați posibilitățile de extindere a algoritmilor 5.2 și 5.3 la cazul sistemelor generalizate (de tip descriptor, vezi exercițiul 5.6).

E 5.14 Formulați și rezolvați exercițiul 5.12 în cazul sistemelor generalizate.

E 5.15 Discutați controlabilitatea și observabilitatea unui sistem $S = (A, b, c^T)$ scris în forma standard controlabilă. Idem, pentru forma standard observabilă.

E 5.16 Se dă o funcție de transfer discretă $T(z) = N(z)/p(z)$. Să se scrie o realizare de stare $S = (A, b, c^T)$ astfel încât problema (5.26) a determinării șirului de comenzi de efort minim să fie cât mai simplă. Scrieți algoritmul de calcul corespunzător.

E 5.17 Se dă o funcție de tranfer discretă $T(z) = N(z)/p(z)$. Să se scrie o realizare de stare $S = (A, b, c^T)$ astfel încât problema (5.36) a determinării unei estimări optimale a stării inițiale să fie cât mai simplă. Scrieți algoritmul de calcul corespunzător.

E 5.18 Analizați posibilitățile de extindere a soluțiilor exercițiilor 5.15 și 5.16 în cazul unei matrice de transfer $T(z) \in \mathcal{R}^{l \times m}(z)$. Considerați separat cazurile $m = 1, l = 1$ și $m, l > 1$.

E 5.19 Scrieți o realizare de stare $S = (A, b, c^T)$ a funcției de transfer $T(z) = N(z)/p(z)$ astfel încât $R(A, b) = I$. Idem, astfel încât $Q(c^T, A) = I$. În ce condiții putem avea simultan $R(A, b) = Q(c^T, A) = I$?

E 5.20 Formulați analogul algoritmilor 5.4 și 5.5 în cazul observabilității (fără a recurge la schema de dualizare din §4).

E 5.21 Formulați și discutați realizarea numerică a testelor elementare de observabilitate, bazate pe duala propoziției 5.1.

E 5.22 Se dau perechea observabilă (C, A) și un vector $\mathbf{y} \in \mathcal{R}^{nl}$. Scrieți algoritmul de calcul al pseudosoluției $x^* = x^*(0)$ a sistemului de ecuații (5.36), utilizând transformări ortogonale.

E 5.23 Se dau perechea controlabilă (A, B) și un vector $x = x(n) \in \mathcal{R}^n$. Scrieți algoritmul de calcul al soluției normale \mathbf{u}^* a sistemului de ecuații (5.26), utilizând transformări ortogonale.

E 5.24 Obțineți condiții suficiente simple de controlabilitate și observabilitate a conexiunilor serie, paralel și în circuit închis. Ce puteți afirma despre stabilitatea acestor conexiuni ?

E 5.25 Extindeți propoziția 5.1 precum și duala sa în cazul sistemelor generalizate (vezi exercițiul 5.6). Discutați realizarea numerică a testelor corespunzătoare.

E 5.26 Se dă un triplet $S = (A, B, C)$. Indicați (cel puțin două) proceduri de calcul a unei baze ortogonale pentru *subspațiul controlabil* $\mathcal{R} \stackrel{\text{def}}{=} \text{Im}R$ al perechii (A, B) . Idem, pentru *subspațiul neobservabil* $\mathcal{N} \stackrel{\text{def}}{=} \text{Ker}Q$.

E 5.27 Formulați problemele de calcul (5.26) și (5.36) în cazul unui sistem S nu neapărat controlabil și observabil.

E 5.28 Scrieți programul MATLAB care implementează algoritmul 5.6 și testați-l considerând perechi (A, b) adecvat alese.

Indicație. Scrieți o funcție MATLAB care generează o pereche (\tilde{A}, \tilde{b}) de ordin n dat, nu neapărat controlabilă, utilizând forma (5.47), în care perechea controlabilă (A_R, b_R) de ordin $r \leq n$ este în FSC, iar $A_{R\bar{R}}$ și $A_{\bar{R}}$ sunt generate aleator. Perechea dorită (A, b) se obține aplicând lui (\tilde{A}, \tilde{b}) o transformare de asemănare oarecare, i.e. $(A, b) = (S\tilde{A}S^{-1}, S\tilde{b})$, unde matricea nesingulară S poate fi, de asemenea, aleasă aleator.

