

**Bor Plestenjak**

## **Numerične metode 2 (finančna matematika)**

**delovna verzija**

verzija: 13. julij 2010

# Kazalo

<b>1 Nesimetrični problem lastnih vrednosti</b>	<b>3</b>
1.1 Uvod . . . . .	3
1.2 Schurova forma . . . . .	4
1.3 Teorija motenj . . . . .	7
1.4 Potenčna metoda . . . . .	11
1.5 Obratna napaka in izračunljive ocene . . . . .	14
1.6 Inverzna iteracija . . . . .	15
1.7 Ortogonalna iteracija . . . . .	16
1.8 QR iteracija . . . . .	17
1.8.1 Redukcija na Hessenbergovo obliko . . . . .	18
1.8.2 Premiki . . . . .	20
1.9 Implicitna QR metoda . . . . .	22
<b>2 Simetrični problem lastnih vrednosti</b>	<b>26</b>
2.1 Uvod . . . . .	26
2.2 Rayleighova iteracija . . . . .	30
2.3 QR iteracija za simetrični lastni problem . . . . .	32
2.4 Sturmovo zaporedje . . . . .	33
2.5 Deli in vladaj . . . . .	36
2.6 Jacobijeva metoda . . . . .	39
2.7 Relativno robustne reprezentacije . . . . .	41
2.7.1 Zasukani razcep . . . . .	42
2.7.2 Relativno robustna reprezentacija matrike . . . . .	44
2.7.3 MRRR . . . . .	46
<b>3 Posplošitve problema lastnih vrednosti</b>	<b>48</b>
3.1 Posplošeni problem lastnih vrednosti . . . . .	48
3.2 QZ algoritem . . . . .	50
3.3 Definiten matrični šop . . . . .	53
3.4 Nelinearni problem lastnih vrednosti . . . . .	53
3.4.1 Newtonova metoda in inverzna iteracija . . . . .	55
3.4.2 Zaporedne linearne aproksimacije . . . . .	56
3.5 Polinomski (kvadratni) problem lastnih vrednosti . . . . .	57
3.5.1 Hermitski kvadratni problem lastnih vrednosti . . . . .	58
3.6 Računanje singularnega razcepa . . . . .	60
3.7 QR iteracija za računanje singularnega razcepa . . . . .	63
3.8 dqds metoda . . . . .	64
3.9 Jacobijeva metoda za singularni razcep . . . . .	66
3.9.1 Lastni razcep simetrične pozitivno definitne matrike . . . . .	71

<b>4</b>	<b>Enakomerna aproksimacija</b>	<b>73</b>
4.1	Aproksimacija . . . . .	73
4.2	Enakomerna aproksimacija s polinomi . . . . .	74
4.3	Ekonomizacija Čebiševa . . . . .	76
<b>5</b>	<b>Interpolacija</b>	<b>81</b>
5.1	Uvod . . . . .	81
5.2	Interpolacijski polinom . . . . .	81
5.3	Deljene difference . . . . .	84
5.4	Končne difference . . . . .	88
5.5	Interpolacija s kosoma polinomskimi funkcijami . . . . .	90
5.6	Beziérove krivulje . . . . .	95
5.7	Interpolacija v Matlabu . . . . .	96
<b>6</b>	<b>Numerično odvajanje</b>	<b>98</b>
6.1	Uvod . . . . .	98
6.2	Drugi načini izpeljave . . . . .	99
6.3	Celotna napaka . . . . .	100
<b>7</b>	<b>Numerično integriranje</b>	<b>103</b>
7.1	Kvadraturene formule . . . . .	103
7.2	Newton–Cotesova pravila . . . . .	104
7.3	Neodstranljiva napaka . . . . .	106
7.4	Sestavljene formule . . . . .	107
7.5	Peanov izrek . . . . .	108
7.6	Richardsonova ekstrapolacija . . . . .	110
7.7	Adaptivne metode . . . . .	110
7.8	Rombergova metoda . . . . .	112
7.9	Gaussove kvadraturene formule . . . . .	114
7.10	Izlimitirani integrali . . . . .	117
7.11	Večdimenzionalni integrali . . . . .	119
7.12	Metoda Monte Carlo . . . . .	120
7.13	Numerično integriranje v Matlabu . . . . .	121
7.14	Numerično seštevanje vrst . . . . .	122
<b>8</b>	<b>Diferencialne enačbe</b>	<b>124</b>
8.1	Uvod . . . . .	124
8.2	Obstoj rešitve, občutljivost in stabilnost . . . . .	126
8.2.1	Občutljivost rešitve začetnega problema . . . . .	126
8.2.2	Stabilnost rešitve začetnega problema . . . . .	127
8.3	Enokoračne metode . . . . .	127
8.3.1	Eulerjeva metoda . . . . .	127
8.3.2	Taylorjeva vrsta . . . . .	129
8.3.3	Runge-Kutta metode . . . . .	131
8.3.4	Adaptivna ocena koraka . . . . .	132
8.4	Stabilnost in konvergenca enokoračnih metod . . . . .	134
8.5	Večkoračne metode . . . . .	136
8.6	Inherentna in inducirana nestabilnost . . . . .	138
8.7	Začetni problemi drugega reda . . . . .	140
8.8	Reševanje začetnih diferencialnih enačb v Matlabu . . . . .	141

8.9	Implicitno podane diferencialne enačbe . . . . .	144
8.10	Metoda zveznega nadaljevanja . . . . .	144
8.11	Robni problemi drugega reda . . . . .	145
8.11.1	Linearni robni problem . . . . .	146
8.11.2	Nelinearni robni problem . . . . .	148
8.11.3	Prevedba na variacijski problem . . . . .	149
8.12	Reševanje robnih problemov v Matlabu . . . . .	150

# Poglavje 1

## Nesimetrični problem lastnih vrednosti

### 1.1 Uvod

Dana je matrika  $A \in \mathbb{R}^{n \times n}$ . Če neničelen vektor  $x \in \mathbb{C}^n$  in skalar  $\lambda \in \mathbb{C}$  zadoščata enačbi  $Ax = \lambda x$ , potem pravimo, da je  $\lambda$  *lastna vrednost*,  $x$  pa (*desni*) *lastni vektor* matrike  $A$ . Neničelen vektor  $y \in \mathbb{C}^n$ , pri katerem je  $y^H A = \lambda y^H$ , je *levi lastni vektor* matrike  $A$ . Levi in desni lastni vektorji, ki pripadajo različnim lastnim vrednostim, so med seboj ortogonalni, saj velja naslednja lema.

**Lema 1.1** *Naj bosta  $\lambda$  in  $\mu$  različni lastni vrednosti matrike  $A$ . Če je  $x$  desni lastni vektor, ki pripada  $\lambda$ ,  $y$  pa levi lastni vektor, ki pripada  $\mu$ , potem sta  $x$  in  $y$  ortogonalna.*

*Dokaz.* Enakost  $Ax = \lambda x$  pomnožimo z leve z  $y^H$ , enakost  $y^H A = \mu y^H$  pa z desne z  $x$ . Dobimo

$$y^H Ax = \lambda y^H x = \mu y^H x,$$

to pa je zaradi  $\lambda \neq \mu$  lahko izpolnjeno le v primeru, ko je  $y^H x = 0$ . ■

V posebnem primeru, ko je matrika  $A$  simetrična, iz zgornje leme sledi, da so lastni vektorji, ki pripadajo različnim lastnim vrednostim, paroma ortogonalni. To je posledica tega, da so za simetrično matriko desni lastni vektorji enaki levim.

Pri algoritmih za računanje lastnih vektorjev zadošča, da obravnavamo le desne lastne vektorje. Namreč, če je  $y$  levi lastni vektor matrike  $A$  za lastno vrednost  $\lambda$ , potem je  $y$  desni lastni vektor matrike  $A^H$  za lastno vrednost  $\bar{\lambda}$ .

Pravimo, da se da matriko  $A$  diagonalizirati, če obstajata nesingularna matrika  $X = [x_1 \cdots x_n]$  in diagonalna matrika  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ , da je  $A = X\Lambda X^{-1}$ . V tem primeru je  $Ax_i = \lambda_i x_i$  za  $i = 1, \dots, n$ , kar pomeni, da so stolpci matrike  $X$  lastni vektorji, diagonalna matrike  $\Lambda$  pa je sestavljena iz lastnih vrednosti.

Če je  $S$  nesingularna matrika, potem pravimo, da sta matriki  $A$  in  $B = S^{-1}AS$  *podobni*. V tem primeru imata  $A$  in  $B$  iste lastne vrednosti. Velja, da je  $x$  desni lastni vektor za matriko  $A$  natanko tedaj, ko je  $S^{-1}x$  desni lastni vektor za matriko  $B$ .

Lastne vrednosti matrike  $A$  so ničle karakterističnega polinoma  $p(\lambda) = \det(A - \lambda I)$ . Vsaka  $n \times n$  matrika  $A$  ima tako  $n$  lastnih vrednosti  $\lambda_1, \dots, \lambda_n$ , pri čemer večkratne ničle štejemo

večkrat. Če je  $\lambda$  lastna vrednost matrike  $A$ , potem je njena *algebraična večkratnost* enaka večkratnosti  $\lambda$  kot ničle karakterističnega polinoma matrike  $A$ , *geometrijska večkratnost* pa je enaka dimenziji podprostora  $\ker(A - \lambda I)$  in je vedno manjša ali enaka algebraični večkratnosti. Če je algebraična večkratnost lastne vrednosti ena, potem pravimo, da je lastna vrednost enostavna,

Ker se ničel splošnega polinoma stopnje pet ali več ne da izračunati drugače kot numerično, to velja tudi za lastne vrednosti in direktne metode za računanje lastnih vrednosti ne obstajajo. Lastne vrednosti vedno računamo z iterativnimi postopki.

Računanje lastnih vrednosti preko ničel eksplicitno izračunanega karakterističnega polinoma v splošnem ni priporočljivo. Prva težava je, da algoritmi za računanje koeficientov karakterističnega polinoma niso stabilni. Če to združimo z dejstvom, da so ničle polinoma lahko zelo občutljive na motnje koeficientov, kot kaže npr. Wilkinsonov primer. iz zgloda ??.

**Zgled 1.1** Pri eksplicitni uporabi karakterističnega polinoma lahko izgubimo natančnost. Naj bo

$$A = \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1 \end{bmatrix}$$

za  $\epsilon = \sqrt{u}/2$ , kjer je  $u$  osnovna zaokrožitvena napaka. Če izračunamo karakteristični polinom, potem zaradi zaokroževanja dobimo  $\det(A - \lambda I) = \lambda^2 - 2\lambda + 1$ . Na ta način bi izračunali dvojno lastno vrednost 1, pravi lastni vrednosti matrike  $A$  pa sta  $1 - \epsilon$  in  $1 + \epsilon$ .  $\square$

Računanje lastnih vrednosti preko eksplicitno izračunanega karakterističnega polinoma torej ni numerično stabilno. V nadaljevanju bomo spoznali več različnih stabilnih algoritmov za izračun lastnih vrednosti in vektorjev.

## 1.2 Schurova forma

Vemo, da za vsako  $n \times n$  matriko  $A$  obstajata taka nesingularna matrika  $X$  in bločno diagonalna matrika  $J = \text{diag}(J_1, \dots, J_k)$ , ki ji pravimo *Jordanova forma*, da je  $X^{-1}AX = J$ , kjer je

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}$$

*Jordanova kletka* velikosti  $m_i \times m_i$  za  $i = 1, \dots, k$  in  $n = m_1 + \dots + m_k$ .

Jordanova forma pove veliko o matriki  $A$ . Iz nje lahko npr. preberemo vse lastne vrednosti matrike in njihove algebraične in geometrične večkratnosti. Zelo pomembna je tudi za računanje vrednosti funkcij matrik. Tako ima npr. rešitev sistema diferencialnih enačb s konstantnimi koeficienti  $y'(t) = Ay(t)$  obliko  $y(t) = y(t_0)e^{A(t-t_0)}$ . Za rešitev je torej potrebno znati izračunati matriko  $e^A$ , ki je definirana z razvojem v vrsto

$$e^A = \sum_{j=0}^{\infty} \frac{1}{j!} A^j.$$

Če je  $f$  definirana in dovoljkrat zvezno odvedljiva funkcija v lastnih vrednostih matrike  $A$ , potem s pomočjo Jordanove forme  $A = XJX^{-1}$  definiramo  $f(A) = Xf(J)X^{-1}$ , kjer je  $f(J) = \text{diag}(f(J_1), \dots, f(J_k))$  in

$$f(J_i) = \begin{bmatrix} f(\lambda_i) & f'(\lambda_i) & \cdots & \frac{f^{(m_i-1)}(\lambda_i)}{(m_i-1)!} \\ & f(\lambda_i) & \ddots & \vdots \\ & & \ddots & f'(\lambda_i) \\ & & & f(\lambda_i) \end{bmatrix}$$

za  $i = 1, \dots, k$ .

Na žalost je Jordanova forma zelo občutljiva in neprimerna za numerično računanje. Ker ni zvezna funkcija elementov matrike  $A$ , jo lahko v primeru večkratnih lastnih vrednosti majhne motnje popolnoma spremenijo. Zato tudi računanje funkcij matrik preko Jordanove forme ni stabilno.

### Zgled 1.2 Matrika

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

je že kar v Jordanovi formi z eno samo kletko  $n \times n$ , za poljuben  $\epsilon > 0$  pa ima Jordanova forma matrike

$$A(\epsilon) = \begin{bmatrix} 0 & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ \epsilon & & & 0 \end{bmatrix}$$

$n$  kletk velikosti  $1 \times 1$  z lastnimi vrednostmi  $\sqrt[n]{\epsilon}$ . □

Računanje Jordanove forme tudi ni obratno stabilno. Denimo, da smo za  $A$  numerično izračunali  $\tilde{X}$  in  $\tilde{J}$ . Sedaj nas zanima, ali je  $\tilde{X}^{-1}(A + \delta A)\tilde{X} = \tilde{J}$  za neko matriko  $\delta A$ , ki je blizu 0. Pri tem predpostavimo, da je  $\tilde{X}$  točna, za  $\tilde{J}$  pa velja  $\tilde{J} = J + \delta J$ . Iz  $X^{-1}(A + \delta A)X = J + \delta J$  sledi  $X^{-1}\delta AX = \delta J$ , od tod pa lahko ocenimo

$$\|\delta A\| \leq \|X\| \cdot \|X^{-1}\| \cdot \|\delta J\| = \kappa(X)\|\delta J\|.$$

Ker je v splošnem občutljivost  $\kappa(X)$  lahko poljubno velika, računanje Jordanove forme ni obratno stabilno.

Za stabilnost bi bilo bolje, če bi bila prehodna matrika  $X$  kar unitarna. Izkaže se, da z unitarno podobnostno transformacijo lahko na stabilen način matriko transformiramo v zgornjo trikotno matriko.

**Izrek 1.2** Za vsako matriko  $A$  obstajata unitarna matrika  $Q$  in zgornja trikotna matrika  $T$ , da je  $Q^H A Q = T$  (Schurova forma).

*Dokaz.* Naredimo indukcijo po  $n$ . Za  $n = 1$  izrek očitno velja, saj vzamemo  $Q = [1]$  in  $T = A$ .

Predpostavimo sedaj, da izrek velja za  $n - 1$  in ga dokažimo za  $n$ . Naj bo  $\lambda$  lastna vrednost matrike  $A$  in  $x$  njen normiran lastni vektor. Obstaja taka unitarna matrika  $U$ , da je  $Ue_1 = x$  (skonstruiramo jo lahko npr. kar s kompleksnim Householderjevim zrcaljenjem). Matrika  $B = U^H AU$  ima obliko

$$B = \begin{matrix} & & 1 & & n-1 \\ & & & & \\ & 1 & & & \\ n-1 & & \begin{bmatrix} \lambda & \times & \cdots & \times \\ & C & & \end{bmatrix} & \end{matrix},$$

saj je  $Be_1 = U^H AUe_1 = U^H Ax = U^H \lambda x = \lambda e_1$ . Po indukcijski predpostavki obstaja Schurova forma za  $(n - 1) \times (n - 1)$  matriko  $C$ . Torej obstaja taka unitarna matrika  $V_1$ , da je  $V_1^H C V_1 = T_1$  zgornja trikotna matrika. Sedaj je

$$\underbrace{\begin{bmatrix} 1 & 0 \\ 0 & V_1^H \end{bmatrix}}_{V^H} B \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & V_1 \end{bmatrix}}_V = \begin{bmatrix} \lambda & \times & \cdots & \times \\ 0 & & & T_1 \end{bmatrix}$$

in zgornja trikotna matrika  $\underbrace{V^H U^H A}_{Q^H} \underbrace{UV}_Q$  je Schurova forma matrike  $A$ . ■

Schurova forma ni enolična, saj je npr. vrstni red diagonalnih elementov  $\lambda_i$  poljuben.

V primeru realne matrike se lahko izognemo kompleksni aritmetiki, če dopustimo, da namesto zgornje trikotne matrike dobimo *kvazi zgornjo trikotno matriko*, ki ima na diagonalni lahko bloke  $2 \times 2$ . V teh diagonalnih blokkih  $2 \times 2$  so skriti konjugirani pari kompleksnih lastnih vrednosti.

**Izrek 1.3** *Za vsako realno matriko  $A$  obstajata ortogonalna matrika  $Q$  in kvazi zgornja trikotna matrika  $T$ , da je  $Q^T A Q = T$  (realna Schurova forma).*

*Dokaz.* Spet uporabimo indukcijo. Če je  $\lambda \in \mathbb{R}$ , potem lahko nadaljujemo tako kot pri dokazu izreka 1.2, zato predpostavimo, da velja  $\lambda \notin \mathbb{R}$ . Potem imamo poleg lastnega para  $(\lambda, x)$  tudi lastni par  $(\bar{\lambda}, \bar{x})$ . Če definiramo realna vektorja

$$x_R = \frac{1}{2}(x + \bar{x}), \quad x_I = \frac{1}{2i}(x - \bar{x}),$$

potem obstaja taka ortogonalna matrika

$$U = \begin{bmatrix} 2 & n-2 \\ U_1 & U_2 \end{bmatrix},$$

da je  $\text{Lin}(U_1) = \text{Lin}(\{x_R, x_I\}) = \text{Lin}(\{x, \bar{x}\})$ . Ker je  $\text{Lin}(U_1)$  invarianten podprostor za  $A$ , velja

$$U^T A U = \begin{matrix} & & 2 & & n-2 \\ & & & & \\ & 2 & & & \\ n-2 & & \begin{bmatrix} B & \times & \cdots & \times \\ & C & & \end{bmatrix} & \end{matrix},$$

kjer sta  $\lambda$  in  $\bar{\lambda}$  lastni vrednosti bloka  $B$ , preostale lastne vrednosti pa so v matriki  $C$ . Nadaljujemo podobno kot pri dokazu izreka 1.2. ■



### 1.3 Teorija motenj

Zanima nas, kaj se dogaja z lastnimi vrednostmi (in vektorji), ko zmotimo matriko. Vemo, da so lastne vrednosti zvezne funkcije elementov matrike, saj so ničle karakterističnega polinoma, ničle polinoma pa so zvezne funkcije koeficientov polinoma. Če se da matriko diagonalizirati, potem naslednji izrek pove, da je sprememba lastnih vrednosti omejena z občutljivostjo matrike lastnih vektorjev.

**Izrek 1.4 (Bauer–Fike)** Naj se da matriko  $A$  diagonalizirati kot  $A = X\Lambda X^{-1}$ , kjer je  $\Lambda$  diagonalna matrika lastnih vrednosti  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Potem vse lastne vrednosti matrike  $A + \epsilon E$  ležijo v uniji  $n$  krogov

$$K_i = \{z \in \mathbb{C} : |z - \lambda_i| \leq \epsilon \kappa(X) \|E\|\}, \quad i = 1, \dots, n,$$

kjer za normo lahko vzamemo katerokoli izmed norm  $\|\cdot\|_1$ ,  $\|\cdot\|_2$  ali  $\|\cdot\|_\infty$ .

*Dokaz.* Naj bo  $\lambda(\epsilon)$  lastna vrednost  $A + \epsilon E$ . Predpostavimo lahko, da se  $\lambda(\epsilon)$  razlikuje od vseh lastnih vrednosti  $\lambda_1, \dots, \lambda_n$ , saj sicer izrek očitno drži. Naj bo torej matrika  $A + \epsilon E - \lambda(\epsilon)I$  singularna. Zapišemo lahko

$$\begin{aligned} X^{-1}(A + \epsilon E - \lambda(\epsilon)I)X &= \Lambda - \lambda(\epsilon)I + \epsilon X^{-1}EX \\ &= (\Lambda - \lambda(\epsilon)I) \left( I + \epsilon (\Lambda - \lambda(\epsilon)I)^{-1} X^{-1}EX \right). \end{aligned}$$

Ker je  $\Lambda - \lambda(\epsilon)I$  nesingularna matrika, mora biti  $I + \epsilon (\Lambda - \lambda(\epsilon)I)^{-1} X^{-1}EX$  singularna matrika. To pomeni, (uporabimo lemo ??), da je

$$1 \leq \|\epsilon (\Lambda - \lambda(\epsilon)I)^{-1} X^{-1}EX\| \leq \epsilon \|(\Lambda - \lambda(\epsilon)I)^{-1}\| \|X^{-1}\| \|E\| \|X\|.$$

Iz

$$\|(\Lambda - \lambda(\epsilon)I)^{-1}\| = \frac{1}{\min_{i=1, \dots, n} |\lambda_i - \lambda(\epsilon)|}$$

sledi

$$\min_{i=1, \dots, n} |\lambda_i - \lambda(\epsilon)| \leq \epsilon \kappa(X) \|E\|. \quad \blacksquare$$

**Opomba 1.1** Če unija krogov razpade na povezane komponente, potem iz zveznosti lastnih vrednosti sledi, da vsaka komponenta vsebuje natanko toliko lastnih vrednosti, kolikor krogov jo sestavlja.

**Posledica 1.5** Če sta matriki  $A$  in  $E$  simetrični, potem pri predpostavkah izreka 1.4 za vsako lastno vrednost  $\lambda(\epsilon)$  matrike  $A + \epsilon E$  velja

$$\min_{i=1, \dots, n} |\lambda(\epsilon) - \lambda_i| \leq \epsilon \|E\|_2.$$

*Dokaz.* Za simetrično matriko velja, da je matrika lastnih vektorjev  $X$  ortogonalna matrika, torej  $\kappa_2(X) = 1$ . ■

Bauer–Fikeov izrek je ponavadi preveč splošen, saj za motnje vseh lastnih vrednosti uporabi isto zgornjo mejo. V resnici so lahko posamezne lastne vrednosti bolj občutljive od ostalih. Naslednji izrek pove, od česa je odvisna občutljivost enostavne lastne vrednosti.

**Izrek 1.6** Naj bo  $\lambda_i$  enostavna lastna vrednost matrike  $A$  z normiranima levim in desnim lastnim vektorjem  $y_i$  in  $x_i$ . Če je  $\lambda_i + \delta\lambda_i$  ustrezna lastna vrednost zmotene matrike  $A + \delta A$  in je  $x_i + \delta x_i$  pripadajoči lastni vektor, potem velja

$$\lambda_i + \delta\lambda_i = \lambda_i + \frac{y_i^H \delta A x_i}{y_i^H x_i} + \mathcal{O}(\|\delta A\|^2).$$

Z izrazom ustrezna lastna vrednost je mišljeno, da si izberemo tisto lastno vrednost zmotene matrike, za katero velja  $\lim_{\|\delta A\| \rightarrow 0} (\delta\lambda_i) = 0$ .

*Dokaz.* Velja  $Ax_i = \lambda_i x_i$  in  $(A + \delta A)(x_i + \delta x_i) = (\lambda_i + \delta\lambda_i)(x_i + \delta x_i)$ . Če zanemarimo člene (od tod na koncu dobimo  $\mathcal{O}(\|\delta A\|^2)$ ), ki vsebujejo produkte dveh majhnih popravkov in pomnožimo enačbo z leve z  $y_i^H$ , dobimo

$$\delta\lambda_i = \frac{y_i^H \delta A x_i}{y_i^H x_i}. \quad \blacksquare$$

**Definicija 1.7** Naj bo  $\lambda_i$  enostavna lastna vrednost,  $x_i$  in  $y_i$  pa pripadajoča desni in levi lastni vektor. Če definiramo

$$s_i := \frac{y_i^H x_i}{\|x_i\|_2 \|y_i\|_2},$$

potem je  $|s_i|^{-1}$  občutljivost enostavne lastne vrednosti  $\lambda_i$ . Če je  $\lambda_i$  večkratna lastna vrednost, je njena občutljivost neskončna.

Če je matrika  $A$  simetrična, potem imajo vse enostavne lastne vrednosti najmanjšo možno občutljivost 1, saj so levi lastni vektorji enaki desnim.

**Zgled 1.3** Matrika

$$A_1 = \begin{bmatrix} -149 & -50 & -154 \\ 537 & 180 & 546 \\ -27 & -9 & -25 \end{bmatrix}$$

ima eksaktne lastne vrednosti 1, 2, 3. V Matlabu z ukazom `eig(A)`, ki uporablja obratno stabilen algoritem, ter z računanjem v dvojni natančnosti, dobimo

$$\begin{aligned} \hat{\lambda}_1 &= 1.00000000010722 \\ \hat{\lambda}_2 &= 1.99999999991790 \\ \hat{\lambda}_3 &= 2.99999999997399. \end{aligned}$$

Matrika Godunova

$$A_2 = \begin{bmatrix} 289 & 2064 & 336 & 128 & 80 & 32 & 16 \\ 1152 & 30 & 1312 & 512 & 288 & 128 & 32 \\ -29 & -2000 & 756 & 384 & 1008 & 224 & 48 \\ 512 & 128 & 640 & 0 & 640 & 512 & 128 \\ 1053 & 2256 & -504 & -384 & -756 & 800 & 208 \\ -287 & -16 & 1712 & -128 & 1968 & -30 & 2032 \\ -2176 & -287 & -1565 & -512 & -541 & -1152 & -289 \end{bmatrix}$$

ima eksaktne lastne vrednosti  $-4, -2, -1, 0, 1, 2, 4$ . V tem primeru Matlab izračuna naslednje približke za lastne vrednosti:

$$\begin{aligned}\widehat{\lambda}_1 &= 5.0968 + 1.8932i \\ \widehat{\lambda}_2 &= 5.0968 - 1.8932i \\ \widehat{\lambda}_3 &= 1.2157 + 4.3784i \\ \widehat{\lambda}_4 &= 1.2157 - 4.3784i \\ \widehat{\lambda}_5 &= -5.6738 \\ \widehat{\lambda}_6 &= -3.4756 + 3.4463i \\ \widehat{\lambda}_7 &= -3.4756 - 3.4463i.\end{aligned}$$

V obeh primerih so izračunane lastne vrednosti eksaktne lastne vrednosti malo zmotene matrike, saj Matlab uporablja obratno stabilen algoritem. Pri matriki  $A_2$  pride do velikih razlik zaradi tega, ker so lastne vrednosti te matrike zelo občutljive. Namreč, njihove občutljivosti so velikostnega reda  $10^{13}$ . Če to primerjamo z matriko  $A_1$ , kjer imajo vse lastne vrednosti občutljivost velikostnega reda  $10^2$ , je očitno, da občutljivost lastne vrednosti pomembno vpliva na natančnost izračunanih približkov.  $\square$

Če se da matriko diagonalizirati, potem lahko govorimo tudi o občutljivosti lastnih vektorjev.

**Izrek 1.8** Naj bo  $A = X\Lambda X^{-1} = Y^H\Lambda Y^{-H}$ , kjer je  $X = [x_1 \cdots x_n]$  matrika normiranih desnih lastnih vektorjev,  $Y = [y_1 \cdots y_n]$  matrika normiranih levih lastnih vektorjev in  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  diagonalna matrika lastnih vrednosti. Če je  $\lambda_i$  enostavna lastna vrednost in  $\lambda_i + \delta\lambda_i$  ustrezna lastna vrednost zmotene matrike  $A + \delta A$  s pripadajočim lastnim vektorjem  $x_i + \delta x_i$ , potem velja

$$x_i + \delta x_i = x_i + \sum_{\substack{j=1 \\ j \neq i}}^n \frac{y_j^H \delta A x_i}{(\lambda_i - \lambda_j) s_j} x_j + \mathcal{O}(\|\delta A\|^2).$$

*Dokaz.* Zmoteni lastni vektor lahko izrazimo v bazi lastnih vektorjev matrike  $A$  kot

$$x_i + \delta x_i = x_i + \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_j x_j.$$

Če v enakost  $(A + \delta A)(x_i + \delta x_i) = (\lambda_i + \delta\lambda_i)(x_i + \delta x_i)$  vstavimo zgornji razvoj, jo pomnožimo z leve z  $y_k^H$ , kjer je  $k \neq i$ , in zanemarimo člene (od tod na koncu spet  $\mathcal{O}(\|\delta A\|^2)$ ), ki vsebujejo produkte dveh majhnih popravkov, dobimo

$$\alpha_k = \frac{y_k^H \delta A x_i}{(\lambda_i - \lambda_k) s_k}.$$

Pri tem smo upoštevali, da je levi lastni vektor  $y_k$  ortogonalen na vse desne lastne vektorje  $x_j$ , kjer je  $k \neq j$ .  $\blacksquare$

**Lema 1.9** Če je  $\lambda_i$  enostavna lastna vrednost, potem je  $s_i \neq 0$ .

*Dokaz.* Brez škode za splošnost lahko privzamemo, da je  $i = 1$ . Naj bo sedaj  $s_1 = 0$ ,  $\|x_1\|_2 = \|y_1\|_2 = 1$  pa naj bosta ustrezni levi in desni lastni vektor.  $U$  naj bo taka unitarna matrika, da je

$Ue_1 = x_1$ . Potem je matrika  $B = U^H A U$  oblike (glej dokaz izreka 1.2)

$$B = \begin{matrix} & & 1 & n-1 \\ & 1 & \left[ \begin{array}{c} \lambda \\ \times \cdots \times \\ 0 \end{array} \right] \\ n-1 & & & C \end{matrix}.$$

Iz enakosti  $Ax_1 = \lambda_1 x_1$ ,  $y_1^H A = \lambda_1 y_1^H$  in  $s_1 = y_1^H x_1 = 0$  dobimo

$$\begin{aligned} U^H A U e_1 &= U^H \lambda_1 U e_1 = \lambda_1 e_1, \\ (y_1^H U) U^H A U &= \lambda_1 (y_1^H U), \end{aligned} \tag{1.1}$$

$$y_1^H U e_1 = 0. \tag{1.2}$$

Iz (1.2) sledi, da je  $y_1^H U$  oblike  $[0 \ z_1^H]$ , ko pa to vstavimo v (1.1), dobimo  $[0 \ z_1^H] B = \lambda_1 [0 \ z_1^H]$  oziroma  $z_1^H C = \lambda_1 z_1^H$ . Ker je  $\lambda_1$  lastna vrednost matrike  $C$ , ima matrika  $A$  vsaj dvojno lastno vrednost  $\lambda_1$ . ■

V primeru večkratne lastne vrednosti lahko vektorja  $x_i$  in  $y_i$  določimo tako, da bosta ortogonalna, ni pa to nujno res za poljubno izbrana lastna vektorja večkratne lastne vrednosti.

**Izrek 1.10** Naj bo  $A = X \Lambda X^{-1} = Y^H \Lambda Y^{-H}$ , kjer je  $X = [x_1 \ \cdots \ x_n]$  matrika normiranih desnih lastnih vektorjev,  $Y = [y_1 \ \cdots \ y_n]$  matrika normiranih levih lastnih vektorjev in  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  diagonalna matrika lastnih vrednosti. Potem velja

$$X^{-1} = \begin{bmatrix} \frac{1}{s_1} y_1^H \\ \vdots \\ \frac{1}{s_n} y_n^H \end{bmatrix}.$$

*Dokaz.*  $Y^H A = \Lambda Y^H$ , po drugi strani pa  $X^{-1} A = \Lambda X^{-1}$ . Od tod sledi, da so stolpci  $X^{-1}$  levi lastni vektorji. Torej je

$$X^{-1} = \begin{bmatrix} c_1 y_1^H \\ \vdots \\ c_n y_n^H \end{bmatrix}$$

za neke konstante  $c_1, \dots, c_n$ . Če želimo  $XX^{-1} = I$ , mora veljati  $c_i = \frac{1}{s_i}$ . ■

**Lema 1.11** Matrika ne more imeti natanko ene zelo občutljive lastne vrednosti.

*Dokaz.* Naj bodo vse lastne vrednosti različne, sicer že imamo občutljiv par. Naj bo  $A = X \Lambda X^{-1} = Y^H \Lambda Y^{-H}$ , kjer je  $X = [x_1 \ \cdots \ x_n]$  matrika normiranih desnih lastnih vektorjev,  $Y = [y_1 \ \cdots \ y_n]$  matrika normiranih levih lastnih vektorjev in  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  diagonalna matrika lastnih vrednosti. Po izreku 1.10 je

$$X \begin{bmatrix} \frac{1}{s_1} y_1^H \\ \vdots \\ \frac{1}{s_n} y_n^H \end{bmatrix} = I$$

oziroma

$$\sum_{i=1}^n \frac{x_i y_i^H}{s_i} = I.$$

Denimo, da je  $|s_1|^{-1}$  največja občutljivost. Ocenimo lahko

$$\frac{1}{|s_1|} \leq 1 + \sum_{j=2}^n \frac{1}{|s_j|},$$

od tod pa je očitno, da mora biti vsaj ena izmed preostalih občutljivosti tudi zelo velika. ■

Zgornjo lemo si lahko razlagamo tudi tako, da velika občutljivost lastne vrednosti pomeni, da je blizu večkratne lastne vrednosti, to pa seveda pomeni, da obstaja vsaj še ena taka lastna vrednost.

Pri občutljivosti linearnega sistema smo videli, da je recipročna vrednost občutljivosti matrike enaka oddaljenosti od najbližjega singularnega sistema. Podobno tudi sedaj velja, da je občutljivost enostavne lastne vrednosti povezana z oddaljenostjo od najbližje matrike z večkratno lastno vrednostjo. Dokaz naslednjega izreka lahko najdete npr. v [8].

**Izrek 1.12** Naj bo  $\lambda$  enostavna lastna vrednost matrike  $A$  z normiranima levim in desnim lastnim vektorjem  $y$  in  $x$  in naj bo  $|s| < 1$ , kjer je  $s = y^H x$ . Potem obstaja matrika  $A + \delta A$  z večkratno lastno vrednostjo  $\lambda$  in

$$\frac{\|\delta A\|_2}{\|A\|_2} \leq \frac{|s|}{\sqrt{1 - |s|^2}}.$$

Velika občutljivost lastne vrednosti pomeni, da je blizu večkratne lastne vrednosti, to pa seveda pomeni, da obstaja vsaj še ena taka lastna vrednost. Zato velja, da matrika ne more imeti natanko ene zelo občutljive lastne vrednosti.

Kaj se zgodi z občutljivostjo lastne vrednosti pri podobnostnih transformacijah? Naj bo  $\lambda$  enostavna lastna vrednost matrike  $A$  z normiranima levim in desnim lastnim vektorjem  $y$  in  $x$ . Naj bo  $B = S^{-1}AS$ , kjer je  $S$  nesingularna matrika. Če je  $\tilde{s}$  občutljivost  $\lambda$  kot lastne vrednosti matrike  $B$ , potem lahko izpeljemo zvezo

$$\frac{|s|}{|\tilde{s}|} = \frac{\|S^H y\| \|S^{-1} x\|}{\|y\| \|x\|}$$

in ocenimo lahko

$$\frac{1}{\kappa(S)} |\tilde{s}| \leq |s| \leq \kappa(S) |\tilde{s}|.$$

Občutljivost se torej lahko v najboljšem primeru zmanjša za  $\kappa(S)$ , v najslabšem primeru pa se za isti faktor poveča. Če je  $S$  ortogonalna matrika, potem se občutljivost ne spremeni, zato so tovrstne transformacije stabilne.

## 1.4 Potenčna metoda

Pri prvem algoritmu za računanje lastnih vrednosti in vektorjev ne potrebujemo drugega kot množenje z matriko  $A$ . Denimo, da izberemo normiran začetni vektor  $z_0$  in nato za  $k = 0, 1, \dots$

generiramo vektorje po naslednjem predpisu:

$$y_{k+1} = Az_k, \quad z_{k+1} = \frac{y_{k+1}}{\|y_{k+1}\|}. \quad (1.3)$$

Dobimo zaporedje normiranih vektorjev  $z_k$ , za katere se izkaže, da ko gre  $k$  proti neskončno, skonvergirajo proti lastnemu vektorju matrike  $A$ .

**Izrek 1.13** Naj bo  $\lambda_1$  dominantna lastna vrednost matrike  $A$ , kar pomeni

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|.$$

Potem, ko gre  $k$  proti neskončnosti, vektorji  $z_k$ , izračunani po predpisu (1.3), po smeri konvergirajo proti lastnemu vektorju za  $\lambda_1$ .

*Dokaz.* Izrek sicer velja za splošno matriko, dokazali pa ga bomo le za primer, ko se da matriko diagonalizirati. Naj velja  $A = X\Lambda X^{-1}$ , kjer je  $X = [x_1 \ \dots \ x_n]$  in  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Začetni vektor  $z_0$  lahko razvijemo po lastnih vektorjih kot

$$z_0 = \sum_{i=1}^n \alpha_i x_i.$$

Potem pri pogoju  $\alpha_1 \neq 0$  velja

$$z_k = \frac{A^k z_0}{\|A^k z_0\|} = \frac{\alpha_1 x_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k x_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1}\right)^k x_n}{\|\alpha_1 x_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k x_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1}\right)^k x_n\|}$$

in  $z_k$  konvergira v smeri proti  $x_1$ , ko gre  $k$  proti neskončnosti. ■

Vidimo, da se da vektor  $z_k$  izraziti s produktom  $k$ -te potence matrike  $A$  z vektorjem  $z_0$ . Zaradi tega metodo, ki generira vektorje po predpisu (1.3), imenujemo *potenčna metoda*.

V zadnjem dokazu smo predpostavili:

- a)  $\alpha_1 \neq 0$ ,
- b)  $A$  ima enostavno dominantno lastno vrednost.

Pri numeričnem računanju je predpostavka a) v praksi vedno izpolnjena, saj zaokrožitvene napake povzročijo, da je  $\alpha_1 \neq 0$ . Priporočljivo je tudi, da je začetni vektor izbran naključno. Pri točki b) se da pokazati, da izrek velja tudi, ko je  $\lambda_1$  večkratna lastna vrednost. Metoda se da ustrezno predelati (če si zapomnimo in hkrati gledamo zadnje tri približke) tudi za primera:

- $|\lambda_1| = |\lambda_2| > |\lambda_3| \geq \dots$  in  $\lambda_1 = -\lambda_2$ ,
- $|\lambda_1| = |\lambda_2| > |\lambda_3| \geq \dots$  in  $\lambda_1 = \overline{\lambda_2}$ .

Vektor  $z_k$  po smeri konvergira proti lastnemu vektorju za  $\lambda_1$ . Zaradi tega normiranje vektorja v vsakem koraku v bistvu ni potrebno, izvajamo ga le zato, da pri numeričnem računanju ne pride do prekoračitve (v primeru  $|\lambda_1| > 1$ ) ali podkoračitve (v primeru  $|\lambda_1| < 1$ ).

Kako ugotovimo, kdaj je postopek že skonvergiral dovolj blizu rešitve? Ker vektor konvergira po smeri, ne pa tudi po komponentah, pogoj  $\|z_{k+1} - z_k\| \leq \epsilon$  ni dober. Potrebujemo kriterij, ki nam za dani približek za lastni vektor pove, kako dober približek je to. Tega pa se ne da ugotoviti drugače kot da za približek za lastni vektor poiščemo še približek za lastno vrednost in potem skupaj pogledamo kako dober približek za lastni par imamo.

Denimo, da imamo približek  $x$  za lastni vektor in iščemo lastno vrednost. Najboljši približek je  $\lambda$ , ki minimizira

$$\|Ax - \lambda x\|_2,$$

rešitev pa je (uporabimo normalni sistem za predoločeni sistem  $x\lambda = Ax$ ) *Rayleighov kvocient*

$$\rho(x, A) = \frac{x^H Ax}{x^H x},$$

definiran za  $x \neq 0$ .

Za Rayleighov kvocient velja  $\rho(x, A) = \rho(ax, A)$  za  $a \neq 0$ . Rayleighov kvocient je torej odvisen le od smeri, ne pa tudi od norme vektorja. Očitno je tudi, da iz  $Ax = \lambda x$  sledi  $\rho(x, A) = \lambda$ .

Pravilen zaustavitveni kriterij za potenčno metodo je, da za vsak vektor  $z_k$  izračunamo Rayleighov kvocient  $\rho_k = \rho(z_k, A)$ , potem pa pogledamo normo ostanka  $\|Az_k - \rho_k z_k\|_2$ . Če je norma dovolj majhna, je  $\rho_k, z_k$  dober približek za lastni par.

---

**Algoritem 1.1** Osnovna varianta potenčne metode. Začetni podatki so matrika  $A$  (zadošča funkcija, ki zna za dani vektor  $x$  izračunati produkt  $Ax$ ), normiran vektor  $z_0$  in toleranca  $\epsilon$ .

---

```

 $y_1 = Az_0, \quad \rho_0 = z_0^H y_1, \quad k = 0$ 
while  $\|y_{k+1} - \rho_k z_k\|_2 \geq \epsilon$ 
   $k = k + 1$ 
   $z_k = \frac{1}{\|y_k\|_2} y_k$ 
   $y_{k+1} = Az_k$ 
   $\rho_k = z_k^H y_{k+1}$ 

```

---

Vidimo, da računanje Rayleighovega kvocienta in preverjanje konvergence ne vplivata bistveno na časovno zahtevnost. Glavna operacija, ki jo v vsakem koraku algoritma izvedemo enkrat, je množenje z matriko  $A$ . Ta ima v primeru splošne matrike zahtevnost  $\mathcal{O}(n^2)$ , ostale operacije pa imajo zahtevnost  $\mathcal{O}(n)$ .

Konvergenca potenčne metode je linearna. Hitrost je odvisna od razmerja  $\left| \frac{\lambda_2}{\lambda_1} \right|$ . Če je blizu 1, bo konvergenca počasna, blizu 0 pa hitrejša. Pomagamo si lahko s premiki. Matrika  $A - \sigma I$  ima enake lastne vektorje kot  $A$ , lastne vrednosti pa so  $\lambda_i - \sigma$ ,  $i = 1, \dots, n$ . Sedaj je konvergenca lahko hitrejša.

Tako dobimo dominantno lastno vrednost  $\lambda_1$ . Naj bo  $x_1$  pripadajoči normiran lastni vektor. Za ostale lastne pare lahko naredimo redukcijo:

a) *Hotellingova redukcija* za  $A = A^T$ . Definiramo

$$B = A - \lambda_1 x_1 x_1^T.$$

Hitro lahko preverimo, da velja  $Bx_1 = 0$  in  $Bx_k = \lambda_k x_k$  za  $k \neq 1$ . Če uporabimo potenčno metodo na matriki  $B$ , bomo tako dobili drugo dominantno lastno vrednost matrike  $A$

Matrike  $B$  nam ni potrebno eksplicitno izračunati. Produkt matrike  $B$  s poljubnim vektorjem  $z$  lahko namreč izračunamo kot  $Bz = Az - \lambda_1(x_1^T z)x_1$ . Na ta način tudi za izračun naslednje lastne vrednosti še vedno zadošča, da poznamo le funkcijo, ki izračuna produkt matrike  $A$  z danim vektorjem.

- b) *Householderjeva redukcija* za splošno matriko. Poiščemo unitarno matriko  $U$ , da je  $Ux_1 = e_1$ , pri čemer lahko seveda uporabimo Householderjeva zrcaljenja. Potem ima matrika  $B = UAU^H$  obliko (glej dokaz izreka 1.2)

$$B = \begin{bmatrix} \lambda_1 & b^T \\ 0 & C \end{bmatrix}.$$

Preostale lastne vrednosti matrike  $A$  se ujemajo z lastnimi vrednostmi matrike  $C$ .

Tudi v primeru Householderjeve redukcije ekspliciten izračun matrike  $C$  ni potreben. Pri potenčni metodi moramo znati izračunati produkt  $Cw$  za vektor  $w \in \mathbb{C}^{n-1}$ . Pri tem si pomagamo z zvezo

$$UAU^H \begin{bmatrix} 0 \\ w \end{bmatrix} = \begin{bmatrix} \lambda_1 & b^T \\ 0 & C \end{bmatrix} \begin{bmatrix} 0 \\ w \end{bmatrix} = \begin{bmatrix} b^T w \\ Cw \end{bmatrix}.$$

Potrebujemo torej množenje z matriko  $A$  in dve množenji z ortogonalno matriko  $Q$ , ki ju lahko v primeru Householderjevega zrcaljenja izvedemo z zahtevnostjo  $\mathcal{O}(n)$ .

Če iščemo lastno vrednost nesingularne matrike  $A$ , ki je najmanjša po absolutni vrednosti, delamo potenčno metodo za  $A^{-1}$ , saj ima  $A^{-1}$  lastne vrednosti  $\lambda^{-1}, \dots, \lambda_n^{-1}$ . V algoritmu namesto množenja  $y_{k+1} = A^{-1}z_k$  rešujemo sistem  $Ay_{k+1} = z_k$ .

## 1.5 Obratna napaka in izračunljive ocene

Denimo, da smo numerično, npr. s potenčno metodo, izračunali približek  $(\hat{\lambda}, \hat{x})$  za lastni par matrike  $A$ . Radi bi ocenili, za koliko se  $\hat{\lambda}$  razlikuje od prave lastne vrednosti.

**Definicija 1.14** Po normi relativna obratna napaka približka  $(\hat{\lambda}, \hat{x})$  za lastni par je

$$\eta(\hat{\lambda}, \hat{x}) = \min \left\{ \epsilon > 0 : (A + \delta A)\hat{x} = \hat{\lambda}\hat{x}, \|\delta A\|_2 \leq \epsilon \|A\|_2 \right\}.$$

Če gledamo samo približek za lastno vrednost  $\hat{\lambda}$ , definiramo obratno napako kot  $\eta(\hat{\lambda}) = \min_{\hat{x} \neq 0} \eta(\hat{\lambda}, \hat{x})$ , podobno za obratno napako približka za lastni vektor vzamemo  $\eta(\hat{x}) = \min_{\hat{\lambda}} \eta(\hat{\lambda}, \hat{x})$ .

**Lema 1.15** Velja

- 1)  $\eta(\hat{\lambda}, \hat{x}) = \frac{\|A\hat{x} - \hat{\lambda}\hat{x}\|_2}{\|A\|_2 \|\hat{x}\|_2},$
- 2)  $\eta(\hat{\lambda}) = \frac{\sigma_{\min}(A - \hat{\lambda}I)}{\|A\|_2},$



$$3) \eta(\hat{x}) = \frac{\|A\hat{x} - \rho(\hat{x}, A)\hat{x}\|_2}{\|A\|_2 \|\hat{x}\|_2}.$$

*Dokaz.* Za točko 1) označimo  $r = A\hat{x} - \hat{\lambda}\hat{x}$ . Iz  $(A + \delta A)\hat{x} = \hat{\lambda}\hat{x}$  sledi  $-\delta A\hat{x} = r$ , od tod pa  $\|r\|_2 \leq \|\delta A\|_2 \|\hat{x}\|_2$ .

Če vzamemo  $\delta A = -r\hat{x}^H / \|\hat{x}\|_2^2$ , potem je  $(A + \delta A - \hat{\lambda}I)\hat{x} = 0$  in  $\|\delta A\|_2 = \|r\|_2 / \|\hat{x}\|_2$ , kar pomeni, da je minimum res dosežen in točka 1) drži.

Pri točki 2) upoštevamo dejstvo, da je  $\min_{x \neq 0} \|Mx\|_2 = \sigma_{\min}(M)$ , pri točki 3) pa lastnost Rayleighovega kvocienta, da je minimum  $\|Ax - \tau x\|_2$  dosežen pri  $\tau = \rho(x, A)$ . ■

Iz zgornje leme sledi, da je potenčna metoda obratno stabilna, saj vrne tak približek  $(\hat{\lambda}, \hat{x})$  za lastni par, za katerega velja  $\|A\hat{x} - \hat{\lambda}\hat{x}\|_2 \leq \epsilon$ .

Če bi radi ocenili, za koliko se približek  $\hat{\lambda}$  razlikuje od točne lastne vrednosti, potrebujemo poleg ocene za obratno napako še oceno za občutljivost lastne vrednosti. Za oceno potrebujemo tudi približek za levi lastni vektor. Tega bodisi izračunamo v algoritmu ali pa uporabimo npr. inverzno iteracijo, ki jo bomo spoznali v naslednjem razdelku.

Naj bosta torej  $\hat{x}$  in  $\hat{y}$  približka za desni in levi lastni vektor, ki pripadata  $\hat{\lambda}$ . Potem velja ocena

$$|\hat{\lambda} - \lambda| \leq \frac{\|r\|_2}{|\hat{s}|},$$

kjer je

$$\hat{s} = \frac{-\hat{y}^H \hat{x}}{\|\hat{y}\|_2 \|\hat{x}\|_2}.$$

## 1.6 Inverzna iteracija

Rayleighov kvocient vrne najboljši približek za lastno vrednost, ki ustreza danemu vektorju. Kaj pa obratno? Denimo, da smo izračunali približek za lastno vrednost  $\sigma$ , sedaj pa potrebujemo še lastni vektor. Tu si lahko pomagamo z *inverzno iteracijo*, ki je v grobem predstavljena v algoritmu 1.2.

---

**Algoritem 1.2** Osnovna verzija inverzne iteracije. Začetni podatki so matrika  $A$ , približek za lastno vrednost  $\sigma$  in normiran vektor  $z_0$ .

---

$$\begin{aligned} k &= 0, 1, \dots \\ &\text{reši sistem } (A - \sigma I)y_{k+1} = z_k \\ z_{k+1} &= \frac{1}{\|y_{k+1}\|} y_{k+1} \end{aligned}$$


---

Naj za približek  $\sigma$  velja, da mu je najbližja lastna vrednost  $\lambda_i$  in velja  $|\lambda_i - \sigma| \ll |\lambda_j - \sigma|$  za  $j \neq i$ . Inverzna iteracija v bistvu ni nič drugega kot potenčna metoda za matriko  $(A - \sigma I)^{-1}$ , zato jo imenujemo tudi *inverzna potenčna metoda*. Od tod vemo, da vektor  $z_k$  po smeri konvergira proti lastnemu vektorju, ki pripada dominantni lastni vrednosti matrike  $(A - \sigma I)^{-1}$ .

Lastne vrednosti matrike  $(A - \sigma I)^{-1}$  so  $(\lambda_j - \sigma)^{-1}$  za  $j = 1, \dots, n$ . Ker velja  $|\lambda_i - \sigma| \ll |\lambda_j - \sigma|$  za  $j \neq i$ , je  $|(\lambda_i - \sigma)^{-1}| \gg |(\lambda_j - \sigma)^{-1}|$  za  $j \neq i$ . Boljši, ko je približek  $\sigma$ , dominantnejša je lastna vrednost  $(\lambda_i - \sigma)^{-1}$  in hitrejša je konvergenca.

Inverzno iteracijo ponavadi uporabljamo zato, da dobimo lastni vektor za numerično izračunano lastno vrednost. V tem primeru je  $\sigma$  kar lastna vrednost matrike  $A$ , izračunana z natančnostjo  $\mathcal{O}(u)$ . V praksi zato potrebujemo le en do dva koraka inverzne iteracije, da iz poljubnega začetnega vektorja izračunamo pripadajoči lastni vektor.

Če je  $\sigma$  res lastna vrednost matrike  $A$ , potem je matrika  $A - \sigma I$  singularna in v algoritmu lahko pričakujemo težave pri reševanju sistema s to matriko. Izkaže se, da nam zaokrožitvene napake pomagajo do tega, da v praksi ne pride do deljenja z nič, zato je inverzna iteracija zelo učinkovita metoda za računanje lastnih vektorjev za lastne vrednosti.

## 1.7 Ortogonalna iteracija

Pravimo, da je  $\mathcal{N}$  invarianten podprostor za matriko  $A$ , če velja  $A\mathcal{N} \subset \mathcal{N}$ . Naj ima matrika  $S_1$   $p$  linearno neodvisnih stolpcev. Če jo dopolnimo do nesingularne matrike  $S = [S_1 \ S_2]$  in je

$$B = S^{-1}AS = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix},$$

potem lahko hitro preverimo, da stolpci  $S_1$  razpenjajo invariantni podprostor natanko takrat, ko je  $B_{21} = 0$ . V tem primeru so lastne vrednosti matrike  $A$  unija lastnih vrednosti matrik  $B_{11}$  in  $B_{22}$ . Če lastne vrednosti matrike  $A$  lahko uredimo po absolutni vrednosti tako, da velja  $|\lambda_1| \geq \dots \geq |\lambda_p| > |\lambda_{p+1}| \geq \dots \geq |\lambda_n|$ , potem stolpci matrike  $S_1$  razpenjajo dominantni invariantni podprostor takrat, ko so lastne vrednosti matrike  $B_{11}$  enake  $\lambda_1, \dots, \lambda_p$ .

Če se da matriko diagonalizirati, potem bazo za dominantni invariantni podprostor dimenzije  $p$  lahko sestavimo iz  $p$  lastnih vektorjev, ki pripadajo po absolutni vrednosti  $p$  največjim lastnim vrednostim.

Za izračun baze za dominantni invariantni podprostor imamo na voljo *ortogonalno iteracijo*.

---

**Algoritem 1.3** Osnovna verzija ortogonalne iteracije. Začetni podatki so matrika  $A$  velikosti  $n \times n$  in matrika  $Z_0$  velikosti  $n \times p$ ,  $p \leq n$ , z ortonormiranimi stolpci.

---

$k = 0, 1, \dots$

$Y_{k+1} = AZ_k$

izračunaj QR razcep  $Y_{k+1} = QR$  in vzemi  $Z_{k+1} = Q$

---

Opazimo lahko, da je pri  $p = 1$  to kar potenčna metoda. Potenčna metoda skonvergira proti lastnemu vektorju, linearni podprostor, ki ga razpenja lastni vektor, pa je očitno dominanten invarianten podprostor dimenzije 1.

**Izrek 1.16** Naj velja  $A = X\Lambda X^{-1}$ , kjer je  $X = [x_1 \ \dots \ x_n]$  in  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Lastne vrednosti naj bodo urejene po absolutni vrednosti in naj velja  $|\lambda_p| > |\lambda_{p+1}|$ . Potem matrika  $Z_k$  iz ortogonalne iteracije konvergira proti ortonormirani bazi za invariantni podprostor  $\text{Lin}(\{x_1, \dots, x_p\})$ .

*Dokaz.* Očitno je  $\text{Lin}(Z_{k+1}) = \text{Lin}(Y_{k+1}) = \text{Lin}(AZ_k)$ , od tod pa sledi  $\text{Lin}(Z_k) = \text{Lin}(A^k Z_0)$ . Ker je  $A^k = X\Lambda^k X^{-1}$ , velja

$$A^k Z_0 = X\Lambda^k X^{-1} Z_0 = \lambda_p^k X \begin{bmatrix} (\lambda_1/\lambda_p)^k & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & (\lambda_n/\lambda_p)^k \end{bmatrix} X^{-1} Z_0.$$

Ko gre  $k$  proti neskončno, gre  $A^k Z_0$  proti  $X \cdot \begin{matrix} p \\ n-p \\ 0 \end{matrix} \begin{pmatrix} \times \\ \times \\ 0 \end{pmatrix}$ , to pa pomeni, da  $\text{Lin}(Z_k)$  konvergira proti  $\text{Lin}(\{x_1, \dots, x_p\})$ . ■

Podobno v primeru, ko je  $|\lambda_r| > |\lambda_{r+1}|$ , za prvih  $r < p$  stolpcev velja, da  $\text{Lin}(Z_k(:, 1:r))$  konvergira proti  $\text{Lin}(\{x_1, \dots, x_r\})$ .

Vzemimo kar  $p = n$  in poljubno nesingularno matriko  $Z_0$ . V tem primeru seveda ne računamo invariantnega podprostora dimenzije  $n$ , saj je to kar celotni prostor. Pri predpostavki, da so absolutne vrednosti  $\lambda_i$  paroma različne (to hkrati pomeni, da so vse realne), lahko pokažemo, da matrika  $A_k := Z_k^T A Z_k$  konvergira proti Schurovi formi.

$A_k$  in  $Z_k^T A Z_k$  sta podobni matriki, saj je  $Z_k$  ortogonalna matrika. Naj bo  $Z_k = [Z_{k1} \ Z_{k2}]$ , kjer ima  $Z_{k1}$   $p$  stolpcev. Potem je

$$Z_k^T A Z_k = \begin{bmatrix} Z_{k1}^T A Z_{k1} & Z_{k1}^T A Z_{k2} \\ Z_{k2}^T A Z_{k1} & Z_{k2}^T A Z_{k2} \end{bmatrix}.$$

Ker  $\text{Lin}(Z_{k1})$  konvergira proti invariantnemu podprostoru  $\text{Lin}(\{x_1, \dots, x_p\})$ , enako velja tudi za  $\text{Lin}(AZ_{k1})$ , to pa pomeni, da  $Z_{k2}^T A Z_{k1}$  konvergira proti 0, saj je  $Z_{k2}^T Z_{k1} = 0$ . Ker to velja za vsak  $p = 1, \dots, n$ , sledi, da matrika  $Z_k^T A Z_k$  res konvergira proti zgornji trikotni matriki, torej proti Schurovi formi.

Poddiagonalni elementi matrike  $A_k$  konvergirajo proti 0 z linearno konvergenco, hitrost konvergence  $(i, j)$ -tega elementa, kjer je  $i > j$ , pa je odvisna od razmerja  $|\lambda_j|/|\lambda_i|$ .

## 1.8 QR iteracija

V prejšnjem razdelku smo videli, da z ortogonalno iteracijo lahko izračunamo Schurovo formo in s tem vse lastne vrednosti matrike. V tem razdelku pa bomo spoznali algoritem, ki zna to narediti na ekonomičnejši način. Gre za *QR iteracijo*<sup>1</sup>, ki je trenutno najboljša numerična metoda za izračun vseh lastnih vrednosti splošne nesimetrične matrike. Osnova verzija je zapisana v algoritmu 1.4.

V vsakem koraku izračunamo QR razcep matrike in faktorja v zamenjanem vrstnem redu zmnožimo v novo matriko. Iz  $A_{k+1} = R_k Q_k = Q_k^T A_k Q_k$  sledi, da sta si matriki  $A_{k+1}$  in  $A_k$

<sup>1</sup>Metodo sta neodvisno leta 1961 odkrila angleški računalnikar John G. F. Francis (1934–) in ruska matematičarka Vera N. Kublanovskaja (1920–). Francis se je leta 1962 nehal ukvarjati z numerično matematiko in se nato do leta 2007 sploh ni zavedal, kakšen vpliv ima njegov algoritem na numerično matematiko. Po oceni, ki sta jo naredila Jack Dongarra in Francis Sullivan leta 2000, spada QR iteracija med 10 algoritmov iz 20. stoletja, ki so najbolj vplivali na razvoj znanosti in tehnike [10].

---

**Algoritem 1.4** Osnovna verzija QR iteracije. Začetni podatek je matrika  $A$ .

---

$$\begin{aligned} A_0 &= A \\ k &= 0, 1, \dots \\ A_k &= Q_k R_k \text{ (izračunaj QR razcep)} \\ A_{k+1} &= R_k Q_k \end{aligned}$$


---

ortogonalno podobni, torej je  $A_k$  ortogonalno podobna začetni matriki  $A$ . Izkaže se, da je QR iteracija povezana z ortogonalno iteracijo, od tod pa sledi, da  $A_k$  konvergira proti Schurovi formi.

**Lema 1.17** Za matriko  $A_k$  iz QR iteracije velja  $A_k = Z_k^T A Z_k$ , kjer je  $Z_k$  matrika, ki jo dobimo v ortogonalni iteraciji iz  $Z_0 = I$ . V primeru, ko imajo lastne vrednosti paroma različne absolutne vrednosti,  $A_k$  konvergira proti Schurovi formi.

*Dokaz.* Uporabimo indukcijo po  $k$ .

Na začetku je  $A_0 = Z_0^T A Z_0$ , saj je  $Z_0 = I$ . Denimo, da je  $A_k = Z_k^T A Z_k$ . Potem je

$$Z_k^T A Z_k = Z_k^T \left( \underbrace{\begin{matrix} Z_{k+1} & S_{k+1} \\ \text{ort.} & \text{zg. trik.} \end{matrix}}_{\text{QR razcep } AZ_k} \right) = \underbrace{Z_k^T Z_{k+1}}_{\text{ort.}} \underbrace{S_{k+1}}_{\text{zg. trik.}} = Q_k R_k,$$

zaradi enoličnosti pa je to kar QR razcep matrike  $A_k$ . Sledi

$$A_{k+1} = R_k Q_k = S_{k+1} Z_k^T Z_{k+1} = \underbrace{Z_{k+1}^T A Z_k}_{S_{k+1}} Z_k^T Z_{k+1} = Z_{k+1}^T A Z_{k+1}. \quad \blacksquare$$

V primeru, ko ima matrika  $A$  tudi kompleksne lastne vrednosti, matrika  $A_k$  skonvergira proti realni Schurovi formi.

Če pogledamo zahtevnost enega koraka QR iteracije, vidimo, da ima časovno zahtevnost  $\mathcal{O}(n^3)$ , saj moramo v vsakem koraku izračunati QR razcep matrike. Hitrost konvergence je odvisna od razmerja med lastnimi vrednostmi. Če sta dve lastni vrednosti zelo blizu, lahko pričakujemo, da bo metoda potrebovala veliko korakov, preden bo skonvergirala do Schurove forme.

Da pridemo do uporabne verzije QR iteracije, moramo vpeljati še nekaj izboljšav.

### 1.8.1 Redukcija na Hessenbergovo obliko

En korak osnovne QR iteracije porabi  $\mathcal{O}(n^3)$  operacij, kar ni najbolj ekonomično. Zahtevnost enega koraka lahko močno zmanjšamo, če matriko  $A$  predhodno reduciramo na zgornjo Hessenbergovo obliko.

**Definicija 1.18** Pravimo, da je matrika  $A$  zgornja Hessenbergova, če velja  $a_{ij} = 0$  za  $i > j + 1$ .

Zgornja Hessenbergova matrika ima torej le zgornji trikotnik in eno poddiagonalo. Od Schurove forme jo loči le poddiagonala. Izkaže se, da se zgornja Hessenbergova oblika ohranja med QR iteracijo.

**Trditev 1.19** Če je  $A$  zgornja Hessenbergova, se oblika med QR iteracijo ohranja.

*Dokaz.* Pri QR razcepu matrike  $A = [a_1 \ \cdots \ a_n]$  dobimo zgornjo Hessenbergovo matriko  $Q$  in zgornjo trikotno matriko  $R$ . Pri  $Q = [q_1 \ \cdots \ q_n]$  je oblika razvidno iz dejstva, da je  $q_i$  linearna kombinacija stolpcev  $a_1, \dots, a_i$ . Hitro lahko preverimo, da je produkt zgornje trikotne in zgornje Hessenbergove matrike spet zgornja Hessenbergova matrika. ■

Vsako realno matriko  $A$  lahko z ortogonalno podobnostno transformacijo spremenimo v zgornjo Hessenbergovo matriko. Za splošno matriko uporabimo Householderjeva zrcaljenja, če pa ima matrika  $A$  v spodnjem trikotniku že veliko ničel, so lahko Givensove rotacije še bolj učinkovite.

**Zgled 1.4** Na zgledu matrike velikosti  $5 \times 5$  pogledimo, kako matriko z ortogonalnimi podobnostnimi transformacijami spremenimo v zgornjo Hessenbergovo matriko. Naj bo

$$A = \begin{bmatrix} \times & \times & \times & \times & \times \\ (\times) & \times & \times & \times & \times \\ (\times) & \times & \times & \times & \times \\ (\times) & \times & \times & \times & \times \\ (\times) & \times & \times & \times & \times \end{bmatrix}.$$

Elementi v oklepajih označujejo elemente vektorja, ki določa Householderjevo zrcaljenje v naslednjem koraku. Zrcaljenje določimo tako, da se označeni vektor prezrcali v smer prvega enotskega vektorja.

Najprej poiščemo tako ortogonalno matriko  $Q_1$ , da je

$$Q_1 A = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \end{bmatrix}, \quad A_1 = Q_1 A Q_1^T = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & (\times) & \times & \times & \times \\ 0 & (\times) & \times & \times & \times \\ 0 & (\times) & \times & \times & \times \end{bmatrix}.$$

Nato poiščemo ortogonalno matriko  $Q_2$ , da je

$$Q_2 A_1 = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & \times & \times & \times \end{bmatrix}, \quad A_2 = Q_2 A_1 Q_2^T = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & (\times) & \times & \times \\ 0 & 0 & (\times) & \times & \times \end{bmatrix},$$

na koncu pa še ortogonalno matriko  $Q_3$ , da je

$$Q_3 A_2 = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}, \quad H = Q_3 A_2 Q_3^T = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}.$$

Tako dobimo zgornjo Hessenbergovo matriko  $H = Q_3 Q_2 Q_1 A (Q_3 Q_2 Q_1)^T$ . □

Algoritem za splošno matriko ima naslednjo obliko.

---

**Algoritem 1.5** Redukcija na zgornjo Hessenbergovo obliko preko Householderjevih zrcaljenj. Začetni podatek je  $n \times n$  matrika  $A$ . Algoritem vrne zgornjo Hessenbergovo matriko  $H$  in po potrebi tudi ortogonalno matriko  $Q$ , da je  $A = Q^T H Q$ .

---

$$Q = I \quad (*)$$

$$i = 1, \dots, n-2$$

določi  $w_i \in \mathbb{R}^{n-i}$  za Householderjevo zrcaljenje  $P_i$ , ki prezrcali  $A(i+1:n, i)$  v  $\pm ke_1$

$$A(i+1:n, i:n) = P_i A(i+1:n, i:n)$$

$$A(1:n, i+1:n) = A(1:n, i+1:n) P_i$$

$$Q(i+1:n, i:n) = P_i Q(i+1:n, i:n) \quad (*)$$


---

Korake označene z (\*) izvedemo le v primeru, če potrebujemo tudi prehodno matriko  $Q$ .

Število operacij je  $\frac{10}{3}n^3 + \mathcal{O}(n^2)$  oziroma  $\frac{14}{3}n^3 + \mathcal{O}(n^2)$  če potrebujemo tudi  $Q$ .

Če na začetku matriko  $A$  reduciramo na Hessenbergovo obliko, porabimo potem za en korak QR iteracije le še  $\mathcal{O}(n^2)$  namesto  $\mathcal{O}(n^3)$  operacij. Matrika  $Q_k$  iz QR razcepa zgornje Hessenbergove matrike  $A_k$  je namreč produkt  $n-1$  Givensovih rotacij, za eno množenje matrike z Givensovo rotacijo pa vemo, da ima zahtevnost  $\mathcal{O}(n)$ .

**Definicija 1.20** Hessenbergova matrika  $H$  je ireducibilna, če so vsi njeni subdiagonalni elementi  $h_{i+1,i}$  neničelni.

Če je  $H$  reducibilna, je npr.

$$H = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}$$

in problem lastnih vrednosti razpade na dva ločena problema. Zaradi tega lahko vedno predpostavimo, da je  $H$  ireducibilna.

Pri numeričnem računanju subdiagonalni element  $a_{i+1,i}^{(k)}$  matrike  $A_k$  postavimo na 0, kadar je dovolj majhen v primerjavi s sosednjima diagonalnima elementoma. To pomeni, da zadošča kriteriju

$$|a_{i+1,i}^{(k)}| < \epsilon(|a_{ii}^{(k)}| + |a_{i+1,i+1}^{(k)}|),$$

kjer je  $\epsilon = \mathcal{O}(u)$  izbrana toleranca.

## 1.8.2 Premiki

Z redukcijo na Hessenbergovo obliko smo zmanjšali zahtevnost posameznega koraka QR iteracije, samo število potrebnih korakov pa se ni zmanjšalo, saj je hitrost konvergence odvisna od razmerja med lastnimi vrednostmi.

Konvergenco lahko pospešimo z vpeljavo premikov:

Naslednja lema nam zagotavlja, da je tudi po vpeljavi premika matrika  $A_k$  še vedno ortogonalno podobna začetni matriki  $A$ .

---

**Algoritem 1.6** QR iteracija s premiki. Začetni podatek je matrika  $A$ .

---

$$\begin{aligned} A_0 &= A \\ k &= 0, 1, \dots \\ &\text{izberi premik } \sigma_k \\ A_k - \sigma_k I &= Q_k R_k \text{ (izračunaj QR razcep)} \\ A_{k+1} &= R_k Q_k + \sigma_k I \end{aligned}$$


---

**Lema 1.21** Matriki  $A_k$  in  $A_{k+1}$  pri QR iteraciji s premikom sta ortogonalno podobni.

*Dokaz.*

$$A_{k+1} = R_k Q_k + \sigma_k I = Q_k^T (Q_k R_k + \sigma_k I) Q_k = Q_k^T A_k Q_k. \quad \blacksquare$$

Za hitro konvergenco moramo za premik izbrati čim boljši približek za lastno vrednost. Če bi za premik izbrali kar lastno vrednost, potem iz naslednje leme sledi, da bi se v enem koraku QR iteracije iz matrike izločila ta lastna vrednost in bi računanje lahko nadaljevali na manjši matriki.

**Lema 1.22** Naj bo  $\sigma$  lastna vrednost ireducibilne zgornje Hessenbergove matrike  $A$ . Če je QR razcep  $A - \sigma I = QR$  in  $B = RQ + \sigma I$ , potem je  $b_{n,n-1} = 0$  in  $b_{nn} = \sigma$ .

*Dokaz.* Ker je  $A$  ireducibilna, je prvih  $n - 1$  stolpcev matrike  $A - \sigma I$  linearno neodvisnih. V razcepu  $A - \sigma I = QR$  zato velja  $r_{ii} \neq 0$  za  $i = 1, \dots, n - 1$ . Ker pa je  $A - \sigma I$  singularna, mora biti  $r_{nn} = 0$ . To pa pomeni, da je zadnja vrstica v matriki  $RQ$  enaka 0, torej v matriki  $B = RQ + \sigma I$  velja  $b_{n,n-1} = 0$  in  $b_{nn} = \sigma$ .

Preostale lastne vrednosti lahko sedaj izračunamo iz matrike  $B(1 : n - 1, 1 : n - 1)$ . ■

Potrebujemo čim boljši približek za lastno vrednost  $A_k$ , zato imamo na voljo različne premike:

- a) *Enojni premik:* za  $\sigma_k$  izberemo  $a_{nn}^{(k)}$ . V tem primeru imamo kvadratično konvergenco v bližini lastne vrednosti, vendar premik ni dober za kompleksne lastne vrednosti.

---

Motivacija, da za premik izberemo element v spodnjem desnem kotu je, da naj bi bil to dober približek za po absolutni vrednosti najmanjšo lastno vrednost  $\lambda_n$  matrike  $A$ .

Naj bo  $y_n$  levi lastni vektor za  $\lambda_n$ . Z malce računanja lahko pokažemo, da pri QR algoritmu brez premikov velja

$$A^k = \tilde{Q}_{k-1} \tilde{R}_{k-1},$$

kjer je  $\tilde{Q}_{k-1} = Q_0 \cdots Q_{k-1}$  in  $\tilde{R}_{k-1} = R_{k-1} \cdots R_0$ . Od tod sledi  $A^{-k} = \tilde{R}_{k-1}^{-1} \tilde{Q}_{k-1}^T$ . Iz zadnje enakosti vidimo, da je zadnja vrstica matrike  $\tilde{Q}_{k-1}^T$  proporcionalna  $e_n^T A^{-k}$ , to je zadnji vrstici matrike  $A^{-k}$ . Razen v izjemnem in pri numeričnem računanju malo verjetnem primeru, ko  $e_n$  nima nobene komponente v smeri  $y_n$ , bo vrstica  $e_n^T A^{-k}$  po smeri konvergirala proti  $y_n^T$ . Od tod sledi, da zadnji stolpec matrike  $\tilde{Q}_{k-1}$  konvergira proti levemu lastnemu vektorju  $y_n$ . Iz zveze  $A_k = \tilde{Q}_{k-1}^T A \tilde{Q}_{k-1}$  je tako razvidno, da  $a_{nn}^{(k)} = e_n^T A_k e_n = (\tilde{Q}_{k-1} e_n)^T A \tilde{Q}_{k-1} e_n$  konvergira proti lastni vrednosti  $\lambda_n$ .

---

b) Dvojni oz. Francisov premik: vzamemo podmatriko

$$A_k(n-1:n, n-1:n) = \begin{bmatrix} a_{n-1, n-1}^{(k)} & a_{n-1, n}^{(k)} \\ a_{n, n-1}^{(k)} & a_{nn}^{(k)} \end{bmatrix},$$

ki ima lastni vrednosti  $\sigma_1^{(k)}, \sigma_2^{(k)}$  (lahko sta tudi kompleksni). Sedaj naredimo dva premika v enem koraku:

$$\begin{aligned} A_k - \sigma_1^{(k)} I &= Q_k R_k \text{ (izračunaj QR razcep)} \\ A'_k &= R_k Q_k + \sigma_1^{(k)} I \\ A'_k - \sigma_2^{(k)} I &= Q'_k R'_k \text{ (izračunaj QR razcep)} \\ A_{k+1} &= R'_k Q'_k + \sigma_2^{(k)} I \end{aligned}$$

Izkaže se, da lahko en korak QR z dvojnimi premiki izvedemo brez kompleksne aritmetike, saj velja:

$$\begin{aligned} Q_k Q'_k R'_k R_k &= Q_k (A'_k - \sigma_2^{(k)} I) R_k \\ &= Q_k Q_k^H (A_k - \sigma_2^{(k)} I) Q_k R_k = (A_k - \sigma_2^{(k)} I) Q_k R_k \\ &= (A_k - \sigma_2^{(k)} I) (A_k - \sigma_1^{(k)} I) \\ &= A_k^2 - (\sigma_1^{(k)} + \sigma_2^{(k)}) A_k + \sigma_1^{(k)} \sigma_2^{(k)} I =: N_k \end{aligned}$$

in  $(Q_k Q'_k)(R'_k R_k)$  je QR razcep realne matrike  $N_k$ . Ker po lemi 1.21 velja tudi

$$A_{k+1} = Q_k^H A'_k Q'_k = Q_k^H Q_k^H A_k Q_k Q'_k,$$

potrebujemo le realni QR razcep realne matrike  $N_k$ .

Na prvi pogled nam zgornja ugotovitev ne pomaga dosti, saj v formuli za  $N_k$  nastopa matrika  $A_k^2$ , za izračun le te pa potrebujemo  $\mathcal{O}(n^3)$  operacij. Kot bomo videli v naslednjem razdelku, pa se izkaže, da v resnici potrebujemo le prvi stolpec matrike  $N_k$ .

## 1.9 Implicitna QR metoda

**Izrek 1.23 (Implicitni Q)** Če je  $Q = [q_1 \cdots q_n]$  taka ortogonalna matrika, da je  $Q^T A Q = H$  ireducibilna Hessenbergova matrika, potem so stolpci  $q_2, \dots, q_n$  do predznaka natančno določeni s  $q_1$ .

*Dokaz.* Denimo, da je  $V^T A V = G$ , kjer je  $V = [v_1 \cdots v_n]$  ortogonalna matrika,  $G$  ireducibilna zgornja Hessenbergova matrika in  $q_1 = v_1$ .

Potem je  $W = V^T Q$  ortogonalna matrika. Če zapišemo  $W = [w_1 \cdots w_n]$ , potem je  $w_1 = e_1$ . Velja

$$G W = G V^T Q = V^T A Q = V^T Q H = W H.$$

Iz te zveze sledi  $G w_i = \sum_{j=1}^{i+1} h_{ji} w_j$  oziroma

$$h_{i+1, i} w_{i+1} = G w_i - \sum_{j=1}^i h_{ji} w_j.$$



Ker je  $w_1 = e_1$  in ima  $Gw_i$  en neničelni element več od  $w_i$ , sledi  $w_i \in \text{Lin}(\{e_1, \dots, e_i\})$ . To pomeni, da je  $W$  zgornja trikotna matrika. Ker pa je  $W$  hkrati ortogonalna, je edina možnost  $W = \text{diag}(1, \pm 1, \dots, \pm 1)$ , torej  $v_i = \pm q_i$  za  $i = 2, \dots, n$ . ■

Posledica je, da če v QR algoritmu  $A_k = Q_k R_k$ ,  $A_{k+1} = R_k Q_k = Q_k^T A_k Q_k$  poznamo prvi stolpec matrike  $Q_k$ , potem lahko matriko  $A_{k+1}$  izračunamo brez tega, da bi računali celoten QR razcep matrike  $A_k$ . Tako dobimo *implicitno QR iteracijo*.

Najprej pogledjmo implicitno QR iteracijo z enojnim premikom. Vemo, da je prvi stolpec  $Q_k$  enak normiranemu prvemu stolpcu  $A_k - \sigma_k I$ . Če uspemo poiskati tako ortogonalno matriko  $Q_k$ , ki bo imela za prvi stolpec normirani prvi stolpec  $A_k - \sigma_k I$  in bo  $Q_k^T A_k Q_k$  zgornja Hessenbergova matrika, potem je po izreku o implicitnem Q matrika  $Q_k^T A_k Q_k$  enaka matriki iz naslednjega koraka QR metode.

Matriko  $Q_k$  poiščemo kot produkt Givensovih rotacij  $Q_k = R_{12} R_{23} \cdots R_{n-1,n}$ . Prva rotacija  $R_{12}$  je že določena s prvim stolpcem  $A_k - \sigma_k I$ , ostale pa določimo tako, da bo  $Q_k^T A_k Q_k$  zgornja Hessenbergova matrika. Če je namreč

$$R_{12} = \begin{bmatrix} c_1 & s_1 & & & \\ -s_1 & c_1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix},$$

potem je

$$Q_k = R_{12} R_{23} \cdots R_{n-1,n} = \begin{bmatrix} c_1 & \times & \cdots & \cdots & \times \\ -s_1 & \times & \cdots & \cdots & \times \\ & \times & & & \vdots \\ & & \ddots & & \vdots \\ & & & \times & \times \end{bmatrix}.$$

Algoritem lahko označimo kot *premikanje grbe*. Pogledjmo si ga na primeru matrike velikosti  $5 \times 5$ . Po prvem koraku dobimo

$$R_{12}^T A_k R_{12} = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ + & \times & \times & \times & \times \\ & & \times & \times & \times \\ & & & \times & \times \end{bmatrix}.$$

Novi neničelni element, označen s +, je grba, ki jo z naslednjimi rotacijami pomikamo navzdol ob diagonali. Tako po vrsti poiščemo  $R_{23}$ ,  $R_{34}$  in  $R_{45}$ , da je

$$R_{23}^T R_{12}^T A_k R_{12} R_{23} = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ & \times & \times & \times & \times \\ + & \times & \times & \times & \times \\ & & & \times & \times \end{bmatrix},$$

$$R_{34}^T R_{23}^T R_{12}^T A_k R_{12} R_{23} R_{34} = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ & \times & \times & \times & \times \\ & & \times & \times & \times \\ & & & + & \times & \times \end{bmatrix}$$

in

$$R_{45}^T R_{34}^T R_{23}^T R_{12}^T A_k R_{12} R_{23} R_{34} R_{45} = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ & \times & \times & \times & \times \\ & & \times & \times & \times \\ & & & \times & \times \end{bmatrix}.$$

Dobimo zgornjo Hessenbergovo matriko, ki je po izreku o implicitnem Q enaka matriki  $A_{k+1}$ .

Še bolj kot pri enojnem premiku nam izrek o implicitnem Q pride prav pri dvojnem premiku. Vemo, da je  $A_{k+1} = U_k^T A_k U_k$ , kjer je  $U_k$  ortogonalna matrika iz QR razcepa matrike  $N_k = A_k^2 - (\sigma_1^{(k)} + \sigma_2^{(k)})A_k + \sigma_1^{(k)}\sigma_2^{(k)}I$ . Dovolj je poznati le prvi stolpec matrike  $N_k$ . Le ta ima obliko

$$\begin{bmatrix} a_{11}^2 + a_{12}a_{21} - sa_{11} + t \\ a_{21}(a_{11} + a_{22} - s) \\ a_{21}a_{32} \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

kjer sta

$$\begin{aligned} s &= a_{n-1,n-1} + a_{nn} \\ t &= a_{n-1,n-1}a_{nn} - a_{n-1,n}a_{n,n-1}. \end{aligned}$$

Sedaj najprej poiščemo Householderjevo zrcaljenje oblike

$$P_1 = \begin{bmatrix} \times & \times & \times & & & \\ \times & \times & \times & & & \\ \times & \times & \times & & & \\ & & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix},$$

ki ima prvi stolpec enak normiranemu prvemu stolpcu matrike  $N_k$ . Po množenju s  $P_1$  dobimo grbo velikosti  $2 \times 2$ , ki jo premikamo navzdol s Householderjevimi zrcaljenji. Poglejmo si, kako to naredimo v primeru matrike velikosti  $6 \times 6$ .

$$P_1 A_k P_1 = \begin{bmatrix} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ + & \times & \times & \times & \times & \times \\ + & + & \times & \times & \times & \times \\ & & & \times & \times & \times \\ & & & & \times & \times \end{bmatrix},$$

$$P_2 = \begin{bmatrix} 1 & & & & & \\ & \times & \times & \times & & \\ & \times & \times & \times & & \\ & \times & \times & \times & & \\ & & & & 1 & \\ & & & & & 1 \end{bmatrix}, \quad P_2 P_1 A_k P_1 P_2 = \begin{bmatrix} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ & \times & \times & \times & \times & \times \\ & + & \times & \times & \times & \times \\ & + & + & \times & \times & \times \\ & & & & \times & \times \end{bmatrix},$$

$$\begin{aligned}
 P_3 &= \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & \times & \times & \times & & \\ & & \times & \times & \times & & \\ & & \times & \times & \times & & \\ & & & & & & 1 \end{bmatrix}, & P_3 P_2 P_1 A_k P_1 P_2 P_3 &= \begin{bmatrix} \times & \times & \times & \times & \times & \times & \\ \times & \times & \times & \times & \times & \times & \\ & \times & \times & \times & \times & \times & \\ & & \times & \times & \times & \times & \\ & & & + & \times & \times & \\ & & & + & + & \times & \times \end{bmatrix}, \\
 P_4 &= \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & & \\ & & & \times & \times & \times & \\ & & & \times & \times & \times & \\ & & & \times & \times & \times & \end{bmatrix}, & P_4 P_3 P_2 P_1 A_k P_1 P_2 P_3 P_4 &= \begin{bmatrix} \times & \times & \times & \times & \times & \times & \\ \times & \times & \times & \times & \times & \times & \\ & \times & \times & \times & \times & \times & \\ & & \times & \times & \times & \times & \\ & & & \times & \times & \times & \\ & & & & + & \times & \times \end{bmatrix}, \\
 P_5 &= \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & & \\ & & & 1 & & & \\ & & & & \times & \times & \\ & & & & \times & \times & \end{bmatrix}, & P_5 \cdots P_1 A_k P_1 \cdots P_5 &= \begin{bmatrix} \times & \times & \times & \times & \times & \times & \\ \times & \times & \times & \times & \times & \times & \\ & \times & \times & \times & \times & \times & \\ & & \times & \times & \times & \times & \\ & & & \times & \times & \times & \\ & & & & \times & \times & \end{bmatrix}.
 \end{aligned}$$

V vsakem koraku grbo velikosti  $2 \times 2$  premaknemo za eno mesto navzdol, dokler je v zadnjem koraku ne izločimo iz matrike in nam ostane  $A_{k+1}$ . S pomočjo izreka o implicitnem Q lahko tako en korak QR z dvojnim premikom izvedemo v  $\mathcal{O}(n^2)$  operacijah. Podrobneje, en korak stane  $10n^2$  operacij, če računamo le  $A_{k+1}$ , in še dodatnih  $10n^2$  operacij, če posodobimo še matriko Q.

## Dodatna literatura

Za računanje lastnih vrednosti in vektorjev je na voljo obsežna literatura. V slovenščini je na voljo knjiga [8]. Kar se tiče tuje literature, lahko vse algoritme, ki smo jih navedli v tem poglavju, skupaj s potrebno analizo, najdete v [9]. Uporabna sta tudi učbenika [7] in [12].

## Poglavje 2

# Simetrični problem lastnih vrednosti

### 2.1 Uvod

Naj bo  $A$  simetrična matrika. Vemo, da so vse lastne vrednosti realne. Schurova forma za simetrično matriko je diagonalna matrika. Lastne vrednosti lahko uredimo tako, da velja

$$\lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_1,$$

pripadajoče lastne vektorje  $x_1, \dots, x_n$  pa lahko izberemo tako, da tvorijo ortonormirano bazo. Hitro vidimo, da velja

$$\|A\|_2 = \max(|\lambda_1|, |\lambda_n|).$$

Ker imamo opravka le z realnimi lastnimi vektorji, tudi *Rayleighov kvocient* računamo le za realne vektorje. Za  $x \neq 0$  je

$$\rho(x, A) = \frac{x^T A x}{x^T x}.$$

Že iz prejšnjega poglavja vemo, da za Rayleighov kvocient veljajo naslednje lastnosti:

- Za  $\alpha \neq 0$  je  $\rho(x, A) = \rho(\alpha x, A)$ .
- Za lastni vektor  $x_i$  je  $\rho(x_i, A) = \lambda_i$ .
- Če je  $x$  približek za lastni vektor, je  $\rho(x, A)$  najboljša aproksimacija za lastno vrednost v smislu, da je  $\min_{\sigma \in \mathbb{R}} \|Ax - \sigma x\|_2$  dosežen pri  $\sigma = \rho(x, A)$ .

Za simetrične matrike velja še naslednja lastnost.

**Lema 2.1** Če je  $A$  simetrična matrika z lastnimi vrednostmi  $\lambda_n \leq \dots \leq \lambda_1$ , potem za vsak  $x \neq 0$  velja

$$\lambda_n \leq \rho(x, A) \leq \lambda_1.$$

*Dokaz.* Če vektor  $x$  razvijemo po bazi lastnih vektorjev kot  $x = \sum_{i=1}^n \alpha_i x_i$ , dobimo

$$\rho(x, A) = \frac{\sum_{i=1}^n \alpha_i^2 \lambda_i}{\sum_{i=1}^n \alpha_i^2},$$

kar lahko očitno ocenimo navzgor in navzdol z  $\lambda_1$  oziroma  $\lambda_n$ . ■

Iz zgornje leme sledi, da bi lahko največjo in najmanjšo lastno vrednost izrazili z maksimumom oziroma minimumom Rayleighovega kvocienta po vseh neničelnih vektorjih. Kot pravi naslednji izrek, lahko podobno izrazimo tudi vse preostale lastne vrednosti. Rezultat je temelj za številne pomembne teoretične rezultate.

**Izrek 2.2 (Courant–Fischerjev minimaks izrek)** Če je  $A$  simetrična matrika z lastnimi vrednostmi  $\lambda_n \leq \dots \leq \lambda_1$ , potem za  $i = 1, \dots, n$  velja

$$\lambda_i = \min_{\substack{S \subset \mathbb{R}^n \\ \dim(S)=n-i+1}} \max_{\substack{x \in S \\ x \neq 0}} \rho(x, A) = \max_{\substack{R \subset \mathbb{R}^n \\ \dim(R)=i}} \min_{\substack{x \in R \\ x \neq 0}} \rho(x, A). \quad (2.1)$$

*Dokaz.* Za poljubna podprostor  $S, R \subset \mathbb{R}^n$ ,  $\dim(R) = i$  in  $\dim(S) = n - i + 1$ , obstaja  $x_{RS} \in R \cap S$ ,  $x_{RS} \neq 0$ , saj je  $\dim(R) + \dim(S) = n + 1$ . Očitno velja

$$\min_{\substack{x \in R \\ x \neq 0}} \rho(x, A) \leq \rho(x_{RS}, A) \leq \max_{\substack{x \in S \\ x \neq 0}} \rho(x, A).$$

Ker to velja za vsak par  $R, S$  velja tudi za par  $\tilde{R}, \tilde{S}$ , kjer je dosežen minimum oz. maksimum v izrazu (2.1). Torej

$$\max_{\substack{R \subset \mathbb{R}^n \\ \dim(R)=i}} \min_{\substack{x \in R \\ x \neq 0}} \rho(x, A) \leq \rho(x_{\tilde{R}\tilde{S}}, A) \leq \min_{\substack{S \subset \mathbb{R}^n \\ \dim(S)=n-i+1}} \max_{\substack{x \in S \\ x \neq 0}} \rho(x, A). \quad (2.2)$$

Po drugi strani pa za par  $\hat{R} = \text{Lin}(x_1, \dots, x_i)$  in  $\hat{S} = \text{Lin}(x_i, \dots, x_n)$  velja

$$\min_{\substack{x \in \hat{R} \\ x \neq 0}} \rho(x, A) = \lambda_i = \max_{\substack{x \in \hat{S} \\ x \neq 0}} \rho(x, A).$$

Od tod sledi

$$\max_{\substack{R \subset \mathbb{R}^n \\ \dim(R)=i}} \min_{\substack{x \in R \\ x \neq 0}} \rho(x, A) \geq \min_{\substack{S \subset \mathbb{R}^n \\ \dim(S)=n-i+1}} \max_{\substack{x \in S \\ x \neq 0}} \rho(x, A), \quad (2.3)$$

saj je

$$\max_{\substack{R \subset \mathbb{R}^n \\ \dim(R)=i}} \min_{\substack{x \in R \\ x \neq 0}} \rho(x, A) \geq \min_{\substack{x \in \hat{R} \\ x \neq 0}} \rho(x, A)$$

in podobno

$$\min_{\substack{S \subset \mathbb{R}^n \\ \dim(S)=n-i+1}} \max_{\substack{x \in S \\ x \neq 0}} \rho(x, A) \leq \max_{\substack{x \in \hat{S} \\ x \neq 0}} \rho(x, A).$$

Iz (2.2) in (2.3) sledi (2.1). ■

**Posledica 2.3** Če sta  $A$  in  $E$  simetrični matriki in so  $\lambda_n \leq \dots \leq \lambda_1$  lastne vrednosti matrike  $A$ ,  $\hat{\lambda}_n \leq \dots \leq \hat{\lambda}_1$  pa lastne vrednosti matrike  $A + E$ , potem za  $i = 1, \dots, n$  velja

$$\lambda_i + \lambda_n(E) \leq \hat{\lambda}_i \leq \lambda_i + \lambda_1(E).$$

*Dokaz.* Iz  $\rho(x, A + E) = \rho(x, A) + \rho(x, E)$  ocenimo

$$\rho(x, A) + \lambda_n(E) \leq \rho(x, A + E) \leq \rho(x, A) + \lambda_1(E)$$

in uporabimo Courant–Fischerjev minimaks izrek. ■

Od tod sledi naslednji, t.i. Weylov izrek. Če se vam zdi znan, to ni nič nenavadnega, saj smo ga kot posledico Bauer–Fikeovega izreka zapisali že v posledici 1.5.

**Posledica 2.4 (Weylov izrek)** Če sta  $A$  in  $E$  simetrični matriki in so  $\lambda_n \leq \dots \leq \lambda_1$  lastne vrednosti matrike  $A$ ,  $\hat{\lambda}_n \leq \dots \leq \hat{\lambda}_1$  pa lastne vrednosti matrike  $A + E$ , potem za  $i = 1, \dots, n$  velja

$$|\lambda_i - \hat{\lambda}_i| \leq \|E\|_2.$$

**Izrek 2.5 (Cauchyjev izrek o prepletanju)** Če je  $A$  simetrična matrika in je  $A_r$  njena vodilna podmatrika velikosti  $r \times r$  za  $r = 1, \dots, n - 1$ , potem velja

$$\lambda_{r+1}(A_{r+1}) \leq \lambda_r(A_r) \leq \lambda_r(A_{r+1}) \leq \dots \leq \lambda_2(A_{r+1}) \leq \lambda_1(A_r) \leq \lambda_1(A_{r+1}).$$

*Dokaz.* Dovolj je dokazati primer  $r = n - 1$ . Če je  $x' \in \mathbb{R}^{n-1}$  in  $x = \begin{bmatrix} x' \\ 0 \end{bmatrix} \in \mathbb{R}^n$ , potem je  $\rho(x', A_{n-1}) = \rho(x, A)$ .

Po Courant–Fischerjevem minimaks izreku velja

$$\lambda_k(A_{n-1}) = \min_{\substack{S' \subset \mathbb{R}^{n-1} \\ \dim(S')=n-k}} \max_{\substack{x' \in S' \\ x' \neq 0}} \rho(x', A_{n-1}).$$

Vsak podprostor lahko podamo z njegovim ortogonalnim komplementom, zato lahko pišemo

$$\lambda_k(A_{n-1}) = \min_{p'_1, \dots, p'_{k-1} \in \mathbb{R}^{n-1}} \max_{\substack{x' \in \mathbb{R}^{n-1}, x' \neq 0 \\ x' \perp p'_i, i=1, \dots, k-1}} \rho(x', A_{n-1}).$$

Pogoj, da morajo biti vektorji  $p'_1, \dots, p'_{k-1}$  linearno neodvisni, lahko izpustimo, saj je minimum očitno dosežen pri linearno neodvisnih vektorjih. Sedaj lahko pišemo

$$\begin{aligned} \lambda_k(A_{n-1}) &= \min_{p_1, \dots, p_{k-1} \in \mathbb{R}^{n-1}} \max_{\substack{x \in \mathbb{R}^n, x \neq 0, x \perp e_n \\ x \perp p_i, i=1, \dots, k-1}} \rho(x, A) \\ &\leq \min_{p_1, \dots, p_{k-1} \in \mathbb{R}^n} \max_{\substack{x \in \mathbb{R}^n, x \neq 0 \\ x \perp p_i, i=1, \dots, k-1}} \rho(x, A) = \lambda_k(A). \end{aligned}$$

Tako smo pokazali  $\lambda_k(A_{n-1}) \leq \lambda_k(A)$ . Po drugi strani pa velja

$$\begin{aligned} \lambda_k(A_{n-1}) &= \min_{p_1, \dots, p_{k-1} \in \mathbb{R}^{n-1}} \max_{\substack{x \in \mathbb{R}^n, x \neq 0, x \perp e_n \\ x \perp p_i, i=1, \dots, k-1}} \rho(x, A) \\ &= \min_{\substack{p_1, \dots, p_{k-1}, p_k \in \mathbb{R}^n \\ p_k = e_n}} \max_{\substack{x \in \mathbb{R}^n, x \neq 0 \\ x \perp p_i, i=1, \dots, k}} \rho(x, A) \\ &\geq \min_{p_1, \dots, p_{k-1}, p_k \in \mathbb{R}^n} \max_{\substack{x \in \mathbb{R}^n, x \neq 0 \\ x \perp p_i, i=1, \dots, k}} \rho(x, A) = \lambda_{k+1}(A). \end{aligned} \quad \blacksquare$$

Naslednji izrek pove, da lahko iz norme ostanka  $Ax - \beta x$  približka  $\beta$  za lastno vrednost in približka  $x$  za lastni vektor dobimo zelo dobro oceno, kako natančno  $\beta$  aproksimira točno lastno vrednost matrike  $A$ .

**Izrek 2.6** Če je  $A$  simetrična matrika,  $\|x\|_2 = 1$  in  $\beta$  približek za lastno vrednost, potem ima matrika  $A$  vsaj eno lastno vrednost  $\lambda_i$ , ki zadošča  $|\beta - \lambda_i| \leq \|Ax - \beta x\|_2$ .

*Dokaz.* Naj bo  $r = Ax - \beta x$ . Dokaz je podoben kot pri Bauer–Fikeovemu izreku. Predpostavimo lahko, da se  $\beta$  razlikuje od vseh lastnih vrednosti, kar pomeni, da je  $A - \beta I$  nesingularna. Matrika  $A$  se da diagonalizirati kot  $A = XDX^T$ , torej  $A - \beta I = X(D - \beta I)X^T$ . Sedaj je  $x = (A - \beta I)^{-1}r$ , kar pomeni  $1 \leq \|(D - \beta I)^{-1}\|_2 \|r\|_2$ , odtod pa dobimo

$$\min_i |\lambda_i - \beta| \leq \|r\|_2. \quad \blacksquare$$

**Definicija 2.7** Vsaki simetrični matriki  $A$  priredimo trojico celih števil  $(\nu, \zeta, \pi)$ , kjer je  $\nu$  število negativnih,  $\zeta$  število ničelnih in  $\pi$  število pozitivnih lastnih vrednosti matrike  $A$ . Omenjeno trojico imenujemo inercija matrike  $A$ .