E 5.29 Utilizând programul MATLAB ce implementează algoritmul 5.4, bazat pe construcția matricei de controlabilitate, ilustrați fenomenul de necontrolabilitate numerică a perechilor (A, b) cu o singură intrare, urmărind variația numărului de condiționare $\text{cond}(R)$ în raport cu creșterea

ordinului perechii. Comentați comparativ rezultatele numerice obținute utilizând procedura din exercițiul 5.28.

E 5.30 Scrieți funcțiile MATLAB ce implementează procedura de deflație controlabilă descrisă în propoziția 5.4 utilizând **DVS**, procedura de triangularizare ortogonală cu pivotarea coloanelor și, respectiv, procedura de eliminare gaussiană cu pivotare completă.

E 5.31 Scrieți programul MATLAB care implementează algoritmul 5.7 și testați-l considerând perechi (A, B) adecvat alese.

Indicație. Generarea perechilor controlabile (A_R, B_R) cu $m \geq 1$ are la bază utilizarea FSC cu intrări decuplate descrise în capitolul 3, § 3. Indicele de controlabilitate ν al perechii astfel obținute coincide cu cel mai mare indice de controlabilitate al perechilor diagonale din FSC.

E 5.32 Considerați un sistem $S = (A, B, C)$, nu neapărat controlabil sau observabil, având perechea (A, B) de forma (5.65), unde (A_R, B_R) este controlabilă și fie $C_R = C(:, 1:r)$. Arătați că S și $S_R \stackrel{\text{def}}{=} (A_R, B_R, C_R)$ au aceeași matrice de transfer, mai precis avem

$$T(s) \stackrel{\text{def}}{=} C(sI - A)^{-1}B = C_R(sI - A_R)^{-1}B_R. \quad (5.119)$$

Egalitatea (5.119) arată că polii fiși $\lambda(A_R)$ corespunzători părții necontrolabile a perechii (A, B) sunt *poli simplificabili* ai lui $T(s)$. Prin dualitate, aceeași proprietate aparține polilor fiși $\lambda(A_Q)$ corespunzători părții neobservabile a perechii (C, A) .

E 5.33 Arătați că în cazul SISO, sistemul $S = (A, b, c^T)$ este o realizare minimală a funcției de transfer ireductibile $T(s) = N(s)/p(s)$ dacă și numai dacă ordinul matricei A coincide cu gradul numitorului $p(s)$. Ce se întâmplă în cazul sistemelor cu mai multe intrări sau/și cu mai multe ieșiri ?

E 5.34 O matrice de transfer $T(s)$ se numește stabilă dacă polii tuturor elementelor lui $T(s)$ sunt stabili, i.e. au partea reală strict negativă, respectiv au modulul strict subunitar. Arătați că $T(s)$ este stabilă dacă și numai dacă orice realizare minimală $S_m = (A_m, B_m, C_m)$ a lui $T(s)$ este stabilă. Scrieți algoritmul de testare a stabilității unei matrice de transfer $T(s)$ date.

E 5.35 Enumerați câteva proprietăți semnificative ale realizărilor stabilizabile sau/și detectabile.

E 5.36 Indicați câteva metode de calcul a unei baze a *subspațiului stabilizabil* $\mathcal{S} = \mathcal{R} + \mathcal{X}^-(A)$. Idem, pentru *subspațiul nedetectabil* $\mathcal{T} = \mathcal{N} \cap \mathcal{X}^+(A)$.