Če je matrika  $X$  nesingularna, potem pravimo, da sta matriki  $A$  in  $X^TAX$  kongruentni. Izkaže se, da imajo kongruentne simetrične matrike isto inercijo.

**Izrek 2.8 (Sylvester)** Naj bo matrika  $A$  simetrična. Za poljubno nesingularno matriko  $X$  imata  $A$  in  $X^TAX$  isto inercijo.

*Dokaz.* Naj bo matrika  $A$  velikosti  $n \times n$ . Denimo, da se inerciji matrik  $A$  in  $X^TAX$  razlikujeta tako, da ima  $A$   $\nu$  negativnih,  $X^TAX$  pa  $\nu' < \nu$  negativnih lastnih vrednosti.

Naj bo  $\mathcal{N}$  invariantni podprostor dimenzije  $\nu$ , ki ga razpenjajo vsi lastni vektorji matrike  $A$ , ki pripadajo negativnim lastnim vrednostim. Podobno naj bo  $\mathcal{P}$  podprostor dimenzije  $n - \nu'$ , ki ga razpenjajo vsi lastni vektorji matrike  $X^TAX$ , ki pripadajo nenegativnim lastnim vrednostim. Zaradi nesingularnosti matrike  $X$  je podprostor  $X\mathcal{P}$  tudi dimenzije  $n - \nu'$ .

Ker je  $n - \nu' + \nu > n$ , obstaja neničelni vektor  $x \in \mathcal{N} \cap X\mathcal{P}$ . Zanj zaradi  $x \in \mathcal{N}$  velja  $x^T Ax < 0$ , po drugi strani pa zaradi  $x \in X\mathcal{P}$  obstaja tak  $y \in \mathcal{P}$ , da je  $x = Xy$ , torej  $x^T Ax = y^T (X^TAX)y \geq 0$ . Prišli smo do protislovja, kar pomeni, da je  $\nu = \nu'$ .

Podobno pokažemo, da velja tudi  $\pi = \pi'$ , od tod pa že sledi, da se inerciji ujemata. ■

S pomočjo Sylvestrovega izreka lahko dokažemo reletivni Weylov izrek, s katerim lahko bolj natančno omejimo spremembe lastnih vrednosti simetrične matrike  $A$ .

**Izrek 2.9 (Relativni Weylov izrek)** Naj bo  $A$  simetrična matrika z lastnimi vrednostmi  $\lambda_n \leq \dots \leq \lambda_1$ . Če je  $X$  nesingularna matrika in so lastne vrednosti matrike  $X^TAX$  enake  $\hat{\lambda}_n \leq \dots \leq \hat{\lambda}_1$ , potem za  $i = 1, \dots, n$  velja

$$|\hat{\lambda}_i - \lambda_i| \leq |\lambda_i| \epsilon,$$

kjer je  $\epsilon = \|X^T X - I\|_2$ .

*Dokaz.* Vemo, da je  $i$ -ta lastna vrednost matrike  $A - \lambda_i I$  enaka 0. Po Sylvestrovem izreku o ohranjanju inercije mora tudi  $i$ -ta lastna vrednost matrike  $X^T(A - \lambda_i I)X$  biti enaka 0. Če zapišemo

$$X^T(A - \lambda_i I)X = X^TAX - \lambda_i I + \lambda_i(I - X^T X),$$

potem iz Weylovega izreka sledi, da se  $i$ -ta lastna vrednost  $X^TAX - \lambda_i I$  razlikuje od  $i$ -te lastne vrednosti  $A - \lambda_i I$  za manj kot  $\|\lambda_i(I - X^T X)\|_2$ , to pa je ravno  $|\hat{\lambda}_i - \lambda_i| \leq |\lambda_i| \epsilon$ . ■

## 2.2 Rayleighova iteracija

Če inverzno iteracijo kombiniramo z Rayleighovim kvocientom, dobimo naslednji algoritem:

$$\begin{aligned} & \text{izberi } z_0 \neq 0 \\ & k = 0, 1, \dots \\ & \sigma_k = \rho(z_k, A) \\ & \text{reši } (A - \sigma_k I)y_{k+1} = z_k \\ & z_{k+1} = \frac{y_{k+1}}{\|y_{k+1}\|_2} \end{aligned}$$

Namesto fiksnega premika  $\sigma$  pri inverzni iteraciji uporabljamo Rayleighov kvocient, ki je najboljši približek za lastno vrednost danega vektorja.

Metoda ni omejena le na simetrične matrice, lahko jo uporabimo tudi za nesimetrične. Pokazati se da, da je konvergenca Rayleighove iteracije v bližini enostavne lastne vrednosti kvadratična, v primeru simetrične matrice pa celo kubična.

**Lema 2.10** Naj bo normiran vektor  $z_0$  približek za lastni vektor  $x_1$ . Če za simetrično matrico  $A$  z lastnimi vrednostmi  $|\lambda_n| \leq \dots \leq |\lambda_2| < |\lambda_1|$  izvedemo en korak potenčne metode z začetnim vektorjem  $z_0$ , potem velja

$$\|z_1 \pm x_1\|_2 \leq \frac{|\lambda_2|}{|\lambda_1|} \|z_0 - x_1\|_2,$$

kjer predznak  $\pm$  izberemo tako, da je norma razlike manjša.

*Dokaz.* Brez škode za splošnost lahko predpostavimo, da je  $\lambda_1 > 0$ . Če vektor  $z_0$  razvijemo po lastnih vektorjih, dobimo  $z_0 = \sum_{i=1}^n \alpha_i x_i$ , kjer je  $\alpha_1 \approx 1$ . Od tod sledi  $z_1 = \sum_{i=1}^n \beta_i x_i$ , kjer je

$$\beta_1 = \frac{\alpha_1 \lambda_1}{\left(\sum_{i=1}^n \alpha_i^2 \lambda_i^2\right)^{1/2}} = \frac{1}{\left(1 + \sum_{i=2}^n \left(\frac{\alpha_i}{\alpha_1}\right)^2 \left(\frac{\lambda_i}{\lambda_1}\right)^2\right)^{1/2}}.$$

Velja  $\|z_0 - x_1\|_2^2 = 2(1 - \alpha_1)$  in podobno  $\|z_1 \pm x_1\|_2^2 = 2(1 - \beta_1)$ . Sedaj lahko ocenimo

$$\beta_1 \approx 1 - \frac{1}{2} \sum_{i=2}^n \left(\frac{\alpha_i}{\alpha_1}\right)^2 \left(\frac{\lambda_i}{\lambda_1}\right)^2,$$

torej

$$1 - \beta_1 \approx \frac{1}{2} \sum_{i=2}^n \left(\frac{\alpha_i}{\alpha_1}\right)^2 \left(\frac{\lambda_i}{\lambda_1}\right)^2 \leq \frac{1}{2} \left(\frac{\lambda_2}{\lambda_1}\right)^2 \sum_{i=2}^n \alpha_i^2 = \frac{1}{2} \left(\frac{\lambda_2}{\lambda_1}\right)^2 (1 - \alpha_1^2).$$

Tako dobimo

$$\|z_1 \pm x_1\|_2^2 = 2(1 - \beta_1) \leq \left(\frac{\lambda_2}{\lambda_1}\right)^2 (1 + \alpha_1)(1 - \alpha_1) \leq \left(\frac{\lambda_2}{\lambda_1}\right)^2 2(1 - \alpha_1) = \left(\frac{\lambda_2}{\lambda_1}\right)^2 \|z_0 - x_1\|_2^2$$

in lema je dokazana. ■



**Posledica 2.11** Naj bo  $A$  simetrična matrika in  $|\lambda_i - \sigma| \ll |\lambda_j - \sigma|$  za  $j \neq i, j = 1, \dots, n$ . Če izvedemo en korak inverzne iteracije z začetnim vektorjem  $z_0$ , potem velja

$$\|z_1 - x_i\|_2 = \mathcal{O}(|\lambda_i - \sigma| \cdot \|z_0 - x_i\|_2).$$

**Lema 2.12** Naj bo  $A$  simetrična matrika in naj bo normiran vektor  $z$  približek za lastni vektor  $x_k$ . Potem velja

$$|\lambda_k - \rho(z, A)| \leq 2\|A\|_2 \|z - x_k\|_2^2.$$

*Dokaz.* Brez škode za splošnost lahko predpostavimo, da je  $k = 1$ . Če vektor  $z$  razvijemo po lastnih vektorjih kot  $z = \sum_{i=1}^n \alpha_i x_i$ , potem je  $\|z - x_1\|_2^2 = 2(1 - \alpha_1)$ . Za razliko Rayleighovega kvocienta dobimo

$$\lambda_1 - \rho(z, A) = \lambda_1 \sum_{i=1}^n \alpha_i^2 - \sum_{i=1}^n \lambda_i \alpha_i^2 = \sum_{i=2}^n (\lambda_1 - \lambda_i) \alpha_i^2$$

in ocenimo

$$|\lambda_1 - \rho(z, A)| \leq 2\|A\|_2 \sum_{i=2}^n \alpha_i^2 = 2\|A\|_2 (1 - \alpha_1^2) \leq 4\|A\|_2 (1 - \alpha_1) = 2\|A\|_2 \|z - x_1\|_2^2. \quad \blacksquare$$

**Izrek 2.13** Naj bo  $A$  simetrična matrika. Potem ima Rayleighova iteracija v bližini enostavne lastne vrednosti kubični red konvergence.

*Dokaz.* Naj Rayleighova iteracija konvergira k enostavni lastni vrednosti  $\lambda_k$  s pripadajočim lastnim vektorjem  $x_k$ . Po enem koraku Rayleighove iteracije po posledici 2.11 velja

$$\|z_1 - x_k\|_2 = \mathcal{O}(|\lambda_k - \sigma_0| \cdot \|z_0 - x_k\|_2).$$

Ker po lemi 2.12 velja tudi

$$|\lambda_k - \rho(z_0, A)| \leq 2\|A\|_2 \|z_0 - x_k\|_2^2,$$

oceni skupaj data  $\|z_1 - x_k\|_2 = \mathcal{O}(\|z_0 - x_k\|_2^3)$  in konvergenca je res kubična.  $\blacksquare$

**Zgled 2.1** Naj bo  $A = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$ ,  $\lambda_1 > \lambda_2$  in  $z_r = \begin{bmatrix} c_r \\ s_r \end{bmatrix}$ ,  $\|z_r\|_2^2 = c_r^2 + s_r^2 = 1$ . Pri enem koraku Rayleighove iteracije dobimo

$$\sigma_r = z_r^T A z_r = c_r^2 \lambda_1 + s_r^2 \lambda_2.$$

Iz sistema

$$\begin{bmatrix} \lambda_1 - c_r^2 \lambda_1 - s_r^2 \lambda_2 & 0 \\ 0 & \lambda_2 - c_r^2 \lambda_1 - s_r^2 \lambda_2 \end{bmatrix} y_{r+1} = z_r$$

dobimo

$$y_{r+1} = \frac{1}{(\lambda_1 - \lambda_2) c_r^2 s_r^2} \begin{bmatrix} c_r^3 \\ -s_r^3 \end{bmatrix},$$

od tod pa

$$z_{r+1} = \frac{1}{\sqrt{c_r^6 + s_r^6}} \begin{bmatrix} c_r^3 \\ -s_r^3 \end{bmatrix}.$$

Od tod v primeru  $s_r \neq c_r$  sledi kubična konvergenca proti  $e_1$  oziroma  $e_2$ .  $\square$

## 2.3 QR iteracija za simetrični lastni problem

V primeru simetrične matrike je zgornja Hessenbergova matrika tridiagonalna. Za redukcijo na tridiagonalno obliko sicer zaradi simetrije porabimo približno polovico toliko operacij kot za nesimetrično matriko, a to še vedno pomeni  $\mathcal{O}(n^3)$  operacij. Med samo QR iteracijo pa je razlika večja. Zaradi tridiagonalne oblike lahko en korak QR iteracije sedaj izvedemo z zahtevnostjo  $\mathcal{O}(n)$ , medtem ko imamo pri zgornji Hessenbergovi obliki, ki nastopa pri nesimetričnem primeru,  $\mathcal{O}(n^2)$  operacij.

Pri QR iteraciji torej najprej poiščemo ortogonalno matriko  $Q$ , da je  $T = QAQ^T$  tridiagonalna, potem pa delamo QR z enojnim premikom:

$$\begin{aligned} T_0 &= T \\ k &= 0, 1, \dots \\ &\text{izberi premik } \sigma_k \\ T_k - \sigma_k I &= Q_k R_k \text{ (izračunaj QR razcep)} \\ T_{k+1} &= R_k Q_k + \sigma_k I \end{aligned}$$

Naj bo

$$T_k = \begin{bmatrix} a_1^{(k)} & b_1^{(k)} & & & & \\ b_1^{(k)} & a_2^{(k)} & b_2^{(k)} & & & \\ & \ddots & \ddots & \ddots & & \\ & & b_{n-2}^{(k)} & a_{n-1}^{(k)} & b_{n-1}^{(k)} & \\ & & & b_{n-1}^{(k)} & a_n^{(k)} & \end{bmatrix}.$$

Kako izberemo premik:

- *Rayleighov premik*: vzamemo  $\sigma_k = a_n^{(k)} = \rho(e_n, T_k)$ . V tem primeru imamo za skoraj vse matrike zagotovljeno kubično konvergenco, a vseeno obstajajo primeri, ko metoda ne konvergira, če začetni približek ni dovolj dober.

**Zgled 2.2** Primer matrike, za katero QR iteracija z Rayleighovim premikom ne konvergira, je

$$T_0 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}. \text{ V tem primeru je } \sigma_0 = 0 \text{ in } T_1 = T_0. \quad \square$$

Kubično konvergenco v bližini enostavnih lastnih vrednosti nam zagotavlja izrek 2.14, ki povezuje Rayleighovo iteracijo in QR iteracijo.

- *Wilkinsonov premik*: za  $\sigma_k$  vzamemo tisto lastno vrednost matrike  $\begin{bmatrix} a_{n-1}^{(k)} & b_{n-1}^{(k)} \\ b_{n-1}^{(k)} & a_n^{(k)} \end{bmatrix}$ , ki je bližja  $a_n^{(k)}$ . Sedaj imamo za vse matrike dokazano vsaj linearno konvergenco, v praksi pa imamo za skoraj vse matrike kubično konvergenco (a brez dokaza).

**Izrek 2.14** Če je  $T$  simetrična in delamo QR iteracijo z Rayleighovimi premiki, potem so premiki  $\sigma_k$  enaki Rayleighovim kvocientom, ki jih dobimo pri Rayleighovi iteraciji za matriko  $T$ , če za začetni vektor vzamemo  $z_0 = e_n$ .

*Dokaz.* Najprej vpeljemo matrike  $\overline{Q}_k = Q_0 Q_1 \cdots Q_k$  in  $\overline{R}_k = R_k R_{k-1} \cdots R_0$ , kjer so  $Q_i$  in  $R_i$  matrike iz QR iteracije. Pokažimo, da za vsak  $k \geq 0$  velja

- a)  $T_{k+1} = \overline{Q}_k^T T \overline{Q}_k$ ,  
 b)  $(T - \sigma_0 I) \cdots (T - \sigma_k I) = \overline{Q}_k \overline{R}_k$ .

Točka a) sledi iz zveze  $T_{k+1} = Q_k^T T_k Q_k$ .

Pri točki b) uporabimo indukcijo. Za  $k = 0$  očitno velja, saj je  $T_0 - \sigma_0 I = Q_0 R_0$ . Pokažimo, da če velja za  $k - 1$ , potem velja tudi za  $k \geq 1$ . Res, produkt  $\overline{Q}_k \overline{R}_k$  lahko pišemo kot

$$\begin{aligned} \overline{Q}_k \overline{R}_k &= \overline{Q}_{k-1} Q_k R_k \overline{R}_{k-1} = \overline{Q}_{k-1} (T_k - \sigma_k I) \overline{R}_{k-1} \\ &= \overline{Q}_{k-1} \overline{Q}_{k-1}^T (T - \sigma_k I) \overline{Q}_{k-1} \overline{R}_{k-1} = (T - \sigma_k I) \overline{Q}_{k-1} \overline{R}_{k-1} \\ &= (T - \sigma_0 I) \cdots (T - \sigma_{k-1} I) (T - \sigma_k I). \end{aligned}$$

Pri izpeljavi smo upoštevali točko a) in dejstvo, da matrike oblike  $T - \sigma_i I$  med seboj komutirajo.

Sedaj lahko dokažemo, da za vsak  $k$  velja  $\rho(e_n, T_k) = \rho(z_k, T)$ . To je očitno res za  $k = 0$ . Vektor  $z_{k+1}$  iz Rayleighove iteracije zadošča enačbi  $(T - \sigma_k I) z_{k+1} = \alpha_k z_k$ , kjer skalar  $\alpha_k$  izberemo tako, da bo  $\|z_{k+1}\|_2 = 1$ . Od tod rekurzivno sledi

$$(T - \sigma_0 I) \cdots (T - \sigma_k I) z_{k+1} = \beta_k e_n$$

za primerno izbran skalar  $\beta_k$ . Po točki b) to lahko zapišemo kot

$$\overline{R}_k^T \overline{Q}_k^T z_{k+1} = \beta_k e_n,$$

kjer smo upoštevali, da zaradi simetrije matrike  $T$  velja  $\overline{Q}_k \overline{R}_k = \overline{R}_k^T \overline{Q}_k^T$ . To pa pomeni, da je

$$z_{k+1} = \beta_k \overline{Q}_k \overline{R}_k^{-T} e_n.$$

Ker je matrika  $\overline{R}_k^{-T}$  spodnja trikotna, ima vektor  $\overline{R}_k^{-T} e_n$  smer  $e_n$ , od koder sledi  $z_{k+1} = \overline{Q}_k e_n$ . Od tod iz točke a) sledi

$$\rho(z_{k+1}, T) = \rho(\overline{Q}_k e_n, T) = \rho(e_n, \overline{Q}_k^T T \overline{Q}_k) = \rho(e_n, T_{k+1})$$

in izrek je dokazan. ■

## 2.4 Sturmovo zaporedje

Naj bo

$$T = \begin{bmatrix} a_1 & b_1 & & & & \\ b_1 & a_2 & b_2 & & & \\ & \ddots & \ddots & \ddots & & \\ & & b_{n-2} & a_{n-1} & b_{n-1} & \\ & & & b_{n-1} & a_n & \end{bmatrix}$$

ireducibilna tridiagonalna simetrična matrika (torej  $b_i \neq 0$  za vsak  $i$ ). Če s  $T_r$  označimo njeno vodilno  $r \times r$  podmatriko in definiramo  $f_r(\lambda) = \det(T_r - \lambda I)$ , potem z razvijanjem po zadnji vrstici pridemo do rekurzivne formule

$$f_{r+1}(\lambda) = (a_{r+1} - \lambda)f_r(\lambda) - b_r^2 f_{r-1}(\lambda) \quad (2.4)$$

za  $r = 0, \dots, n-1$ , ki se začne z  $f_0(\lambda) \equiv 1$  in  $f_1(\lambda) = a_1 - \lambda$ .

**Izrek 2.15** *Polinomi  $f_0, \dots, f_n$  tvorijo Sturmovo zaporedje, kar pomeni, da zadoščajo naslednjim trem točkam:*

- 1)  $f_0(\lambda) \neq 0$  za vsak  $\lambda$ .
- 2) Če je  $f_r(\lambda_0) = 0$  za  $r < n$ , potem je  $f_{r-1}(\lambda_0)f_{r+1}(\lambda_0) < 0$ .
- 3) Če je  $f_n(\lambda_0) = 0$ , potem je  $f_{n-1}(\lambda_0)f'_n(\lambda_0) < 0$ .

*Dokaz.* Točka 1) je očitna. Pri točki 2) iz rekurzivne formule sledi  $f_{r+1}(\lambda_0) = -b_r^2 f_{r-1}(\lambda_0)$ . Obe vrednosti sta neničelni, saj bi sicer veljalo  $f_i(\lambda_0) = 0$  za  $i = 0, \dots, n$ , to pa je v protislovju s točko 1).

Pri točki 3) definiramo  $\Delta_r(\lambda_0) = f_r(\lambda_0)f'_{r-1}(\lambda_0) - f_{r-1}(\lambda_0)f'_r(\lambda_0)$ . Z računanjem lahko hitro preverimo, da velja

$$\Delta_{r+1}(\lambda_0) = f_r^2(\lambda_0) + b_r^2 \Delta_r(\lambda_0).$$

Ker je  $\Delta_1(\lambda_0) = 1 > 0$ , je  $\Delta_r(\lambda_0) > 0$  za vsak  $r$ . Torej tudi  $\Delta_n(\lambda_0) = -f_{n-1}(\lambda_0)f'_n(\lambda_0) > 0$ . ■

**Posledica 2.16** *Ireducibilna tridiagonalna simetrična matrika ima enostavne lastne vrednosti.*

*Dokaz.* To sledi iz točke 3) zadnjega izreka, saj v primeru  $f_n(\lambda_0) = 0$  velja  $f'_n(\lambda_0) \neq 0$ . ■

Pri fiksnem  $\lambda_0$  označimo z  $u(\lambda_0)$  število ujemanj predznaka v zaporedju  $f_0(\lambda_0), \dots, f_n(\lambda_0)$ . Pri tem vsako notranjo ničlo štejemo za eno ujemanje, ničlo na koncu pa ne. Primeri:

$$\begin{aligned} u(+ + - - +) &= 2, \\ u(+ + 0 - +) &= 2, \\ u(+ + 0 - + 0) &= 2. \end{aligned}$$

**Lema 2.17** *Število  $u(\lambda_0)$  je enako številu lastnih vrednosti matrike  $T$ , ki so strogo večje od  $\lambda_0$ .*

*Dokaz.* Naj  $\lambda$  teče od  $-\infty$  do  $\infty$ . Pri  $\lambda = -\infty$  imamo očitno zaporedje  $+ + + \dots$ , pri  $\lambda = \infty$  pa  $+ - + - \dots$ . Tako je  $u(-\infty) = n$  in  $u(\infty) = 0$ .

Pokazali bomo, da se število  $u(\lambda)$  lahko spremeni le, če prečkamo ničlo polinoma  $f_n$ , ne pa tudi, če prečkamo ničlo polinoma  $f_r$ ,  $r < n$ .

Naj bo  $f_r(\lambda_0) = 0$ ,  $r < n$ . Potem je iz tabele

	$\lambda_0 - \epsilon$	$\lambda_0$	$\lambda_0 + \epsilon$
$f_{r-1}$	$\pm$	$\pm$	$\pm$
$f_r$	$\pm$	0	$\mp$
$f_{r+1}$	$\mp$	$\mp$	$\mp$

razvidno, da pri zadosti majhnem  $\epsilon > 0$  velja  $u(\lambda_0 - \epsilon) = u(\lambda_0 + \epsilon)$ .

V primeru  $f_n(\lambda_0) = 0$  pa iz tabele

	$\lambda_0 - \epsilon$	$\lambda_0$	$\lambda_0 + \epsilon$
$f_{n-1}$	$\pm$	$\pm$	$\pm$
$f_n$	$\pm$	0	$\mp$

vidimo, da pri zadosti majhnem  $\epsilon > 0$  velja  $u(\lambda_0 - \epsilon) = u(\lambda_0 + \epsilon) + 1$ . ■

Sedaj lahko z bisekcijo ali kakšno drugo metodo poiščemo  $k$ -to lastno vrednost. Iščemo točko  $\lambda_k$ , za katero velja  $u(\lambda_k - \epsilon) = k$  in  $u(\lambda_k + \epsilon) = k - 1$  za dovolj majhen  $\epsilon > 0$ . Če nimamo boljše ocene za začetni interval, lahko vzamemo  $[-\|T\|, \|T\|]$  za poljubno normo matrike  $T$ .

Metoda je uporabna tudi, če nas zanimajo samo lastne vrednosti na določenem intervalu ali pa nekaj lastnih vrednosti. Tako je npr. razlika  $u(\beta) - u(\alpha)$  enaka število lastnih vrednosti matrike  $T$ , ki so na intervalu  $(\alpha, \beta]$ .

Za izračun  $u(\lambda)$  preko tričlenske rekurzivne formule potrebujemo  $4n + \mathcal{O}(1)$  operacij, pri čemer smo predpostavili, da kvadrate obdiagonalnih elementov izračunamo vnaprej. Izkaže pa se, da lahko naredimo še bolje. Če definiramo kvociente  $d_i(\lambda) = f_i(\lambda)/f_{i-1}(\lambda)$ , potem zanje velja  $d_1(\lambda) = a_1 - \lambda$  in

$$d_{r+1}(\lambda) = a_{r+1} - \lambda - \frac{b_r^2}{d_r(\lambda)}. \quad (2.5)$$

Za izračun te formule potrebujemo le  $3n + \mathcal{O}(1)$  operacij, namesto zaporednih ujemanj predznakov pa zdaj preštejemo, koliko vrednosti  $d_1(\lambda), \dots, d_n(\lambda)$  je pozitivnih.

Matrika  $T - \lambda I$  je tridiagonalna. Denimo, da za to matriko obstaja razcep  $LDL^T$ , kjer je  $L$  spodnja trikotna bidiagonalna matrika z enicami na diagonali oblike

$$L = \begin{bmatrix} 1 & & & & \\ l_1 & 1 & & & \\ & \ddots & \ddots & & \\ & & & l_{n-1} & 1 \end{bmatrix},$$

matrika  $D = \text{diag}(d_1, \dots, d_n)$  pa je diagonalna. S primerjanjem elementov matrik  $T - \lambda I$  in  $LDL^T$  pridemo do naslednjega algoritma za izračun matrik  $L$  in  $D$ .

---

**Algoritem 2.1**  $LDL^T$  razcep matrike  $T - \lambda I$ .

---

$$d_1 = a_1 - \lambda$$

$$i = 1, \dots, n - 1$$

$$l_i = b_i / d_i$$

$$d_{i+1} = a_{i+1} - \lambda - l_{i-1}^2 d_{i-1}$$


---

Iz zgornjega algoritma dobimo formulo (2.5), če formulo za  $l_i$  vstavimo v formulo za  $d_{i+1}$ . Ker ne pivotiramo, se lahko v (2.5) pojavi deljenje z 0. Če računamo v aritmetiki, ki zadošča standardom IEEE, to ni nobena težava. Sploh je računanje po formuli (2.5) zelo stabilno, saj velja naslednji izrek, katerega dokaz lahko najdete npr. v [8].

**Lema 2.18** Če ne pride do prekoračitve ali podkoračitve, se predznaki numerično izračunanih diagonalnih elementov  $\hat{d}_1, \dots, \hat{d}_n$  ujemajo s predznaki točnih  $d_1, \dots, d_n$  iz  $LDL^T$  razcepa matrike  $T + \delta T$ , kjer je  $\delta a_k = 0$  in  $|\delta b_k| \leq (5/2)|b_k|u$  za  $k = 1, \dots, n$ .

## 2.5 Deli in vladaj

To je trenutno najhitrejša metoda za ireducibilno simetrično tridiagonalno matriko. Naj bo

$$T = \begin{bmatrix} a_1 & b_1 & & & \\ b_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & b_{n-1} & a_n \\ & & & & \end{bmatrix}.$$

Za  $m < n$  razdelimo matriko  $T$  kot

$$T = \begin{bmatrix} T_1 & 0 \\ 0 & T_2 \end{bmatrix} + b_m v v^T,$$

kjer je

$$T_1 = \begin{bmatrix} a_1 & b_1 & & & \\ b_1 & \ddots & \ddots & & \\ & \ddots & \ddots & a_{m-1} & b_{m-1} \\ & & & b_{m-1} & a_m - b_m \end{bmatrix}, \quad T_2 = \begin{bmatrix} a_{m+1} - b_m & b_{m+1} & & & \\ b_{m+1} & a_{m+2} & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & b_{n-1} & a_n \end{bmatrix}$$

in

$$v = e_m + e_{m+1} = [0 \ \dots \ 0 \ 1 \ 1 \ 0 \ \dots \ 0]^T.$$

$T_1$  in  $T_2$  sta simetrični tridiagonalni matriki, zato obstajata ortogonalni matriki  $Q_1, Q_2$  in diagonalni matriki  $D_1, D_2$ , da je  $T_1 = Q_1 D_1 Q_1^T$  in  $T_2 = Q_2 D_2 Q_2^T$ . Potem je

$$T = \begin{bmatrix} Q_1 & 0 \\ 0 & Q_2 \end{bmatrix} \left( \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix} + b_m u u^T \right) \begin{bmatrix} Q_1^T & 0 \\ 0 & Q_2^T \end{bmatrix},$$

kjer je

$$u = \begin{bmatrix} Q_1^T & 0 \\ 0 & Q_2^T \end{bmatrix} v = \begin{bmatrix} Q_1(m, :)^T \\ Q_2(1, :)^T \end{bmatrix}.$$

Lastne vrednosti  $T$  so tako enake lastnim vrednostim  $D + \rho u u^T$ , kjer je  $D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$  in  $\rho = b_m$ . Če izračunamo lastne vrednosti in vektorje  $D + \rho u u^T = Q' \Lambda Q'^T$ , potem so na diagonali  $\Lambda$  lastne vrednosti matrike  $T$ ,  $Q = \begin{bmatrix} Q_1 & 0 \\ 0 & Q_2 \end{bmatrix} Q'$  pa so ustrezni lastni vektorji.

Osnovna verzija metode deli in vladaj je zapisana v algoritmu 2.5. Zapisati moramo še, kako lahko ekonomično in natančno izračunamo lastne vrednosti in vektorje matrike  $D + \rho u u^T$ . Za samo računanje uredimo diagonalne elemente matrike  $D$  tako, da je  $d_1 \geq d_2 \geq \dots \geq d_n$ . Seveda moramo ustrezno preurediti tudi elemente vektorja  $u$ .

Najprej izločimo vse lastne vrednosti in vektorje, kjer lahko uporabimo naslednjo lemo.

**Lema 2.19** Naj bo  $A = D + \rho u u^T$ , kjer je  $D = \text{diag}(d_1, \dots, d_n)$ ,  $d_1 \geq d_2 \geq \dots \geq d_n$  in  $u = [u_1 \ \dots \ u_n]^T$ .

- a) Če je  $d_i = d_{i+1}$ , je  $d_i$  lastna vrednost matrike  $A$ , lastni vektor pa je  $[0 \ \dots \ 0 \ -u_{i+1} \ u_i \ 0 \ \dots \ 0]^T$ .

---

**Algoritem 2.2** Osnovna varianta metode deli in vladaj. Začetni podatek je ireducibilna tridiagonalna simetrična matrika  $T$ . Algoritem vrne diagonalno matriko  $\Lambda$  in ortogonalno matriko  $Q$  da je  $T = Q\Lambda Q^T$ .

---

$[Q, \Lambda] = \text{deliinvladaj}(T)$

če je  $T$  velikosti  $1 \times 1$ , potem vrni  $Q = 1$ ,  $\Lambda = T$

sicer:

razdeli  $T = \begin{bmatrix} T_1 & 0 \\ 0 & T_2 \end{bmatrix} + b_m v v^T$ .

$[Q_1, D_1] = \text{deliinvladaj}(T_1)$

$[Q_2, D_2] = \text{deliinvladaj}(T_2)$

iz  $Q_1, Q_2, D_1, D_2$  izračunaj  $D + \rho u u^T$

izračunaj lastne vrednosti  $\Lambda$  in vektorje  $Q'$  za  $D + \rho u u^T$

$Q = \begin{bmatrix} Q_1 & 0 \\ 0 & Q_2 \end{bmatrix} Q'$ .

---

b) Če je  $u_i = 0$ , je  $d_i$  lastna vrednost matrike  $A$ , lastni vektor pa je  $e_i$ .

V nadaljevanju lahko predpostavimo, da so vsi  $d_i$  različni in da je  $u_i \neq 0$  za  $i = 1, \dots, n$ . Predpostavimo še, da je  $\lambda$  lastna vrednost matrike  $D + \rho u u^T$ , matrika  $D - \lambda I$  pa je nesingularna. Potem je

$$\det(D + \rho u u^T - \lambda I) = \det\left((D - \lambda I)(I + \rho(D - \lambda I)^{-1} u u^T)\right),$$

od tod pa sledi

$$\det(I + \rho(D - \lambda I)^{-1} u u^T) = 0.$$

Ker vemo, da je  $\det(I + x y^T) = 1 + y^T x$ , so lastne vrednosti  $D + \rho u u^T$  ničle sekularne enačbe  $f(\lambda) = 0$ , kjer je

$$\begin{aligned} f(\lambda) &= \det(I + \rho(D - \lambda I)^{-1} u u^T) = 1 + \rho u^T (D - \lambda I)^{-1} u \\ &= 1 + \rho \sum_{i=1}^n \frac{u_i^2}{d_i - \lambda}. \end{aligned}$$

Kako zgleda graf  $f(\lambda)$ ? Asimptota je  $y = 1$ . Ker je  $f'(\lambda) = \rho \sum_{i=1}^n \frac{u_i^2}{(d_i - \lambda)^2}$ , je za  $\rho > 0$  funkcija strogo naraščajoča (med poli), sicer pa padajoča. Ničle ležijo med poli, ena pa desno od zadnjega pola (pri  $\rho > 0$ ) ali levo od prvega pola (pri  $\rho < 0$ ).

Denimo, da z neko numerično metodo, npr. Newtonovo, poiščemo lastne vrednosti. Za lastne vektorje potem velja:

**Lema 2.20** Če je  $\alpha$  lastna vrednost  $D + \rho u u^T$ , je  $(D - \alpha I)^{-1} u$  ustrezní lastni vektor.

*Dokaz.*

$$\begin{aligned} (D + \rho u u^T)(D - \alpha I)^{-1} u &= (D - \alpha I + \alpha I + \rho u u^T)(D - \alpha I)^{-1} u \\ &= u + \alpha(D - \alpha I)^{-1} u + u(\rho u^T (D - \alpha I)^{-1} u) \\ &= u + \alpha(D - \alpha I)^{-1} u - u \end{aligned}$$

$$= \alpha(D - \alpha I)^{-1}u,$$

saj iz  $f(\alpha) = 1 + \rho u^T(D - \alpha I)^{-1}u = 0$  sledi  $\rho u^T(D - \alpha I)^{-1}u = -1$ . ■

Ker rešujemo diagonalni sistem, lahko vsak lastni vektor izračunamo v času  $\mathcal{O}(n)$ .

Za konec si podrobno pogledjmo, kako rešujemo sekularno enačbo. Navadne tangentne metode ne moremo uporabiti, saj so ničle lahko zelo blizu polov in bi potrebovali zelo dobre približke. Namesto aproksimacije funkcije s tangento zato raje uporabimo preprosto racionalno funkcijo, ki se prilega funkciji  $f$ .

Denimo, da iščemo rešitev na intervalu  $(d_{i+1}, d_i)$ , začetni približek pa je  $x_r$ . Sedaj poiščemo racionalno funkcijo oblike

$$h(\lambda) = \frac{c_1}{d_i - \lambda} + \frac{c_2}{d_{i+1} - \lambda} + c_3,$$

za katero velja  $h(x_r) = f(x_r)$  in  $f'(x_r) = h'(x_r)$ . Zaradi stabilnosti razdelimo  $f$  na dva dela kot

$$f(\lambda) = 1 + \rho \sum_{k=1}^i \frac{u_k^2}{d_k - \lambda} + \rho \sum_{k=i+1}^n \frac{u_k^2}{d_k - \lambda} =: 1 + \psi_1(\lambda) + \psi_2(\lambda).$$

To naredimo zato, da v vsoti za  $\psi_1(\lambda)$  oziroma  $\psi_2(\lambda)$  seštevamo enako predznačene člene. Sedaj določimo  $c_1, c'_1$  tako, da za

$$h_1(\lambda) = \frac{c_1}{d_i - \lambda} + c'_1$$

velja  $h_1(x_r) = \psi_1(x_r)$  in  $h'_1(x_r) = \psi'_1(x_r)$ . Podobno določimo  $c_2, c'_2$  tako, da za

$$h_2(\lambda) = \frac{c_2}{d_{i+1} - \lambda} + c'_2$$

velja  $h_2(x_r) = \psi_2(x_r)$  in  $h'_2(x_r) = \psi'_2(x_r)$ . Sedaj je  $h(\lambda) = 1 + h_1(\lambda) + h_2(\lambda)$  iskana racionalna funkcija. Enačba  $h(\lambda) = 0$  ima dve rešitvi, za  $x_{r+1}$  pa vzamemo tisto, ki leži znotraj  $(d_{i+1}, d_i)$ . Konvergenca je zelo hitra, saj  $h$  zelo dobro aproksimira  $f$  na intervalu  $(d_{i+1}, d_i)$ .

Pri numeričnem računanju se izkaže, da so lastni vektorji, ki jih izračunamo preko leme 2.20 bolj slabo ortogonalni, zato je potrebno algoritem popraviti. Pri popravku vektor  $u$  nadomestimo z bližnjim  $\hat{u}$ , za katerega je potem vse v redu, vse pa temelji na naslednjem izreku.

**Izrek 2.21 (Löwner)** Naj bo  $D = \text{diag}(d_1, \dots, d_n)$ , kjer je  $d_n < d_{n-1} < \dots < d_1$ , in naj bodo  $\alpha_n < \alpha_{n-1} < \dots < \alpha_1$  dana števila, ki se prepletajo s števili  $d_i$ :

$$d_n < \alpha_n < \dots < d_1 < \alpha_1.$$

Potem za vektor  $\hat{u}$ , podan z

$$|\hat{u}_i| = \left( \frac{\prod_{j=1}^n (\alpha_j - d_i)}{\prod_{j=1, j \neq i}^n (d_j - d_i)} \right)^{1/2} \quad (2.6)$$

velja, da so  $\alpha_i$  točne lastne vrednosti matrike  $\hat{D} = D + \hat{u}\hat{u}^T$ .

*Dokaz.* Po eni strani je karakteristični polinom matrike  $\hat{D}$  enak  $\det(\hat{D} - \lambda I) = \prod_{j=1}^n (\alpha_j - \lambda)$ , po drugi strani pa zaradi  $\hat{D} - \lambda I = (D - \lambda I)(I + (D - \lambda I)^{-1}\hat{u}\hat{u}^T)$  velja

$$\det(\hat{D} - \lambda I) = \prod_{j=1}^n (d_j - \lambda) \left( 1 + \sum_{j=1}^n \frac{\hat{u}_j^2}{d_j - \lambda} \right).$$



Ko vstavimo  $\lambda = d_i$  in izenačimo izraza, dobimo

$$\prod_{j=1}^n (\alpha_j - d_i) = \hat{u}_i^2 \prod_{\substack{j=1 \\ j \neq i}}^n (d_j - d_i).$$

Od tod sledi

$$\hat{u}_i^2 = (\alpha_i - d_i) \prod_{\substack{j=1 \\ j \neq i}}^n \left( \frac{\alpha_j - d_i}{d_j - d_i} \right).$$

Zaradi prepletanja se predznak  $\alpha_j - d_i$  ujema s predznakom  $d_j - d_i$  za  $j \neq i$ . Ker velja še  $\alpha_i > d_i$ , je izraz na desni strani res pozitiven in ga lahko korenimo. ■

**Opomba 2.1** Zaradi enostavnosti smo v zadnjem izreku privzeli, da je  $\rho = 1$ . V primeru, ko je  $\rho > 0$ , lahko  $u$  in  $\rho$  vedno tako normiramo, da je to res. Na podoben način bi lahko pokazali, da formula (2.6) drži tudi za negativni  $\rho$ , ko lahko privzamemo, da je  $\rho = -1$  in so  $\alpha_1, \dots, \alpha_n$  točne lastne vrednosti matrike  $\hat{D} = D - \hat{u}\hat{u}^T$ . V tem primeru se števila prepletajo kot  $\alpha_n < d_n < \dots < \alpha_1 < d_1$ .

Stabilno računanje lastnih vrednosti in vektorjev matrike  $D + \rho uu^T$ , kjer je  $\rho \pm 1$ , poteka na naslednji način:

- Iz sekularne enačbe izračunamo lastne vrednosti  $\alpha_1, \dots, \alpha_n$  matrike  $D + \rho uu^T$ .
- Po Löwnerjevem izreku izračunamo vektor  $\hat{u}$ , da so  $\alpha_1, \dots, \alpha_n$  točne lastne vrednosti matrike  $D + \rho \hat{u}\hat{u}^T$ . Predznake izberemo tako, da se predznak  $\hat{u}_i$  ujema s predznakom  $u_i$  za  $i = 1, \dots, n$ .
- Za lastne vektorje vzamemo vektorje  $(D - \alpha_i I)^{-1} \hat{u}$  za  $i = 1, \dots, n$ .

Naj bo  $T(n)$  število operacij, ki jih porabi algoritem za matriko velikosti  $n$ . Potem velja

$$\begin{aligned} T(n) &= 2T(n/2) && \text{(rekurzivni klic dveh podproblemov)} \\ &+ \mathcal{O}(n^2) && \text{(izračun lastnih vrednosti } D + \rho uu^T) \\ &+ \mathcal{O}(n^2) && \text{(izračun lastnih vektorjev } D + \rho uu^T) \\ &+ n^3. && \text{(izračun produkta za matriko } Q) \end{aligned}$$

Od tod sledi, da je časovna zahtevnost algoritma za izračun vseh lastnih vrednosti in vektorjev enaka  $T(n) = (4/3)n^3 + \mathcal{O}(n^2)$ , v praksi pa je še manjša, saj lahko velikokrat lastne vektorje in vrednosti izračunamo preko leme 2.19, kar nam močno olajša izračun produkta za matriko  $Q$ .

## 2.6 Jacobijeva metoda

Pri tej metodi matrike na začetku ne reduciramo na tridiagonalno obliko. Ideja je, da matriko  $A$  z množenji z Givensovimi rotacijami z leve in desne poskusimo spraviti čim bližje diagonalni matriki.<sup>1</sup>

<sup>1</sup>To je najstarejša numerična metoda za računanje lastnih vrednosti. Pruski matematik Carl Gustav Jacob Jacobi (1804-1851) jo je razvil leta 1846. Ob pojavi prvih računalnikov se je izkazalo, da je metodo zelo preprosto implementirati, zato je bila na začetku to glavna metoda za računanje lastnih vrednosti simetričnih matrik, potem pa jo je izpodrinila hitrejša QR iteracija.

Denimo, da bi radi v matriki  $A$  uničili element na mestu  $pq$ . Najprej poiščemo rotacijo  $R = \begin{bmatrix} c & -s \\ s & c \end{bmatrix}$ , da bo veljalo

$$R^T \begin{bmatrix} a_{pp} & a_{pq} \\ a_{pq} & a_{qq} \end{bmatrix} R = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a_{pp} & a_{pq} \\ a_{pq} & a_{qq} \end{bmatrix} \begin{bmatrix} c & -s \\ s & c \end{bmatrix} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}. \quad (2.7)$$

Iz zveze  $(a_{qq} - a_{pp})sc + a_{pq}(c^2 - s^2) = 0$  dobimo

$$\tau := \frac{\cos 2\varphi}{\sin 2\varphi} = \frac{c^2 - s^2}{2sc} = \frac{a_{pp} - a_{qq}}{2a_{pq}}.$$

Če definiramo  $t := \frac{s}{c} = \tan \varphi$ , potem velja (uporabimo formule za tangens dvojnega kota)

$$t^2 + 2\tau t - 1 = 0.$$

Rešitev je

$$t = \frac{\text{sign}(\tau)}{|\tau| + \sqrt{1 + \tau^2}}, \quad c = \sqrt{\frac{1}{1 + t^2}}, \quad s = ct.$$

Tako smo dobili formule za izračun Jacobijeve rotacije  $\text{jac}(A, p, q)$ , ki v  $A$  uniči element  $a_{pq}$ .

Algoritem za eno rotacijo je

$$\begin{aligned} A &= \text{jac}(A, p, q) \\ \tau &= \frac{a_{pp} - a_{qq}}{2a_{pq}} \\ t &= \frac{\text{sign}(\tau)}{|\tau| + \sqrt{1 + \tau^2}} \\ c &= \sqrt{\frac{1}{1 + t^2}} \\ s &= ct \\ A &= R_{pq}(\varphi)^T A R_{pq}(\varphi) \\ J &= J R_{pq}(\varphi) \quad (\text{če potrebujemo tudi lastne vektorje}) \end{aligned}$$

Po rotaciji  $(p, q)$  se v  $A$  spremenita vrstica  $p$  in  $q$  ter stolpca  $p$  in  $q$ . To pomeni, da se ničle, ki smo jih naredili v prejšnjih korakih, lahko s tem zrcaljenjem pokvarijo. Kljub temu pa se norma izvediodagonalnih elementov z rotacijami zmanjšuje.

**Definicija 2.22** Za  $n \times n$  matriko  $A$  definiramo

$$\text{off}(A) = \sqrt{\sum_{\substack{j,k=1 \\ j \neq k}}^n |a_{jk}|^2}.$$

Vidimo, da je  $\text{off}(A)$  v bistvu Frobeniusova norma matrike  $A$  brez diagonale.

**Lema 2.23** Če  $A'$  dobimo iz  $A$  z Jacobijevo rotacijo  $\text{jac}(A, p, q)$ , potem velja

$$\text{off}(A')^2 = \text{off}(A)^2 - 2a_{pq}^2.$$

*Dokaz.* Iz  $A' = R_{pq}^T A R_{pq}$  sledi  $\|A'\|_F = \|A\|_F$ . Za diagonalne elemente  $A'$  velja  $a'_{ii} = a_{ii}$  za  $i \neq p, q$ , za preostala dva elementa pa zaradi (2.7) velja  $a'_{pp}{}^2 + a'_{qq}{}^2 = a_{pp}{}^2 + a_{qq}{}^2 + 2a_{pq}{}^2$ . Torej mora za izvendiagonalne elemente veljati  $\text{off}(A')^2 = \text{off}(A)^2 - 2a_{pq}{}^2$ . ■

Z vsako Jacobijevo rotacijo se tako zmanjša  $\text{off}(A)$ . Postopek ponavljamo, dokler ni  $\text{off}(A) \leq \epsilon$ .

Kako uničujemo elemente:

- *klasična varianta*: v vsakem koraku poiščemo po absolutni vrednosti največji izvendiagonalni element in ga uničimo. Sicer imamo po številu rotacij res najhitrejšo konvergenco, a imamo veliko primerjanj, ki povečajo časovno zahtevnost metode. Če si shranjujemo po absolutni vrednosti največje elemente v vsakem stolpcu, potem lahko največji izvendiagonalni element poiščemo v času  $\mathcal{O}(n)$ , prav toliko pa tudi porabimo v vsakem koraku, da posodobimo seznam največjih elementov. Iskanje torej poveča časovno zahtevnost, a ima isti red kot sicer porabimo za en korak Jacobijeve metode.
- *ciklična varianta*: v vedno enakem vrstnem redu gremo skozi vse elemente. Tu nimamo primerjanj, lahko pa zaradi uničevanja elementov, ki so majhni po absolutni vrednosti porabimo veliko korakov.
- *pragovna varianta*: v vedno enakem vrstnem redu gremo skozi vse elemente, a uničimo le tiste elemente, ki so po absolutni vrednosti čez neko mejo, ki jo zmanjšamo v vsakem sprehodu.

Jacobijeva metoda porabi več operacij kot QR ali deli in vladaj, njena prednost pa je, da za simetrične pozitivno definitne matrike lastne vrednosti blizu 0 izračuna relativno bolj natančno od ostalih metod.

## 2.7 Relativno robustne reprezentacije

Denimo, da imamo tridiagonalno simetrično matriko  $T$ , za katero bi radi izračunali vse lastne pare. S pomočjo algoritmov, ki smo jih že spoznali, to najhitreje izvedemo tako, da najprej izračunamo samo lastne vrednosti preko QR iteracije, za kar porabimo  $\mathcal{O}(n^2)$  operacij. Potem lahko preko inverzne iteracije in  $\mathcal{O}(n^2)$  operacij izračunamo enega po enega še vse lastne vektorje. Končni algoritem bi tako porabil le  $\mathcal{O}(n^2)$  operacij. Vendar, tako izračunani lastni vektorji niso tako ortogonalni, kot bi morali biti. Lahko je potrebna naknadna ortogonalizacija, ki pa ima zahtevnost  $\mathcal{O}(n^3)$  operacij.

V tem razdelku bomo na grobo predstavili algoritem, ki zna s pomočjo RRR (relativno robustnih reprezentacij) ta problem rešiti v  $\mathcal{O}(n^2)$  operacij.<sup>2</sup>

Začnimo s problemom, kako za tridiagonalno ireducibilno matriko natančno izračunati lastni vektor za znano lastno vrednost. Naj bo  $\lambda$  lastna vrednost ireducibilne simetrične tridiagonalne matrike

$$T = \begin{bmatrix} a_1 & b_1 & & & \\ b_1 & \ddots & \ddots & & \\ & \ddots & \ddots & & \\ & & & b_{n-1} & \\ & & & b_{n-1} & a_n \end{bmatrix}.$$

<sup>2</sup>Algoritem je predstavil indijski matematik Inderjit S. Dhillon v svojem doktoratu leta 1997.

Naj bo  $z$  lastni vektor, ki ga iščemo in zadošča  $(T - \lambda I)z = 0$ . Ker je  $T$  ireducibilna, sta prvi in zadnji element vektorja  $z$  neničelna, torej  $z_1 \neq 0$  in  $z_n \neq 0$ . Lastni vektor bi lahko določili tako, da fiksiramo  $z_1 = 1$  in opustimo zadnjo enačbo, ki ji z avtomatično zadošča. Algoritem je

$$\begin{aligned} z_1 &= 1 \\ z_2 &= -(a_1 - \lambda)/b_1 \\ k &= 2, \dots, n-1 \\ z_{k+1} &= -(b_{k-1}z_{k-1} + (a_k - \lambda)z_k)/b_k \end{aligned}$$

Če namesto točne lastne vrednosti  $\lambda$  uporabimo približek  $\tilde{\lambda}$ , potem tako izračunani vektor  $z$  ne zadošča zadnji enačbi in velja  $(T - \tilde{\lambda}I)z = \delta_n e_n$  za nek  $\delta_n$ .

Podobno, če opustimo  $k$ -to enačbo, za ostanek velja  $(T - \tilde{\lambda}I)z = \delta_k e_k$ . Opustitev  $k$ -te enačbe je ekvivalentna temu, da kot približek za lastni vektor vzamemo  $z = (T - \tilde{\lambda}I)^{-1}e_k$ . Vprašanje je, katero enačbo je najbolje opustiti oziroma kateri vektor  $(T - \tilde{\lambda}I)^{-1}e_i$ , kjer je  $i = 1, \dots, n$ , je najboljši približek za lastni vektor.

Naj bodo  $x_1, \dots, x_n$  lastni vektorji matrike  $T$ ,  $\tilde{\lambda}$  pa naj bo približek za lastno vrednost  $\lambda_j$ . Iz razvoja  $e_k = \sum_{i=1}^n (x_i)_k x_i$  sledi

$$z^{(k)} = (T - \tilde{\lambda}I)^{-1}e_k = \frac{(x_j)_k}{\lambda_j - \tilde{\lambda}} \left( x_j + \sum_{i \neq j} \frac{(x_i)_k}{(x_j)_k} \cdot \frac{\lambda_j - \tilde{\lambda}}{\lambda_i - \tilde{\lambda}} x_i \right).$$

Vidimo, da bo  $z^{(k)}$  dober približek za lastni vektor  $x_j$ , če

- opustimo  $k$ -to enačbo, kjer ima lastni vektor  $x_j$  maksimalno komponento,
- velja  $|\lambda_j - \tilde{\lambda}| \ll |\lambda_i - \tilde{\lambda}|$  za  $i \neq j$ .

### 2.7.1 Zasukani razcep

V tem podrazdelku bomo obravnavali zasukani razcep, s katerim lahko dobro ocenimo, katera komponenta lastnega vektorja je maksimalna.

Za matriko  $T - \lambda I$  lahko izračunamo razcepa  $T - \lambda I = LD^+L^T = UD^-U^T$ , kjer je

$$\begin{aligned} L &= \begin{bmatrix} 1 & & & & \\ l_1 & 1 & & & \\ & \ddots & \ddots & & \\ & & l_{n-1} & 1 & \\ & & & & 1 \end{bmatrix}, & D^+ &= \begin{bmatrix} d_1^+ & & & & \\ & d_2^+ & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & d_n^+ \end{bmatrix}, \\ U &= \begin{bmatrix} 1 & u_1 & & & \\ & \ddots & \ddots & & \\ & & 1 & u_{n-1} & \\ & & & & 1 \end{bmatrix}, & D^- &= \begin{bmatrix} d_1^- & & & & \\ & d_2^- & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & d_n^- \end{bmatrix}. \end{aligned}$$

Ustrezna algoritma sta



---

**Algoritem 2.5** Reševanje  $(T - \tilde{\lambda}I)z = \gamma_k e_k$  preko zasukanega razcepa  $T - \tilde{\lambda}I = N_k D_k N_k^T$ .

---

$$\begin{aligned} z_k &= 1 \\ j &= k-1, \dots, 1 \\ z_j &= -l_j z_{j+1} && (\text{v primeru } z_{j+1} = 0 \text{ pa } z_j = -b_{j+1} z_{j+2} / b_j) \\ j &= k+1, \dots, n \\ z_j &= -u_{j-1} z_{j-1} && (\text{v primeru } z_{j-1} = 0 \text{ pa } z_j = -b_{j-2} z_{j-2} / b_{j-1}) \end{aligned}$$


---

Iz zasukanega razcepa lahko tudi enostavno izračunamo približek za lastni vektor. Če je  $T - \tilde{\lambda}I = N_k D_k N_k^T$ , potem ima sistem  $(T - \tilde{\lambda}I)z = \gamma_k e_k$  enolično rešitev, ki jo dobimo s postopkom, opisanim v algoritmu 2.5.

Naj bo  $\tilde{\lambda}$  približek za lastno vrednost matrike  $T$ . Pripadajoči lastni vektor  $z$  izračunamo po naslednjem postopku, katerega zahtevnost je  $12n + \mathcal{O}(1)$  operacij.

---

**Algoritem 2.6** Izračun lastnega vektorja preko zasukanega razcepa.

---

- 1) Izračunaj razcepa  $T - \tilde{\lambda}I = LD^+L^T$  in  $T - \tilde{\lambda}I = UD^-U^T$ .
  - 2) Izračunaj  $\gamma_1, \dots, \gamma_n$  za zasukane razcepe in poišči minimalni  $|\gamma_k|$ .
  - 3) Reši sistem  $(T - \tilde{\lambda}I)z = \gamma_k e_k$  z algoritmom 2.5 in normiraj  $z$ .
- 

Opisani postopek lahko izračuna zelo natančne približke za lastne vektorje, če je lastna vrednost dobro separirana od ostalih. Če to ni res, potem lahko z ustreznim premikom  $\sigma$  poskrbimo, da bo v matriki  $T - \sigma I$  lastna vrednost  $\lambda_j - \sigma$  dobro separirana od ostalih. Tu pa je pomembno, da imamo matriko predstavljeno tako, da se s takim premikom ne pokvari relativna natančnost lastnih vrednosti. Kako to naredimo, bomo obravnavali v naslednjem podrazdelku.

### 2.7.2 Relativno robustna reprezentacija matrike

Tridiagonalno matriko  $T$  lahko predstavimo z vektorjema diagonalnih in obdiagonalnih elementov  $a$  in  $b$ . Majhne relativne spremembe  $a$  in  $b$  lahko povzročijo velike relativne spremembe lastnih vrednosti in vektorjev, zato ta predstavitev ni *relativno robustna*.

Če je  $T$  pozitivno definitna, jo lahko predstavimo z vektorjema diagonalnih in obdiagonalnih elementov faktorja Choleskega. Majhne relativne spremembe teh elementov povzročijo majhne relativne spremembe lastnih parov, zato je takšna predstavitev relativno robustna. Točna definicija relativno robustne reprezentacije je naslednja:

**Definicija 2.26** Množica števil  $\{p_i\}$ , ki določa matriko  $T$ , je relativno robustna reprezentacija (RRR), če se lastni pari matrike  $T + \delta T$ , ki jo določa množica zmotenih elementov  $\{p_i(1 + \epsilon_i)\}$ , relativno malo razlikujejo od lastnih parov matrike  $T$ .

Če reprezentacija z visoko relativno natančnostjo določa le lastne vrednosti  $(\lambda_j, \dots, \lambda_k)$ , pravimo da je delna RRR( $j, \dots, k$ ).

Če je matrika  $T$  pozitivno definitna, potem lahko namesto razcepa Choleskega uporabimo  $LDL^T$  razcep, ki je prav tako RRR. S tem se izognemo kvadratnim korenem, ki so računsko zahtevnejši od osnovnih štirih operacij.

Če je  $T$  nedefinitna, lahko uporabimo razcep  $T - \sigma I = LDL^T$  za ustrezní premik  $\sigma$ . Ob primerno izbranem  $\sigma$  dosežemo, da je ta razcep delna RRR za iskane lastne vrednosti.

Medtem ko lastne vrednosti lahko izračunamo z visoko relativno natančnostjo npr. z bisekcijo, pa za natančne lastne vektorje potrebujemo, da je lastna vrednost relativno dobro izolirana, kar lahko dosežemo s premiki.

Denimo, da poznamo RRR matrike  $T + \tau I$  oblike  $LDL^T$ , radi pa bi izračunali razcepa  $L^+D^+L^{+T}$  in  $U^-D^-U^{-T}$  za premaknjeno matriko  $T + \tau I - \mu I = LDL^T - \mu I$ . Če to naredimo za naslednjimi algoritmi, potem ohranimo relativno robustno reprezentacijo.

---

### Algoritem 2.7 Prehod RRR.

---

#### dstqds algoritem:

$$\begin{aligned} s_1 &= -\mu \\ i &= 1, \dots, n-1 \\ d_i^+ &= s_i + d_i \\ l_i^+ &= l_i d_i / d_i^+ \\ s_{i+1} &= l_i^+ l_i s_i - \mu \\ d_n^+ &= s_n + d_n \end{aligned}$$

#### dqds algoritem:

$$\begin{aligned} p_n &= d_n - \mu \\ i &= n-1, \dots, 1 \\ d_{i+1}^- &= d_i l_i^2 + p_{i+1} \\ t &= d_i / d_{i+1}^- \\ u_i^- &= l_i t \\ p_i &= p_{i+1} t - \mu \\ d_1^- &= p_1 \end{aligned}$$

#### izračun $\gamma$ :

$$\begin{aligned} i &= 1, \dots, n \\ \gamma_i &= s_i + p_i + \mu \end{aligned}$$


---

Naj bo  $\tilde{\lambda}$  približek za lastno vrednost matrike  $LDL^T$ , ki je RRR. Potem lastni vektor dobimo z naslednjim algoritmom.

---

### Algoritem 2.8 GetVec.

---

- 1) Preko dstqds izračunaj razcep  $LDL^T - \tilde{\lambda}I = L^+D^+L^{+T}$
  - 2) Preko dqds izračunaj razcep  $LDL^T - \tilde{\lambda}I = U^-D^-U^{-T}$
  - 3) Izračunaj  $\gamma_1, \dots, \gamma_n$  za zasukane razcepe in poišči minimalni  $|\gamma_k|$
  - 4) Reši sistem  $N_k D_k N_k^T z = \gamma_k e_k$  in normiraj  $z$
- 

Naj bo  $\tilde{\lambda}$  približek za lastno vrednost  $\lambda_j$  matrike  $T = LDL^T$ . Za lastni vektor  $z$ , izračunan preko

inverzne iteracije, velja

$$|\sin \angle(z, x_j)| \leq \frac{\mathcal{O}(nu\|T\|)}{\text{gap}(\tilde{\lambda})},$$

kjer je  $\text{gap}(\tilde{\lambda}) = \min_{i \neq j} |\tilde{\lambda} - \lambda_i|$ . Za lastni vektor, izračunan preko RRR, pa velja

$$|\sin \angle(z, x_j)| \leq \frac{\mathcal{O}(nu)}{\text{relgap}(\tilde{\lambda})},$$

kjer je  $\text{relgap}(\tilde{\lambda}) = \frac{\text{gap}(\tilde{\lambda})}{|\tilde{\lambda}|}$ . Kadar je  $|\tilde{\lambda}|$  dosti manjša kot  $\|T\|$ , je RRR lahko veliko natančnejši.

### 2.7.3 MRRR

Končni algoritem so večkratne relativno robustne reprezentacije. S premiki in preračuni RRR dosežemo, da bodo lastne vrednosti, za katere računamo lastne vrednosti, dobro relativno separirane od ostalih.

---

#### Algoritem 2.9 MRRR

---

- 1) Poišči premik  $\tau$  da bo  $LDL^T$  razcep matrike  $T + \tau I$  RRR.
  - 2) Izračunaj  $T + \tau I = L_0 D_0 L_0^T$
  - 3) Relativno natančno izračunaj lastne vrednosti  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$  matrike  $L_0 D_0 L_0^T$ .
  - 4)  $l = 1, m = n$
  - 5) Razdeli lastne vrednosti  $\tilde{\lambda}_l, \dots, \tilde{\lambda}_m$  na izolirane in gruče.
  - 6) Za vsako izolirano lastno vrednost  $\tilde{\lambda}_j$  izračunaj  $z_j$  preko GetVec.
  - 7) Za vsako gručo  $\tilde{\lambda}_j, \dots, \tilde{\lambda}_{j+k-1}$ 
    - Preko dstqds izračunaj  $L_0 D_0 L_0^T - \tau_s I = L_S D_S L_S^T$ , kjer  $\tau_s$  izbereš tako, da
      - bo v  $L_S D_S L_S^T$  vsaj ena lastna vrednost izolirana,
      - bo  $L_S D_S L_S^T$  delna RRR( $j, \dots, j+k-1$ ).
    - Relativno natančno izračunaj lastne vrednosti  $\mu_j, \dots, \mu_{j+k-1}$  matrike  $L_S D_S L_S^T$ .
    - Določi  $\tilde{\lambda}_i = \mu_i$  za  $i = j, \dots, j+k-1$
    - $l = j, m = j+k-1, L_0 = L_S, D_0 = D_S$ , rekurzivno se vrni na korak 5).
- 

Za lastno vrednost  $\tilde{\lambda}_j$  pravimo, da je relativno izolirana, če velja

$$\text{relgap}(\tilde{\lambda}_j) \geq \delta.$$

Lastne vrednosti  $\tilde{\lambda}_j, \dots, \tilde{\lambda}_{j+k-1}$  tvorijo gručo, če velja

$$\text{reldis}(\tilde{\lambda}_i, \tilde{\lambda}_{i+1}) < \delta \quad \text{za } i = j, \dots, j+k-2,$$



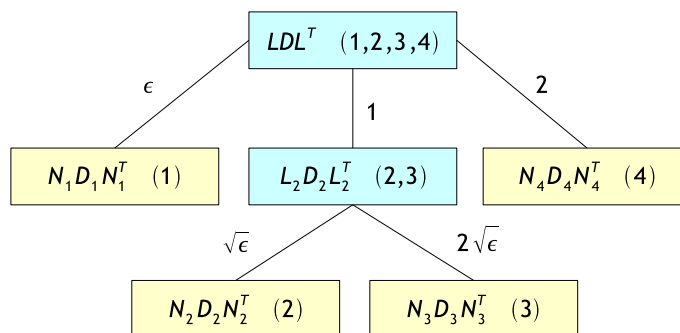
$$\begin{aligned} \text{reldis}(\tilde{\lambda}_{j-1}, \tilde{\lambda}_j) &\geq \delta, \\ \text{reldis}(\tilde{\lambda}_{j+k-1}, \tilde{\lambda}_{j+k}) &\geq \delta, \end{aligned}$$

kjer je

$$\text{reldis}(\lambda, \mu) = \frac{|\lambda - \mu|}{|\lambda|}.$$

Za  $\delta$  vzamemo npr.  $10^{-3}$ .

**Zgled 2.3** Denimo, da ima matrika lastne vrednosti  $\epsilon$ ,  $1 + \sqrt{\epsilon}$ ,  $1 + 2\sqrt{\epsilon}$ , 2. Slika 2.1 predstavlja, kako poteka v tem primeru izračun lastnih vektorjev preko MRRR. Vsak list, ki jih je toliko kot lastnih vrednosti, predstavlja en izračun lastnega vektorja po metodi GetVec. Vsaka povezava predstavlja prehod na delno RRR premaknjene matrike, kjer uporabimo algoritem 2.7, povezava pa je tudi označena s premikom, ki ga uporabimo.  $\square$



Slika 2.1: Zgled predstavitevne diagrama metode MRRR.

Zahtevnost končne metode za izračun  $k$  lastnih parov je  $\mathcal{O}(kn)$ . Predvideva se, da bo MRRR sčasoma postal osnovna metoda za reševanje simetričnega problema lastnih vrednosti. Metoda je večinoma hitrejša od ostalih in potrebuje najmanj dodatnega spomina.

## Dodatna literatura

Za računanje lastnih vrednosti je na voljo obsežna literatura. V slovenščini lahko skoraj celotno snov najdete v knjigi [8]. Kar se tiče tuje literature, lahko skoraj vse algoritme, ki smo jih navedli, skupaj s potrebno analizo, najdete v [9]. Uporabna sta tudi učbenika [7] in [12].

Relativno robustne reprezentacije so lepo opisane v diplomskem delu [20].

## Poglavje 3

# Posplošitve problema lastnih vrednosti

### 3.1 Posplošeni problem lastnih vrednosti

**Definicija 3.1** Dani sta kvadratni matriki  $A$  in  $B$ . Množico vseh matrik oblike  $A - \lambda B$ , kjer je  $\lambda \in \mathbb{C}$ , imenujemo matrični šop in označimo z  $(A, B)$  ali  $A - \lambda B$ . Karakteristični polinom matričnega šopa  $(A, B)$  je  $p(\lambda) = \det(A - \lambda B)$ . Če njegov karakteristični polinom ni identično enak 0, potem je matrični šop regularen, sicer pa singularen.

Če je matrični šop  $(A, B)$  regularen in je

$$Ax = \lambda Bx$$

za neničelni vektor  $x$ , potem pravimo, da je  $\lambda$  (končna) lastna vrednost in  $x$  (desni) lastni vektor. Podobno je neničelni vektor  $y$  levi lastni vektor za  $\lambda$ , če je  $y^H A = \lambda y^H B$ .

Problemu iskanja lastnih vrednosti matričnega šopa pravimo *posplošeni problem lastnih vrednosti*. Če je  $(A, B)$  regularen matrični šop, potem so njegove končne lastne vrednosti ničle karakterističnega polinoma  $\det(A - \lambda B)$ , ki je stopnje  $m \leq n$ . V primeru  $m < n$  ima šop še lastno vrednost  $\infty$  z večkratnostjo  $n - m$ .

**Zgled 3.1** V primeru

$$A - \lambda B = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 0 \end{bmatrix} - \lambda \begin{bmatrix} 2 & & \\ & 0 & \\ & & 1 \end{bmatrix}$$

dobimo  $p(\lambda) = (2\lambda - 1)\lambda$ , torej so lastne vrednosti  $\lambda_1 = 1/2$ ,  $\lambda_2 = 0$  in  $\lambda_3 = \infty$ . □

Neskončne lastne vrednosti se pojavijo natanko takrat, ko je matrika  $B$  singularna. Vsak neničelni vektor iz  $\ker(B)$  je desni lastni vektor za lastno vrednost  $\infty$ .

**Izrek 3.2** Za regularen matrični šop  $(A, B)$  velja:

- 1) Če je matrika  $B$  nesingularna, so vse lastne vrednosti šopa  $(A, B)$  končne in enake lastnim vrednostim matrike  $B^{-1}A$  ali  $AB^{-1}$ .

- 2) Če je matrika  $B$  singularna, ima šop  $(A, B)$  lastno vrednost  $\infty$  z večkratnostjo  $n - \text{rang}(B)$ .
- 3) Če je matrika  $A$  nesingularna, so lastne vrednosti šopa  $(A, B)$  recipročne lastne vrednosti matrike  $A^{-1}B$  oziroma  $BA^{-1}$ , kjer lastna vrednost  $0$  ustreza neskončni lastni vrednosti  $(A, B)$ .

**Definicija 3.3** Če sta matriki  $U$  in  $V$  nesingularni, sta matrična šopa  $(A, B)$  in  $(UAV, UBV)$  ekvivalentna.

**Izrek 3.4** Ekvivalentna regularna matrična šopa  $(A, B)$  in  $(UAV, UBV)$  imata iste lastne vrednosti. Za lastne vektorje velja:

- a)  $x$  je desni lastni vektor za  $(A, B)$  natanko tedaj, ko je  $V^{-1}x$  desni lastni vektor za  $(UAV, UBV)$ ,
- b)  $y$  je levi lastni vektor za  $(A, B)$  natanko tedaj, ko je  $U^{-H}y$  levi lastni vektor za  $(UAV, UBV)$ .

Posplošitev Jordanove forme za regularne matrične šope je Weierstrassova forma. Za vsak regularen matrični šop  $(A, B)$  obstajata nesingularni matriki  $U$  in  $V$ , da je

$$U(A - \lambda B)V = \text{diag}(J_{n_1}(\lambda_1) - \lambda I_{n_1}, \dots, J_{n_k}(\lambda_k) - \lambda I_{n_k}, N_{m_1}, \dots, N_{m_l}),$$

kjer je

$$J_{n_i}(\lambda_i) = \begin{bmatrix} \lambda_i & 1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & \lambda_i \end{bmatrix}, \quad N_{m_i} = \begin{bmatrix} 1 & -\lambda & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & -\lambda \\ & & & & 1 \end{bmatrix}.$$

Blok  $J_{n_i}(\lambda_i)$  pripada končni lastni vrednosti  $\lambda_i$ , blok  $N_{m_i}$  pa neskončni lastni vrednosti. Opazimo lahko, da je  $N_{m_i} = I_{m_i} - \lambda J_{m_i}(0)$ . Tako kot Jordanova forma je tudi Weierstrassova neprimerna za numerično računanje. Stabilnejša je posplošitev Schurove forme, o kateri govori naslednji izrek.

**Izrek 3.5** Za vsak regularen matrični šop  $A - \lambda B$  obstajata unitarni matriki  $Q$  in  $Z$ , da je

$$Q^H(A - \lambda B)Z = S - \lambda T,$$

kjer sta matriki  $S$  in  $T$  zgoraj trikotni matriki. Lastne vrednosti so potem kvocienti  $\lambda_i = t_{ii}/s_{ii}$  za  $s_{ii} \neq 0$  in  $\infty$  v primeru  $s_{ii} = 0$ .

*Dokaz.* Tako, kot pri dokazu o obstoju navadne Schurove forme, uporabimo indukcijo. Za matrike velikosti  $1 \times 1$  je forma trivialna. Denimo, da obstaja za vse matrike velikosti  $(n-1) \times (n-1)$ , matriki  $A$  in  $B$  pa sta velikosti  $n \times n$ .

Ker je šop  $(A, B)$  regularen, ima vsaj eno lastno vrednost  $\lambda$ , pri čemer je možno tudi  $\lambda = \infty$ . Naj bo  $x$  normiran desni lastni vektor za  $\lambda$ . Iz  $Ax = \lambda Bx$  sledi, da sta vektorja  $Ax$  in  $Bx$  kolinearna. Ker oba hkrati ne moreta biti enaka nič, obstaja tak normiran vektor  $y$ , da  $Ax$  in  $Bx$  ležita v podprostoru, ki ga razpenja  $y$ .

Naj bo  $X = [x \ X_1]$  taka unitarna matriki, da je njen prvi stolpec enak vektorju  $x$ , podobno naj velja za matriko  $Y = [y \ Y_1]$  in vektor  $y$ . Potem velja

$$Y^H A X = \begin{bmatrix} \alpha & a^T \\ 0 & A_1 \end{bmatrix}, \quad Y^H B X = \begin{bmatrix} \beta & b^T \\ 0 & B_1 \end{bmatrix}.$$

Če je lastna vrednost  $\lambda$  končna, potem je  $\lambda = \alpha/\beta$ , sicer pa je  $\beta = 0$ . Ker po indukcijski predpostavki za šop  $(A_1, B_1)$  obstaja posplošena Schurova forma, nam to omogoča zapisati Schurovo formo za šop  $(A, B)$ . ■

Situacija  $s_{ii} = t_{ii} = 0$  je možna le, če je matrični šop  $(A, B)$  singularen.

Če sta matriki  $A$  in  $B$  realni, potem obstaja *realna posplošena Schurova forma*, kjer sta matriki  $Q$  in  $Z$  ortogonalni,  $S$  je kvazi zgornja trikotna,  $T$  pa zgornja trikotna matrika. Tako se v primeru realnih matrik tudi pri posplošenem problemu lastnih vrednosti lahko izognemo kompleksni aritmetiki.

Kaj lahko povemo o občutljivosti lastnih vrednosti posplošenega problema lastnih vrednosti? Da lahko enakovredno v analizo vključimo še neskončne lastne vrednosti, uporabljamo ločno razdaljo, definirano z

$$\chi(\alpha, \beta) = \frac{|\alpha - \beta|}{\sqrt{1 + |\alpha|^2} \sqrt{1 + |\beta|^2}},$$

za poljubni kompleksni števili  $\alpha, \beta$ . V limiti dobimo

$$\chi(\alpha, \infty) = \frac{1}{\sqrt{1 + |\alpha|^2}}.$$

**Izrek 3.6** Naj bo  $\lambda$  enostavna lastna vrednost šopa  $(A, B)$  z normiranim desnim lastnim vektorjem  $x$  in levim  $y$ . Če je  $\tilde{\lambda}$  ustrezna lastna vrednost zmotenega šopa  $(\tilde{A}, \tilde{B})$ , kjer je  $\|A - \tilde{A}\|_2 \leq \epsilon$ ,  $\|B - \tilde{B}\|_2 \leq \epsilon$ , potem je

$$\chi(\lambda, \tilde{\lambda}) \leq \frac{\epsilon}{|y^H A x|^2 + |y^H B x|^2} + \mathcal{O}(\epsilon^2).$$

## 3.2 QZ algoritem

Za numerično računanje posplošene Schurove forme imamo na voljo QZ algoritem, ki je izpeljanka QR metode za nesimetričen problem lastnih vrednosti. Na začetku šop  $(A, B)$  z ortogonalnimi ekvivalentnimi transformacijami reduciramo na šop  $(\tilde{A}, \tilde{B}) = (\tilde{Q}A\tilde{Z}, \tilde{Q}B\tilde{Z})$ , kjer je prva matrika zgornja Hessenbergova, druga pa zgornja trikotna.

Poglejmo, kako naredimo to redukcijo v primeru, ko sta matriki velikosti  $4 \times 4$ . Najprej poiščemo tako ortogonalno matriko  $P$ , ki pretvori matriko  $B$  v zgornjo trikotno obliko, torej

$$A_1 = PA = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}, \quad B_1 = PB = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times \end{bmatrix}.$$

Matriko  $P$  lahko najenostavneje skonstruiramo s pomočjo Householderjevih zrcaljenj. Nato z Givensovimi rotacijami enega po enega uničujemo elemente pod subdiagonalo matrike  $A$  in hkrati popravljamo matriko  $B$  nazaj na trikotno obliko. Najprej poiščemo rotacijo  $R_{34}$ , da je

$$A_2 = R_{34}^T A_1 = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, \quad B_2 = R_{34}^T B_1 = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}.$$

Element v matriki  $B_2$  uničimo z Givensovo rotacijo z desne, ki ne spremeni oblike matrike  $A_2$ :

$$A_3 = A_2 \tilde{R}_{34} = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, \quad B_3 = B_2 \tilde{R}_{34} = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times \end{bmatrix}.$$

Podobno uničimo naslednji element v matriki  $A$  in nato popravimo matriko  $B$ , dobimo

$$A_4 = R_{23}^T A_3 = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, \quad B_4 = R_{23}^T B_3 = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & 0 & \times \end{bmatrix},$$

$$A_5 = A_4 \tilde{R}_{23} = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, \quad B_5 = B_4 \tilde{R}_{23} = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times \end{bmatrix}.$$

V matriki  $A$  gremo na drugi stolpec in uničimo še element na mestu  $(4, 2)$ .

$$A_6 = R_{34}^T A_5 = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}, \quad B_6 = R_{34}^T B_5 = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix},$$

$$\tilde{A} = A_6 \tilde{R}_{34} = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}, \quad \tilde{B} = B_6 \tilde{R}_{34} = \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times \end{bmatrix}.$$

Tako smo s ortogonalnimi transformacijami pretvorili matriki  $A$  in  $B$  v zgornjo Hessenbergovo in trikotno obliko. Za splošni  $n \times n$  matriki je časovna zahtevnost  $8n^3 + \mathcal{O}(n^2)$  operacij za izračun  $\tilde{A}$  in  $\tilde{B}$ , še dodatnih  $7n^3 + \mathcal{O}(n^2)$  pa potrebujemo za izračun prehodnih ortogonalnih matrik  $\tilde{Q}$  in  $\tilde{Z}$ , ki jih dobimo iz produktov uporabljenih zrcaljenj in rotacij.

Za nadaljevanje predpostavimo, da smo redukcijo že naredili in je torej matrika  $A$  zgornja Hessenbergova, matrika  $B$  pa zgornja trikotna. Pri QZ iteraciji v bistvu izvajamo implicitno QR metodo na matriki  $C = AB^{-1}$ , ki je zgornja Hessenbergova.

Začnemo z  $(A_0, B_0) = (A, B)$ , potem pa v vsakem koraku posodobimo matrični šop  $(A_k, B_k)$  v

$$(A_{k+1}, B_{k+1}) = (Q_k A_k Z_k, Q_k B_k Z_k),$$

kjer ortogonalni matriki  $Q_k$  in  $Z_k$  določimo tako, da je  $A_{k+1}$  zgornja Hessenbergova in  $B_{k+1}$  zgornja trikotna, matrika  $A_{k+1} B_{k+1}^{-1}$  pa se ujema z matriko, ki bi jo dobili, če bi izvedli en korak QR iteracije za matriko  $A_k B_k^{-1}$ .

Velja  $A_{k+1} B_{k+1}^{-1} = Q_k (A_k B_k^{-1}) Q_k^T$ . Če je  $A_{k+1}$  zgornja Hessenbergova matrika,  $B_{k+1}$  zgornja trikotna matrika in se prvi stolpec  $Q_k$  ujema s prvim stolpcem ustrezne matrike pri QR iteraciji za matriko  $A_k B_k^{-1}$ , potem nam izrek o implicitnem Q zagotavlja, da je to ekvivalentno metodi QR na  $AB^{-1}$ .

Naj bo  $C_k = A_k B_k^{-1}$ . Za izračun Francisovega premika  $\sigma_1, \sigma_2$  potrebujemo  $2 \times 2$  podmatriko  $R = C(n-1:n, n-1:n)$ , ki jo lahko izračunamo v  $\mathcal{O}(1)$  operacij. Naj bo  $v_1$  prvi stolpec matrike  $(C_k - \sigma_1 I)(C_k - \sigma_2 I)$ . Tudi to lahko izračunamo v  $\mathcal{O}(1)$  operacij. Sedaj poiščemo Householderjevo zrcaljenje  $P_0$ , ki vektor  $v_1$  prezrcali v smer  $e_1$ . Če matriki  $A_k$  in  $B_k$  z leve pomnožimo s  $P_0$ , dobimo grbo, ki jo podobno kot pri QR iteraciji za nesimetrično matriko z naslednjimi ortogonalnimi transformacijami premakamo navzdol in na koncu spravimo iz matrik.

Poglejmo, kako grbo premaknemo v primeru, ko imamo matriki velikosti  $5 \times 5$ . Najprej dobimo

$$A_k^{(1)} = P_0 A_k = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}, \quad B_k^{(1)} = P_0 B_k = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & \times \end{bmatrix}.$$

Z dvema zrcaljenjima z desne spravimo  $B$  nazaj v trikotno obliko. Najprej z  $Z_1$  uničimo poddiagonalne elemente v tretji vrstici, nato pa z  $Z_2$  še v drugi vrstici. Pri tem se v  $A$  pojavijo novi neničelni elementi v četrti vrstici. Dobimo

$$A_k^{(2)} = A_k^{(1)} Z_1 = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}, \quad B_k^{(2)} = B_k^{(1)} Z_1 = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & \times \end{bmatrix},$$

$$A_k^{(3)} = A_k^{(2)} Z_2 = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}, \quad B_k^{(3)} = B_k^{(2)} Z_2 = \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & \times \end{bmatrix}.$$

Z zrcaljenjem, ki deluje od druge do četrte vrstice, uničimo elementa v tretji in četrti vrstici prvega stolpca matrike  $A$ . Tako se je grba premaknila za eno mesto navzdol.

$$A_k^{(4)} = P_1 A_k^{(3)} = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times \end{bmatrix}, \quad B_k^{(4)} = P_1 B_k^{(3)} = \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times \end{bmatrix}.$$

Postopek ponavljamo, dokler grbo ne izločimo iz matrik in nam na koncu ne ostane par zgornje Hessenbergove in zgornje trikotne matrike.

Ko izvajamo QZ algoritem, lahko izvedemo deflacijo in nadaljujemo z matriko manjše velikosti, če velja  $a_{k+1,k} = 0$  oziroma  $b_{kk} = 0$ . Pri situaciji  $a_{k+1,k} = 0$  problem preprosto razdelimo na dva manjša splošena problema lastnih vrednosti velikosti  $k \times k$  in  $(n-k) \times (n-k)$ .

V primeru  $b_{kk} = 0$  pa z ustreznimi ortogonalnimi transformacijami problem prevedemo na podobnega, kjer velja  $a_{n,n-1} = 0$  in  $b_{nn} = 0$ , potem pa lahko nadaljujemo samo z vodilnima podmatrikama velikosti  $(n-1) \times (n-1)$ . V tem primeru smo izločili lastno vrednost  $\infty$ .

S QZ algoritmom na koncu izračunamo splošeno Schurovo formo. Če nas zanimajo samo lastne vrednosti, potem ni potrebno shranjevati produktov ortogonalnih transformacij. Zadoščata nam končni trikotni matriki.

Naj bo  $\sigma$  izračunana lastna vrednost šopa  $(A, B)$ , bodisi z metodo QZ ali s kakšno drugo metodo, za katero potrebujemo pripadajoči lastni vektor. Izračunamo ga lahko s s posplošitvijo inverzne iteracije.

---

**Algoritem 3.1** Inverzna iteracija za posplošeni problem lastnih vrednosti
 

---

izberi začetni vektor  $q_0$

$k = 1, 2, \dots$

reši  $(A - \sigma B)z_k = Bq_{k-1}$

$q_k = z_k / \|z_k\|$

---

Če je matrika  $B$  nesingularna, potem je zgornja metoda ekvivalentna inverzni iteraciji za  $B^{-1}A$ .

### 3.3 Definiten matrični šop

Denimo, da iščemo lastne vrednosti matričnega para  $(A, B)$ , kjer sta obe matriki simetrični, matrika  $B$  pa je še pozitivno definitna. V tem primeru pravimo, da je matrični šop *definiten*.

Zaradi ohranjanja simetrije takega problema ni stabilno reševati preko matrike  $C = AB^{-1}$ . Bolje je, če najprej matriko  $B$  razcepimo po Choleskem v  $B = VV^T$  in dobimo,

$$\begin{aligned} Ax &= \lambda VV^T x \\ V^{-1}Ax &= \lambda V^T x \\ V^{-1}AV^{-T}V^T x &= \lambda V^T x \end{aligned}$$

Dobimo simetričen lastni problem  $Cy = \lambda y$ , kjer je  $C = V^{-1}AV^{-T}$  in  $y = V^T x$ . Matrika  $C$  je simetrična in od tod sledi, da ima definiten matrični šop sam realne lastne vrednosti. V nadaljevanju lahko uporabimo algoritme za simetrične matrike, kar ne bi bilo možno, če bi delali direktno z matriko  $AB^{-1}$ , saj ta ni simetrična.

Za definiten matrični šop  $(A, B)$  lahko definiramo posplošeni Rayleighov kvocient

$$\rho(x, A, B) = \frac{x^T Ax}{x^T Bx}.$$

Podobno kot za simetričen problem velja  $\lambda_n \leq \rho(x, A, B) \leq \lambda_1$ .

**Lema 3.7** Posplošeni Rayleighov kvocient  $\rho(x, A, B)$  vrne skalar  $\lambda$ , ki minimizira

$$\|Ax - \lambda Bx\|_{B^{-1}},$$

kjer je  $\|z\|_{B^{-1}}^2 = z^T B^{-1}z$ .

Za iskanje posameznih lastnih parov lahko uporabimo tudi posplošeno Rayleighovo iteracijo.

### 3.4 Nelinearni problem lastnih vrednosti

Naj bo  $T(\lambda)$  matrika velikosti  $n \times n$ , katere elementi so funkcije parametra  $\lambda$ . Za vsak  $\lambda \in \mathbb{C}$  je torej  $T(\lambda)$  kompleksna matrika iz  $\mathbb{C}^{n \times n}$ . Za elemente lahko predpostavimo, da so zadosti krat zvezno odvedljivi.

**Algoritem 3.2** Posplošena Rayleighova iteracijaizberi  $x_0 \neq 0$  $k = 0, 1, \dots$ 

$$\rho_k = \frac{x_k^T A x_k}{x_k^T B x_k}$$

reši  $(A - \rho_k B)y_{k+1} = Bx_k$ 

$$x_{k+1} = y_{k+1} / \|y_{k+1}\|$$

Če obstajata tak skalar  $\lambda \in \mathbb{C}$  in neničelni vektor  $x \in \mathbb{C}^n$ , da je

$$T(\lambda)x = 0,$$

potem je  $\lambda$  lastna vrednost,  $x$  pa (desni) lastni vektor. Podobno je neničelni vektor  $y \in \mathbb{C}^n$  levi lastni vektor, če je  $y^H T(\lambda) = 0$ . Pravimo, da je  $(\lambda, x, y)$  lastna trojica za  $T$ .

Posebni primeri nelinearnih problemov lastnih vrednosti so:

- Navadni linearni problem lastnih vrednosti (EP). Če za matriko  $A \in \mathbb{C}^{n \times n}$  definiramo

$$T(\lambda) = \lambda I - A,$$

potem je očitno  $T(\lambda)x = 0$  natanko tedaj, ko je  $Ax = \lambda x$ .

- Posplošeni problem lastnih vrednosti (GEP). Če za matriki  $A, B \in \mathbb{C}^{n \times n}$  vzamemo

$$T(\lambda) = \lambda B - A,$$

potem je  $T(\lambda)x = 0$  ekvivalentno  $Ax = \lambda Bx$ .

Pri GEP se lahko zgodi, da sta matriki  $A$  in  $B$  izbrani tako, da je  $\det(T(\lambda)) \equiv 0$ . V tem primeru pravimo, da je problem *singularen*. Nasprotno za problem  $T(\lambda)x = 0$  pravimo, da je *regularen*, če je  $\det(T(\lambda)) \not\equiv 0$ . V primeru, ko je problem regularen, so lastne vrednosti rešitve karakteristične enačbe  $\det(T(\lambda)) = 0$ .

- Kvadratni problem lastnih vrednosti (QEP). Za matrike  $M, C, K \in \mathbb{C}^{n \times n}$  definiramo

$$T(\lambda) = \lambda^2 M + \lambda C + K.$$

V primeru, ko je matrika  $M$  nesingularna, imamo  $2n$  lastnih vrednosti, saj je  $\det(T(\lambda))$  v tem primeru polinom stopnje  $2n$  z vodilnim koeficientom  $\det(M)$ .

Problem je sicer nelinearen, a se da linearizirati. Ena izmed možnosti je, da ga prevedemo na GEP oblike

$$\left( \lambda \begin{bmatrix} M & 0 \\ 0 & I \end{bmatrix} + \begin{bmatrix} C & K \\ -I & 0 \end{bmatrix} \right) \begin{bmatrix} \lambda x \\ x \end{bmatrix} = 0.$$

- Polinomski problem lastnih vrednosti (PEP). Tu imamo

$$T(\lambda) = \lambda^m A_m + \lambda^{m-1} A_{m-1} + \dots + \lambda A_1 + A_0.$$

Podobno kot kvadratni problem lastnih vrednosti lahko tudi polinomski problem lastnih vrednosti prevedemo na GEP z matrikama velikosti  $(mn) \times (mn)$ . Če je  $A_m$  nesingularna matrika, potem ima PEP  $mn$  lastnih vrednosti.



- Racionalni problem lastnih vrednosti (REP), npr.:

$$T(\lambda) = -K + \lambda M + \sum_{k=1}^m \frac{\lambda}{\sigma_k - \lambda} C_k.$$

Obstaja več vrst REP. Vse lahko prevedemo na polinomski problem lastnih vrednosti, če izraz spravimo na skupni imenovalac. Torej bi tudi REP lahko linearizirali v GEP.

- Iracionalni problem lastnih vrednosti, npr.:

$$T(\lambda) = K - \lambda M + i \sum_{k=1}^m \sqrt{\lambda - \sigma_k^2} W_k.$$

Zgornji problem, kjer so  $M, K, W_1, \dots, W_m$  simetrične matrice, pri čemer je  $K$  nenegativno,  $M$  pa pozitivno definitna,  $\sigma_1, \dots, \sigma_m$  pa so nenegativni skalarji, nastopa pri razvoju linearnih pospeševalnikov.

Iracionalnega problema lastnih vrednosti ne moremo linearizirati in zanj lahko rečemo, da je pristno nelinearen.

- Pri študiju diferencialnih enačb s časovnimi zamiki nastopajo nelinearni problemi lastnih vrednosti oblike

$$T(\lambda) = -\lambda I + A_0 + \sum_{k=1}^m A_j e^{-h_k \lambda},$$

kjer so  $h_1, \dots, h_m$  pozitivni premiki,  $A_0, \dots, A_m$  pa so realne matrice.

Poglejmo si dve numerični metodi, ki ju lahko uporabimo za splošen nelinearni problem lastnih vrednosti.

### 3.4.1 Newtonova metoda in inverzna iteracija

Če izberemo vektor  $v \in \mathbb{C}^n$  in dodamo k enačbi  $T(\lambda)x = 0$  še pogoj  $v^H x = 1$ , potem lahko zapišemo NEP v obliki nelinearnega sistema  $n + 1$  enačb za  $n + 1$  neznank:

$$F(x, \lambda) := \begin{bmatrix} T(\lambda)x \\ v^H x - 1 \end{bmatrix} = 0.$$

Za reševanje zgornjega sistema uporabimo Newtonovo metodo. Če je  $(x_k, \lambda_k)$  tekoči približek za lastni par, potem po Newtonovi metodi naslednji približek  $(x_{k+1}, \lambda_{k+1})$  dobimo kot

$$\begin{bmatrix} x_{k+1} \\ \lambda_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ \lambda_k \end{bmatrix} - JF(x_k, \lambda_k)^{-1} F(x_k, \lambda_k),$$

kjer je Jacobijeva matrika enaka

$$JF(x_k, \lambda_k) = \begin{bmatrix} T(\lambda_k) & T'(\lambda_k)x_k \\ v^H & 0 \end{bmatrix}.$$

Tako pridemo do sistema

$$\begin{bmatrix} T(\lambda_k) & T'(\lambda_k)x_k \\ v^H & 0 \end{bmatrix} \begin{bmatrix} x_{k+1} - x_k \\ \lambda_{k+1} - \lambda_k \end{bmatrix} = - \begin{bmatrix} T(\lambda_k)x_k \\ v^H x_k - 1 \end{bmatrix},$$

oziroma

$$\begin{aligned} T(\lambda_k)x_{k+1} &= -(\lambda_{k+1} - \lambda_k)T'(\lambda_k)x_k \\ v^H x_{k+1} &= 1. \end{aligned} \quad (3.1)$$

Iz zveze (3.1) sledi

$$x_{k+1} = -(\lambda_{k+1} - \lambda_k)T(\lambda_k)^{-1}T'(\lambda_k)x_k. \quad (3.2)$$

To pomeni, da iz zgornje enačbe lahko določimo pravo smer vektorja  $x_{k+1}$ . Več kot smer tudi ne potrebujemo, saj lahko potem  $x_{k+1}$  enostavno pomnožimo s pravim skalarjem, da bo  $v^H x_{k+1} = 1$ . Za naslednji korak potrebujemo še nov približek za lastno vrednost  $\lambda_{k+1}$ . Če definiramo

$$u_{k+1} := T(\lambda_k)^{-1}T'(\lambda_k)x_k$$

in enačbo (3.2) skalarno pomnožimo z  $v$ , potem dobimo

$$v^H x_{k+1} = -(\lambda_{k+1} - \lambda_k)v^H u_{k+1}.$$

Ob predpostavki, da je prejšnji približek za lastni vektor bil normiran, torej  $v^H x_k = 1$ , dobimo

$$\lambda_{k+1} = \lambda_k - \frac{v^H x_k}{v^H u_{k+1}}.$$

Tako pridemo do prvega algoritma. Čeprav gre za Newtonovo metodo, metodo zaradi podobnosti s podobnim algoritmom za EP imenujemo *inverzna iteracija*. V nadaljevanju bomo videli, kako so v standardnem primeru EP povezane inverzna iteracija, Newtonova metoda in Rayleighova iteracija. Podobne povezave potem veljajo tudi za NEP.

---

### Algoritem 3.3 Inverzna iteracija

---

izberi vektor  $v$  in tak začetni približek  $(\lambda_0, x_0)$  za lastni par, da je  $v^H x_0 = 1$   
 $k = 0, 1, \dots$

- a) reši linearni sistem  $T(\lambda_k)u_{k+1} = T'(\lambda_k)x_k$
  - b)  $\lambda_{k+1} = \lambda_k - \frac{v^H x_k}{v^H u_{k+1}}$
  - c)  $x_{k+1} = \frac{1}{v^H u_{k+1}}u_{k+1}$
- 

### 3.4.2 Zaporedne linearne aproksimacije

Naj bo preslikava  $T : \mathbb{C} \rightarrow \mathbb{C}^{n \times n}$  spet dvakrat zvezno odvedljiva. Denimo, da že imamo približek  $\lambda_k$  za lastno vrednost. Iščemo popravek  $\Delta\lambda_k$  in vektor  $x \neq 0$ , da bo

$$T(\lambda_k - \Delta\lambda_k)x = 0.$$

Če zgornjo enačbo razvijemo v vrsto, dobimo

$$(T(\lambda_k) - \Delta\lambda_k T'(\lambda_k) + \mathcal{O}(|\Delta\lambda_k|^2))x = 0.$$

Če zanemarimo kvadratne člene, dobimo, da je popravek  $\Delta\lambda_k$  lastna vrednost posplošenega problema lastnih vrednosti

$$T(\lambda_k)x = \Delta\lambda_k T'(\lambda_k)x. \quad (3.3)$$

Pri metodi zaporednih linearnih aproksimacij za  $\Delta\lambda_k$  vzamemo po absolutni vrednosti najmanjšo lastno vrednost (3.3). V bližini lastne vrednosti je metoda podobna inverzni iteraciji, saj dobimo

$$\frac{1}{\Delta\lambda_k}x = T(\lambda_k)^{-1}T'(\lambda_k)x.$$

---

### Algoritem 3.4 Zaporedne linearne aproksimacije

---

izberi začetni približek  $\lambda_0$  za lastno vrednost

$k = 0, 1, \dots$

za  $\Delta\lambda_k$  vzemi po absolutni vrednosti najmanjšo lastno vrednost GEP  $T(\lambda_k)x = \theta T'(\lambda_k)x$

$\lambda_{k+1} = \lambda_k - \Delta\lambda_k$

---

Dela v enem koraku je več kot pri prejšnji metodi, saj moramo v vsakem koraku rešiti GEP, je pa konvergenca zato ponavadi hitrejša.

## 3.5 Polinomski (kvadratni) problem lastnih vrednosti

Pri polinomskem problemu lastnih vrednosti (PEP) je dan matrični polinom  $P(\lambda) = A_0 + \lambda A_1 + \dots + \lambda^m A_m$ , kjer so  $A_0, \dots, A_m$  matrice velikosti  $n \times n$ . Iščemo tak skalar  $\lambda_0 \in \mathbb{C}$  in neničelni vektor  $x_0 \in \mathbb{C}^n$ , da je  $P(\lambda_0)x_0 = 0$ .