Indicație. În cazul în care coloanele matricelor $S_k \in \mathcal{R}^{n \times n_k}$ constituie sisteme de generatori pentru subspațiile $\mathcal{X}_k \subset \mathcal{R}^n$, adică $\mathcal{X}_k = \text{Im } S_k$, $k = 1 : 2$, atunci coloanele matricei bloc

$$S = [S_1 \ S_2]$$

constituie un sistem de generatori ai subspațiului $\mathcal{X}_1 + \mathcal{X}_2$. Prin dualitate, dacă liniile matricelor $T_k \in \mathcal{R}^{n_k \times n}$ constituie un sistem de anulatori pentru subspațiile $\mathcal{X}_k \subset \mathcal{R}^n$, adică $\mathcal{X}_k = \text{Ker } T_k$, $k = 1 : 2$, atunci liniile matricei

$$T = \begin{bmatrix} T_1 \\ T_2 \end{bmatrix}$$

constituie un sistem de anulatori ai subspațiului $\mathcal{X} = \mathcal{X}_1 \cap \mathcal{X}_2$. În ambele cazuri, pentru a obține o bază a lui \mathcal{X} , se utilizează **DVS**.

E 5.37 Analizați posibilitatea combinării algoritmilor de dublare și de rădăcină pătrată din §9, utilizând schema "candidat"

$$\begin{bmatrix} L \\ 0 \end{bmatrix} \leftarrow U \begin{bmatrix} L \\ LA \end{bmatrix}, \quad A \leftarrow A^2.$$

E 5.38 Formulați versiunile algoritmilor 5.9 și 5.10 corespunzătoare rezolvării DEMAL (5.85) fără a recurge la schema de dualizare expusă în finalul §4.

E 5.39 Arătați că dacă X, Y sunt soluțiile ecuațiilor (5.87), (5.88) (respectiv (5.85), (5.86) în cazul discret) atunci avem $\text{tr}(B^T X B) = \text{tr}(C^T Y C)$, unde $\text{tr}(A)$ este *urma* matricei pătrate A .

Numărul $\delta \geq 0$ introdus prin egalitatea de mai sus coincide cu norma H_2 a funcției de transfer $T(s)$, definită prin

$$\|T\|_2 = \frac{1}{2\pi} \int_{-\infty}^{\infty} \text{tr} T^H(i\omega) T(i\omega) d\omega.$$

(În cazul discret integrala se efectuează pe cercul unitate).

Bibliografie

- [1] **Ionescu V., Popeea C.** CONDOCEREA STRUCTURALĂ A SISTEMELOR LINIARE, ET București, 1986.
- [2] **Ionescu V.** INTRODUCERE ÎN TEORIA STRUCTURALĂ A SISTEMELOR LINIARE, EA București, 1975.
- [3] **Aplevich J.D.** *Direct Computation of Canonical Forms for Linear Systems by Elementary Matrix Operations*, IEEE Trans. AC - 19, No. 2, pp. 124 – 126, 1974.
- [4] **Aplevich J.D.** *Tableau Methods for Analysis and Design of Linear Systems*, Automatica, Vol. 15, pp. 419 – 429, 1979.
- [5] **Varga A.** *Numerically Reliable Algorithm to Test Controllability*, Electron. Lett. 15, pp. 454 – 455, 1981.
- [6] **Varga A.** *Numerically Stable Algorithm for Standard Controllability Form*, Electron. Lett. 17, pp. 74 – 75, 1981.
- [7] **Klema V., Laub A.J.** *The Singular Value Decomposition: Its Computation and Some Applications*, IEEE Trans. AC - 25, No. 2, pp. 164 – 176, 1980.
- [8] **Moore B.C.** *Principal Component Analysis in Linear Systems Controllability, Observability and Model Reduction*, IEEE Trans. AC - 26, No. 2, pp. 17 – 32, 1981.
- [9] **Laub A.J., Heath M.T., Paige C.C., Ward R.C.** *Computation of System Balancing Transformations*, IEEE Trans. AC - 32, No. 2, 1987.
- [10] **Uhlig F.** *Simultaneous Block - Diagonalization of Two Real Symmetric Matrices*. Linear Algebra and Its Applications, Vol. 7, pp. 281 – 289, 1973.
- [11] **Heath M.T., Laub A.J., Paige C.C., Ward R.C.** *Computing SVD of a Product of Two Matrices*, SIAM J. Sci. Stat. Comp., Vol. 7, pp. 1147 – 1159, 1986.
- [12] **Hammarling S.J.** *Numerical Solution of the Stable, Non-negative Definite Lyapunov Equations*, IMA J. numer. Anal., Vol. 2, pp. 303 – 323, 1982.