PEP  $P$  je regularen, če njegov karakteristični polinom, definiran z  $g(\lambda) := \det(P(\lambda))$ , ni identično enak 0. Ničle polinoma  $g$  so končne lastne vrednosti PEP  $P$ . Če jih je manj kot  $mn$ , jih do  $mn$  dopolnimo z neskončnimi lastnimi vrednostmi. Neskončne lastne vrednosti ustrezajo ničelnim lastnim vrednostim *zvratnega polinoma*  $P_R(\lambda) := \lambda^m P(1/\lambda) = \lambda^m A_0 + \lambda^{m-1} A_1 + \dots + A_m$ , pojavijo pa se le, če je matrika  $A_m$  nesingularna.

Regularni PEP  $P$  ima tako  $mn$  (končnih ali neskončnih) lastnih vrednosti.

**Zgled 3.2** Naslednji primer kvadratnega problema lastnih vrednosti je povzet iz [23].

Če vzamemo  $Q(\lambda) = \lambda^2 M + \lambda C + K$ , kjer so

$$M = \begin{bmatrix} 0 & 6 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & -6 & 0 \\ 2 & -7 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad K = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

potem je  $\det Q(\lambda) = -6\lambda^5 + 11\lambda^4 - 12\lambda^3 + 12\lambda^2 - 6\lambda + 1$  in problem je regularen. Lastni pari so

$k$	1	2	3	4	5	6
$\lambda_k$	1/3	1/2	1	$i$	$-i$	$\infty$
$x_k$	$\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

Pet lastnih vrednosti je končnih, ena pa neskončna. Opazimo, da imata različni lastni vrednosti lahko isti lastni vektor, kar pri standardnem in posplošenem lastnem problemu seveda ni možno.  $\square$

Lahko se zgodi, da si največ  $m$  različnih lastnih vrednosti deli isti lastni vektor. Namreč, če je  $x$  lastni vektor, potem je vsaj ena izmed rešitev enačbe  $x^H P(\lambda)x = 0$ , ki ima največ  $m$  rešitev, enaka lastni vrednosti.

Dan je QEP  $Q(\lambda) = \lambda^2 M + \lambda C + K$ . Standardni postopek za numerično reševanje QEP je *linearizacija*, kjer problem prevedemo na navadni ali posplošni problem lastnih vrednosti reda  $2n$ . Možnih linearizacij je več, primera sta:

- a) prevedba na standardni nesimetrični lastni problem (če je  $M$  nesingularna)

$$Au = \lambda u$$

kjer sta

$$A = \begin{bmatrix} 0 & I \\ -M^{-1}K & -M^{-1}C \end{bmatrix}, \quad u = \begin{bmatrix} x \\ \lambda x \end{bmatrix},$$

- b) prevedba na posplošeni simetrični problem lastnih vrednosti

$$Az = \lambda Bz$$

kjer so

$$A = \begin{bmatrix} C & K \\ K & 0 \end{bmatrix}, \quad z = \begin{bmatrix} x \\ \lambda x \end{bmatrix}, \quad B = \begin{bmatrix} -M & 0 \\ 0 & K \end{bmatrix}.$$

Podobno lahko splošni polinomski problem lastnih vrednosti  $P(\lambda) = \lambda^m A_m + \dots + \lambda A_1 + A_0$  lineariziramo kot GEP z matrikami velikosti  $mn \times mn$ . Primer je t.i. *prva spremljevalna forma*, ki ima obliko

$$C_1(\lambda) = \lambda \begin{bmatrix} I & & & \\ & I & & \\ & & \ddots & \\ & & & A_m \end{bmatrix} + \begin{bmatrix} 0 & -I & & \\ & \ddots & \ddots & \\ & & 0 & -I \\ A_0 & A_1 & \cdots & A_{m-1} \end{bmatrix}.$$

V tem primeru je  $(\lambda, x)$  lastni par PEP natanko tedaj, ko je  $\left( \lambda, \begin{bmatrix} x \\ \lambda x \\ \vdots \\ \lambda^{m-1}x \end{bmatrix} \right)$  lastni par GEP.

### 3.5.1 Hermitski kvadratni problem lastnih vrednosti

V tem podrazdelku bomo obravnavali poseben razred kvadratnih problemov lastnih vrednosti, ki se v praksi velikokrat pojavi. Matrike  $M, K, C$  so hermitske, matrika  $M$  pa je še pozitivno definitna. Za poljuben neničelen  $x \in \mathbb{C}^n$  lahko definiramo

$$m(x) = x^H M x, \quad k(x) = x^H K x, \quad c(x) = x^H C x.$$

**Definicija 3.8** Dan je kvadratni problem lastnih vrednosti  $Q(\lambda) = \lambda^2 M + \lambda C + K$ , kjer so matrike  $M, K, C$  hermitske, matrika  $M$  pa je pozitivno definitna. Problem  $Q$  je

- hiperboličen, če za vsak  $x \neq 0$  velja  $c(x)^2 > 4m(x)k(x)$ ,

- eliptičen, če za vsak  $x \neq 0$  velja  $c(x)^2 < 4m(x)k(x)$ ,
- nadkritično dušen, če je hiperboličen in je matrika  $C$  pozitivno definitna, matrika  $K$  pa nenegativno definitna.

Če je  $x$  lastni vektor, potem je vsaj ena izmed rešitev kvadratne enačbe  $x^H Q(\lambda)x = 0$  enaka lastni vrednosti, ki pripada  $x$ . Od tod sledi, da so vse lastne vrednosti hiperboličnega QEP realne, medtem ko eliptičen QEP nima nobene realne lastne vrednosti. Podobno vidimo, da so vse lastne vrednosti nadkritično dušenega QEP negativne.

Dan je hiperboličen QEP  $Q(\lambda) = \lambda^2 M + \lambda C + K$ , kjer so  $M, K, C$  hermitske,  $M$  je pozitivno definitna in za vsak  $x \neq 0$  velja  $c(x)^2 > 4m(x)k(x)$ . Vse lastne vrednosti so realne in zanje velja, da jih lahko uredimo, da velja

$$\lambda_{2n} \leq \dots \leq \lambda_{n+1} < \lambda_n \leq \dots \leq \lambda_1.$$

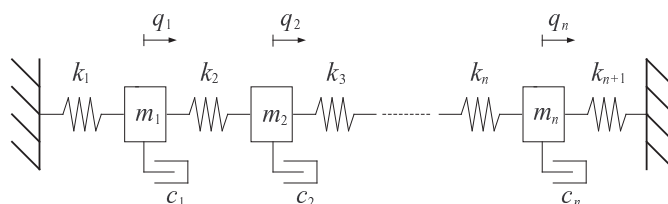
Definiramo lahko t.i. *primarni in sekundarni Rayleighov funkcional*

$$p_1(x) = \frac{-c(x) + d(x)}{2m(x)}, \quad p_2(x) = \frac{-c(x) - d(x)}{2m(x)},$$

kjer je  $d(x) = \sqrt{c(x)^2 - 4m(x)k(x)}$ . Za vsak  $x \neq 0$  velja  $p_2(x) < p_1(x)$ .

Lastni vektorji, ki pripadajo t.i. primarnim lastnim vrednostim  $\lambda_1, \dots, \lambda_n$ , so linearno neodvisni, enako velja za lastne vektorje, ki pripadajo sekundarnim lastnim vrednostim  $\lambda_{n+1}, \dots, \lambda_{2n}$ .

**Zgled 3.3** Za zgled uporabe kvadratnega problema lastnih vrednosti vzemimo nihanje dušenega sistema mas in vzmeti.



Če predpostavimo  $q_0 = q_{n+1} = 0$ , potem iz Newtonovega zakona dobimo enačbe

$$m_i \ddot{q}_i(t) = -k_i (q_i(t) - q_{i-1}(t)) - k_{i+1} (q_i(t) - q_{i+1}(t)) - c_i \dot{q}_i(t),$$

$i = 1, \dots, n$ , iz katerih sestavimo

$$M \ddot{q}(t) + C \dot{q}(t) + K q(t) = 0,$$

kjer je

$$M = \begin{bmatrix} m_1 & & \\ & \ddots & \\ & & m_n \end{bmatrix}, \quad C = \begin{bmatrix} c_1 & & \\ & \ddots & \\ & & c_n \end{bmatrix},$$

$$K = \begin{bmatrix} k_1 + k_2 & -k_2 & & & \\ & -k_2 & \ddots & \ddots & \\ & & \ddots & \ddots & -k_n \\ & & & -k_n & k_n + k_{n+1} \end{bmatrix}.$$

$M$  je masna matrika,  $C$  matrika dušenja,  $K$  pa togostna matrika.

V primeru, ko so vse lastne vrednosti enostavne, ima splošna rešitev homogene diferencialne enačbe

$$M\ddot{q}(t) + C\dot{q}(t) + Kq(t) = 0,$$

obliko

$$q(t) = \sum_{k=1}^{2n} \alpha_k e^{\lambda_k t} x_k,$$

kjer so  $\alpha_1, \dots, \alpha_{2n}$  poljubne konstante,  $\lambda_1, \dots, \lambda_n$  so lastne vrednosti,  $x_1, \dots, x_{2n}$  pa lastni vektorji za

$$\lambda^2 Mx + \lambda Cx + Kx = 0.$$

Konstante  $\alpha_1, \dots, \alpha_{2n}$  določimo iz začetnih odmikov  $q(0)$  in hitrosti  $\dot{q}(0)$ .

Sistem je stabilen, če velja  $\operatorname{Re}(\lambda_k) < 0$  za vse  $k$  oziroma šibko stabilen, če velja  $\operatorname{Re}(\lambda_k) \leq 0$  za vse  $k$ , če pa je  $\operatorname{Re}(\lambda_k) = 0$  je  $\lambda_k$  enostavna lastna vrednost.

V primeru vsiljenega nihanja imamo enačbo

$$M\ddot{q}(t) + C\dot{q}(t) + Kq(t) = f(t).$$

Če je desna stran oblike  $f(t) = e^{i\omega_0 t} f_0$ , kar ustreza vsiljenemu nihanju, potem je partikularna rešitev

$$q_p(t) = e^{i\omega_0 t} \sum_{k=1}^{2n} \frac{y_k^H f_0}{i\omega_0 - \lambda_k} x_k,$$

kjer so  $y_1, \dots, y_{2n}$  levi lastni vektorji.

Če se  $i\omega_0$  približa eni izmed lastnih vrednosti  $\lambda_k$ , se lahko pojavi resonanca. □

### 3.6 Računanje singularnega razcepa

Za matriko  $A \in \mathbb{R}^{m \times n}$  bi radi izračunali singularni razcep  $A = U\Sigma V^T$ . Ker velja

$$A^T A = V\Sigma^T \Sigma V^T,$$

dobimo preprosto idejo:

- 1) izračunamo  $A^T A$ ,
- 2) rešimo lastni problem  $A^T A = V\Lambda V^T$ , odtod dobimo  $V$ ,
- 3) za  $\Sigma$  vzamemo  $m \times n$  matriko, ki ima v zgornjem bloku kvadratni koren  $\Lambda$ ,

4) rešimo sistem  $U\Sigma = AV$  za matriko  $U$ .

**Zgled 3.4** Računanje na ta način ni numerično stabilno. Če vzamemo npr.  $A = \begin{bmatrix} 1 & 1 \\ 0 & \sqrt{\eta} \end{bmatrix}$ , kjer je  $\eta = \frac{1}{2}u$ , potem je  $A^T A = \begin{bmatrix} 1 & 1 \\ 1 & 1 + \eta \end{bmatrix}$ , v premični piki pa se to zaokroži v  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ . Numerično izračunane singularne vrednosti so zato  $\sqrt{2}$  in 0, točne pa so  $\sqrt{2}$  in  $\sqrt{\eta/2}$ . Pri majhnih singularnih vrednostih lahko torej pride do velike relativne napake. ■

Računanje  $A^T A$  tako ni najboljša ideja in se mu poskušamo izogniti. Pri računanju SVD za  $A$  (razen pri Jacobijevi metodi) matriko najprej reduciramo na bidiagonalno obliko. Postopek je:

- 1) poišči ortogonalni matriki  $U_1, V_1$ , da bo  $A = U_1 B V_1^T$  in  $B$  bidiagonalna,
- 2) izračunaj SVD za  $B$ :  $B = U_2 \Sigma V_2^T$ ,
- 3) SVD razcep za  $A$  je  $A = (U_1 U_2) \Sigma (V_1 V_2)^T$ .

Redukcijo na bidiagonalno obliko naredimo z množenji s Householderjevimi zrcaljenji z leve in z desne. V primeru matrike  $6 \times 4$  dobimo:

$$\begin{aligned}
 A &= \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}, & P_1 A &= \begin{bmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, \\
 P_1 A P_1' &= \begin{bmatrix} \times & \times & 0 & 0 \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & \times & \times & \times \end{bmatrix}, & P_2 P_1 A P_1' &= \begin{bmatrix} \times & \times & 0 & 0 \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}, \\
 P_2 P_1 A P_1' P_2' &= \begin{bmatrix} \times & \times & 0 & 0 \\ 0 & \times & \times & 0 \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}, & P_4 P_3 P_2 P_1 A P_1' P_2' &= \begin{bmatrix} \times & \times & 0 & 0 \\ 0 & \times & \times & 0 \\ 0 & 0 & \times & \times \\ 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = B.
 \end{aligned}$$

Nadaljujemo z računanjem singularnega razcepa za  $B$ , pri čemer lahko predpostavimo, da je  $B$  kvadratna matrika  $n \times n$ .

**Lema 3.9** Naj bo  $B$  kvadratna matrike  $n \times n$  s singularnimi vrednostmi in vektorji  $Bv_i = \sigma_i u_i$ ,  $i = 1, \dots, n$ . Lastne vrednosti matrike

$$C = \begin{bmatrix} 0 & B^T \\ B & 0 \end{bmatrix}$$

so  $\pm\sigma_i$ , ustrezni lastni vektorji pa  $\frac{1}{\sqrt{2}} \begin{bmatrix} v_i \\ \pm u_i \end{bmatrix}$ ,  $i = 1, \dots, n$ .





Za večje singularne vrednosti je to v redu, za manjše pa ne. Tako lahko pri  $\sigma_n$  pričakujemo napako velikosti  $\mathcal{O}(u\|A\|\kappa_2(A))$ .

V nadaljevanju si bomo pogledali nekaj specialnih metod za računanje singularnega razcepa matrike  $B$ .

### 3.7 QR iteracija za računanje singularnega razcepa

Naj bo

$$B = \begin{bmatrix} a_1 & b_1 & & & \\ & \ddots & \ddots & & \\ & & a_{n-1} & b_{n-1} & \\ & & & & a_n \end{bmatrix}.$$

Mislimo si, da računamo implicitno QR iteracijo z Wilkinsonovimi premiki za  $B^T B$ , a  $B^T B$  nikoli ne izračunamo eksplicitno. Za Wilkinsonov premik potrebujemo lastne vrednosti desne spodnje  $2 \times 2$  podmatrike  $B^T B$

$$\begin{bmatrix} a_{n-1}^2 + b_{n-2}^2 & a_{n-1}b_{n-1} \\ a_{n-1}b_{n-1} & a_n^2 + b_{n-1}^2 \end{bmatrix},$$

za  $\sigma^2$  pa vzamemo bližnjo lastno vrednost  $a_n^2 + b_{n-1}^2$ . Nato izračunamo Givensovo rotacijo  $R = \begin{bmatrix} c & -s \\ s & c \end{bmatrix}$ , da bo (glede na prvi stolpec  $B^T B$ )

$$R \begin{bmatrix} a_1^2 - \sigma^2 \\ a_1 b_2 \end{bmatrix} = \begin{bmatrix} \times \\ 0 \end{bmatrix}$$

in vzamemo  $R'_{12} = \begin{bmatrix} R & 0 \\ 0 & I \end{bmatrix}$ . Ko izračunamo  $BR'_{12}$  (to je ekvivalentno  $R'^T_{12} B^T BR'_{12}$  pri QR za  $B^T B$ ), dobimo (v primeru  $5 \times 5$ )

$$BR'_{12} = \begin{bmatrix} \times & \times & & & \\ + & \times & \times & & \\ & & \times & \times & \\ & & & \times & \times \\ & & & & \times \end{bmatrix}.$$

Dodatni neničelni element  $+$  uničimo tako, da ga z množenjem z ustreznimi rotacijami pogajamo navzdol z leve na desno. Tako dobimo

$$R_{12}BR'_{12} = \begin{bmatrix} \times & \times & + & & \\ 0 & \times & \times & & \\ & & \times & \times & \\ & & & \times & \times \\ & & & & \times \end{bmatrix}, \quad R_{12}BR'_{12}R'_{23} = \begin{bmatrix} \times & \times & 0 & & \\ & \times & \times & & \\ & + & \times & \times & \\ & & & \times & \times \\ & & & & \times \end{bmatrix}.$$

Na koncu je  $(R_{45}R_{34}R_{23}R_{12})B(R'_{12}R'_{23}R'_{34}R'_{45})$  bidiagonalna matrika za naslednji korak QR metode.

Za en korak enostranske QR iteracije porabimo  $30n + \mathcal{O}(1)$  osnovnih operacij in  $2n$  kvadratnih korenov. Kvadratni koren je zahtevnejši od osnovnih operacij in zanj lahko porabimo npr. tudi 10 ali več osnovnih operacij.

Kriterija za deflacijo sta  $|b_i| \leq \epsilon(|a_i| + |a_{i+1}|)$  in  $|a_i| \leq \epsilon\|B\|$ .

### 3.8 dqds metoda

Če je  $T$  simetrična pozitivno definitna matrika, lahko delamo naslednji postopek, ki ga imenujemo *LR metoda*:

$$\begin{aligned} T_0 &= T \\ i &= 0, 1, 2, \dots: \\ T_i &= V_i^T V_i \text{ (razcep Choleskega)} \\ T_{i+1} &= V_i V_i^T \end{aligned}$$

Presenetljivo ugotovimo, da  $T_i$  konvergira proti diagonalni matriki lastnih vrednosti matrike  $T$ . Za začetek je očitno, da sta matriki  $T_i$  in  $T_{i+1}$  podobni, saj velja

$$T_{i+1} = V_i V_i^T = V_i V_i^T V_i V_i^{-1} = V_i T_i V_i^{-1}.$$

Torej, če  $T_i$  konvergira proti diagonalni matriki, mora le ta na diagonalni imeti lastne vrednosti  $T$ .

Če je  $T$  tridiagonalna matrika, se da hitro preveriti, da so vse  $T_i$  tridiagonalne.

Ker  $V_i$  ni ortogonalna matrika, je metoda LR na prvi pogled potencialno nestabilna, a je skrb odveč, saj nam stabilnost in konvergenco zagotavlja naslednji presenetljivi izrek, ki povezuje QR in LR metodo.

**Lema 3.11** *Dva koraka LR algoritma ustrezata enemu koraku QR algoritma brez premikov.*

*Dokaz.* En korak QR algoritma je:

$$T_0 = QR, \quad T' = RQ.$$

Dva koraka LR metode sta

$$T_0 = V_0^T V_0, \quad T_1 = V_0 V_0^T = V_1^T V_1, \quad T_2 = V_1 V_1^T.$$

Pokazati moramo, da je  $T' = T_2$ .

Ker je  $T_0$  simetrična je

$$T_0^2 = T_0^T T_0 = (QR)^T QR = R^T R.$$

To je razcep Choleskega za s.p.d.  $T_0^2$ . Po drugi strani je

$$T_0^2 = T_0 T_0 = V_0^T V_0 V_0^T V_0 = V_0^T V_1^T V_1 V_0 = (V_1 V_0)^T (V_1 V_0).$$

Ker je to tudi razcep Choleskega, ta pa je enoličen, velja  $R = V_1 V_0$ .

To pa pomeni

$$\begin{aligned} T' &= RQ = RQR^{-1} = RT_0R^{-1} = (V_1V_0)(V_0^TV_0)V_0^{-1}V_1^{-1} \\ &= V_1V_0V_0^TV_1^{-1} = V_1V_1^TV_1V_1^{-1} = V_1V_1^T = T_2. \quad \blacksquare \end{aligned}$$

Namesto enega koraka QR metode lahko tako naredimo dva koraka LR metode, kar je ekonomičneje. Konvergenco pospešimo s premiki, pri čemer pa moramo sedaj paziti, da bo premik manjši od najmanše lastne vrednosti, saj sicer izgubimo pozitivno definitnost.

$$\begin{aligned} T_0 &= T \\ i &= 0, 1, 2, \dots: \\ &\text{izberi premik } \tau_i^2 > 0, \text{ tako da je } \tau_i^2 \leq \lambda_n(T) \\ T_i - \tau_i^2 I &= V_i^T V_i \text{ (razcep Choleskega)} \\ T_{i+1} &= V_i V_i^T + \tau_i^2 I \end{aligned}$$

Matriki  $T_{i+1}$  in  $T_i$  sta podobni, saj velja

$$T_{i+1} = V_i V_i^T + \tau_i^2 I = V_i^{-T} V_i^T (V_i V_i^T + \tau_i^2 I) = V_i^{-T} T_i V_i^T.$$

Pri *dqds metodi* delamo LR metodo s premikom za  $B^T B$ , pri čemer pa te matrike ne računamo eksplicitno. Z ekonomičnim algoritmom iz bidiagonalne matrike  $B_i$  dobimo bidiagonalno  $B_{i+1}$ , pri tem pa velja, da če bi naredili LR s premikom za  $B_i^T B_i$ , bi dobili  $B_{i+1}^T B_{i+1}$ .

Označimo

$$B_i = \begin{bmatrix} a_1 & b_1 & & & \\ & \ddots & \ddots & & \\ & & a_{n-1} & b_{n-1} & \\ & & & a_n & \\ & & & & \end{bmatrix} \quad \text{in} \quad B_{i+1} = \begin{bmatrix} \tilde{a}_1 & \tilde{b}_1 & & & \\ & \ddots & \ddots & & \\ & & \tilde{a}_{n-1} & \tilde{b}_{n-1} & \\ & & & \tilde{a}_n & \\ & & & & \end{bmatrix}.$$

S primerjavo elementov  $T_i = B_i^T B_i + \tau_i^2 I$  in  $T_{i+1} = B_{i+1}^T B_{i+1} + \tau_{i+1}^2 I = B_i B_i^T + \tau_i^2 I$  pridemo do osnovnega *dqds* algoritma, kjer je  $\delta = \tau_{i+1}^2 - \tau_i^2 \geq 0$  in privzamemo  $b_0 = \tilde{b}_0 = 0$ .

$$\begin{aligned} k &= 1, \dots, n-1 \\ \tilde{a}_k^2 &= a_k^2 + b_k^2 - \tilde{b}_{k-1}^2 - \delta \\ \tilde{b}_k^2 &= a_{k+1}^2 b_k^2 / \tilde{a}_k^2 \\ \tilde{a}_n^2 &= a_n^2 - \tilde{b}_{n-1}^2 - \delta \end{aligned}$$

Ves čas lahko računamo s kvadrati in korenimo le na koncu. Če označimo  $c_k = a_k^2$  in  $d_k = b_k^2$ , dobimo novo različico algoritma:

$$\begin{aligned} k &= 1, \dots, n-1 \\ \tilde{c}_k &= c_k + d_k - \tilde{d}_{k-1} - \delta \\ d_k &= c_{k+1} d_k / \tilde{c}_k \\ \tilde{c}_n &= c_n - \tilde{d}_{n-1} - \delta \end{aligned}$$

Zadnji algoritem porabi le  $5n + \mathcal{O}(1)$  operacij za en korak LR algoritma. Ta algoritm lahko še izboljšamo, da bo bolj natančen (ob enakem številu operacij) in tako končno pridemo do *dqds* algoritma.

Če definiramo  $g_k = c_k - \tilde{d}_{k-1} - \delta$ , potem velja

$$g_k = c_k - \frac{c_k d_{k-1}}{\tilde{c}_{k-1}} - \delta = c_k \frac{\tilde{c}_{k-1} - d_{k-1}}{\tilde{c}_{k-1}} - \delta = c_k \frac{c_{k-1} - \tilde{d}_{k-2} - \delta}{\tilde{c}_{k-1}} - \delta = \frac{c_k}{\tilde{c}_{k-1}} g_{k-1} - \delta.$$

Končni dqds algoritem je

$$\begin{aligned} g &= c_1 - \delta \\ k &= 1, \dots, n-1 \\ \tilde{c}_k &= g + d_j \\ t &= c_{k+1} / \tilde{c}_k \\ \tilde{d}_k &= d_k \cdot t \\ g &= g \cdot t - \delta \\ \tilde{c}_n &= g \end{aligned}$$

Za razliko od prejšnjega algoritma ima eno množenje namesto seštevanja. Pokazati se da, da nam ravno to zagotavlja visoko obratno stabilnost.

**Izrek 3.12** Če je  $B$  bidiagonalna nesingularna matrika s singularnimi vrednostmi  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$ , potem za izračunane singularne vrednosti  $\tilde{\sigma}_1 \geq \dots \tilde{\sigma}_n > 0$  po dqds algoritmu velja (če ne uporabljamo premikov)

$$\frac{|\tilde{\sigma}_i - \sigma_i|}{\sigma_i} \leq (10n - 5)u + O(u^2). \quad \blacksquare$$

### 3.9 Jacobijeva metoda za singularni razcep

Spet delamo s polno matriko  $A$ . Mislimo si, da delamo Jacobijevo metodo za računanje lastnih vrednosti na  $A^T A$ , pri čemer spet ne računamo  $A^T A$ .

Če delamo  $\text{jac}(A^T A, p, q)$  to pomeni, da iz  $A^T A$  dobimo  $R_{pq}^T A^T A R_{pq}$ . Zato lahko delamo z  $A$  in v enem koraku iz  $A$  dobimo  $AR_{pq}$ .

Algoritem za enostransko Jacobijevo metodo je:

$$\begin{aligned} A &= \text{oneside\_jac}(A, p, q) \\ \text{izračunaj } b_{pp} &= (A^T A)_{pp}, b_{pq} = (A^T A)_{pq}, b_{qq} = (A^T A)_{qq} \\ \text{če velja } |b_{pq}| &> \epsilon \sqrt{b_{pp} b_{qq}}, \text{ potem} \\ \tau &= \frac{b_{pp} - b_{qq}}{2b_{pq}} \\ t &= \frac{\text{sign}(\tau)}{|\tau| + \sqrt{1 + \tau^2}} \\ c &= \frac{1}{1 + t^2} \\ s &= ct \\ A &= AR_{pq}(\varphi) \\ J &= JR_{pq}(\varphi) \text{ (če potrebujemo tudi singularne vektorje)}. \end{aligned}$$

Po rotaciji  $(p, q)$  se v  $A$  spremenita stolpca  $p$  in  $q$ .

V poštev pride le pragovna varianta, saj bi morali pri klasični izračunati vse elemente  $A^T A$ , da bi lahko poiskali maksimalni element in ustrezno rotacijo. Končamo, ko velja  $|b_{pq}| \leq \epsilon \sqrt{b_{pp}b_{qq}}$  za vse  $p \neq q$ .

Na koncu dobimo:

- $\sigma_i = \|A(:, i)\|_2$ , (to je v bistvu  $\sqrt{b_{ii}}$ ),
- $U = [u_1 \ \cdots \ u_n]$ , kjer je  $u_i = \frac{1}{\sigma_i} A(:, i)$ ,
- $V = J$  (če smo shranjevali produkte).

**Posledica 3.13 (Relativni Weylov izrek za singularne vrednosti)** Naj bodo  $\sigma_1 \geq \cdots \geq \sigma_n$  singularne vrednosti matrike  $A$  in  $\tilde{\sigma}_1 \geq \cdots \geq \tilde{\sigma}_n$  singularne vrednosti matrike  $\tilde{A} = Y^T A X$ . Potem za  $i = 1, \dots, n$  velja

$$|\hat{\sigma}_i - \sigma_i| \leq |\sigma_i| \epsilon,$$

kjer je  $\epsilon = \max(\|X^T X - I\|_2, \|Y^T Y - I\|_2)$ .

*Dokaz.* Če definiramo matriki

$$C = \begin{bmatrix} 0 & A^T \\ A & 0 \end{bmatrix}, \quad Z = \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix},$$

potem velja

$$Z^T C Z = \begin{bmatrix} 0 & \tilde{A}^T \\ \tilde{A} & 0 \end{bmatrix}.$$

Od tod za  $i = 1, \dots, n$  iz leme 3.9 in relativnega Weylovega izreka za lastne vrednosti simetrične matrike sledi  $|\hat{\sigma}_i - \sigma_i| \leq |\sigma_i| \epsilon$ , kjer je

$$\epsilon = \|Z^T Z - I\|_2 = \left\| \begin{bmatrix} X^T X - I & 0 \\ 0 & Y^T Y - I \end{bmatrix} \right\|_2 = \max(\|X^T X - I\|_2, \|Y^T Y - I\|_2). \quad \blacksquare$$

Potrebujemo še nekaj pomožnih rezultatov, ki nam povejo, kako je omejena obratna napaka, če matriko množimo z osnovnimi ortogonalnimi transformacijami.

**Izrek 3.14** Naj bo  $R$  točna Givensova rotacija (Householderjevo zrcaljenje) in  $\tilde{R}$  njena aproksimacija v plavajoči vejici. Potem je

$$\begin{aligned} fl(\tilde{R}A) &= R(A + E), & \|E\|_2 &= \mathcal{O}(u)\|A\|_2 \\ fl(A\tilde{R}) &= (A + F)R, & \|F\|_2 &= \mathcal{O}(u)\|A\|_2. \end{aligned}$$

*Dokaz.* Izrek bomo dokazali le za primer Givensove rotacije. Dokaz za Householderjevo zrcaljenje lahko najdete npr. v [14].

S podrobno analizo (ki je opravljena npr. v [25]), se da pokazati, da za izračunana parametra Givensove rotacije velja  $\hat{c} = c(1 + \delta_1)$  in  $\hat{s} = s(1 + \delta_2)$ , kjer je  $|\delta_1|, |\delta_2| = \mathcal{O}(u)$ . Od tod sledi

$\|\widehat{R} - R\| = \mathcal{O}(u)$ , pomembno pa je tudi, da ima matrika  $\widehat{R} - R$  lahko največ štiri neničelne elemente.

Naj bo

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

Z analizo zaokrožitvenih napak dobimo

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \end{bmatrix} = \begin{bmatrix} \widehat{c}x_1(1 + \epsilon_1) + \widehat{s}x_2(1 + \epsilon_2) \\ -\widehat{s}x_1(1 + \epsilon_3) + \widehat{c}x_2(1 + \epsilon_4) \end{bmatrix},$$

kjer je  $|\epsilon_i| \leq 2u$  za  $i = 1, \dots, 4$ . Iz

$$\begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} + \begin{bmatrix} c\delta_1 & s\delta_2 \\ -s\delta_2 & c\delta_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} \widehat{c}x_1\epsilon_1 + \widehat{s}x_2\epsilon_2 \\ -\widehat{s}x_1\epsilon_3 + \widehat{c}x_2\epsilon_4 \end{bmatrix}$$

sledi

$$\left\| \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\|_2 = \mathcal{O}(u) \left\| \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|_2.$$

Denimo, da  $R$  predstavlja rotacijo v  $p$ -ti in  $q$ -ti vrstici. Potem pri računanju produkta  $\widehat{R}A$  lahko pride do napak le v  $p$ -ti in  $q$ -ti vrstici. Za  $p$ -to in  $q$ -to vrstico matrike  $RA$  potem dobimo

$$\|fl(\widehat{R}A) - RA\|_2 \leq \|fl(\widehat{R}A) - RA\|_F = \mathcal{O}(u) \|A([p \ q], :)\|_F \leq \mathcal{O}(u) \sqrt{2} \|A\|_2. \quad \blacksquare$$

**Izrek 3.15** *Vzemimo zaporedje ortogonalnih transformacij (rotacij ali zrcaljenj)  $P_1, \dots, P_j$  in  $Q_1, \dots, Q_j$ . Za numerično izračunano matriko  $B = P_j \cdots P_1 A Q_1 \cdots Q_j$  velja*

$$\widetilde{B} = P_j \cdots P_1 (A + E_j) Q_1 \cdots Q_j,$$

kjer je  $\|E_j\|_2 = 2j\mathcal{O}(u)\|A\|_2$ .

*Dokaz.* Uporabimo indukcijo. Za  $j = 0$  očitno drži.

Predpostavimo, da izrek velja za  $j - 1$  množenj z ortogonalnimi transformacijami. Naj bo  $B_{j-1} = P_{j-1} \cdots P_1 A Q_1 \cdots Q_{j-1}$ . Če definiramo  $P^{(k)} = P_k \cdots P_1$  in  $Q^{(k)} = Q_1 \cdots Q_k$ , potem je predpostavka, da velja  $\widehat{B}_{j-1} = P^{(j-1)}(A + E_{j-1})Q^{(j-1)}$ , kjer je  $\|E_{j-1}\|_2 = 2(j-1)\mathcal{O}(u)\|A\|_2$ .

Pri analizi računanja  $B_j = P_j B_{j-1} Q_j$  dobimo  $\widehat{B}_{j-1/2} = P_j(\widehat{B}_{j-1} + F_1)$  in  $\widehat{B}_j = (\widehat{B}_{j-1/2} + F_2)Q_j$ , kjer je  $\|F_1\|_2, \|F_2\|_2 = \mathcal{O}(u)\|A\|_2$ . Tako pridemo do

$$\begin{aligned} \widehat{B}_j &= (P_j(\widehat{B}_{j-1} + F_1) + F_2)Q_j \\ &= P_j \widehat{B}_{j-1} Q_j + P_j F_1 Q_j + F_2 Q_j \\ &= P_j \left( P^{(j-1)}(A + E_{j-1})Q^{(j-1)} \right) + P_j F_1 Q_j + F_2 Q_j \\ &= P_j [A + E_{j-1} + P^{(j-1)T} F_1 Q^{(j-1)T} + P^{(j)T} F_2 Q^{(j-1)T}] Q^{(j)}. \end{aligned}$$

Sledi  $E_j = E_{j-1} + P^{(j-1)T} F_1 Q^{(j-1)T} + P^{(j)T} F_2 Q^{(j-1)T}$  in ocena  $\|E_j\|_2 \leq 2j\mathcal{O}(u)\|A\|_2$ . \blacksquare

**Izrek 3.16** Naj bo  $A = DX$   $n \times n$  matrika, kjer je  $D$  nesingularna diagonalna matrika,  $X$  pa nesingularna matrika. Če  $\hat{A}_m$  dobimo po  $m$  korakih enostranske Jacobijeve metode in so  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$  singularne vrednosti  $A$ ,  $\hat{\sigma}_1 \geq \dots \geq \hat{\sigma}_n$  pa singularne vrednosti  $\hat{A}_m$ , potem velja

$$\frac{|\hat{\sigma}_i - \sigma_i|}{\sigma_i} \leq O(mu)\kappa(X). \quad \blacksquare$$

*Dokaz.* Uporabimo indukcijo. Najprej pogledjmo, kaj dobimo v primeru  $m = 1$ , ko matriko  $A$  množimo le z eno Givensovo rotacijo  $R_1$ .

Naj bo  $F_1 = \hat{A}_1 - AR_1$ . Če uporabimo izrek 3.14 za  $i$ -to vrstico matrike  $A_1 = AR_1$ , potem dobimo

$$\|F_1(i, :)\|_2 = \|\hat{A}_1(i, :) - A(i, :)R_1\|_2 = \mathcal{O}(u)\|A(i, :)\|_2 = \mathcal{O}(u)\|d_{ii}X((i, :))\|_2.$$

To pomeni

$$\left\| \frac{1}{d_{ii}} F_1(i, :) \right\|_2 = \mathcal{O}(u)\|X(i, :)\|_2,$$

in  $\|D^{-1}F_1\|_2 = \mathcal{O}(u)\|X\|_2$ . Iz zveze

$$\hat{A}_1 = AR_1 + F_1 = AR_1(I + R_1^T A^{-1} F_1) = AR_1(I + R_1^T X^{-1} D^{-1} F_1)$$

sledi  $\hat{A}_1 = AR_1(1 + E_1)$ , kjer lahko ocenimo  $\|E_1\|_2 = \mathcal{O}(u)\|X^{-1}\|_2\|X\|_2$ .

Denimo, da za  $A_{m-1} = AR_1 \cdots R_{m-1}$  za numerično izračunano matriko velja  $\hat{A}_{m-1} = A_{m-1}(I + E_{m-1})$ , kjer je  $\|E_{m-1}\|_2 = \mathcal{O}((m-1)u)\kappa_2(X)$ . Vemo, da velja  $\hat{A}_m = \hat{A}_{m-1}R_m(I + E)$ , kjer je  $\|E\|_2 = \mathcal{O}(u)\kappa_2(X)$ . Če razpišemo

$$\begin{aligned} \hat{A}_m &= A_{m-1}R_1 \cdots R_{m-1}(I + E_{m-1})R_m(I + E) \\ &= A_{m-1}R_1 \cdots R_m(I + R_m^T E_{m-1} R_m + E + R_m^T E_{m-1} R_m E), \end{aligned}$$

dobimo

$$E_m = R_m^T E_{m-1} R_m + E + R_m^T E_{m-1} R_m E.$$

Če zanemarimo  $R_m^T E_{m-1} R_m E$ , dobimo oceno  $\|E_m\|_2 = \mathcal{O}(mu)\kappa_2(X)$ .

Sedaj uporabimo relativni Weylov izrek in dobimo končno oceno  $|\hat{\sigma}_i - \sigma_i| \leq \sigma_i \epsilon$ , kjer je

$$\epsilon = \|(I + E_m)^T(I + E_m) - I\|_2 \leq 3\|E_m\|_2 \leq 3\mathcal{O}(mu)\kappa_2(X),$$

za  $i = 1, \dots, n$ . ■

**Izrek 3.17** Če Jacobijevo enostransko metodo zaustavimo, ko za vse  $j \neq k$  velja

$$|g_{jk}| \leq \epsilon \sqrt{g_{jj}g_{kk}},$$

kjer je  $G = A^T A$ , in so  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$  singularne vrednosti  $A$ ,  $\alpha_1^2 \geq \dots \geq \alpha_n^2$  pa diagonalni elementi  $G$ , potem je

$$|\sigma_i - \alpha_i| \leq n\epsilon|\alpha_i|. \quad \blacksquare$$

*Dokaz.* Pišemo lahko  $G = D\tilde{G}D$ , kjer je  $D = \text{diag}(\sqrt{g_{11}}, \dots, \sqrt{g_{nn}})$  in  $|\tilde{g}_{jk}| \leq \epsilon$  za vse  $j \neq k$ . Matrika  $\tilde{G}$  je simetrična in pozitivno definitna, zato obstaja razcep Choleskega  $\tilde{G} = LL^T$ , kjer je  $L$  spodnja trikotna matrika. Tako dobimo  $G = A^T A = DLL^T D$ .

Ker na koncu kot približke za singularne vrednosti vzamemo diagonalne elemente matrike  $G$ , bi radi ocenili razliko med singularnimi vrednostmi matrik  $A$  in  $D$ . Iz  $A^T A = DLL^T D$  sledi

$$L^{-1}D^{-1}A^T A D^{-1}L^{-T} = (AD^{-1}L^{-T})^T (AD^{-1}L^{-T}) = I,$$

to pa pomeni, da je matrika  $AD^{-1}L^{-T}$  ortogonalna. Torej obstaja taka ortogonalna matrika  $Q$ , da je  $L^T D = QA$ . Ker imata matriki  $A$  in  $QA$  iste singularne vrednosti, iz relativnega Weylovega izreka sledi  $|\sigma_i - \alpha_i| \leq \alpha_i \theta$ , kjer je

$$\theta = \|L^T L - I\|_2 = \|LL^T - I\|_2 \leq \|\tilde{G} - I\|_F \leq n\epsilon. \quad \blacksquare$$

Jacobijeva metoda torej izračuna singularne vrednosti (in vektorje) z visoko relativno natančnostjo, če lahko zapišemo  $A = DX$ , kjer je  $X$  dobro pogojena matrika,  $D$  pa nesingularna diagonalna matrika.

Jacobijeva metoda tako ponavadi deluje bolje od drugih v primeru, ko ima  $D$  elemente, ki se po velikosti bistveno razlikujejo. Za zgled si oglejmo naslednji primer iz [8].

**Zgled 3.5** Naj bo

$$A = \begin{bmatrix} \delta & 1 & 1 & 1 \\ \delta & \delta & 0 & 0 \\ \delta & 0 & \delta & 0 \\ \delta & 0 & 0 & \delta \end{bmatrix},$$

kjer je  $\delta = 10^{-20}$ . Velja

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

Po prvem koraku bidiagonalizacije dobimo

$$\begin{bmatrix} -2\delta & -1/2 - \delta/2 & -1/2 - \delta/2 & -1/2 - \delta/2 \\ 0 & -1/2 + 5\delta/6 & -1/2 - \delta/6 & -1/2 - \delta/6 \\ 0 & -1/2 - \delta/6 & -1/2 + 5\delta/6 & -1/2 - \delta/6 \\ 0 & -1/2 - \delta/6 & -1/2 - \delta/6 & -1/2 + 5\delta/6 \end{bmatrix},$$

kar se zaokroži v

$$\begin{bmatrix} -2\delta & -1/2 & -1/2 & -1/2 \\ 0 & -1/2 & -1/2 & -1/2 \\ 0 & -1/2 & -1/2 & -1/2 \\ 0 & -1/2 & -1/2 & -1/2 \end{bmatrix}$$

in v naslednjem koraku dobimo

$$\begin{bmatrix} -2\delta & \sqrt{3}/2 & 0 & 0 \\ 0 & 3/2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

Singularne vrednosti dobljene bidiagonalne matrike so  $\sqrt{3}, \sqrt{3}\delta, 0, 0$ , točne singularne vrednosti pa so  $\sqrt{3}, \sqrt{3}\delta, \delta, \delta$ .  $\blacksquare$



Za bidiagonalno matriko je dqds obratno stabilen algoritem, problem pa je, če imamo na začetku polno matriko. Natančnost lahko izgubimo pri redukciji na bidiagonalno obliko, pri Jacobijevi metodi pa se to ne zgodi, saj matrike na začetku ne reduciramo.

### 3.9.1 Lastni razcep simetrične pozitivno definitne matrike

Lastni razcep simetrične pozitivno definitne matrike  $A$  lahko izračunamo na naslednji način:

- izračunamo razcep Choleskega  $A = LL^T$ ,
- izračunamo singularni razcep  $L = U\Sigma V^T$ ,
- $A = U\Sigma^2U^T$ .

Če za točko b) uporabimo enostransko Jacobijevo metodo in je  $L = DX$ , kjer je  $D$  diagonalna matrika in  $X$  dobro pogojena, imajo izračunane singularne vrednosti relativno napako omejeno z  $\mathcal{O}(u)\kappa(X)$ .

Če upoštevamo še obratno stabilnost razcepa Choleskega v točki a), celotni postopek vrne singularne vrednosti z relativno napako  $\mathcal{O}(u)\kappa(X)^2$ .

V primeru, ko je matrika  $X$  dobro pogojena, lahko s tem postopkom vse lastne vrednosti simetrične pozitivno definitne matrike  $A$  izračunamo z visoko relativno natančnostjo. Za zgled si oglejmo naslednji primer iz [8].

**Zgled 3.6** *Vzemimo*

$$A = \begin{bmatrix} 1 & \sqrt{\delta} & \sqrt{\delta} \\ \sqrt{\delta} & 1 & 10\delta \\ \sqrt{\delta} & 10\delta & 100\delta \end{bmatrix},$$

kjer je  $\delta = 10^{-20}$ . Če  $A$  reduciramo na tridiagonalno matriko, je točen rezultat

$$T = \begin{bmatrix} 1 & \sqrt{2\delta} & \\ \sqrt{2\delta} & 1/2 + 60\delta & 1/2 - 50\delta \\ & 1/2 - 50\delta & 1/2 + 40\delta \end{bmatrix},$$

pri numeričnem računanju v dvojni natančnosti pa dobimo

$$\tilde{T} = \begin{bmatrix} 1 & \sqrt{2\delta} & \\ \sqrt{2\delta} & 1/2 & 1/2 \\ & 1/2 & 1/2 \end{bmatrix},$$

ki ni simetrična pozitivno definitna, pri čemer smo izgubili relativno natančnost najmanjše lastne vrednosti.

Razcep Choleskega in nato enostranski Jacobi izračunata vse lastne vrednosti s polno natančnostjo:  $\lambda_1 = 1 + \sqrt{\delta}$ ,  $\lambda_2 = 1 - \sqrt{\delta}$ ,  $\lambda_3 = 99\delta$ .  $\square$

## Dodatna literatura

Za računanje singularnega razcepa je na voljo obsežna literatura. V slovenščini lahko skoraj celotno snov najdete v knjigi [8]. Kar se tiče tuje literature, lahko skoraj vse algoritme, ki smo jih navedli, skupaj s potrebno analizo, najdete v [9]. Uporabna sta tudi učbenika [7] in [12].

Več o kvadratnem problemu lastnih vrednosti lahko najdete v preglednem članku [23].

Več o regularizaciji s pomočjo singularnega razcepa lahko najdete npr. v [11].

## Poglavje 4

# Enakomerna aproksimacija

### 4.1 Aproksimacija

Denimo, da imamo podano funkcijo  $f$ , ki je zvezna na intervalu  $[a, b]$ . Radi bi jo čim boljše aproksimirali s kakšno preprosto funkcijo  $g$ , ki bi bila lažje izračunljiva in bi se jo dalo po potrebi preprosto odvajati oziroma integrirati.

Ponavadi za aproksimacijo uporabljamo polinome, odvisno od oblike funkcije  $f$  in namena aproksimacije pa uporabljamo tudi druge funkcije, npr. trigonometrične polinome, racionalne funkcije in zlepke. Polinomi imajo zelo lepe lastnosti, saj je računanje vrednosti enostavno, prav tako pa jih je enostavno odvajati in integrirati. Pomembno je tudi, da množica polinomov stopnje kvečjemu  $k$  tvori vektorski prostor, ki ga bomo označili s  $P_k$ . Kadar iščemo aproksimacijsko funkcijo, ki se jo da zapisati kot linearno kombinacijo baznih funkcij, govorimo o linearni aproksimaciji.

Kakovost aproksimacije merimo z normo ostanka  $\|f - g\|$ . Različne norme pripeljejo do različnih aproksimacijskih funkcij. Najpogostejše izbire so:

- a) *diskretna aproksimacija po metodi najmanjših kvadratov*, kjer iščemo  $g$ , ki minimizira

$$\sum_{i=1}^m |f(x_i) - g(x_i)|^2,$$

- b) *zvezna aproksimacija po metodi najmanjših kvadratov*, kjer iščemo  $g$ , ki minimizira

$$\|f - g\|_2 := \left( \int_a^b |f(x) - g(x)|^2 dx \right)^{1/2},$$

- c) *diskretna enakomerna aproksimacija*, kjer iščemo  $g$ , ki minimizira

$$\max_{i=1, \dots, m} |f(x_i) - g(x_i)|,$$

- d) *zvezna enakomerna aproksimacija*, kjer iščemo  $g$ , ki minimizira

$$\|f - g\|_\infty := \max_{x \in [a, b]} |f(x) - g(x)|.$$

Z aproksimacijo po metodi najmanjših kvadratov smo se že srečali v ?? poglavju. Izkazalo se je, da v primeru linearnega modela problem lahko prevedemo na reševanje linearnega sistema in tako dobimo aproksimacijsko funkcijo z direktno metodo. V tem poglavju si bomo podrobneje pogledali enakomerno aproksimacijo, kjer ne obstajajo direktne metode za izračun najboljše aproksimacijske funkcije.

Enakomerno aproksimacijo srečamo med drugim pri računanju vrednosti elementarnih funkcij. Če npr. od računalnika zahtevamo, da nam vrne vrednost funkcije  $\sin$  pri izbranem argumentu  $x$ , bomo kot rezultat dobili numerično izračunano vrednost  $g(x)$ , kjer je  $g$  polinom (ali racionalna funkcija) nizke stopnje, izbran tako, da se izračunani  $g(x)$  od  $\sin(x)$  ne razlikuje za več kot osnovno zaokrožitveno napako. Za elementarne funkcije zadošča, če poznamo tak polinom na nekem zaprtem intervalu  $[a, b]$ , saj lahko izračun  $f(x)$  prevedemo na ta interval. Npr., za izračun  $f(x) = \sin(x)$  očitno zadošča, če znamo aproksimirati vrednosti funkcije na intervalu  $[0, \pi]$ .

## 4.2 Enakomerna aproksimacija s polinomi

Vsako zvezno funkcijo se da na zaprtem intervalu poljubno dobro aproksimirati s polinomom. To nam zagotavlja znani Weierstrassov izrek, katerega dokaz lahko najdete npr. v [17].

**Izrek 4.1 (Weierstrass)** Naj bo  $f$  zvezna funkcija na  $[a, b]$ . Za poljuben  $\epsilon > 0$  obstaja polinom  $p$ , da je  $|f(x) - p(x)| \leq \epsilon$  za vsak  $x \in [a, b]$ .

Seveda nam zgornji izrek nič ne pove o stopnji polinoma  $p$ , ki dovolj dobro aproksimira funkcijo  $f$ . Če se omejimo le na polinome stopnje kvečjemu  $n$ , potem lahko definiramo

$$\text{dist}(f, P_n) = \min_{p \in P_n} \|f - p\|_\infty.$$

Weierstrassov izrek nam pove, da za vsako zvezno funkcijo  $f$  velja  $\lim_{n \rightarrow \infty} \text{dist}(f, P_n) = 0$ , nas pa zanima, kako poiščemo polinom najboljše enakomerne aproksimacije stopnje kvečjemu  $n$ .

Izkaže se, da je zadosten pogoj za to, da je  $p$  polinom najboljše enakomerne aproksimacije za  $f$ , ta, da na  $[a, b]$  razlika  $f - p$  alternirajoče zavzame maksimum po absolutni vrednosti v  $n + 2$  točkah.

**Izrek 4.2** Naj bo  $f$  zvezna funkcija na  $[a, b]$ . Če za polinom  $p \in P_n$  na  $[a, b]$  obstaja  $n + 2$  takih točk  $x_0 < x_1 < \dots < x_{n+1}$ , da za razliko  $r = f - p$  velja, da v točkah  $x_0, \dots, x_{n+1}$  doseže  $|r|$  svoj maksimum, pri čemer predznak  $r$  med točkami z zaporednim indeksom alternira, potem je  $p$  polinom najboljše enakomerne aproksimacije za  $f$  na  $[a, b]$ .

*Dokaz.* Denimo, da  $p$  ni polinom najboljše aproksimacije. Potem obstaja tak polinom  $q \in P_n$ , da velja

$$|f(x_i) - q(x_i)| < |f(x_i) - p(x_i)|, \quad i = 0, \dots, n + 1.$$

Iz zgornje neenakosti sledi, da je predznak  $(f(x_i) - p(x_i)) - (f(x_i) - q(x_i))$  enak predznaku  $f(x_i) - p(x_i)$  za  $i = 0, \dots, n + 1$ . Toda,

$$(f - p) - (f - q) = q - p$$

je polinom stopnje kvečjemu  $n$ , ki ne more alternirati v  $n + 2$  točkah, saj bi potem imel  $n + 1$  ničel. Torej tak polinom  $q$  ne obstaja in je  $p$  res polinom najboljše enakomerne aproksimacije. ■

Dokazati se da tudi, da tak polinom res obstaja, dokaz lahko najdete npr. v [17]. To dejstvo se potem uporabi v Remezovem<sup>1</sup> postopku za iskanje polinoma najboljše enakomerne aproksimacije, ki je opisan v algoritmu 4.1. Postopek skonstruira zaporedje polinomov, ki konvergirajo proti polinomu  $p \in P_n$ , ki zadošča predpostavkam izreka 4.2 in je tako polinom najboljše enakomerne aproksimacije.

---

**Algoritem 4.1** Remezov postopek. Začetni podatki so funkcija  $f$ , interval  $[a, b]$ , stopnja  $n$  in toleranca  $\epsilon$ . Algoritem vrne polinom najboljše enakomerne aproksimacije stopnje kvečjemu  $n$  za  $f$  na  $[a, b]$ .

---

izberi delilne točke  $x_0 < x_1 < \dots < x_{n+1}$  na intervalu  $[a, b]$   
ponavljaj

- a) določi polinom  $p \in P_n$ , za katerega velja  $f(x_i) - p(x_i) = (-1)^i r$  za nek  $r$
  - b) določi točko  $y$ , kjer  $|f - p|$  doseže maksimum na  $[a, b]$
  - c) če je  $|f(y) - p(y)| - |r| < \epsilon$ , končaj in vrni  $p$
  - č) poišči delilno točko  $z$  za zamenjavo:
    - $x_k \leq y \leq x_{k+1}$ : če je predznak  $f - p$  v  $x_k$  in  $y$  enak, potem  $z = x_k$ , sicer  $z = x_{k+1}$ ;
    - $a \leq y \leq x_0$ : če je predznak  $f - p$  v  $x_0$  in  $y$  enak, potem  $z = x_0$ , sicer  $z = x_{n+1}$ ;
    - $x_{n+1} \leq y \leq b$ : če je predznak  $f - p$  v  $x_{n+1}$  in  $y$  enak, potem  $z = x_{n+1}$ , sicer  $z = x_0$
  - d) zamenjaj točko  $z$  z  $y$  in po potrebi preštevilči delilne točke po velikosti
- 

Poglejmo si podrobneje, kako izvedemo posamezne korake algoritma 4.1. Pri točki a) moramo rešiti linearni sistem. Če npr. polinom  $p$  zapišemo v standardni bazi kot  $p(x) = a_0 + \dots + a_n x^n$ , dobimo sistem

$$\begin{aligned} a_0 + a_1 x_0 + \dots + a_n x_0^n - (-1)^0 r &= f(x_0) \\ &\vdots \\ a_0 + a_1 x_{n+1} + \dots + a_n x_{n+1}^n - (-1)^{n+1} r &= f(x_{n+1}) \end{aligned}$$

za koeficiente  $a_0, \dots, a_n$  in  $r$ .

Pri točki b) moramo uporabiti eno izmed numeričnih metod za iskanje lokalnih ekstremov. Če lahko  $f$  odvajamo, potem lahko npr. z Newtonovo metodo pridemo do približkov za lokalne ekstreme.

Naslednja lema nam zagotavlja, da se v primeru izpolnjene točke c) polinom  $p$  od polinoma najboljše enakomerne aproksimacije na celem intervalu  $[a, b]$  razlikuje za manj kot  $\epsilon$ .

**Lema 4.3** Naj bo  $f$  zvezna funkcija na  $[a, b]$ . Če za polinom  $p \in P_n$  velja, da razlika  $f - p$  alternira v

---

<sup>1</sup>Evgenij Jakovlevič Remez (1896-1975) je bil ukrajinski matematik, rojen v Belorusiji. Postopek konstrukcije polinoma najboljše enakomerne aproksimacije je objavil leta 1936 v francoščini kot Eugene Remes, zato je postopek znan tudi kot Remesov. Tako kot se danes veliko ukrajinskih in ruskih matematičnih člankov objavi v angleščini, da so dostopni širšemu občinstvu, tako so v obdobju med prvo in drugo svetovno vojno za to uporabljali francoščino.

$n + 2$  točkah  $x_0 < x_1 < \dots < x_{n+1}$  z intervala  $[a, b]$ , potem velja

$$\min_{i=0, \dots, n+1} |f(x_i) - p(x_i)| \leq \text{dist}(f, P_n) \leq \|f - p\|_\infty.$$

*Dokaz.* Desna neenakost  $\text{dist}(f, P_n) \leq \|f - p\|_\infty$  očitno velja za vsak polinom  $p \in P_n$ .

Naj bo  $q$  polinom najboljše enakomerne aproksimacije za  $f$  na  $[a, b]$ . Denimo, da leva neenakost ne velja in je torej  $|f(x_i) - p(x_i)| > \text{dist}(f, P_n)$  za  $i = 0, \dots, n + 1$ . Od tod tako kot v dokazu izreka 4.2 sledi, da je v vseh točkah  $x_i$ ,  $i = 0, \dots, n + 1$ , predznak  $(f(x_i) - p(x_i)) - (f(x_i) - q(x_i))$  enak predznaku  $f(x_i) - p(x_i)$ . To pa pomeni, da polinom  $q$  ne more obstajati, kar je seveda protislovje. ■

**Zgled 4.1** Poglejmo si preprost primer uporabe Remezovega postopka. Funkcijo  $e^x$  bi radi na intervalu  $[0, 1]$  po metodi najboljše enakomerne aproksimacije čim boljše aproksimirali s premico oblike  $y = a_0 + a_1 x$ .

Denimo, da smo na začetku izbrali točke  $x_0 = 0$ ,  $x_1 = 1/2$ ,  $x_2 = 1$ . Sedaj nastavimo sistem  $f(x_i) - p(x_i) = (-1)^i r$  za  $i = 0, 1, 2$ . Enačbam

$$\begin{aligned} 1 - a_0 &= r \\ e^{1/2} - a_0 - 1/2 a_1 &= -r \\ e - a_0 - a_1 &= r \end{aligned}$$

ustreza rešitev  $a_0 = \frac{1}{4}(1 - e^{1/2})(3 + e^{1/2})$ ,  $a_1 = e - 1$ ,  $r = \frac{1}{4}(e^{1/2} - 1)^2$ .

Z odvajanjem  $f - p$  ugotovimo, da ima  $|f - p|$  maksimum pri  $y = \ln(e - 1) = 0.54132$ . Ker velja  $f(y) - p(y) > 0$ , v naslednjem koraku zamenjamo točko  $x_1$  z  $y$ . Nove točke za naslednji korak so tako  $x_0 = 0$ ,  $x_1 = \ln(e - 1)$ ,  $x_2 = 1$ .

Z novimi točkami ponovimo postopek od začetka in dobimo  $a_0 = 1/2(e + (1 - e) \ln(e - 1))$ ,  $a_1 = e - 1$ ,  $r = 1 - a_0$ . Ker se v naslednjem koraku izkaže, da  $|f - p|$  doseže maksimum na  $[0, 1]$  ravno v delilnih točkah, se postopek konča in našli smo polinom najboljše enakomerne aproksimacije. □

Zgornji zgled ni povsem tipičen, saj se ponavadi Remezov postopek ne konča v končnem številu korakov. Običajno postopek končamo, ko je pogoj v točki c) izpolnjen za izbrano toleranco  $\epsilon$ , saj  $|r|$  konvergira proti  $\text{dist}(f, P_n)$ . Pokazati se da (glej npr. [17]), da je konvergenca linearna.

Omenimo še, da obstaja tudi različica Remezovega postopka, kjer v vsakem koraku zamenjamo vse točke s točkami, kjer ima razlika  $f - p$  lokalne ekstreme. Pri tej metodi imamo sicer več dela z iskanjem ekstremov, ima pa zato v bližini rešitve kvadratično konvergenco.

### 4.3 Ekonomizacija Čebiševa

Remezov postopek ni preprost in velikokrat bi nam zadoščalo, če bi znali funkcijo  $f$  na kak enostavnejši način aproksimirati s polinomom nizke stopnje  $p$ , za katerega bi veljalo  $\|f - p\|_\infty \leq \epsilon$  za izbrano konstanto  $\epsilon > 0$ .

Tako npr. vemo, da lahko vsako dovolj gladko funkcijo aproksimiramo s pomočjo Taylorjevega polinoma, ki ga dobimo tako, da odrežemo razvoj v Taylorjevo vrsto. Iz ocene za napako

Taylorjeve vrste lahko ocenimo, koliko členov potrebujemo, da bo skupna napaka na celem intervalu  $[a, b]$  dovolj majhna. Ker za Taylorjev polinom velja, da dobro aproksimira le v okolici točke, okrog katere smo ga razvili, je smiselno vzeti razvoj v Taylorjevo vrsto okrog točke  $(a + b)/2$ .

Izkaže se, da so polinomi, ki jih tako dobimo, ponavadi občutno višjih stopenj kot pa bi jih lahko dobili preko algoritmov za najboljšo enakomerno aproksimacijo. Stopnjo polinoma lahko zmanjšamo s pomočjo ekonomizacije Čebiševa, ki temelji na lastnostih polinomov Čebiševa<sup>2</sup>

**Definicija 4.4** Polinomi Čebiševa so definirani s  $T_0(x) = 1$ ,  $T_1(x) = x$  in z rekurzivno zvezo

$$T_{m+1}(x) = 2xT_m(x) - T_{m-1}(x)$$

za  $m \geq 1$ .

Preverimo lahko, da za polinome Čebiševa velja:

- $T_m(x) = 2^{m-1}x^m + \mathcal{O}(x^{m-1})$ ,
- $T_m(x) = \begin{cases} \cos(m \arccos x), & |x| \leq 1, \\ \cosh(m \operatorname{arccosh} x), & |x| \geq 1, \end{cases}$
- $|T_m(x)| \leq 1$  za  $|x| \leq 1$ .

Naslednji seznam prikazuje nekaj prvih polinomov Čebiševa in kako se z njimi izražajo polinomi iz standardne baze:

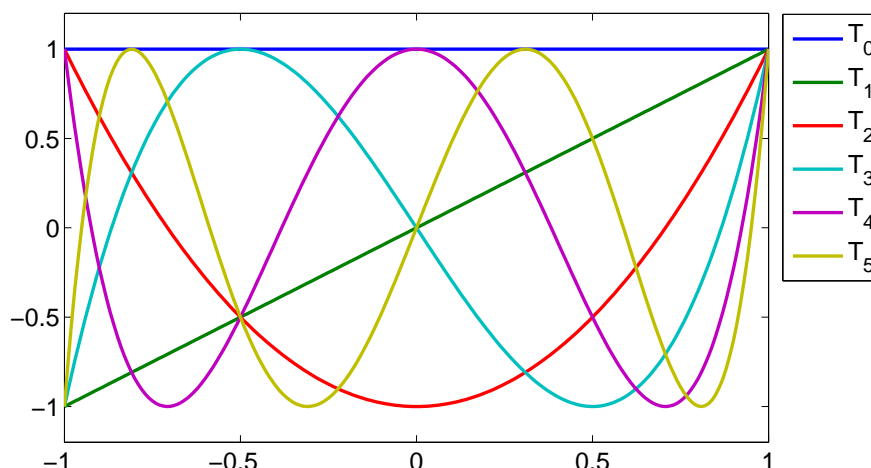
$$\begin{array}{ll} T_0(x) = 1 & 1 = T_0 \\ T_1(x) = x & x = T_1 \\ T_2(x) = 2x^2 - 1 & x^2 = \frac{1}{2}(T_0 + T_2) \\ T_3(x) = 4x^3 - 3x & x^3 = \frac{1}{4}(3T_1 + T_3) \\ T_4(x) = 8x^4 - 8x^2 + 1 & x^4 = \frac{1}{8}(3T_0 + 4T_2 + T_4) \\ T_5(x) = 16x^5 - 20x^3 + 5x & x^5 = \frac{1}{16}(10T_1 + 5T_3 + T_5). \end{array}$$

Razlog, da so polinomi Čebiševa zelo pomembni pri najboljši enakomerni aproksimaciji se skriva v naslednji lemi.

**Lema 4.5** Velja:

- Polinom stopnje  $n - 1$ , ki najboljše enakomerno aproksimira  $x^n$  na intervalu  $[-1, 1]$  je podan s  $q(x) = x^n - 2^{1-n}T_n(x)$ .
- Izmed vseh polinomov stopnje  $n$  z vodilnim koeficientom 1 ima  $2^{1-n}T_n$  najmanjšo neskončno normo na intervalu  $[-1, 1]$ .

<sup>2</sup>Polinomi se imenujejo po ruskem matematiku Pafnutiju Lvoviču Čebiševu (1821-1894), s črko  $T$  pa se označujejo zaradi transkripcije njegovega priimka v Tchebychef. Čebišev je prvi razvil splošno teorijo ortogonalnih polinomov, med drugim pa je znan tudi po pomembnih rezultatih iz aproksimacije in teorije verjetnosti.

Slika 4.1: Nekaj prvih polinomov Čebiševa na intervalu  $[-1, 1]$ .

*Dokaz.* Velja  $2^{1-n}T_n(x) = x^n - q(x)$ . Ker  $2^{1-n}T_n$  na  $[-1, 1]$  zavzame maksimum absolutne vrednosti v  $n + 1$  točkah z alternirajočim predznakom, je po izreku 4.2 potem  $q(x)$  polinom najboljše enakomerne aproksimacije stopnje kvečjemu  $n - 1$  za  $x^n$  na  $[-1, 1]$ .

Poljuben polinom  $p$  stopnje  $n$  z vodilnim koeficientom 1 lahko zapišemo kot  $p(x) = x^n - r(x)$ , kjer je  $r$  polinom stopnje kvečjemu  $n - 1$ . Očitno bo norma  $\|p\|_\infty$  minimalna natanko takrat, ko bo  $r(x)$  polinom najboljše enakomerne aproksimacije za  $x^n$  na  $[-1, 1]$ , to pa pomeni  $r = q$  in  $p = 2^{1-n}T_n$ . ■

Denimo, da na intervalu  $[-1, 1]$  za polinom  $p$  stopnje  $n$  iščemo polinom najboljše enakomerne aproksimacije stopnje kvečjemu  $n - 1$ . Rešitev lahko direktno izračunamo s pomočjo leme 4.5. Če polinom  $p$  zapišemo v bazi polinomov Čebiševa kot

$$p = \sum_{i=0}^n a_i T_i,$$

potem je polinom najboljše aproksimacije stopnje  $n - 1$  kar

$$q = \sum_{i=0}^{n-1} a_i T_i.$$

Za razliko  $p - q = a_n T_n$  namreč velja, da ima na  $[-1, 1]$   $n + 1$  ekstremov, v katerih predznak alternira. Po izreku 4.2 je potem  $q$  polinom najboljše enakomerne aproksimacije iz  $P_{n-1}$  za  $p$  na  $[-1, 1]$ .

Na zgornji ugotovitvi temelji postopek ekonomizacije Čebiševa, zapisan v algoritmu 4.2, s katerim lahko polinom aproksimiramo s polinomom nižje stopnje in tako zmanjšamo število operacij potrebnih za izračun vrednosti aproksimacijskega polinoma. Če nismo na intervalu  $[-1, 1]$  to ni nobena težava, saj uporabimo linearno substitucijo in problem prevedemo na  $[-1, 1]$ .

V primeru, ko stopnjo polinoma zmanjšamo za ena, torej ko je  $k = n - 1$ , je  $p_{n-1}^{\text{econ}}$ , ki ga vrne algoritem 4.2, kar polinom najboljše enakomerne aproksimacije stopnje  $n - 1$  za začetni polinom  $p_n$ . V primeru  $k < n$  pa si lahko predstavljamo, da ekonomizacija po vrsti niža stopnjo poli-



**Algoritem 4.2** Ekonomizacija Čebiševa. Začetni podatki so polinom  $p_n$  stopnje  $n$  na intervalu  $[a, b]$  in  $k < n$ . Algoritem vrne polinom  $p_k^{\text{econ}}$  stopnje  $k$  in tak  $\delta > 0$ , da velja  $\|p_n - p_k^{\text{econ}}\| \leq \delta$ .

a) Polinom  $p_n$  s substitucijo preslikaj v polinom  $q_n$  na intervalu  $[-1, 1]$ , da je

$$g_n(\xi) = p_n \left( a + \frac{b-a}{2}(\xi + 1) \right).$$

b) Polinom  $q_n$  razvij po polinomih Čebiševa kot  $q_n(\xi) = d_0 + d_1 T_1(\xi) + \dots + d_n T_n(\xi)$ .

c) Razvoj  $q_n$  odreži pri  $k < n$  v  $q_k(\xi) = d_0 + d_1 T_1(\xi) + \dots + d_k T_k(\xi)$  in za oceno razlike vzemi  $\delta = |d_{k+1}| + \dots + |d_n|$ .

č) Polinom  $q_k$  izrazi spet v standardni bazi kot  $q_k(\xi) = b_0 + b_1 \xi + \dots + b_k \xi^k$ .

d) Polinom  $q_k$  z obratno substitucijo preslikaj nazaj na intervalu  $[a, b]$  v

$$p_k^{\text{econ}}(x) = q_k \left( -1 + 2 \frac{x-a}{b-a} \right).$$

noma za eno. Če vpeljemo še vmesne polinome  $p_{k+1}^{\text{econ}}, \dots, p_{n-1}^{\text{econ}}$ , potem je vedno  $p_m^{\text{econ}}$  polinom najboljše enakomerne aproksimacije stopnje  $m$  za polinom  $p_{m+1}^{\text{econ}}$  za  $m = n-1, n-2, \dots, k$ .

Vemo, da na  $[a, b]$  velja  $\|p_n - p_k^{\text{econ}}\|_{\infty} \leq \delta$ . Denimo, da smo začetni polinom  $p_n$  skonstruirali tako, da velja  $\|f - p_n\|_{\infty} \leq \epsilon$ . Od tod sledi

$$\|f - p_k^{\text{econ}}\|_{\infty} \leq \|f - p_n\|_{\infty} + \|p_n(x) - p_k^{\text{econ}}\|_{\infty} \leq \epsilon + \delta.$$

**Zgled 4.2** Funkcijo  $f(x) = e^x$  bi radi na intervalu  $[-1/2, 1/2]$  aproksimirali s polinomom stopnje 3 oziroma 4. Če  $f$  aproksimiramo s prvimi 6 členi razvoja v Taylorjevo vrsto okrog 0, dobimo

$$p_5(x) = 1 + x + \frac{1}{2}x^2 + \frac{1}{6}x^3 + \frac{1}{24}x^4 + \frac{1}{120}x^5.$$

Napako Taylorjevega polinoma lahko ocenimo z

$$|f(x) - p_5(x)| \leq \frac{e^{1/2}}{720} = 3.6 \cdot 10^{-5}.$$

Če  $p_5$  preslikamo na  $[-1, 1]$  s substitucijo  $q_5(\xi) = p_5(x/2)$ , dobimo

$$q_5(\xi) = 1 + \frac{1}{2}\xi + \frac{1}{8}\xi^2 + \frac{1}{48}\xi^3 + \frac{1}{384}\xi^4 + \frac{1}{3840}\xi^5.$$

Polinom  $q_5$  se v polinomih Čebiševa izraža kot

$$q_5 = \frac{1089}{1024}T_0 + \frac{3169}{6144}T_1 + \frac{49}{768}T_2 + \frac{65}{12288}T_3 + \frac{1}{3072}T_4 + \frac{1}{61440}T_5.$$

Ker iščemo polinom stopnje 4, odrežemo zadnji člen in dobimo

$$q_4 = \frac{1089}{1024}T_0 + \frac{3169}{6144}T_1 + \frac{49}{768}T_2 + \frac{65}{12288}T_3 + \frac{1}{3072}T_4.$$

Ko polinom preslikamo nazaj na  $[-1/2, 1/2]$  in ga spet izrazimo v standardni bazi, dobimo

$$p_4^{\text{econ}} = 1 + \frac{6143}{6144}x + \frac{1}{2}x^2 + \frac{65}{384}x^3 + \frac{1}{24}x^4.$$

Podobno, če bi odrezali še člen  $T_4$ , bi dobili polinom stopnje 3

$$p_3^{\text{econ}} = \frac{3071}{3072} + \frac{6143}{6144}x + \frac{49}{96}x^2 + \frac{65}{384}x^3.$$

Vemo, da velja  $\|p_5 - p_4^{\text{econ}}\| \leq 1/61400 = 1.6 \cdot 10^{-5}$  oziroma  $\|p_5 - p_3^{\text{econ}}\| \leq 1/61400 + 1/3072 = 3.4 \cdot 10^{-4}$ . Od tod lahko dobimo skupno oceno  $\|f - p_4^{\text{econ}}\| \leq 1.6 \cdot 10^{-5} + 3.6 \cdot 10^{-5} = 5.2 \cdot 10^{-5}$  oziroma  $\|f - p_3^{\text{econ}}\| \leq 3.4 \cdot 10^{-4} + 3.6 \cdot 10^{-5} = 3.8 \cdot 10^{-4}$ .

Če bi izračunali še polinoma najboljše enakomerne aproksimacije stopnje 3 in 4 bi dobili  $\text{dist}(f, P_4) = 1.6 \cdot 10^{-5}$  in  $\text{dist}(f, P_3) = 3.3 \cdot 10^{-4}$ . Od tod lahko ugotovimo:  $\|f - p_4^{\text{econ}}\| \leq 3.3 \cdot \text{dist}(f, P_4)$  in  $\|f - p_3^{\text{econ}}\| \leq 1.2 \cdot \text{dist}(f, P_3)$ , kar pomeni, da smo brez Remezovega postopka dobili zelo dobra približka za polinoma najboljše enakomerne aproksimacije. Polinom stopnje 3 je skoraj optimalen, pri polinomu stopnje 4 pa bi lahko dobili še boljši približek, če bi na začetku uporabili več členov razvoja v Taylorjevo vrsto.  $\square$

Iz samega postopka Čebiševe ekonomizacije tudi ugotovimo, da je namesto ekvidistantnih točk bolje za začetno delitev pri Remezovem postopku uporabiti točke, kjer polinom Čebiševa  $T_{n+1}$ , preslikan na  $[a, b]$ , alternirajoče doseže svoje ekstreme, to pa so

$$x_k = \frac{a+b}{2} + \frac{b-a}{2} \cos\left(\frac{k}{n+1}\right)$$

za  $k = 0, \dots, n+1$ .

## Dodatna literatura

V slovenščini lahko več podrobnosti o problemu najboljše enakomerne aproksimacije najdete v knjigi [17], o ekonomizaciji Čebiševa pa imate na voljo literaturo v angleščini, npr. [6].

## Poglavje 5

# Interpolacija

### 5.1 Uvod

Podane imamo vrednosti funkcije  $f$  v  $n + 1$  paroma različnih točkah  $x_0, \dots, x_n$ , iščemo pa preprostejšo *interpolacijsko funkcijo*  $g$ , ki se v točkah  $x_i$  ujema s funkcijo  $f$ , torej  $g(x_i) = f(x_i)$  za  $i = 0, \dots, n$ . Interpolacijsko funkcijo potrebujemo npr. takrat, kadar imamo funkcijo  $f$  tabelirano, zanima pa nas vrednost funkcije v točki  $x$ , ki v tabeli ni vsebovana. Kot približek za vrednost  $f(x)$  potem vzamemo vrednost interpolacijske funkcije  $g(x)$ . Preprost primer za interpolacijsko funkcijo je npr. kosoma linearna funkcija, kjer vrednost v točki  $x_i < x < x_{i+1}$  dobimo s konveksno kombinacijo vrednosti v  $x_i$  in  $x_{i+1}$ .

Ponavadi za interpolacijsko funkcijo vzamemo polinome. Obstaja natanko en polinom stopnje  $n$  ali manj, ki se v  $n + 1$  paroma različnih točkah  $x_0, \dots, x_n$  ujema s funkcijo  $f$  in tak polinom imenujemo *interpolacijski polinom*.

Polinomsko interpolacijo so v preteklosti uporabljali večinoma za določanje vrednosti tabeliranih funkcij, danes pa se v glavnem uporablja kot orodje pri numeričnem odvajanju, integriranju in reševanju diferencialnih enačb.

Namesto polinomov lahko odvisno od funkcije in potreb uporabljamo tudi druge funkcije, npr. trigonometrične polinome, racionalne funkcije in zlepke. Polinomi imajo zelo lepe lastnosti, saj sta konstrukcija in računanje vrednosti enostavni, prav tako pa jih je zelo preprosto odvajati in integrirati.

Skoraj nikoli ne iščemo enačbe polinoma v standardni bazi. Največkrat zadošča, da znamo izračunati vrednost polinoma v dani točki.

### 5.2 Interpolacijski polinom

Denimo, da imamo paroma različne točke  $x_0, \dots, x_n$  in vrednosti  $y_0, \dots, y_n$ . Iščemo polinom  $I_n$  stopnje  $n$  ali manj, za katerega velja  $I_n(x_i) = y_i$  za  $i = 0, \dots, n$ .

Naivna varianta je, da za koeficiente polinoma

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

nastavimo linearni sistem

$$\begin{array}{cccccc} a_0 & + & a_1 x_0 & + & a_2 x_0^2 & + & \cdots & + & a_n x_0^n & = & y_0 \\ a_0 & + & a_1 x_1 & + & a_2 x_1^2 & + & \cdots & + & a_n x_1^n & = & y_1 \\ & & \vdots & & \vdots & & & & \vdots & & \\ a_0 & + & a_1 x_n & + & a_2 x_n^2 & + & \cdots & + & a_n x_n^n & = & y_n \end{array}$$

Matrika zgornjega sistema je t.i. *Vandermondova matrika*

$$V(x_0, x_1, \dots, x_n) = \begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^n \end{bmatrix}.$$

Zanjo velja, da je v primeru paroma različnih točk nesingularna, saj je

$$\det(V(x_0, x_1, \dots, x_n)) = \prod_{i>j} (x_i - x_j).$$

Izkaže se, da je računanje interpolacijskega polinoma preko reševanja zgornjega linearnega sistema po eni strani zelo zamudno, po drugi strani pa je sistem lahko zelo slabo pogojen. Če npr. interpoliramo na intervalu  $[0, 1]$  in vzamemo ekvidistantne točke  $x_0 = 0, x_1 = 1/n, \dots, x_n = 1$ , potem ima Vandermondova matrika v primeru  $n = 5$  pogojenostno število  $4.9 \cdot 10^3$ , pri  $n = 10$  je  $1.2 \cdot 10^8$ , pri  $n = 20$  pa že kar  $9.7 \cdot 10^{16}$ .

V resnici znamo interpolacijski polinom izračunati na ekonomičnejši način. To izkoriščajo tudi hitre metode za reševanje sistemov z Vandermondovimi matrikami, saj lahko reševanje sistemov te oblike prevedemo na iskanje koeficientov interpolacijskega polinoma.

**Izrek 5.1** Za paroma različne točke  $x_0, \dots, x_n$  in vrednosti  $y_0, \dots, y_n$  obstaja natanko en polinom  $I_n$  stopnje  $n$  ali manj, za katerega velja  $I_n(x_i) = y_i, i = 0, \dots, n$ .

*Dokaz.* Eksistenco pokažemo s konstrukcijo. Definirajmo polinome

$$L_{n,i}(x) = \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)} = \prod_{\substack{k=0 \\ k \neq i}}^n \frac{x - x_k}{x_i - x_k},$$

ki so stopnje  $n$ , imenujemo pa jih *Lagrangeevi koeficienti*. Očitno velja  $L_{n,i}(x_j) = \delta_{ij}$ . Če sedaj definiramo

$$I_n(x) = \sum_{k=0}^n y_k L_{n,k}(x),$$

je to očitno polinom stopnje  $n$  ali manj, prav tako pa velja  $I_n(x_i) = y_i$  za  $i = 0, \dots, n$ .

Enoličnost sledi iz tega, da se dva različna polinoma stopnje  $n$  ali manj ne moreta ujemati v več kot  $n$  različnih točkah. Če bi se namreč ujemale, bi imela njuna razlika več kot  $n$  ničel. Ker pa je razlika spet polinom stopnje največ  $n$ , je to možno le, če je ta polinom kar enak 0. ■

Če vpeljemo polinom  $\omega(x) = (x - x_0)(x - x_1) \cdots (x - x_n)$ , potem velja

$$L_{n,i}(x) = \frac{\omega(x)}{(x - x_i)\omega'(x_i)}.$$

Interpolacijski polinom  $I_n$  se s funkcijo  $f$  ujema v točkah  $x_0, \dots, x_n$ , drugje pa ne nujno. Lahko pa ocenimo razliko.

**Izrek 5.2** Če je  $f$   $(n + 1)$ -krat zvezno odvedljiva funkcija na intervalu  $[a, b]$ , ki vsebuje vse paroma različne točke  $x_0, \dots, x_n$ , potem za vsak  $x \in [a, b]$  obstaja tak  $\xi \in (a, b)$ , da velja

$$f(x) - I_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x) \quad (5.1)$$

in  $\min(x, x_0, \dots, x_n) < \xi < \max(x, x_0, \dots, x_n)$ .

*Dokaz.* Za  $x = x_i$  nimamo kaj dokazovati, saj je na obeh straneh 0. Za ostale pa definiramo funkcijo

$$F(z) := f(z) - I_n(z) - R \cdot \omega(z). \quad (5.2)$$

$F$  je očitno  $(n + 1)$ -krat zvezno odvedljiva in velja  $F(x_i) = 0$  za  $i = 0, \dots, n$ . Za poljubno točko  $x \in [a, b]$ , ki je različna od  $x_0, \dots, x_n$ , lahko konstanto  $R$  določimo tako, da bo  $F(x) = 0$ .

Pri tako izbranem  $R$  ima  $F$  vsaj  $n + 2$  različnih ničel na  $[a, b]$ . Po Rolleovem izreku ima zato  $F'$  vsaj  $n + 1$  različnih ničel na  $(a, b)$ ,  $F''$  vsaj  $n$  ničel, ..., in končno,  $F^{(n+1)}$  ima vsaj eno ničlo na  $(a, b)$ , ki jo označimo s  $\xi$ .

Iz enačbe (5.2) sledi  $0 = F^{(n+1)}(\xi) = f^{(n+1)}(\xi) - R(n+1)!$ , torej

$$R = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

in za izbrani  $x$  smo dokazali (5.1). ■

Lagrangeeva oblika ni najprimernejša za konstrukcijo in računanje vrednosti interpolacijskega polinoma. Prva težava je veliko število operacij, druga težava pa je, da moramo že na začetku določiti stopnjo polinoma. Boljše so npr. zaporedne linearne interpolacije, pa tudi deljene in končne diference.

Označimo z  $I_{i,i+1,\dots,i+k}$  interpolacijski polinom, ki se z  $f$  ujema v točkah  $x_i, x_{i+1}, \dots, x_{i+k}$ . Hitro se lahko prepričamo, da velja

$$I_{i,i+1,\dots,i+k}(x) = \frac{I_{i,i+1,\dots,i+k-1}(x)(x - x_{i+k}) - I_{i+1,i+2,\dots,i+k}(x)(x - x_i)}{x_i - x_{i+k}}$$

oziroma

$$I_{i,i+1,\dots,i+k}(x) = \frac{1}{x_{i+k} - x_i} \begin{vmatrix} x - x_i & I_{i,i+1,\dots,i+k-1}(x) \\ x - x_{i+k} & I_{i+1,i+2,\dots,i+k}(x) \end{vmatrix}.$$

Z zaporednimi linearnimi interpolacijami lahko povečujemo stopnjo interpolacijskega polinoma. Obstaja več postopkov, opisali pa bomo *Nevilleovo shemo*, ki se jo da preprosto implementirati v računalniku. Predstavimo jo s trikotno shemo:

$x_i$	$x - x_i$	$y_i$	$I_{..}(x)$	$I_{...}(x)$	$I_{....}(x)$
$x_0$	$x - x_0$	$y_0$	$I_{01}(x)$		
$x_1$	$x - x_1$	$y_1$	$I_{12}(x)$	$I_{012}(x)$	
$x_2$	$x - x_2$	$y_2$	$I_{23}(x)$	$I_{123}(x)$	$I_{0123}(x)$
$x_3$	$x - x_3$	$y_3$			

**Zgled 5.1** Poišči vrednost interpolacijskega polinoma za podatke  $x : 0, 2, 4$  in  $f(x) : 2, 4, 8$  v točki  $x = 1$ .

Uporabili bomo Nevilleovo shemo:

$x_i$	$x - x_i$	$y_i$	$I_{..}(x)$	$I_{...}(x)$
$x_0 = 0$	$x - x_0 = 1$	$y_0 = 2$	$I_{01}(x) = 3$	
$x_1 = 2$	$x - x_1 = -1$	$y_1 = 4$	$I_{12}(x) = 2$	$I_{012}(x) = \frac{11}{4}$
$x_2 = 4$	$x - x_2 = -3$	$y_2 = 9$		

Do vrednosti smo prišli z naslednjimi računi:

$$I_{01}(x) = \frac{1}{x_1 - x_0} \begin{vmatrix} x - x_0 & I_0(x) \\ x - x_1 & I_1(x) \end{vmatrix} = \frac{1}{2} \begin{vmatrix} 1 & 2 \\ -1 & 4 \end{vmatrix} = 3,$$

$$I_{12}(x) = \frac{1}{x_2 - x_1} \begin{vmatrix} x - x_1 & I_1(x) \\ x - x_2 & I_2(x) \end{vmatrix} = \frac{1}{2} \begin{vmatrix} -1 & 4 \\ -3 & 8 \end{vmatrix} = 2,$$

$$I_{012}(x) = \frac{1}{x_2 - x_0} \begin{vmatrix} x - x_0 & I_{01}(x) \\ x - x_2 & I_{12}(x) \end{vmatrix} = \frac{1}{4} \begin{vmatrix} 1 & 3 \\ -3 & 2 \end{vmatrix} = \frac{11}{4}. \quad \square$$

### 5.3 Deljene difference

**Definicija 5.3** Deljena diferenca  $f[x_0, x_1, \dots, x_k]$  je vodilni koeficient interpolacijskega polinoma stopnje  $k$ , ki se ujema s funkcijo  $f$  v paroma različnih točkah  $x_0, x_1, \dots, x_k$ .

**Izrek 5.4** Za deljene difference velja:

a)

$$P_n(x) = f[x_0] + f[x_0, x_1](x - x_0) + \dots + f[x_0, x_1, \dots, x_n](x - x_0) \dots (x - x_{n-1}) \quad (5.3)$$

je interpolacijski polinom, ki se v točkah  $x_0, \dots, x_n$  ujema z  $f$ .

b)  $f[x_0, x_1, \dots, x_k]$  je simetrična funkcija svojih argumentov.

c)  $(\alpha f + \beta g)[x_0, \dots, x_k] = \alpha f[x_0, \dots, x_k] + \beta g[x_0, \dots, x_k]$ .

d) Velja rekurzivna formula

$$f[x_0, \dots, x_k] = \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0}. \quad (5.4)$$

Dokaz.

a) Uporabimo indukcijo. Očitno se  $f[x_0]$  ujema z  $f(x_0)$ . Naj bo

$$P_i(x) = f[x_0] + f[x_0, x_1](x - x_0) + \dots + f[x_0, x_1, \dots, x_i](x - x_0) \dots (x - x_{i-1}).$$

Polinom, ki interpolira  $f$  v  $x_0, \dots, x_{i+1}$  lahko zapišemo v obliki

$$P_{i+1}(x) = P_i(x) + c(x - x_0) \dots (x - x_i).$$

Po definiciji deljene diference se mora  $c$  ujemati z  $f[x_0, \dots, x_{i+1}]$ .

b), c) Očitno.

d) Naj bo  $p_0$  interpolacijski polinom, ki se ujema z  $f$  v točkah  $x_0, \dots, x_{k-1}$  in  $p_1$  polinom, ki se ujema z  $f$  v točkah  $x_1, \dots, x_k$ . Interpolacijski polinom  $p$ , ki se z  $f$  ujema v vseh točkah, ima obliko

$$p(x) = \frac{x - x_k}{x_0 - x_k} p_0(x) + \frac{x - x_0}{x_k - x_0} p_1(x).$$

Če primerjamo vodilne koeficiente, dobimo zvezo (5.4). ■

Obliko (5.3) imenujemo *Newtonov interpolacijski polinom*.

Za deljeni diferenci na eni in dveh točkah lahko izpeljemo preprosti formuli:

- $f[x_0] = f(x_0)$ ,
- $f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$ .

Kaj se zgodi, če vzamemo interpolacijski polinom na dveh točkah in pošljemo eno točko proti drugi? Iz

$$p(x) = f[x_0] + f[x_0, x_1](x - x_0) = f(x_0) + \frac{f(x_1) - f(x_0)}{x_1 - x_0}(x - x_0)$$

v limiti, ko gre  $x_1$  proti  $x_0$ , dobimo  $p(x) = f(x_0) + f'(x_0)(x - x_0)$ .

To pomeni, da lahko definicijo interpolacijskega polinoma in deljenih diferenc razširimo na polinome, ki se poleg vrednosti ujemajo tudi v odvodih. Če se točke ujemajo, potem to pomeni, da naj bo ujemanje tudi v odvodih. Tako npr. pri točkah  $x_1, x_2, x_2, x_3, x_3$  iščemo polinom  $p$  za katerega velja  $p(x_1) = f(x_1)$ ,  $p(x_2) = f(x_2)$ ,  $p'(x_2) = f'(x_2)$ ,  $p''(x_2) = f''(x_2)$ ,  $p(x_3) = f(x_3)$ ,  $p'(x_3) = f'(x_3)$ .

Če dopuščamo tudi ujemanje točk, potem za deljene diference velja rekurzivna zveza

$$f[x_0, x_1, \dots, x_k] = \begin{cases} \frac{f^k(x_0)}{k!}, & x_0 = x_1 = \dots = x_k \\ \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0}, & \text{sicer} \end{cases}.$$

Če je funkcija  $f$  v točki  $x_0$   $k$ -krat zvezno odvedljiva, potem iz interpolacijskega polinoma

$$p(x) = f[x_0] + f[x_0, x_1](x - x_0) + \dots + f[x_0, x_1, \dots, x_k](x - x_0) \dots (x - x_{k-1}),$$

ko pošljemo vse točke proti  $x_0$ , dobimo Taylorjev polinom razvoja  $f$  okrog točke  $x_0$ :

$$T_k(x) = f[x_0] + \frac{f'(x_0)}{1!}(x - x_0) + \dots + \frac{f^{(k)}(x_0)}{k!}(x - x_0)^k.$$

Deljene diference računamo v trikotni shemi iz katere na koncu hitro preberemo enačbo polinoma oz. izračunamo njegovo vrednost v iskani točki. Konstrukcija je razvidna iz naslednjega zglada.

**Zgled 5.2** Zapiši enačbo polinoma stopnje 5, za katerega velja:  $p(0) = 1$ ,  $p'(0) = 2$ ,  $p''(0) = 3$ ,  $p(1) = -1$ ,  $p'(1) = 3$ ,  $p(2) = 4$ .

$x_i$	$f[.]$	$f[.,.]$	$f[.,.,.]$	$f[.,.,.,.]$	$f[.,.,.,.,.]$	$f[.,.,.,.,.,.]$
0	1					
		2				
0	1		$\frac{3}{2}$			
		2		$-\frac{11}{2}$		
0	1		-4		$\frac{29}{2}$	
		-2		9		$-\frac{79}{8}$
1	-1		5		$-\frac{21}{4}$	
		3		$-\frac{3}{2}$		
1	-1		2			
		5				
2	4					

Vrednosti v tabeli, ki so označene rdeče, smo izračunali direktno iz začetnih podatkih, saj so v teh deljenih diferencah vse točke enake. Ostale vrednosti smo izračunali po rekurzivni zvezi. Tako smo npr. vrednost  $f[0, 0, 0] = 3/2$  dobili iz formule  $f[0, 0, 0] = f''(0)/2$ , vrednost  $f[0, 0, 0, 1] = 29/2$  pa smo dobili iz  $f[0, 0, 0, 1] = (f[0, 0, 1] - f[0, 0, 0]) / (1 - 0)$ .

Koeficiente interpolacijskega polinoma preberemo iz zgornje diagonale in dobimo

$$p(x) = 1 + 2x + \frac{3}{2}x^2 - \frac{11}{2}x^3 + \frac{29}{2}x^3(x - 1) - \frac{79}{8}x^3(x - 1)^2.$$

Namesto tega bi lahko koeficiente vzeli tudi iz spodnje diagonale. V tem primeru interpolacijske točke nastopajo v obratnem vrstnem redu. Interpolacijski polinom je seveda enak kot prej, saj je enoličen, le zapisan je v drugačni bazi. Če vzamemo spodnjo diagonalno, dobimo

$$p(x) = 4 + 5(x - 2) + 2(x - 2)(x - 1) - \frac{3}{2}(x - 2)(x - 1)^2 - \frac{21}{4}(x - 2)(x - 1)x - \frac{79}{8}(x - 2)(x - 1)^2x^2. \quad \square$$

**Izrek 5.5 (Hermite–Genochijeva formula)** Za  $k$ -krat zvezno odvedljivo funkcijo  $f$  velja

$$f[x_0, \dots, x_k] = \int_0^1 dt_1 \int_0^{t_1} dt_2 \int \dots \int_0^{t_{k-1}} f^{(k)} \left( t_k(x_k - x_{k-1}) + \dots + t_1(x_1 - x_0) + x_0 \right) dt_k.$$



*Dokaz.* Uporabimo indukcijo. Če je  $x_0 \neq x_1$ , dobimo

$$\begin{aligned} f[x_0, x_1] &= \int_0^1 f'(t_1(x_1 - x_0) + x_0) dt_1 = \frac{1}{x_1 - x_0} f(t_1(x_1 - x_0) + x_0) \Big|_{t_1=0}^{t_1=1} \\ &= \frac{f(x_1) - f(x_0)}{x_1 - x_0}, \end{aligned}$$

sicer pa

$$f[x_0, x_0] = \int_0^1 f'(x_0) dt_1 = f'(x_0).$$

Sedaj predpostavimo, da izrek velja za  $1, \dots, k-1$ . Če definiramo  $\xi = t_k(x_k - x_{k-1}) + \dots + t_1(x_1 - x_0) + x_0$ , lahko zapišemo integral kot

$$f[x_0, x_1, \dots, x_k] = \int_0^1 dt_1 \int_0^{t_1} dt_2 \int \dots \int_0^{t_{k-1}} f^{(k)}(\xi) dt_k.$$

Sedaj s substitucijo spremenljivk dobimo:

$$d\xi = (x_k - x_{k-1}) dt_k,$$

$$\begin{aligned} t_k = 0 &\implies \xi_0 = t_{k-1}(x_{k-1} - x_{k-2}) + t_{k-2}(x_{k-2} - x_{k-3}) + \dots + t_1(x_1 - x_0) + x_0 \\ t_k = t_{k-1} &\implies \xi_1 = t_{k-1}(x_k - x_{k-2}) + t_{k-2}(x_{k-2} - x_{k-3}) + \dots + t_1(x_1 - x_0) + x_0 \end{aligned}$$

in

$$\int_0^{t_{k-1}} f^{(k)}(\xi) dt_k = \frac{1}{x_k - x_{k-1}} \int_{\xi_0}^{\xi_1} f^{(k)}(\xi) d\xi = \frac{f^{(k-1)}(\xi_1) - f^{(k-1)}(\xi_0)}{x_k - x_{k-1}}.$$

Po indukcijski predpostavki velja

$$\int_0^1 dt_1 \int_0^{t_1} dt_2 \int \dots \int_0^{t_{k-2}} f^{(k-1)}(\xi_1) dt_{k-1} = f[x_0, x_1, \dots, x_{k-2}, x_k]$$

in

$$\int_0^1 dt_1 \int_0^{t_1} dt_2 \int \dots \int_0^{t_{k-2}} f^{(k-1)}(\xi_0) dt_{k-1} = f[x_0, x_1, \dots, x_{k-2}, x_{k-1}],$$

ker pa vemo

$$f[x_0, x_1, \dots, x_k] = \frac{f[x_0, x_1, \dots, x_{k-2}, x_k] - f[x_0, x_1, \dots, x_{k-2}, x_{k-1}]}{x_k - x_{k-1}},$$

je dokaz končan. ■

**Posledica 5.6** Za  $k$ -krat zvezno odvedljivo funkcijo  $f$  velja

$$f[x_0, \dots, x_k] = \frac{1}{k!} f^{(k)}(\xi),$$

kjer je

$$\min_{i=0, \dots, k} (x_i) \leq \xi \leq \max_{i=0, \dots, k} (x_i). \quad \blacksquare$$

**Izrek 5.7** Za  $(n + 1)$ -krat zvezno odvedljivo funkcijo  $f$  in interpolacijski polinom  $I_n$  na točkah  $x_0, \dots, x_n$  velja

$$f(x) = I_n(x) + f[x_0, \dots, x_n, x](x - x_0) \cdots (x - x_n). \quad (5.5)$$

*Dokaz.* Naj polinom  $q$  interpolira funkcijo  $f$  v točkah  $x_0, \dots, x_n, t$ . V Newtonovi obliki lahko  $q$  zapišemo kot

$$q(x) = I_n(x) + f[x_0, \dots, x_n, t](x - x_0) \cdots (x - x_n).$$

Očitno pri izbranem  $t$  potem velja  $q(t) = I_n(t) + f[x_0, \dots, x_n, t](t - x_0) \cdots (t - x_n)$ . Ker to velja za vsak  $t$ , dobimo formulo (5.5), ki jo je bilo potrebno dokazati. ■

**Posledica 5.8** Za  $(n + 1)$ -krat zvezno odvedljivo funkcijo  $f$  in interpolacijski polinom  $I_n$  na točkah  $x_0, \dots, x_n$  velja

$$f(x) - I_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x),$$

kjer je  $\min(x, x_0, \dots, x_n) \leq \xi \leq \max(x, x_0, \dots, x_n)$ . ■

Zgornja ocena se ujema z oceno, ki smo jo zapisali v izreku 5.2 za paroma različne točke. Prednost nove ocene je, da pride v poštev tudi, če vse točke niso medsebojno različne.

---

## 5.4 Končne diference

Če so vse točke *ekvidistantne*, kar pomeni  $x_i = x_0 + ih$ , kjer je  $h$  razmik med sosednjima točkama, lahko formule še poenostavimo. Naj bodo  $y_i = f(x_i)$  vrednosti, iz katerih sestavimo končne diference.

**Definicija 5.9** Končne diference so definirane rekurzivno kot

$$\Delta^m y_k = \begin{cases} y_k, & m = 0 \\ \Delta^{m-1} y_{k+1} - \Delta^{m-1} y_k, & m > 0 \end{cases}.$$

Iz vrednosti  $y_i$  sestavimo *diferenčno tabelo*:

$x_i$	$y_i$	$\Delta$	$\Delta^2$	$\Delta^3$
$x_0$	$y_0$			
$x_1$	$y_1$	$\Delta y_0$		
$x_2$	$y_2$	$\Delta y_1$	$\Delta^2 y_0$	
$x_3$	$y_3$	$\Delta y_2$	$\Delta^2 y_1$	$\Delta^3 y_0$

Za razliko od deljenih diferenc tu nimamo deljenja z razliko interpolacijskih točk, temveč moramo samo odšteti ustrezni diferenci iz levega stolpca. Tako za prve tri diference dobimo naslednje formule:

$$\Delta y_0 = y_1 - y_0,$$

$$\begin{aligned}\Delta^2 y_0 &= \Delta y_1 - \Delta y_0 = (y_2 - y_1) - (y_1 - y_0) = y_2 - 2y_1 + y_0, \\ \Delta^3 y_0 &= \Delta^2 y_1 - \Delta^2 y_0 = (y_3 - 2y_2 + y_1) - (y_2 - 2y_1 + y_0) = y_3 - 3y_2 + 3y_1 - y_0.\end{aligned}$$

Obratno lahko vrednosti izrazimo nazaj iz diferenc. Tako dobimo:

$$\begin{aligned}y_1 &= y_0 + \Delta y_0 = (I + \Delta)y_0 \\ y_2 &= y_1 + \Delta y_1 = (I + \Delta)y_1 = (I + \Delta)^2 y_0\end{aligned}$$

Seveda lahko podobne formule zapišemo tudi za višje diference oziroma točke. Splošni formuli sta zapisani v naslednji lemi, ki se jo da preprosto dokazati z uporabo indukcije.

**Lema 5.10** *Dane so vrednosti  $y_0, \dots, y_m$ . Potem veljata formuli:*

$$\begin{aligned}a) \quad \Delta^m y_0 &= \sum_{k=0}^m (-1)^k \binom{m}{k} y_{m-k}, \\ b) \quad y_m &= (I + \Delta)^m y_0 = \sum_{k=0}^m \binom{m}{k} \Delta^k y_0.\end{aligned}$$

Namesto na vrednostih  $y_0, y_1, \dots$  lahko definiramo končne diference z razmikom  $h$  tudi za funkcijo  $f$  v točki  $x$ . Tako za funkcijo  $f$  definiramo

$$\Delta_h^m f(x) = \begin{cases} f(x), & m = 0 \\ \Delta_h^{m-1} f(x+h) - \Delta_h^{m-1} f(x), & m > 0 \end{cases}.$$

Če definiramo  $y_k = f(x_0 + kh)$  za  $k = 0, 1, \dots$ , potem je očitno  $\Delta_h^k f(x) = \Delta^k y_0$ .

**Lema 5.11** *Če je  $p$  polinom stopnje  $n$  z vodilnim koeficientom  $a_0$ , potem velja  $\Delta_h^n p(x) = n!h^n a_0$  in  $\Delta_h^m p(x) = 0$  za  $m > n$ .*

*Dokaz.* Iz zveze

$$\begin{aligned}\Delta_h p(x) &= p(x+h) - p(x) \\ &= (a_0(x+h)^n + a_1(x+h)^{n-1} + \dots + a_n) - (a_0x^n + a_1x^{n-1} + \dots + a_n) \\ &= nha_0x^{n-1} + \dots\end{aligned}$$

vidimo, da se pri vsaki diferenci stopnja polinoma zmanjša za ena, vodilni koeficient pa pomnoži s stopnjo  $n$  in s  $h$ . ■

Končne diference lahko uporabimo za zapis interpolacijskega polinoma. Naj bo  $f$  dana funkcija, ki jo želimo interpolirati v točkah  $x_0, \dots, x_n$ , kjer je  $x_i = x_0 + ih$  za  $i = 0, \dots, n$ . Če je  $p$  interpolacijski polinom za  $f$  stopnje manjše ali enake  $n$  v točkah  $x_0, \dots, x_n$ , potem velja  $p(x_i) = f(x_i)$  za  $i = 0, \dots, n$ . Od tod očitno sledi, da za končne diference, ki so definirane le z vrednostmi funkcije, velja  $\Delta_h^k f(x_0) = \Delta_h^k p(x_0)$  za  $k = 0, \dots, n$ . V točki  $x_0 + sh$  zapišemo  $p(x_0 + sh) = (I + \Delta_h)^s p(x_0)$  in  $(I + \Delta_h)^s$  simbolno razvijemo po binomski formuli. Tako dobimo

$$p(x_0 + sh) = \sum_{k=0}^{\infty} \binom{s}{k} \Delta_h^k f(x_0).$$

Ker so po lemi 5.11 vse diference polinoma od  $(n+1)$ -ve dalje enake 0, odpadejo vsi členi vsote kjer je  $k > n$ . Da tako res dobimo interpolacijski polinom, nam dokazuje naslednja lema.

**Lema 5.12** Naj bo  $s = (x - x_0)/h$ . Potem je

$$I_n(x) = I_n(x_0 + sh) = \sum_{k=0}^n \binom{s}{k} \Delta_h^k f(x_0) \quad (5.6)$$

interpolacijski polinom za funkcijo  $f$  stopnje kvečjemu  $n$  na točkah  $x_i = x_0 + ih$  za  $i = 0, \dots, n$ .

*Dokaz.* Očitno je  $I_n$  polinom stopnje kvečjemu  $n$ . Sedaj moramo samo še preveriti, da res velja  $f(x_i) = p(x_i)$  za  $i = 0, \dots, n$ , kar pa očitno sledi iz leme 5.10. ■

Formulo (5.6) imenujemo *prema Newtonova interpolacijska formula*.

**Zgled 5.3** Sestavi diferenčno tabelo za naslednje podatke do vključno tretjih diferenc in izračunaj  $I_3(0.02)$ .

$x$	0.0	0.1	0.2	0.3	0.4
$y$	0.00	0.11	0.28	0.57	1.04

$x_i$	$y_i$	$\Delta$	$\Delta^2$	$\Delta^3$
0.0	0.00			
		0.11		
0.1	0.11		0.06	
		0.17		0.06
0.2	0.28		0.12	
		0.29		0.06
0.3	0.57		0.18	
		0.47		
0.4	1.04			

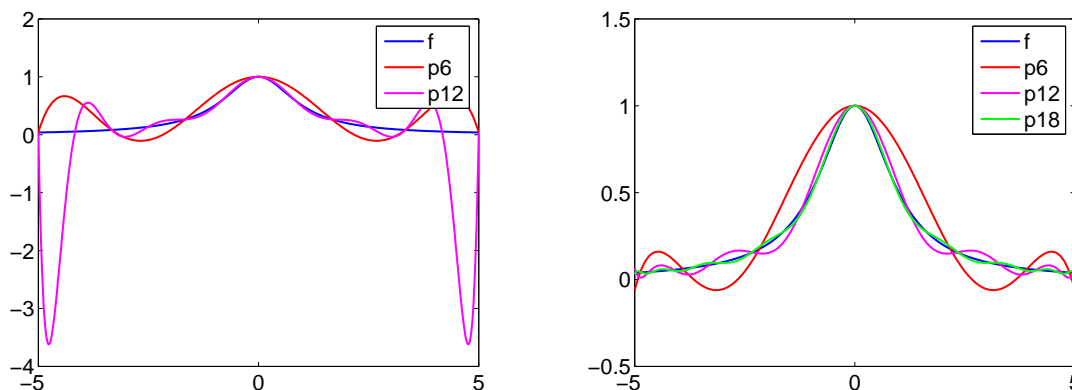
$$h = 0.1, \quad s = \frac{x - x_0}{h} = 0.2$$

$$I_3(0.02) = 0.00 + 0.2 \cdot 0.11 + \frac{0.2 \cdot (-0.8)}{2} 0.06 + \frac{0.2 \cdot (-0.8) \cdot (-1.8)}{6} 0.06 = 0.02008. \quad \square$$

## 5.5 Interpolacija s kosoma polinomskimi funkcijami

Če želimo, da bi se na danem intervalu interpolacijski polinom čim bolj prilegal funkciji, je najbrž prva ideja ta, da povečamo stopnjo interpolacijskega polinoma. To naredimo tako, da povečamo število interpolacijskih točk. Kot kaže naslednji zgled, ki je znan pod imenom *Rungejev protiprimer*, ni nujno, da večanje stopnje interpolacijskega polinoma res izboljša aproksimacijo funkcije s polinomom.

**Zgled 5.4** Denimo, da interpoliramo funkcijo  $f(x) = 1/(1 + x^2)$  na intervalu  $[-5, 5]$ . Če uporabimo ekvidistantne točke, potem se izkaže, da z večanjem stopnje interpolacijski polinom vedno slabše aproksimira  $f$ , kot kaže slika 5.1 (leva slika).



Slika 5.1: Rungejev protiprimer. Če funkcijo  $f(x) = 1/(1+x^2)$  na  $[-5, 5]$  interpoliramo na ekvidistantnih točkah (levo), razlika med funkcijo in interpolacijskim polinomom z večanjem  $n$  narašča. Če pa interpoliramo na Čebiševskih točkah (desno), potem z večanjem števila točk interpolacijski polinom vedno bolje aproksimira  $f$ .

Težava se v tem primeru pojavi zaradi ekvidistantnih točk. Če za interpolacijske točke vzamemo Čebiševske točke, ki so v tem primeru definirane kot  $x_j = 5 \cos(j\pi/n)$  za  $j = 0, \dots, n$ , potem pa z večanjem stopnje interpolacijski polinom vedno bolje aproksimira funkcijo  $f$  (desna slika).  $\square$

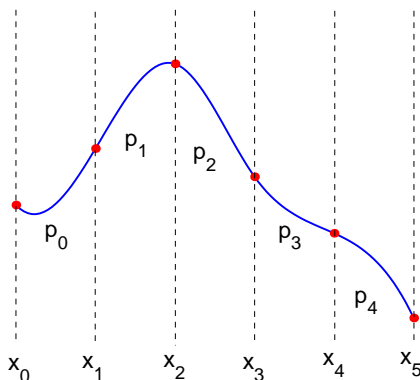
Tu je potrebno še omeniti, da za razliko med interpolacijo z ekvidistantnimi in Čebiševskimi točkami v zadnjem zgledu ni odgovorno numerično računanje in zaokrožitvene napake, saj pride do razlik pri eksaktnem računanju.

Ponavadi se Čebiševske točke sicer res obnašajo bolje kot ekvidistantne točke, a se da pokazati, da za poljubno zaporedje predpisanih interpolacijskih točk obstaja funkcija, pri kateri se zgodi, da z večanjem stopnje interpolacijskega polinoma prihaja do vse večjih razlik med funkcijo in interpolacijskim polinomom.

To ni v nasprotju z Weierstrassovim izrekom iz prejšnjega poglavja. Ta pravi, da je možno zvezno funkcijo poljubno dobro aproksimirati s polinomom dovolj visoke stopnje, a pri tem ni predpisano, v katerih točkah se mora polinom ujemanjati s funkcijo. Če za te točke predpišemo, da morajo biti npr. ekvidistantne, potem izreka seveda ne moremo več uporabiti.

Rungejev protiprimer kaže, da večanje stopnje interpolacijskega polinoma ni dober način za boljšo aproksimacijo funkcije. Rešitev je, da sicer povečamo število interpolacijskih točk, a izberemo drugačno vrsto interpolacijske funkcije. Namesto enega polinoma visoke stopnje interpolacijsko funkcijo zlepimo iz polinomov nizkih stopenj. Tako dobimo kosoma polinomsko funkcijo oziroma *zlepko*. Ker imamo še vedno opravka s polinomi, sta konstrukcija zlepka in izračun vrednosti še vedno preprosta. Primer zlepka je prikazan na sliki 5.2

Denimo, da iščemo interpolacijsko funkcijo, ki bo v točkah  $x_0, \dots, x_{n+1}$  po vrsti imela vrednosti  $y_0, \dots, y_{n+1}$ . Pri polinomskem zlepku interval  $[x_0, x_{n+1}]$  razdelimo na podintervale  $[x_i, x_{i+1}]$  za  $i = 0, \dots, n$ . Na vsakem intervalu  $[x_i, x_{i+1}]$  vzamemo polinom  $p_i$  nizke stopnje, za katerega velja  $p_i(x_i) = y_i$  in  $p_i(x_{i+1}) = y_{i+1}$ . Želimo, da bo sestavljena krivulja čim bolj gladka in da bo čim bolj opisovala obliko začetnih podatkov. Ker gre krivulja čez točke  $(x_i, y_i)$  za  $i = 1, \dots, n$



Slika 5.2: Zgled polinomskega zlepka. Na vsakem intervalu  $[x_i, x_{i+1}]$  imamo definiran polinom  $p_i$  za  $i = 0, \dots, 4$ , vsi pa so povezani v zvezno funkcijo na  $[x_0, x_5]$ .

je vsaj zvezna, saj v notranjih točkah velja

$$p_i(x_{i+1}) = p_{i+1}(x_{i+1}) = y_{i+1}.$$

Osnovna varianta je *kosoma linearna interpolacija*, kjer na vsakem intervalu  $[x_i, x_{i+1}]$  podatke interpoliramo z linearno funkcijo

$$p_i(x) = y_i + \frac{y_{i+1} - y_i}{x_{i+1} - x_i}(x - x_i).$$

Kvadratni zleпки niso najbolj uporabni, podrobneje pa bomo opisali *kubične zlepkе*, kjer je  $p_i$  kubičen polinom na  $[x_i, x_{i+1}]$ . Za kubični interpolacijski zlepek ponavadi zahtevamo še, da je vsaj zvezno odvedljiv.

Posamezni kubični polinom  $p_i$  je natanko določen z vrednostmi in odvodi v krajiščih:

$$\begin{aligned} p_i(x_i) &= y_i, & p_i(x_{i+1}) &= y_{i+1} \\ p_i'(x_i) &= d_i, & p_i'(x_{i+1}) &= d_{i+1}. \end{aligned}$$

Če označimo  $h_i = x_{i+1} - x_i$ ,  $\delta_i = (y_{i+1} - y_i)/h_i$  in uporabimo deljene diference, dobimo

$$\begin{array}{l|l} x_i & y_i \\ & d_i \\ x_i & y_i & \frac{\delta_i - d_i}{h_i} \\ & \delta_i & \frac{d_i + d_{i+1} - 2\delta_i}{h_i^2} \\ x_{i+1} & y_{i+1} & \frac{d_{i+1} - \delta_i}{h_i} \\ & d_{i+1} \\ x_{i+1} & y_{i+1} \end{array}$$

in lahko zapišemo polinom  $p_i$  v obliki

$$p_i(x) = y_i + d_i(x - x_i) + \frac{\delta_i - d_i}{h_i}(x - x_i)^2 + \frac{d_i + d_{i+1} - 2\delta_i}{h_i^2}(x - x_i)^2(x - x_{i+1}).$$

Kubični zlepek je tako določen s parametri  $d_0, \dots, d_{n+1}$ , ki so zaenkrat še prosti, določiti pa jih moramo tako, da bo imel zlepek željene lastnosti. Obstaja veliko možnih izbir, najpogostejše pa so:

- Pri *Hermiteovem kubičnem zlepku* izberemo kar  $d_i = y'_i$ , kjer je  $y'_i$  odvod funkcije, ki jo interpoliramo, v točki  $x_i$  za  $i = 0, \dots, n + 1$ .
- Pri *dvakrat zvezno odvedljivem kubičnem zlepku* parametre  $d_0, \dots, d_{n+1}$  določimo tako, da je zlepek dvakrat zvezno odvedljiv. Velja

$$p''_i(x_{i+1}) = \frac{1}{h_i}(2d_i + 4d_{i+1} - 6\delta_i)$$

in

$$p''_{i+1}(x_{i+1}) = \frac{1}{h_{i+1}}(-4d_{i+1} - 2d_{i+2} + 6\delta_{i+1}).$$

Iz enačb

$$p''_i(x_{i+1}) = p''_{i+1}(x_{i+1}), \quad i = 0, \dots, n - 1$$

dobimo sistem  $n$  linearnih enačb za  $n + 2$  parametrov:

$$h_{i+1}d_i + 2(h_i + h_{i+1})d_{i+1} + h_id_{i+2} = 3h_{i+1}\delta_i + 3h_i\delta_{i+1}.$$

Manjkajoči dve enačbi dobimo:

- *kompletni zlepek*: če poznamo vrednosti odvodov funkcije, ki jo interpoliramo, v začetni in končni točki, potem lahko vzamemo  $d_0 = y'_0$  in  $d_{n+1} = y'_{n+1}$ .
- *naravni zlepek*: s pogojeva  $p''_0(x_0) = 0$  in  $p''_n(x_{n+1}) = 0$  dosežemo, da ima zlepek v začetni in končni točki prevoj.
- *zlepek brez vozlov*: zahtevamo, da je zlepek na  $[x_0, x_2]$  in  $[x_{n-1}, x_n]$  kubični polinom oziroma, da velja  $p'''_0(x_1) = p'''_1(x_1)$  in  $p'''_{n-1}(x_n) = p'''_n(x_n)$ . Tako se v bistvu znebimo vozlov  $x_1$  in  $x_n$ .

V vseh treh primerih moramo na koncu rešiti linearni sistem, da dobimo vrednosti parametrov  $d_0, \dots, d_n$ . Ker je sistem tridiagonalen, lahko to naredimo v časovni zahtevnosti  $\mathcal{O}(n)$ .

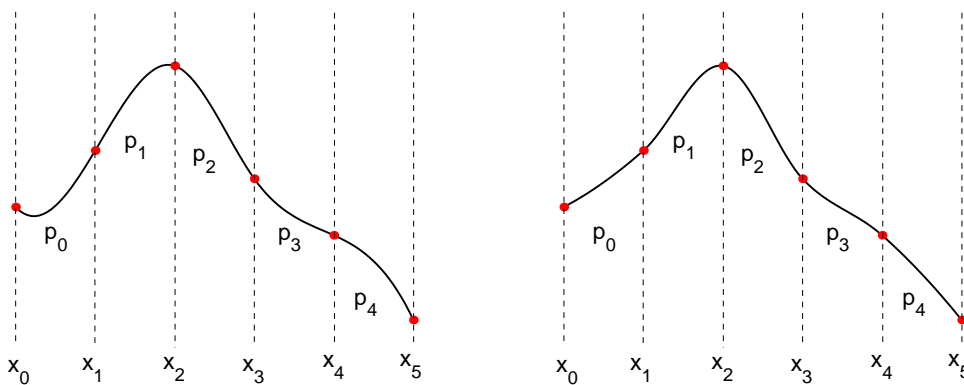
- Pri *kubičnem zlepku, ki ohranja odliko*, naklone  $d_i$  določimo tako, da za  $1/d_i$  vzamemo povprečje recipročnih vrednosti smernih koeficientov premic skozi  $(x_{i-1}, x_i)$  in  $(x_i, x_{i+1})$ . Če sta smerna koeficienta  $\delta_i$  in  $\delta_{i+1}$  enako predznačena, potem pri predpostavki  $h_i = h_{i+1}$  dobimo

$$\frac{1}{d_i} = \frac{1}{2} \left( \frac{1}{\delta_i} + \frac{1}{\delta_{i+1}} \right).$$

Če sta širini intervalov  $[x_{i-1}, x_i]$  in  $[x_i, x_{i+1}]$  različni, se formula za  $d_i$  spremeni.

Če sta smerna koeficienta  $\delta_i$  in  $\delta_{i+1}$  nasprotno predznačena, vzamemo  $d_i = 0$ .

Dobljeni kubični zlepek je enkrat zvezno odvedljiv, preskok v drugem odvodu je lepo razviden iz slike 5.3. Na levi strani je dvakrat zvezno odvedljiv zlepek, na desni strani pa zlepek, ki ohranja obliko.



Slika 5.3: Primerjava kubičnega zleпка, ki ohranja obliko (desno), z dvakrat zvezno odvedljivim kubičnim zleptom (levo).

Pokažimo, da v primeru dvakrat zvezno odvedljivega kubičnega zleпка res velja, da z večanjem števila interpolacijskih točk, zlepek vedno bolj aproksimira funkcijo, ki jo interpolira. Pri večanju števila točk moramo seveda paziti na to, da se manjša maksimalna širina podintervalov

$$h := \max_{i=0, \dots, n} h_i.$$

**Izrek 5.13** Naj bo  $f$  štirikrat zvezno odvedljiva funkcija na intervalu  $[a, b]$ . Če interval razdelimo s točkami  $a = x_0 < x_1 < \dots < x_n < x_{n+1} = b$ , potem za Hermiteov kubični zlepek  $p$  za vsak  $x \in [a, b]$  velja ocena

$$|f(x) - p(x)| \leq \frac{h^4}{384} M_4, \tag{5.7}$$

kjer je  $h = \max_{i=0, \dots, n} (x_{i+1} - x_i)$  in  $M_4 = \max_{x \in [a, b]} |f^{(4)}(x)|$ .

*Dokaz.* Na podintervalu  $[x_i, x_{i+1}]$  funkcijo  $f$  interpoliramo s kubičnim polinomom  $p_i$ , za katerega velja  $p_i(x_i) = f(x_i)$ ,  $p_i'(x_i) = f'(x_i)$ ,  $p_i(x_{i+1}) = f(x_{i+1})$  in  $p_i'(x_{i+1}) = f'(x_{i+1})$ . Za razliko vemo, da velja

$$f(x) - p_i(x) = f[x_i, x_i, x_{i+1}, x_{i+1}, x](x - x_i)^2(x - x_{i+1})^2.$$

Če je  $h_i = x_{i+1} - x_i$ , lahko ocenimo

$$|f(x) - p_i(x)| \leq \frac{M_4}{4!} \cdot \frac{h_i^4}{16}.$$

Od tod sledi ocena (5.7). ■

Podobne ocene se da izpeljati tudi za ostale zleпke. Tako npr. za kompletni kubični zlepek  $p$  in štirikrat zvezno odvedljivo funkcijo  $f$  pri predpostavkah izreka 5.13 veljajo ocene

$$\begin{aligned} |f(x) - p(x)| &\leq \frac{5h^4}{384} M_4, \\ |f'(x) - p'(x)| &\leq \frac{h^3}{24} M_4, \\ |f''(x) - p''(x)| &\leq \frac{3h^2}{8} M_4. \end{aligned}$$



Ne samo, da kompletni kubični zlepek dobro aproksimira funkcijo  $f$ , iz zadnjih dveh ocen sledi, da tudi prvi in drugi odvod zleпка dobro aproksimirata prvi in drugi odvod funkcije  $f$ .

## 5.6 Beziérove krivulje

Pri interpolaciji krivulj v ravnini so pomembno orodje *Beziérove krivulje*. Vsaka Beziérova krivulja je določena s kontrolnimi točkami  $p_i = (x_i, y_i)$ ,  $i = 0, \dots, n$ .

Formula za točke na krivulji je

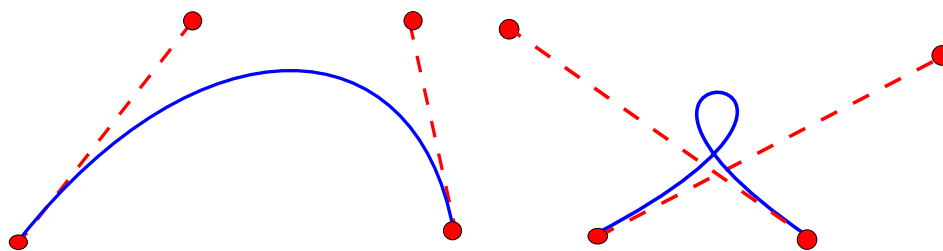
$$P(t) = \sum_{k=0}^n \binom{n}{k} t^k (1-t)^{n-k} p_k,$$

kjer je  $t \in [0, 1]$ . V formuli nastopajo *Bernsteinovi polinomi*

$$B_{n,k}(t) = \binom{n}{k} t^k (1-t)^{n-k}$$

za  $k = 0, \dots, n$ , za katere velja  $\sum_{k=0}^n B_{n,k}(t) = 1$  za vsak  $t \in [0, 1]$ .

Najpogosteje se uporabljajo kubične Beziérove krivulje. Dva zglada sta predstavljena na sliki 5.4.



Slika 5.4: Kubični Beziérovi krivulji.

Nekaj lastnosti Beziérovih krivulj je:

- $P(0) = p_0$  in  $P(1) = p_n$ .
- Krivulja leži znotraj konveksne ogrinjače točk  $p_0, \dots, p_n$ .
- Naklon pri  $t = 0$  je enak naklonu premice skozi  $p_0$  in  $p_1$ , podobno se naklon pri  $t = 1$  ujema z naklonom premice skozi  $p_{n-1}$  in  $p_n$ .

Namesto Beziérovih krivulj na veliko točkah sestavljamo kubične Beziérove krivulje v zlepek. Pogoji za geometrijsko zveznost odvoda je, da zadnje dve točki prve krivulje in prvi dve točki druge krivulje ležijo na isti premici.

Beziérovi kubični zlepek se npr. uporabljajo za zapis računalniških pisav v vektorski obliki. Zapis ne porabi veliko prostora in omogoča, da se lahko velikost pisave poljubno poveča.

Omenimo še nekaj pojmov, ki se pojavijo, kadar pri interpolaciji ravninskih krivulj uporabljamo zlepek:

- *Geometrijska zveznost*: krivulji se dotikata, kar pomeni, da se zadnja točka prve krivulje ujema z začetno točko druge krivulje.
- *Geometrijska zveznost odvodov* ( $G^1$ ): tangenti prve in druge krivulje imata v skupni točki enako smer.
- *Parametrična zveznost odvodov* ( $C^1$ ): prva odvoda prve in druge krivulje v skupni točki se ujemata.

## 5.7 Interpolacija v Matlabu

V Matlabu imamo na voljo naslednje funkcije za delo s polinomi in za konstrukcijo interpolacijskih polinomov in zlepkov:

- `polyval(p, x)`: vrednost polinoma, podanega s koeficienti v  $p$  v točki  $x$
- `roots(p)`: ničle polinoma, podanega s koeficienti v  $p$
- `poly(r)`: vrne koeficiente polinoma z danimi ničlami v  $r$
- `polyfit(x, y, n)`: vrne koeficiente polinoma  $p$  stopnje  $m$ , ki po metodi najmanjših kvadratov najbolje aproksimira točke  $(x_i, y_i)$ , kar pomeni da je vsota

$$\sum_{i=0}^n (p(x_i) - y_i)^2$$

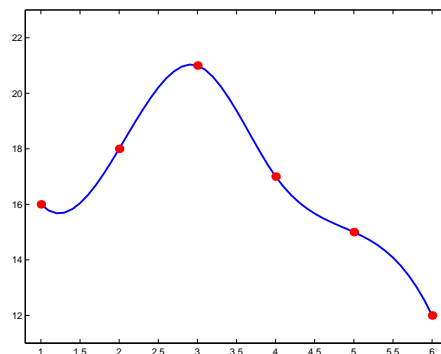
minimalna. Če izberemo  $m \geq n$ , potem dobimo interpolacijski polinom.

- `yi=interp1(x, y, xi)`: vrne vrednost kosoma polinomske interpolacijske funkcije v točki  $x_i$ . Kot četrti argument lahko podamo vrsto interpolacijske funkcije, na voljo so
  - `'nearest'`: kosoma konstantna interpolacija,
  - `'linear'`: kosoma linearna interpolacija,
  - `'spline'`: dvakrat zvezno odvedljiv kubični zlepek,
  - `'cubic'`: kubični zlepek, ki ohranja obliko.

Za interpolacijo v dveh dimenzijah sta na voljo ukaza `interp2` in `griddata`.

**Zgled 5.5** Naslednji ukazi v Matlabu generirajo dvakrat zvezno odvedljivo kubični zlepek in narišejo njegov graf, ki je prikazan na desni strani.

```
x=[1 2 3 4 5 6];
y=[16 18 21 17 15 12];
t=1:0.1:6;
y3=interp1(x,y,t,'spline');
plot(t,y3,'b-', 'LineWidth',2);
hold on
plot(x,y,'r.', 'MarkerSize',27);
axis([0.8 6.2 11 23])
hold off
```



## **Dodatna literatura**

Za interpolacijo je na voljo obsežna literatura, saj je obravnavana skoraj v vseh učbenikih numerične analize. V slovenščini imate na voljo knjigo [17], osnove interpolacije pa najdete tudi v knjigah [3] in [27]. Od tuje literature omenimo npr. knjigi [6] in [16].

## Poglavje 6

# Numerično odvajanje

### 6.1 Uvod

Iščemo odvod funkcije, ki je podana s tabelo vrednosti v točkah  $x_0, \dots, x_n$ . Ideja je, da za približek vzamemo odvod interpolacijskega polinoma. Če je funkcija dovoljkrat zvezno odvedljiva, potem vemo, da je

$$f(x) = I_n(x) + \frac{\omega(x)}{(n+1)!} f^{(n+1)}(\xi),$$

kjer je  $\omega(x) = (x - x_0) \cdots (x - x_n)$ . Z odvajanjem zgornje zveze dobimo

$$f'(x) = I_n'(x) + \underbrace{\frac{\omega'(x)}{(n+1)!} f^{(n+1)}(\xi) + \frac{\omega(x)}{(n+1)!} \cdot \frac{df^{(n+1)}(\xi)}{dx}}_{\text{napaka}}.$$

Izraz za napako ni najlepši, saj ne poznamo odvisnosti  $\xi$  od  $x$ . To nam preprečuje, da bi lahko ocenili maksimalno napako. Če pa računamo odvod v eni izmed točk  $x_0, \dots, x_n$ , zadnji člen odpade in dobimo

$$f'(x_k) = I_n'(x_k) + \frac{\omega'(x_k)}{(n+1)!} f^{(n+1)}(\xi). \quad (6.1)$$

V zgornji formuli je  $I_n'(x_k)$  odvod interpolacijskega polinoma v točki  $x_k$ . Formulo za  $I_n'(x_k)$  lahko izpeljemo preko Lagrangeevih koeficientov. Iz  $I_n(x) = \sum_{i=0}^n f(x_i) L_{n,i}(x)$  sledi  $I_n'(x_k) = \sum_{i=0}^n f(x_i) L'_{n,i}(x_k)$ . Z odvajanjem Lagrangeevega koeficienta

$$L_{n,i}(x) = \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}$$

dobimo

$$\text{a) } i \neq k: \quad L'_{n,i}(x_k) = \frac{\omega'(x_k)}{(x_k - x_i)\omega'(x_i)},$$

$$\text{b) } i = k: \quad L'_{n,k}(x_k) = \sum_{\substack{j=0 \\ j \neq k}}^n \frac{1}{x_k - x_j}.$$

S pomočjo teh formul lahko zapišemo

$$I'_n(x_k) = \omega'(x_k) \sum_{\substack{j=0 \\ j \neq k}}^n \frac{f(x_j)}{(x_k - x_j)\omega'(x_j)} + f(x_k) \sum_{\substack{j=0 \\ j \neq k}}^n \frac{1}{x_k - x_j}.$$

V primeru, ko so točke ekvidistantne in velja  $x_i = x_0 + ih$  za  $i = 0, 1, \dots, n$ , se formula (6.1) še poenostavi. Z malce računanja lahko izpeljemo formulo

$$f'(x_k) = \frac{1}{h} \left( \frac{(-1)^k}{\binom{n}{k}} \sum_{\substack{j=0 \\ j \neq k}}^n \frac{(-1)^j \binom{n}{j} f(x_j)}{k-j} + f(x_k) \sum_{\substack{j=0 \\ j \neq k}}^n \frac{1}{k-j} \right) + \frac{(-1)^{n-k} h^n}{(n+1) \binom{n}{k}} f^{(n+1)}(\xi).$$

Nekaj prvih formul, ki jih dobimo z vstavljanjem  $n$  in  $k$ , je:

- $n = 1$ :

$$\begin{aligned} f'(x_0) &= \frac{1}{h}(f(x_1) - f(x_0)) - \frac{1}{2}hf''(\xi_0) \\ f'(x_1) &= \frac{1}{h}(f(x_1) - f(x_0)) + \frac{1}{2}hf''(\xi_1) \end{aligned}$$

- $n = 2$ :

$$\begin{aligned} f'(x_0) &= \frac{1}{2h}(-3f(x_0) + 4f(x_1) - f(x_2)) + \frac{1}{3}h^2f'''(\xi_0) \\ f'(x_1) &= \frac{1}{2h}(-f(x_0) + f(x_2)) - \frac{1}{6}h^2f'''(\xi_1) \quad (\text{simetrična diferenca}) \\ f'(x_2) &= \frac{1}{2h}(f(x_0) - 4f(x_1) + 3f(x_2)) + \frac{1}{3}h^2f'''(\xi_2) \end{aligned}$$

Če primerjamo formuli pri  $n = 1$  in simetrično diferenco, potem v vseh treh formulah vrednost odvoda aproksimiramo z vrednostnima funkcije v dveh točkah. Pri simetrični diferenci zaradi simetrije pridobimo red  $h^2$  namesto  $h$ , višji red pa imamo v bistvu zaradi tega, ker smo odvajali interpolacijski polinom na treh točkah namesto na dveh.

## 6.2 Drugi načini izpeljave

Formule za numerično odvajanje lahko izpeljujemo tudi iz razvoja v Taylorjevo vrsto. Naj bodo točke ekvidistantne in  $y_i = f(x_i)$ . Iz razvoj

$$\begin{aligned} y_0 &= y_1 - hy'_1 + \frac{1}{2}h^2y''_1 - \frac{1}{6}h^3y'''_1 + \frac{1}{24}h^4f^{(4)}(\xi_0) \\ y_1 &= y_1 \\ y_2 &= y_1 + hy'_1 + \frac{1}{2}h^2y''_1 + \frac{1}{6}h^3y'''_1 + \frac{1}{24}h^4f^{(4)}(\xi_2) \end{aligned}$$

s seštevanjem dobimo formulo  $f''(x_1) = \frac{1}{h^2}(y_0 - 2y_1 + y_2) - \frac{1}{24}h^2(f^{(4)}(\xi_0) + f^{(4)}(\xi_2))$ . Namesto  $f^{(4)}(\xi_0) + f^{(4)}(\xi_2)$  lahko pišemo  $2f^{(4)}(\xi)$  in dobimo končno formulo

$$f''(x_1) = \frac{1}{h^2}(y_0 - 2y_1 + y_2) - \frac{1}{12}h^2f^{(4)}(\xi). \quad (6.2)$$

Formula (6.2) je simetrična diferenca za drugi odvod. Skupaj s simetrično diferenco za prvi odvod se npr. uporabljata za numerično reševanje robnih problemov preko diferenčne metode. Formuli aproksimirata drugi oz. prvi odvod z natančnostjo  $\mathcal{O}(h^2)$ .

Alternativni način izpeljave je *metoda nedoločenih koeficientov*. Pri tej metodi nastavimo sistem

$$f'(x_k) = \sum_{j=0}^n \alpha_j f(x_j) + R(f)$$

in določimo koeficiente  $\alpha_0, \dots, \alpha_n$  tako, da je formula točna za polinome čim višje stopnje. Napako dobimo iz polinoma najnižje stopnje, za katerega formula ni točna.

**Zgled 6.1** Izpeljimo formulo  $f''(x_1) = af(x_0) + bf(x_1) + cf(x_2) + R(f)$ .

Za bazo izberemo  $1, x - x_1, (x - x_1)^2, \dots$ , saj tako dobimo lepši sistem:

$$\left. \begin{array}{l} 1 : 0 = a + b + c \\ x - x_1 : 0 = -ha + hc \\ (x - x_1)^2 : 2 = h^2a + h^2c \end{array} \right\} \Rightarrow a = \frac{1}{h^2}, b = \frac{-2}{h^2}, c = \frac{1}{h^2}.$$

Napaka ima obliko  $R(f) = Ch^p f^{(r)}(\xi)$ , kjer je  $C$  konstanta,  $p$  in  $r$  pa sta stopnji, ki ju moramo določiti. Odvod  $r$  ustreza najnižji stopnji polinoma, za katerega formula ni točna.

Napako dobimo tako, da po vrsti v formulo vstavljamo  $f(x) = x^r$  in poiščemo najnižjo stopnjo, za katero formula ni točna. Ker imamo v formuli tri točke, je formula točna za polinome stopnje 2 ali manj, zato najprej vstavimo  $r = 3$ , ker pa se izkaže, da je formula točna, nadaljujemo z  $r = 4$ .

$$\left. \begin{array}{l} (x - x_1)^3 : 0 = -h^3a + h^3c \\ (x - x_1)^4 : 0 \neq h^4a + h^4c \end{array} \right\} \Rightarrow r = 4, p = 2.$$

Ker za  $f(x) = (x - x_1)^4$  velja  $f^{(4)}(\xi) = 4!$ , sledi, da je napaka enaka  $R(f) = -\frac{1}{12}h^2 f^{(4)}(\xi)$ .  $\square$

### 6.3 Celotna napaka

Naslednji zgled prikazuje težave, ki se pojavijo pri numeričnem računanju odvodov.

**Zgled 6.2** Če po formulah

$$\begin{aligned} f'(0) &= \frac{1}{2h}(f(h) - f(-h)) + \mathcal{O}(h^2) \\ f''(0) &= \frac{1}{h^2}(f(-h) - 2f(0) + f(h)) + \mathcal{O}(h^2) \end{aligned}$$

računamo  $f'(0)$  in  $f''(0)$  za  $f(x) = e^x$ , potem v enojni natančnosti dobimo:

$h$	$f'(0)$	$f''(0)$
0.4	1.0268809	1.0134048
0.04	1.0002673	1.0001659
0.004	1.0000094	1.0021030
0.0004	1.0000169	0.7450581
0.00004	1.0006130	$3.7252907 \cdot 10^1$
0.000004	1.0058284	$3.7252903 \cdot 10^3$

Napaka se z manjšanjem  $h$  nekaj časa manjša, potem pa začne naraščati, čeprav za obe formuli velja, da imata napako reda  $\mathcal{O}(h^2)$ .  $\square$

Do težav, predstavljenih v prejšnjem zgledu, pride zaradi neodstranljive napake. Oglejmo si situacijo na zgledu formule

$$f''(x_1) = \frac{1}{h^2}(y_0 - 2y_1 + y_2) - \frac{1}{12}h^2 f^{(4)}(\xi). \quad (6.3)$$

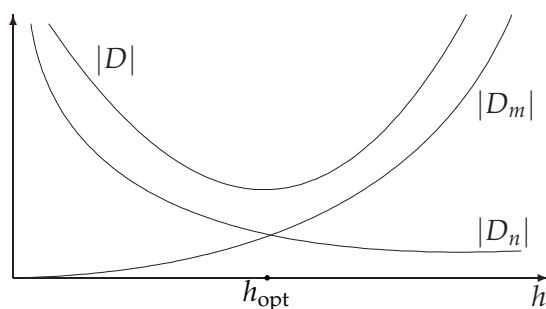
Vemo, da za napako metode  $D_m$  velja  $D_m = -\frac{1}{12}h^2 f^{(4)}(\xi)$ . Od tod lahko ocenimo

$$|D_m| \leq \frac{h^2}{12} \|f^{(4)}\|_\infty.$$

Poleg napake metode pa se pojavi tudi neodstranljiva napaka  $D_n$ , saj namesto s točnimi vrednostmi  $f(x_i)$  računamo s približki  $\tilde{y}_i$ , za katere predpostavimo  $|f(x_i) - \tilde{y}_i| \leq \epsilon$ . Če nič drugega, pride do neodstranljive napake že zaradi tega, ker namesto točnih vrednosti uporabljamo najbližja predstavljava števila. Ocenimo lahko

$$|D_n| \leq \frac{4\epsilon}{h^2}.$$

Pri samem računanju se pojavi še zaokrožitvena napaka  $D_z$ , ki pa smo jo v tokratni analizi zanemarili.



Slika 6.1: Ocene za neodstranljivo napako  $D_n$ , napako metode  $D_m$  in celotno napako  $D$  v odvisnosti od  $h$ .

Ocena za celotno napako (za formulo (6.3)) je

$$|D| \leq |D_m| + |D_n| \leq \frac{h^2}{12} \|f^{(4)}\|_\infty + \frac{4\epsilon}{h^2}.$$

Če poznamo oceni za  $\|f^{(4)}\|_\infty$  in  $\epsilon$ , lahko iz ocen za  $|D_m|$  in  $|D_n|$  določimo optimalni  $h$ , kjer bo ocena za skupno napako najmanjša.

Numerično odvajanje torej ni numerično dobro pogojen problem, saj se pri premajhnem  $h$  zgodi, da se z manjšanjem  $h$  napaka poveča.

## **Dodatna literatura**

Za numerično odvajanje imate v slovenščini na voljo knjige [17], [3] in [27], od tuje literature pa omenimo npr. knjigi [6] in [16].



## Poglavje 7

# Numerično integriranje

### 7.1 Kvadraturene formule

Radi bi izračunali integral

$$I(f) = \int_a^b f(x) dx.$$

Pri tem predpostavimo, da je funkcija  $f$  taka, da integral obstaja, npr., da je kosoma zvezna.

Ideja je, da namesto funkcije  $f$ , ki se je v splošnem primeru ne da analitično integrirati, integriramo kakšno funkcijo, ki se prilega  $f$  in za katero obstaja določeni integral. Prva izbira so interpolacijski polinomi.

Če za paroma različne interpolacijske točke izberemo  $x_0, \dots, x_n$ , potem lahko s pomočjo Lagrangeevih koeficientov zapišemo

$$f(x) = \sum_{i=0}^n f(x_i) L_{n,i}(x) + \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x).$$

Če to integriramo, dobimo

$$\int_a^b f(x) dx = \sum_{i=0}^n f(x_i) \int_a^b L_{n,i}(x) dx + \int_a^b \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x) dx.$$

Od tod sledi, da integral  $I(f)$  lahko aproksimiramo s kvadratureno formulo oblike

$$\int_a^b f(x) dx = \sum_{i=0}^n A_i f(x_i) + R(f), \quad (7.1)$$

kjer paroma različne točke  $x_0, \dots, x_n$  imenujemo *vozli*, koeficiente  $A_0, \dots, A_n$ , za katere velja

$$A_i = \int_a^b L_{n,i}(x) dx$$

za  $i = 0, \dots, n$ , imenujemo *uteži*,  $R(f)$  pa je napaka, ki je enaka

$$R(f) = \int_a^b \frac{f^{(n+1)}(\xi)}{(n+1)!} dx.$$

Formuli pravimo kvadratura zato, ker je integral  $I(f)$  enak ploščini lika, ki ga določa graf funkcije  $f$  na intervalu  $[a, b]$ .

Alternativno bi lahko izpeljali kvadraturno formulo (7.1) tudi tako, da bi najprej izbrali vozle  $x_0, \dots, x_n$ , potem pa določili uteži  $A_0, \dots, A_n$  tako, da je formula točna za polinome čim višje stopnje. Pri tem bi podobno kot pri numeričnem odvajanju lahko uporabili metodo nedoločnih koeficientov.

Za kvadraturno formulo (7.1) pravimo, da je reda  $k$ , če je točna za vse polinome stopnje manjše ali enake  $k$ , ni pa točna za polinome stopnje  $k + 1$ . V obeh zgornjih primerih dobimo kvadraturno formulo, ki je očitno reda vsaj  $n$ , saj je pravilna za polinome stopnje manjše ali enake  $n$ , ne glede na to, kako smo izbrali  $n + 1$  vozlov. Kot bomo videli kasneje, pa lahko s primerno izbiro vozlov dosežemo, da bo formula točna tudi za polinome višjih stopenj.

## 7.2 Newton–Cotesova pravila

Pri *Newton–Cotesovih pravilih* so vozli ekvidistantni:  $a = x_0, b = x_n, h = (b - a)/n, x_i = x_0 + ih$  za  $i = 0, \dots, n$ . Naj bo  $y_k = f(x_k)$ . Ločimo dva tipa Newton–Cotesovih pravil:

a) *zaprti tip*: upoštevamo tudi krajišča: 
$$\int_a^b f(x)dx = \sum_{k=0}^n A_k y_k + R_n(f);$$

b) *odprti tip*: brez krajišč: 
$$\int_a^b f(x)dx = \sum_{k=1}^{n-1} B_k y_k + R_n(f).$$

Nekaj osnovnih zaprtih Newton–Cotesovih pravil je:

- Pri  $n = 1$  dobimo *trapezno pravilo*

$$\int_{x_0}^{x_1} f(x)dx = \frac{h}{2}(y_0 + y_1) - \frac{h^3}{12}f''(\xi).$$

Kot se da sklepati že iz imena, funkcijo interpoliramo s premico in površino trapeza (glej sliko 7.1) vzamemo za približek za integral funkcije.

Koeficiente lahko izračunamo preko integralov Lagrangeevih koeficientov. Tako dobimo

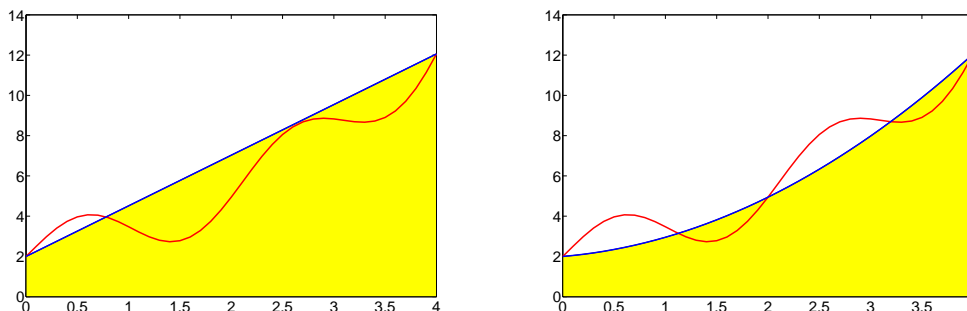
$$A_0 = \int_{x_0}^{x_1} \frac{x - x_1}{x_0 - x_1} dx = \frac{h}{2},$$

$$A_1 = \int_{x_0}^{x_1} \frac{x - x_0}{x_1 - x_0} dx = \frac{h}{2}.$$

Za napako velja

$$R_1(f) = \int_{x_0}^{x_1} \frac{f''(\xi_x)}{2}(x - x_0)(x - x_1)dx = \frac{f''(\xi)}{2} \int_{x_0}^{x_1} (x - x_0)(x - x_1)dx = -\frac{h^3}{12}f''(\xi).$$

Pri izpeljavi napake smo uporabili izrek o povprečni vrednosti, saj je  $(x - x_0)(x - x_1)$  konstantnega predznaka na  $[x_0, x_1]$ . Pravilo je točno za polinome stopnje manjše ali enake 1, kar je očitno že iz same interpolacije s premico.



Slika 7.1: Primerjava trapeznega pravila (levo) in Simpsonovega pravila (desno) na integralu  $\int_0^4 (x^2 - x + 2 + 3 \sin(2x) \cos x) dx$ .

- Pri  $n = 2$  dobimo *Simpsonovo pravilo*

$$\int_{x_0}^{x_2} f(x) dx = \frac{h}{3} (y_0 + 4y_1 + y_2) - \frac{h^5}{90} f^{(4)}(\xi).$$

Kot je razvidno tudi iz slike 7.1, funkcijo interpoliramo s parabolo in integral parabole vzamemo za približek za integral funkcije. Podobno kot pri trapezni metodi bi tudi sedaj lahko uteži določili z integriranjem Lagrangeevih koeficientov. Velja

$$\begin{aligned} A_0 &= \int_{x_0}^{x_2} \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} dx = \frac{h}{3}, \\ A_1 &= \int_{x_0}^{x_2} \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} dx = \frac{4h}{3}, \\ A_2 &= \int_{x_0}^{x_2} \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} dx = \frac{h}{3}. \end{aligned}$$

Iz interpolacije s parabolo je očitno, da je pravilo točno za polinome stopnje 2 ali manj. Izkaže pa se, da je pravilo točno tudi za kubične polinome, saj za  $f(x) = x^3$  velja

$$\int_{x_0}^{x_2} \frac{f^{(3)}(\xi)}{6} (x - x_0)(x - x_1)(x - x_2) dx = 0.$$

Podobno za vsa Newton–Cotesova pravila z lihim številom točk (sodi  $n$ ) zaradi simetrije velja, da so točne tudi za polinome stopnje  $n + 1$ .

Tu sicer ne moremo več uporabiti izreka o povprečni vrednosti, a na srečo za Newton–Cotesova pravila velja, da napako za dovoljkrat zvezno odvedljivo funkcijo lahko ugotovimo iz napake polinoma  $x^{n+1}$  ( $x^{n+2}$  za sodi  $n$ ).

- Pri  $n = 3$  dobimo *3/8 pravilo*

$$\int_{x_0}^{x_3} f(x) dx = \frac{3h}{8} (y_0 + 3y_1 + 3y_2 + y_3) - \frac{3h^5}{80} f^{(4)}(\xi).$$

Čeprav vsebuje en vozle več, ima 3/8 pravilo enak red natančnosti kot Simpsonovo pravilo. Zaradi tega se 3/8 pravila večinoma ne uporablja.

- Pri  $n = 4$  dobimo Boolovo pravilo

$$\int_{x_0}^{x_4} f(x)dx = \frac{2h}{45}(7y_0 + 32y_1 + 12y_2 + 32y_3 + 7y_4) - \frac{8h^7}{945}f^{(6)}(\xi).$$

Tako kot pri Simpsonovi metodi, zaradi simetrije pridobimo en red natančnosti.

Nekaj osnovnih odprtih Newton–Cotesovih pravil je:

- Pri  $n = 2$  dobimo *sredinsko pravilo*

$$\int_{x_0}^{x_2} f(x)dx = 2hy_1 + \frac{h^3}{3}f''(\xi).$$

Čeprav uporabimo vrednost v eni sami točki, je zaradi simetrije pravilo točno tudi za polinome stopnje 1.

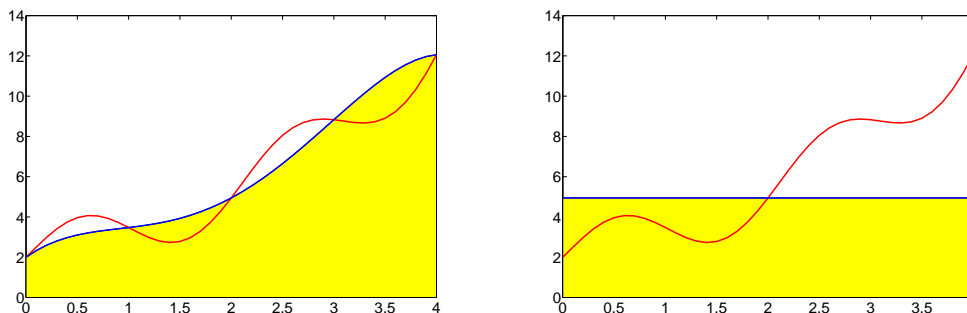
- Pri  $n = 3$  dobimo pravilo

$$\int_{x_0}^{x_3} f(x)dx = \frac{3h}{2}(y_1 + y_2) + \frac{3h^3}{4}f''(\xi).$$

- Pri  $n = 4$  dobimo *Milneovo pravilo*

$$\int_{x_0}^{x_4} f(x)dx = \frac{4h}{3}(2y_1 - y_2 + 2y_3) + \frac{28h^5}{90}f^{(4)}(\xi).$$

To je Newton–Cotesovo pravilo z najmanjšim številom vozlov, ki nima samih uteži s pozitivnim predznakom.



Slika 7.2: Primerjava Boolovega pravila (levo) in sredinskega pravila (desno) na integralu  $\int_0^4 (x^2 - x + 2 + 3 \sin(2x) \cos x) dx$ .

### 7.3 Neodstranljiva napaka

Pri numeričnem odvajanju smo imeli pri numeričnem računanju težave, ko je šel razmik med točkami  $h$  proti 0, saj se je izkazalo, da zaradi neodstranljive napake ne moremo odvoda izračunati poljubno natančno. Glavni problem pri formulah za numerično odvajanje je, da imamo  $h$  v

imenovalcu in ko gre potem  $h$  proti 0, zgoraj v števcu pa je majhna napaka, lahko neodstranljiva napaka naraste preko vseh meja.

Kako pa je s tem pri kvadraturnih formulah? Ker se sedaj  $h$  pojavi v števcu, je na prvi pogled vse v redu, a se izkaže, da imamo tudi tukaj lahko težave.

Vzemimo kvadraturno pravilo

$$\int_a^b f(x)dx = \sum_{i=0}^n A_i f(x_i) + R(f).$$

Mislimo si, da je pri izračunu  $f(x_i)$  absolutna napaka omejena z  $\epsilon > 0$ . Ocena za neodstranljivo napako je potem

$$|D_n| \leq \epsilon \sum_{i=0}^n |A_i|.$$

Očitno velja (sledi iz tega, da je pravilo točno za konstante)

$$\sum_{i=0}^n A_i = b - a.$$

Če so vse uteži pozitivne, od tod sledi  $|D_n| \leq (b - a)\epsilon$  in dobimo lepo oceno za neodstranljivo napako, ki pada proti 0 ko gre  $h$  proti 0. Na žalost pa tako pri zaprtih (za  $n > 8$ ) kot pri odprtih (za  $n > 3$ ) Newton–Cotesovih pravilih dobimo negativne uteži in vsota  $\sum_{i=0}^n |A_i|$  je lahko zelo velika. Tako se tudi tukaj pojavi napaka in višanje stopnje polinoma ni dobra odločitev.

Naslednja tabela prikazuje vsote  $\sum_{i=0}^n |A_i|$ , ki jih dobimo pri zaprtih Newton–Cotesovih pravilih, če integriramo na intervalu  $[0, 1]$ .

$n$	$\sum_{i=0}^n  A_i $
10	3.06
20	$5.44 \cdot 10^2$
30	$2.12 \cdot 10^5$
40	$1.10 \cdot 10^8$
50	$6.70 \cdot 10^{10}$

Iz tabele je očitno, da pri Newton–Cotesovih formulah z višanjem stopnje polinoma ne moremo natančno izračunati integrala. Namesto tega je bolje uporabljati sestavljena pravila.

## 7.4 Sestavljene formule

Pri sestavljenih formulah interval, po katerem integriramo, razdelimo na manjše podintervale. Na vsakem podintervalu uporabimo kvadraturno pravilo nizke stopnje, nato pa rezultate seštejemo.

Denimo, da interval  $[a, b]$  razdelimo z ekvidistantnimi točkami  $a = x_0, b = x_n, h = (b - a)/n, x_i = x_0 + ih$  za  $i = 0, \dots, n$ . Osnovne sestavljene formule so:

- Trapezna formula:

$$\int_{x_0}^{x_n} f(x) dx = \underbrace{\frac{h}{2} (y_0 + 2y_1 + \dots + 2y_{n-1} + y_n)}_{T_h(f)} - \frac{h^2(x_n - x_0)}{12} f''(\xi).$$

Kratka izpeljava:

$$\int_{x_0}^{x_n} f(x) dx = \sum_{k=0}^{n-1} \int_{x_k}^{x_{k+1}} f(x) dx = \sum_{k=0}^{n-1} \left( \frac{h}{2} (y_k + y_{k+1}) - \frac{h^3}{12} f''(\xi_k) \right).$$

Opazimo, da je globalna napaka za en red nižja kot napaka na vsakem intervalu  $[x_i, x_{i+1}]$ . Do tega pride, ker moramo sešteti  $n$  lokalnih napak, pri čemer upoštevamo, da je  $nh = x_n - x_0$ .

- Simpsonova formula:

$$\int_{x_0}^{x_n} f(x) dx = \underbrace{\frac{h}{3} (y_0 + 4y_1 + 2y_2 + 4y_3 + \dots + 2y_{n-2} + 4y_{n-1} + y_n)}_{S_h(f)} - \frac{h^4(x_n - x_0)}{180} f^{(4)}(\xi).$$

- Sredinska formula:

$$\int_{x_0}^{x_n} f(x) dx = \underbrace{2h (y_1 + y_3 + \dots + y_{n-1})}_{U_h(f)} + \frac{h^2(x_n - x_0)}{6} f''(\xi).$$

Pri Simpsonovi in sredinski formuli mora biti  $n = 2m$ .

---

## 7.5 Peanov izrek

**Izrek 7.1 (Peano)** Naj bo  $L$  linearen funkcional oblike

$$\begin{aligned} L(f) &= \int_a^b \left( a_0(x)f(x) + a_1(x)f'(x) + \dots + a_n(x)f^{(n)}(x) \right) dx \\ &+ \sum_{i=0}^{j_0} b_{i0}f(x_{i0}) + \sum_{i=0}^{j_1} b_{i1}f'(x_{i1}) + \dots + \sum_{i=0}^{j_n} b_{in}f^{(n)}(x_{in}), \end{aligned}$$

kjer so funkcije  $a_i$  odsekoma zvezne na  $[a, b]$ , točke  $x_{ij}$  pa ležijo na intervalu  $[a, b]$ . Funkcional  $L$  naj bo za vse polinome stopnje  $n$  ali manj enak 0. Tedaj za vsako funkcijo  $f$ , ki je  $(n+1)$ -krat zvezno odvedljiva na  $[a, b]$ , velja

$$L(f) = \int_a^b f^{(n+1)}(t) K_n(t) dt,$$

kjer je  $K_n$  Peanovo jedro

$$K_n(t) = \frac{1}{n!} L((x-t)_+^n), \quad (x-t)_+^n = \begin{cases} (x-t)^n, & x \geq t, \\ 0, & x < t. \end{cases}$$

*Dokaz.* Taylorjev izrek z ostankom v integralski obliki pravi

$$f(x) = f(a) + f'(a)(x-a) + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + \frac{1}{n!} \int_a^x f^{(n+1)}(t)(x-t)^n dt,$$

kar lahko zapišemo kot

$$f(x) = f(a) + f'(a)(x-a) + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + \frac{1}{n!} \int_a^b f^{(n+1)}(t)(x-t)_+^n dt. \quad (7.2)$$

Če na obeh straneh (7.2) uporabimo funkcional  $L$ , dobimo

$$L(f) = \frac{1}{n!} L \left( \int_a^b f^{(n+1)}(t)(x-t)_+^n dt \right) = \frac{1}{n!} \int_a^b f^{(n+1)}(t) L((x-t)_+^n) dt,$$

saj lahko zamenjamo vrstni red  $L$  in integriranja. ■

**Posledica 7.2** Če je Peanovo jedro  $K_n$  konstantnega predznaka, potem za  $(n+1)$ -krat zvezno odvedljivo funkcijo  $f$  velja

$$L(f) = \frac{f^{(n+1)}(\xi)}{(n+1)!} L(x^{n+1}).$$

*Dokaz.* Po Peanovem izreku velja  $L(x^{n+1}) = \int_a^b (n+1)! K_n(t) dt$ , torej je

$$\int_a^b K_n(t) dt = \frac{1}{(n+1)!} L(x^{n+1}).$$

Ker je  $K_n$  konstantnega predznaka, po izreku o povprečni vrednosti sledi

$$L(f) = \int_a^b f^{(n+1)}(t) K_n(t) dt = f^{(n+1)}(\xi) \int_a^b K_n(t) dt. \quad \blacksquare$$

Peanov izrek uporabimo za računanje napak kvadrature in formul za numerično odvajanje.

**Zgled 7.1** Ocenimo napako trapeznega pravila pri integriranju enkrat zvezno odvedljive funkcije  $f$ . Standardne formule za napako ne moremo uporabiti, saj  $f$  ni dvakrat zvezno odvedljiva. Uporabili bomo Peanov izrek in  $K_0$ .

$$\begin{aligned} K_0(t) &= \int_{x_0}^{x_1} (x-t)_+^0 dx - \frac{h}{2} ((x_0-t)_+^0 + (x_1-t)_+^0) \\ &= \int_t^{x_1} (x-t)^0 dx - \frac{h}{2}(0+1) \\ &= x_1 - t - \frac{h}{2} = \frac{x_0 + x_1}{2} - t. \end{aligned}$$

Peanovo jedro ni konstantnega predznaka, zato lahko le ocenimo

$$|R(f)| \leq \int_{x_0}^{x_1} |f'(t)| \left| \frac{x_0 + x_1}{2} - t \right| dt \leq \frac{h^2}{4} \|f'\|_\infty. \quad \square$$

## 7.6 Richardsonova ekstrapolacija

Z razpolavljanjem  $h$  pridemo pri sestavljenih formulah do natančnejših rezultatov. Iz približkov pri različnih  $h$  lahko ocenimo napako in ugotovimo, če je potrebno še nadaljnje razpolavljanje.

S tovrstnimi ocenami se je prvi ukvarjal Richardson. Naj bo  $S_h(f)$  Simpsonova formula s korakom  $h$  in  $R_h(f)$  napaka Simpsonove formule pri koraku  $h$ . Vemo, da velja

$$I(f) = \int_a^b f(x)dx = S_h(f) + R_h(f) = S_{h/2}(f) + R_{h/2}(f).$$

Za napaki velja

$$R_h(f) = \frac{-(b-a)h^4}{180} f^{(4)}(\xi_1), \quad R_{h/2}(f) = \frac{-(b-a)h^4}{16 \cdot 180} f^{(4)}(\xi_2).$$

Če predpostavimo, da je  $f^{(4)}(\xi_1) \approx f^{(4)}(\xi_2)$ , dobimo

$$R_h(f) \approx 16R_{h/2}(f).$$

Iz  $R_{h/2}(f) = I(f) - S_{h/2}(f) = S_h(f) + R_h(f) - S_{h/2}(f) \approx S_h(f) + 16R_{h/2}(f) - S_{h/2}(f)$  dobimo

$$R_{h/2}(f) \approx \frac{S_{h/2}(f) - S_h(f)}{15}.$$

To formulo lahko uporabimo za oceno napake Simpsonove formule. Ocena sicer ni preveč zanesljiva, saj smo izenačili odvode, a v večini primerov vseeno dobimo spodobne rezultate in lahko ocenimo velikostni razred napake. Poleg ocene lahko iz  $S_{h/2}(f)$  in  $S_h(f)$  z ekstrapolacijo dobimo še boljši približek, saj je

$$I(f) = S_{h/2}(f) + R_{h/2}(f) \approx \frac{16S_{h/2}(f) - S_h(f)}{15}.$$

Podobno lahko naredimo pri trapezni in sredinski formuli, le konstante so drugačne.

## 7.7 Adaptivne metode

Večina metod, ki se jih v praksi uporablja za numerično integriranje, deluje na adaptivnem principu. To pomeni, da metoda sproti ocenjuje napako in temu prilagaja velikost podintervalov. Pri adaptivnih metodah tako na območjih, kjer je funkcija pohlevna, uporabimo večji razmik  $h$ , kjer je njeno obnašanje bolj divje, pa manjši  $h$ .

Rekurzivna adaptivna metoda, ki temelji na Simpsonovem pravilu in Richardsonovi ekstrapolaciji, je predstavljena v algoritmu 7.1.

V algoritmu 7.1 prvi približek  $Q_1$  za integral dobimo tako, da na celem intervalu uporabimo Simpsonovo pravilo, drugi približek  $Q_2$  pa dobimo tako, da interval razdelimo na dva enaka dela in na vsakem uporabimo Simpsonovo pravilo. Potem primerjamo  $Q_1$  in  $Q_2$ .



---

**Algoritem 7.1** Adaptivna Simpsonova metoda. Vhodni podatki so funkcija  $f$ , interval  $[a, b]$  in toleranca  $\delta$ . Metoda vrne vrednost integrala  $Q$  in oceno  $\epsilon \leq \delta$  za napako tega približka.

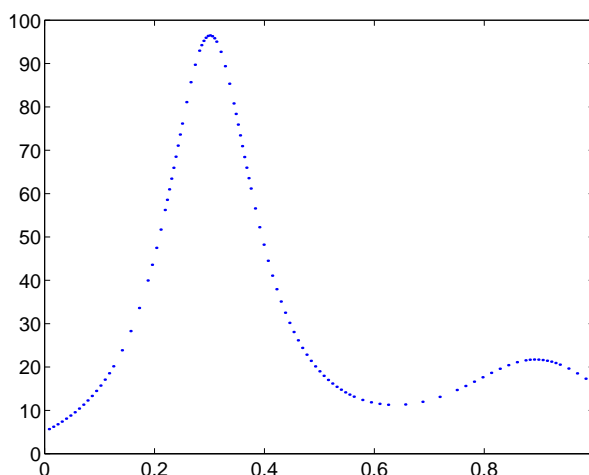
---

```

[Q, ε] =AdaptivniSimpson(f, a, b, δ)
  h = b - a
  c = (a + b)/2
  d = (a + c)/2,  e = (c + b)/2
  Q1 = (h/6)(f(a) + 4f(c) + f(b))
  Q2 = (h/12)(f(a) + 4f(d) + 2f(c) + 4f(e) + f(b))
  ε = |Q1 - Q2|
  if ε ≤ δ
    Q = Q2 + (Q2 - Q1)/15
  else
    [Qa, εa] =AdaptivniSimpson(f, a, c, δ/2)
    [Qb, εb] =AdaptivniSimpson(f, c, b, δ/2)
    Q = Qa + Qb
    ε = εa + εb
  end

```

---



Slika 7.3: Primer uporabe adaptivnega kvadraturega pravila v Matlabu. Vrednosti, označene na grafu, so bile uporabljene za izračun integrala. Vidi se, da točke niso ekvidistantne.

- Če je razlika  $|Q_2 - Q_1|$  večja od zahtevane natančnosti  $\delta$ , potem interval razdelimo na dva dela in na vsakem rekurzivno izračunamo integral z isto metodo, zahtevamo pa, da je sedaj napaka na vsaki polovici pod  $\delta/2$ .
- Če je  $|Q_2 - Q_1| \leq \delta$ , je  $Q_2$  približek za integral, izračunan z zahtevano natančnostjo. Ker pa imamo na voljo še  $Q_1$ , lahko naredimo še korak Richardsonove ekstrapolacije in tako še izboljšamo rezultat.

Namesto Simpsonovega pravila bi lahko uporabili tudi kakšno drugo kvadraturno pravilo. V grobem pri adaptivnem integriranju vedno vrednost integrala ocenimo z dvema kvadraturnima praviloma. Njuna razlika nam služi kot ocena za napako. Če je razlika prevelika, interval

razpolovimo in postopek rekurzivno ponovimo. Da bo postopek ekonomičen, moramo kvadrturni pravili izbrati tako, da se vozli čim bolj prekrivajo in da se vsaj del vozlov prenese v razpolovljeni interval. Tako dosežemo, da je končnih izračunov vrednosti funkcije čim manj, saj pri numeričnem integriranju zahtevnost merimo v tem, koliko vrednosti funkcije moramo izračunati, da pridemo do dovolj dobrega približka.

## 7.8 Rombergova metoda

*Rombergova metoda* na prvi pogled zgleda kot zaporedna uporaba Richardsonove ekstrapolacije za trapezno formulo, a temelji na povsem drugačni osnovi. Medtem, ko Richardsonovo ekstrapolacijo lahko naredimo le enkrat, lahko Rombergovo ekstrapolacijo izvajamo na več nivojih.

**Definicija 7.3** Bernoullijeva števila  $B_k$  so določena z razvojem

$$\frac{x}{e^x - 1} = \sum_{k=0}^{\infty} \frac{B_k}{k!} x^k, \quad |x| < 2\pi.$$

Vsa Bernoullijeva števila so racionalna, vsa z lihimi indeksi razen  $B_1$  pa so enaka 0. Nekaj prvih Bernoullijevih števil je

$$B_0 = 1, \quad B_1 = -\frac{1}{2}, \quad B_2 = \frac{1}{6}, \quad B_4 = -\frac{1}{30}, \quad B_6 = \frac{1}{42}, \quad \dots$$

**Izrek 7.4** Za neskončnokrat zvezno odvedljivo funkcijo  $f$  velja Euler-Maclaurinova sumacijska formula

$$I(f) = \int_a^b f(x) dx = T_h(f) - \sum_{k=1}^{\infty} \frac{B_{2k}}{(2k)!} h^{2k} \left( f^{(2k-1)}(b) - f^{(2k-1)}(a) \right). \quad (7.3)$$

*Dokaz.* Formulo bomo dokazali s simbolnim računanjem. Definiramo operatorje

$$\begin{aligned} Ef(x) &= f(x+h), \\ \Delta f(x) &= f(x+h) - f(x), \\ Df(x) &= f'(x). \end{aligned}$$

Veljajo naslednje formule:

$$\begin{aligned} I + \Delta &= E \\ E &= e^{hD} \quad (\text{razvoj } f(x+h) \text{ v Taylorjevo vrsto}) \\ \Delta &= e^{hD} - I \\ \frac{1}{D} &= \int_a^x f(x) dt + C \quad (\text{nedoločeni integral}) \\ \frac{1}{\Delta} &= \frac{1}{E - I} = \frac{1}{e^{hD} - I} = \frac{1}{hD} - \frac{1}{2}I + \sum_{m=1}^{\infty} B_{2m} \frac{(hD)^{2m-1}}{(2m)!} \quad (\text{razvoj } \frac{z}{e^z - 1}) \end{aligned}$$

Sedaj je

$$(E^n - I) \frac{1}{\Delta} = \frac{E^n - I}{E - I} = E^{n-1} + E^{n-2} + \dots + E + I,$$

to pa pomeni

$$(E^n - I)\frac{1}{\Delta}f(x_0) = \sum_{k=0}^{n-1} f(x_k).$$

Po drugi strani je

$$(E^n - I)\frac{1}{\Delta}f(x_0) = \frac{1}{h} \int_{x_0}^{x_n} f(t)dt - \frac{1}{2}(f(x_n) - f(x_0)) + \sum_{m=1}^{\infty} B_{2m} \frac{h^{2m-1}}{(2m)!} \left( f^{(2m-1)}(x_n) - f^{(2m-1)}(x_0) \right),$$

od tod pa sledi formula (7.3). ■

Vrsta v (7.3) praviloma ni konvergentna, je pa asimptotska, ko gre  $h$  proti 0. Tako dobimo, ko gre  $h$  proti 0,

$$I(f) = T_h(f) - \frac{h^2}{12}(f'(b) - f'(a)) + \frac{h^4}{720}(f'''(b) - f'''(a)) + \dots,$$

kar pomeni

$$I(f) = T_h(f) + \sum_{k=1}^{\infty} a_{k,0}h^{2k},$$

pri čemer so koeficienti  $a_{k,0}$  neodvisni od  $h$ . Če to ponovimo za  $h/2$  in  $h/4$ , dobimo

$$\begin{aligned} I(f) &= T_h(f) + a_{1,0}h^2 + a_{2,0}h^4 + a_{3,0}h^6 + \dots \\ I(f) &= T_{h/2}(f) + a_{1,0} \left(\frac{h}{2}\right)^2 + a_{2,0} \left(\frac{h}{2}\right)^4 + a_{3,0} \left(\frac{h}{2}\right)^6 + \dots \\ I(f) &= T_{h/4}(f) + a_{1,0} \left(\frac{h}{4}\right)^2 + a_{2,0} \left(\frac{h}{4}\right)^4 + a_{3,0} \left(\frac{h}{4}\right)^6 + \dots \end{aligned}$$

Če enačbo s  $T_{h/2}$  pomnožimo s 4 in odštejemo od enačbe s  $T_h$ , se znebimo člena  $h^2$  in dobimo natančnejši približek, podobno pa naredimo s približkoma  $T_{h/4}(f)$  in  $T_{h/2}(f)$ . Velja

$$\begin{aligned} I(f) &= T_{h/2}^{(1)}(f) + a_{2,1}h^4 + a_{3,1}h^6 + \dots \\ I(f) &= T_{h/4}^{(1)}(f) + a_{2,1} \left(\frac{h}{2}\right)^4 + a_{3,1} \left(\frac{h}{2}\right)^6 + \dots, \end{aligned}$$

kjer sta

$$T_{h/2}^{(1)}(f) = \frac{4T_{h/2}(f) - T_h(f)}{3}, \quad T_{h/4}^{(1)}(f) = \frac{4T_{h/4}(f) - T_{h/2}(f)}{3}.$$

Postopek sedaj nadaljujemo v

$$I(f) = T_{h/4}^{(2)}(f) + a_{3,2}h^6 + a_{4,2}h^8 + \dots,$$

kjer je

$$T_{h/4}^{(2)}(f) = \frac{16T_{h/4}^{(1)}(f) - T_{h/2}^{(1)}(f)}{15}.$$

V splošnem postopku tvorimo shemo

napaka	$\mathcal{O}(h^2)$	$\mathcal{O}(h^4)$	$\mathcal{O}(h^6)$	$\mathcal{O}(h^8)$	...
	$T_h^{(0)}(f)$				
	$T_{h/2}^{(0)}(f)$	$T_{h/2}^{(1)}(f)$			
	$T_{h/4}^{(0)}(f)$	$T_{h/4}^{(1)}(f)$	$T_{h/4}^{(2)}(f)$		
	$T_{h/8}^{(0)}(f)$	$T_{h/8}^{(1)}(f)$	$T_{h/8}^{(3)}(f)$	$T_{h/8}^{(4)}(f)$	

kjer je splošna formula

$$T_{h/2^k}^{(j)}(f) = \frac{4^j T_{h/2^k}^{(j-1)}(f) - T_{h/2^{k-1}}^{(j-1)}(f)}{4^j - 1}.$$

**Zgled 7.2** Z Rombergovo metodo bomo izračunali  $\int_1^{2.2} \ln x \, dx = 0.5346062$ . Začeli bomo s  $h = 0.6$  in nato naredili dve razpolavljanji.

Dobimo:

$$\begin{aligned} T_h^{(0)} &= 0.6 \left( \frac{1}{2} \ln 2.2 + \ln 1.6 + \frac{1}{2} \ln 2.2 \right) = 0.5185394 \\ T_{h/2}^{(0)} &= \frac{1}{2} T_h^{(0)} + 0.3(\ln 1.3 + \ln 1.9) = 0.5305351 \\ T_{h/4}^{(0)} &= \frac{1}{2} T_{h/2}^{(0)} + 0.15(\ln 1.15 + \ln 1.45 + \ln 1.75 + \ln 2.05) = 0.5335847 \end{aligned}$$

Pomembno je, da  $T_{h/2^k}$  vedno računamo kot

$$T_{h/2^k} = \frac{1}{2} T_{h/2^{k-1}} + \frac{h}{2^k} (y_1 + y_3 + \dots + y_{2^{k-1}}).$$

Tako vsako funkcijsko vrednost izračunamo enkrat in imamo z Rombergovo ekstrapolacijo zanemarljivo dodatnega dela v primerjavi z računanjem  $T_{h/2^k}^{(0)}$ , rezultat pa je lahko mnogo natančnejši.

Sedaj z Rombergovo ekstrapolacijo dobimo

$$\begin{aligned} T_{h/2}^{(1)} &= \frac{4T_{h/2}^{(0)} - T_h^{(0)}}{3} = 0.5345337 \\ T_{h/4}^{(1)} &= \frac{4T_{h/4}^{(0)} - T_{h/2}^{(0)}}{3} = 0.5346013 \\ T_{h/4}^{(2)} &= \frac{16T_{h/4}^{(1)} - T_{h/2}^{(1)}}{15} = 0.5346058. \end{aligned}$$

□

## 7.9 Gaussove kvadraturene formule

Integral  $\int_a^b f(x)\rho(x)dx$ , kjer je  $\rho$  nenegativna utež, aproksimiramo s kvadraturno formulo

$$\int_a^b f(x)\rho(x)dx = \sum_{i=0}^n A_i^{(n)} f(x_i^{(n)}) + R(f). \quad (7.4)$$

Tako kot prej so koeficienti določeni z vozli, saj velja  $A_i^{(n)} = \int_a^b L_{n,i}(x)\rho(x)dx$ , formula pa je točna za polinome stopnje vsaj  $n$ . S primerno izbiro vozlov lahko dosežemo, da bo formula (7.4) točna za polinome stopnje vsaj  $2n + 1$ , v ozadju pa so ortogonalni polinomi.

To, da je formula (7.4) lahko točna za polinome stopnje  $2n + 1$ , je očitno iz tega, da v njej nastopa  $2n + 2$  parametrov (vozlov in uteži). Po metodi nedoločenih koeficientov bi lahko te parametre nastavili tako, da bi bilo pravilo točno za polinome stopnje  $2n + 1$ . Težava je, ker

dobimo nelinearni sistem, ki ga je težko numerično rešiti. Lažje je, če najprej s pomočjo ortogonalnih polinomov določimo vozle, uteži pa lahko izračunamo z integriranjem Lagrangeevih koeficientov.

Za funkciji  $f$  in  $g$  na  $[a, b]$  definiramo skalarni produkt kot

$$\langle f, g \rangle = \int_a^b f(x)g(x)\rho(x)dx.$$

Funkciji sta si ortogonalni, če je  $\langle f, g \rangle = 0$ . Iz standardne baze polinomov  $1, x, x^2, \dots$  z ortogonalizacijo dobimo ortonormirano bazo  $P_0(x), P_1(x), P_2(x), \dots$ , kjer je  $P_i$  polinom stopnje  $i$  in velja

$$\langle P_i, P_k \rangle = \delta_{ik}.$$

Naj bo sedaj  $P_{n+1}$  tak normiran polinom, da je  $\langle P_{n+1}, q \rangle = 0$  za vsak polinom  $q$  stopnje kvečjemu  $n$  in naj bo  $P_{n+1} = k_{n+1}(x - x_0^{(n)}) \cdots (x - x_n^{(n)})$ . Pri *Gaussovi kvadraturi formuli* za vozle izberemo ničle  $x_0^{(n)}, \dots, x_n^{(n)}$ , torej je  $\omega(x) = (x - x_0^{(n)}) \cdots (x - x_n^{(n)})$ .

Poljuben polinom  $f$  stopnje  $2n + 1$  lahko zapišemo kot  $f(x) = q(x)\omega(x) + r(x)$ , kjer sta  $q$  in  $r$  polinoma stopnje kvečjemu  $n$ . Ker sta si polinoma  $q$  in  $\omega$  ortogonalna, dobimo:

$$\begin{aligned} \int_a^b f(x)\rho(x)dx &= \int_a^b q(x)\omega(x)\rho(x)dx + \int_a^b r(x)\rho(x)dx \\ &= 0 + \sum_{i=0}^n A_i^{(n)} r(x_i^{(n)}) = \sum_{i=0}^n A_i^{(n)} f(x_i^{(n)}). \end{aligned}$$

Torej je pravilo res točno za vse polinome stopnje  $2n + 1$  ali manj.

**Lema 7.5** *Uteži Gaussovih kvadraturnih pravil so pozitivne.*

*Dokaz.* Vzemimo

$$P_i(x) = \frac{\omega^2(x)}{(x - x_i^{(n)})^2}$$

za  $i = 0, \dots, n$ .  $P_i$  je polinom stopnje  $2n$ , torej velja

$$\int_a^b P_i(x)\rho(x)dx = \sum_{k=0}^n A_k P_i(x_k^{(n)}) = A_i P_i(x_i^{(n)}).$$

Ker je  $P_i(x_i^{(n)}) > 0$  in je  $P_i$  nenegativna funkcija, mora biti  $A_i > 0$ . ■

Koeficiente  $A_i$  lahko izračunamo z integriranjem Lagrangeevih koeficientov, še boljše pa je, če uporabimo *Darboux–Cristoffelove formule*:

$$A_i^{(n)} = \frac{1}{\sum_{j=0}^n P_j^2(x_i^{(n)})}, \quad i = 0, \dots, n. \quad (7.5)$$

Iz teorije ortogonalnih polinomov sledi, da so vse ničle  $x_i^{(n)}$  enostavne, realne in leže na  $(a, b)$ .

**Izrek 7.6** Za  $f \in \mathcal{C}^{(2n+2)}[a, b]$  velja

$$\int_a^b f(x)\rho(x) = \sum_{i=0}^n A_i^{(n)} f(x_i^{(n)}) + \frac{f^{(2n+2)}(\xi)}{(2n+2)!k_{n+1}^2}.$$

*Dokaz.* Vzamemo Hermiteov interpolacijski polinom  $H$  v točkah  $x_0^{(n)}, \dots, x_n^{(n)}$ , za katerega velja  $H(x_i^{(n)}) = f(x_i^{(n)})$ ,  $H'(x_i^{(n)}) = f'(x_i^{(n)})$ ,  $i = 0, \dots, n$ . Iz

$$f(x) = H(x) + \frac{f^{(2n+2)}(\xi)}{(2n+2)!}\omega^2(x)$$

sledi

$$\int_a^b f(x)\rho(x)dx = \int_a^b H(x)\rho(x)dx + \int_a^b \frac{f^{(2n+2)}(\xi)}{(2n+2)!}\omega^2(x)dx.$$

Ker je  $H$  polinom stopnje  $2n+1$ , velja

$$\int_a^b H(x)\rho(x)dx = \sum_{k=0}^n A_k H(x_k^{(n)}) = \sum_{k=0}^n A_k f(x_k^{(n)}),$$

za drugi integral pa velja

$$\int_a^b \frac{f^{(2n+2)}(\xi)}{(2n+2)!}\omega^2(x)dx = \int_a^b \frac{f^{(2n+2)}(\xi)}{(2n+2)!} \cdot \frac{P_{n+1}^2(x)}{k_{n+1}^2} dx = \frac{f^{(2n+2)}(\xi)}{(2n+2)!k_{n+1}^2} \int_a^b P_{n+1}^2(x)dx. \quad \blacksquare$$

Najbolj znane družine ortogonalnih polinomov so:

$[-1, 1]$ ,	$\rho(x) = 1$ ,	(Legendre)
$[-1, 1]$ ,	$\rho(x) = (1 - x^2)^{-\frac{1}{2}}$ ,	(Čebišev 1. vrste)
$[-1, 1]$ ,	$\rho(x) = (1 - x^2)^{\frac{1}{2}}$ ,	(Čebišev 2. vrste)
$[-1, 1]$ ,	$\rho(x) = (1 - x)^\alpha(1 + x)^\beta$ , $\alpha, \beta > -1$ ,	(Jacobi)
$[-1, 1]$ ,	$\rho(x) = (1 - x^2)^{\sigma - \frac{1}{2}}$ , $\sigma > \frac{1}{2}$ ,	(Gegenbauer)
$[0, \infty)$ ,	$\rho(x) = x^\sigma e^{-x}$ , $\sigma > -1$	(Laguerre)
$(-\infty, \infty)$ ,	$\rho(x) = e^{-x^2}$ ,	(Hermite).

Za te ortogonalne polinome imamo tabelirane ničle in Gaussove kvadrature formule.

**Zgled 7.3** Gauss–Legendreovi kvadratureni formuli na dveh in treh točkah sta

$$\int_{-1}^1 f(x)dx = f\left(-\sqrt{\frac{1}{3}}\right) + f\left(\sqrt{\frac{1}{3}}\right) + \frac{1}{135}f^{(4)}(\xi),$$

$$\int_{-1}^1 f(x)dx = \frac{5}{9}f\left(-\sqrt{\frac{3}{5}}\right) + \frac{8}{9}f(0) + \frac{5}{9}f\left(\sqrt{\frac{3}{5}}\right) + \frac{1}{15750}f^{(6)}(\xi).$$

Za primerjavo, pri trapeznem in Simpsonovem pravilu dobimo

$$\int_{-1}^1 f(x)dx = f(-1) + f(1) - \frac{2}{3}f^{(2)}(\xi),$$

$$\int_{-1}^1 f(x)dx = \frac{1}{3}f(-1) + \frac{4}{3}f(0) + \frac{1}{3}f(1) + \frac{1}{90}f^{(4)}(\xi). \quad \square$$

Računanje uteži in vozlov za Gaussovo kvadraturno formulo se da prevesti na iskanje lastnih vrednosti in lastnih vektorjev tridiagonalne matrike. Iz teorije ortogonalnih polinomov vemo, da ortonormirani polinomi zadoščajo tričlenski rekurzivni formuli

$$xp_k(x) = b_{k-1}p_{k-1}(x) + a_k p_k(x) + b_k p_{k+1}(x).$$

Če poznamo koeficiente  $a_k$  in  $b_k$ , potem so vozli  $x_0, \dots, x_n$  lastne vrednosti tridiagonalne matrike

$$T_n = \begin{bmatrix} a_0 & b_0 & & & \\ b_0 & a_1 & b_1 & & \\ & \ddots & \ddots & \ddots & \\ & & b_{n-2} & a_{n-1} & b_{n-1} \\ & & & b_{n-1} & a_n \end{bmatrix},$$

uteži pa dobimo iz prvih komponent ortonormiranih lastnih vektorjev:

$$a_k = x_{1k}^2 \int_a^b \rho(x) dx,$$

kjer je  $X$  matrika ortonormiranih lastnih vektorjev matrike  $T_n$ .

Gaussove kvadrature formule lahko posplošimo tako, da so nekateri vozli fiksni, ostale pa določimo tako, da bo formula točna za polinome čim višje stopnje. Tako pri Gauss–Lobattovih formulah fiksiramo obe krajišči, pri Gauss–Radaujevih formulah pa eno krajišče.

Gauss–Lobatto: fiksna vozla sta  $x_0 = a$  in  $x_n = b$ , npr.

$$\int_{-1}^1 f(x) dx = \frac{5}{6}f(-1) + \frac{1}{6}f\left(-\sqrt{\frac{1}{5}}\right) + \frac{1}{6}f\left(\sqrt{\frac{1}{5}}\right) + \frac{5}{6}f(1) + Cf^{(4)}(\xi).$$

Gauss–Radau: fiksni vozli sta  $x_0 = a$ , npr.

$$\int_{-1}^1 f(x) dx = \frac{2}{9}f(-1) + \frac{16 + \sqrt{6}}{18}f\left(\frac{1 - \sqrt{6}}{5}\right) + \frac{16 - \sqrt{6}}{18}f\left(\frac{1 + \sqrt{6}}{5}\right) + Cf^{(5)}(\xi).$$

## 7.10 Izlimitirani integrali

Poglejmo si nekaj možnih pristopov pri računanju izlimitiranih integralov, kjer imamo bodisi pol v krajišču ali pa neskončen interval. Pole in morebitne točke nezveznosti v notranjosti intervala moramo poznati že na začetku in potem interval razdeliti na take podintervale, da v nobenem ni pola ali nezveznosti v notranjosti intervala.

Denimo, da računamo

$$I = \int_a^b \frac{g(x)}{(x-a)^p} dx, \quad 0 < p < 1,$$

kjer je funkcija  $g$  zvezna na  $[a, b]$ . Po definiciji je

$$I = \lim_{\epsilon \rightarrow 0} \int_{a+\epsilon}^b \frac{g(x)}{(x-a)^p} dx.$$

Ta konvergenca je lahko prepočasna, da bi bila uporabna za praktično računanje. Poleg tega lahko pri formulah, kjer vozli vsebujejo tudi krajišča, pričakujemo velike napake pri računanju vrednosti v levem krajišču, ko bo  $\epsilon$  blizu 0. Zato iščemo boljše možnosti.

Varianta 1: Integral  $I$  razdelimo na  $I = I_1 + I_2$ , kjer sta

$$I_1 = \int_a^{a+\epsilon} \frac{g(x)}{(x-a)^p} dx, \quad I_2 = \int_{a+\epsilon}^b \frac{g(x)}{(x-a)^p} dx.$$

Integral  $I_2$  izračunamo s standardnimi metodami, pri računanju integrala  $I_1$  pa funkcijo  $g$  razvijemo v Taylorjevo vrsto okoli točke  $a$ :

$$g(x) = g(a) + (x-a)g'(a) + \frac{(x-a)^2}{2}g''(a) + \dots$$

Tako dobimo

$$I_1 = \epsilon^{1-p} \left( \frac{g(a)}{1-p} + \frac{\epsilon g'(a)}{1!(2-p)} + \frac{\epsilon^2 g''(a)}{2!(3-p)} + \dots \right).$$

Varianta 2: Funkcijo  $g$  razvijemo v Taylorjevo vrsto okoli  $a$ :

$$g(x) = P_s(x) + \text{ostanek},$$

kjer je  $P_s$  polinom stopnje  $s$ . Sedaj  $I$  zapišemo kot  $I = I_1 + I_2$ , kjer sta

$$I_1 = \int_a^b \frac{P_s(x)}{(x-a)^p} dx, \quad I_2 = \int_a^b \frac{g(x) - P_s(x)}{(x-a)^p} dx.$$

Integral  $I_1$  lahko izračunamo eksplicitno, za računanje integrala  $I_2$  pa uporabimo standardne numerične metode.

Varianta 3: Uporabimo substitucijo, npr.

$$(x-a) = t^m, \quad dx = mt^{m-1} dt, \quad m = \frac{k}{1-p}, \quad k \in \mathbb{N}.$$

Dobimo

$$I = m \int_0^{(b-a)^{1/m}} g(a+t^m) t^{k-1} dt$$

in lahko uporabimo standardne metode, saj smo se znebili pola.

Varianta 4: Pomagamo si z Gaussovimi kvadrturnimi formulami. Npr., če imamo

$$I = \int_{-1}^1 \frac{g(x)}{\sqrt{1-x^2}} dx,$$

potem je utež  $\frac{1}{\sqrt{1-x^2}}$  in lahko uporabimo Gauss-Čebiševske kvadrturne formule. Tako lahko npr. pridemo do formule

$$\int_{-1}^1 \frac{g(x)}{\sqrt{1-x^2}} dx = \frac{\pi}{3} \left( g\left(-\frac{\sqrt{3}}{2}\right) + g(0) + g\left(\frac{\sqrt{3}}{2}\right) \right) + Cf^{(6)}(\xi).$$

Pri računanju integrala  $I = \int_a^\infty f(x) dx$  na neskončnem intervalu lahko  $I$  razdelimo na  $I = I_1 + I_2$ , kjer sta

$$I_1 = \int_a^b f(x) dx, \quad I_2 = \int_b^\infty f(x) dx.$$

Integral  $I_1$  izračunamo normalno s standardnimi metodami. Variante za izračun  $I_2$  so



- $I_2$  zanemarimo, če poznamo kakšno dobro oceno za  $|I_2|$ .
- Naredimo substitucijo  $u = 1/x$  in dobimo integral po končnem intervalu

$$I_2 = - \int_{1/b}^0 f(1/u)u^{-2}du.$$

Tudi tu si lahko mogoče pomagamo z Gaussovimi kvadraturnimi formulami.

## 7.11 Večdimenzionalni integrali

Denimo, da računamo integral funkcije  $f$  dveh spremenljivk po pravokotniku  $\Omega = [a, b] \times [c, d]$ . Pišemo lahko

$$\int_{\Omega} f(x, y) dx dy = \int_a^b dx \int_c^d f(x, y) dy.$$

Če interval  $[a, b]$  razdelimo s točkami  $x_i = a + ih, i = 0, \dots, n, h = (b - a)/n$ , interval  $[c, d]$  razdelimo s točkami  $y_j = c + jk, j = 0, \dots, m, k = (d - c)/m$ , in za integrale uporabimo trapezno pravilo, dobimo

$$\int_a^b dx \int_c^d f(x, y) dy = \frac{hk}{4} \sum_{i=0}^n \sum_{j=0}^m A_{ij} f(x_i, y_j) + \mathcal{O}(h^2 + k^2),$$

kjer za koeficiente  $A_{ij}$  velja  $A_{ij} = A_1(i)A_2(j)$  in  $A_1(0) = A_1(n) = 1, A_1(i) = 2$  za  $i \in \{1, \dots, n - 1\}$ , ter  $A_2(0) = A_2(m) = 1, A_2(j) = 2$  za  $j \in \{1, \dots, m - 1\}$ . V bistvu dobimo tenzorski produkt sestavljenega trapeznega pravila, podobno pa lahko naredimo tudi za ostala sestavljena pravila.

**Zgled 7.4** Če imamo npr. mrežo  $5 \times 4$

$(x_0, y_3)$	$(x_1, y_3)$	$(x_2, y_3)$	$(x_3, y_3)$	$(x_4, y_3)$
$(x_0, y_2)$	$(x_1, y_2)$	$(x_2, y_2)$	$(x_3, y_2)$	$(x_4, y_2)$
$(x_0, y_1)$	$(x_1, y_1)$	$(x_2, y_1)$	$(x_3, y_1)$	$(x_4, y_1)$
$(x_0, y_0)$	$(x_1, y_0)$	$(x_2, y_0)$	$(x_3, y_0)$	$(x_4, y_0)$

potem so koeficienti pri trapeznem pravilu

$$\frac{hk}{4} \cdot \begin{array}{|c|c|c|c|c|} \hline 1 & 2 & 2 & 2 & 1 \\ \hline 2 & 4 & 4 & 4 & 2 \\ \hline 2 & 4 & 4 & 4 & 2 \\ \hline 1 & 2 & 2 & 2 & 1 \\ \hline \end{array}.$$

Če pa bi npr. uporabili Simpsonovo metodo za mrežo  $5 \times 5$ , bi dobili

$$\frac{hk}{9} \cdot \begin{array}{|c|c|c|c|c|} \hline 1 & 4 & 2 & 4 & 1 \\ \hline 2 & 8 & 4 & 8 & 2 \\ \hline 4 & 16 & 8 & 16 & 4 \\ \hline 2 & 8 & 4 & 8 & 2 \\ \hline 1 & 4 & 2 & 4 & 1 \\ \hline \end{array}.$$

Če je  $\Omega = [0, 1]^d$  in v vsaki dimenziji porabimo  $n$  točk, uporabimo pa sestavljeno pravilo reda  $k$ , potem za dobljeno tenzorsko pravilo velja, da je napaka reda  $\mathcal{O}(N^{-k/d})$ , kjer je  $N = n^d$  število vseh točk.

Težava je, da pri velikem  $d$  prvič potrebujemo zelo veliko točk ( $n^d$ ), po drugi strani pa napaka prepočasi pada glede na število uporabljenih točk. Zaradi tega se pri velikih  $n$  bolje obnaša metoda Monte Carlo, kjer lahko z znosnim številom točk ponavadi dobimo zadovoljiv približek.

## 7.12 Metoda Monte Carlo

Integral  $I(f) = \int_0^1 f(x)dx$  je enak povprečni vrednosti funkcije  $f$  na intervalu  $[0, 1]$ . To je osnova za *metodo Monte Carlo*.

Če je  $X$  slučajna spremenljivka, enakomerno porazdeljena po  $[0, 1]$ , potem za matematično upanje velja  $E(f(X)) = I(f)$ . Približek za matematično upanje je

$$I_N(f) = \frac{1}{N} \sum_{i=1}^N f(X_i),$$

kjer so  $X_i$  neodvisne naključne vrednosti z  $[0, 1]$ .

Za napako  $\epsilon_N(f) = I(f) - I_N(f)$  velja

$$\epsilon \approx \sigma(f)N^{-1/2},$$

kjer je  $\sigma(f) = \sqrt{D(f)}$  standardna deviacija,

$$D(f) = \int_0^1 (f(x) - I(f))^2 dx$$

pa varianca funkcije  $f$ .

Pri računanju integrala po  $\Omega = [0, 1]^d$  velja enako, le na mestu  $X_i$  sedaj izbiramo vektorje iz  $\Omega$ , kjer je vsak element naključno število med 0 in 1.

Za integral po območju  $I^d = [0, 1]^d$  dobimo

$$I[f] = E[f(\mathbf{X})] = \int_{I^d} f(\mathbf{X}) d\mathbf{X},$$

kjer je  $\mathbf{X}$  vektor neodvisnih slučajnih spremenljivk  $X_1, \dots, X_d$ , enakomerno porazdeljenih na  $[0, 1]$ . Za napako še vedno velja, da pada z  $\mathcal{O}(N^{-1/2})$ , neodvisno od  $d$ .

Če iz kvadrature formule reda  $k$  s tenzorskimi produkti izpeljemo mrežno kvadraturno formulo za integriranje po  $I^d$ , potem pri  $N$  točkah napaka pada z  $\mathcal{O}(N^{-k/d})$ , saj velja  $N = n^d$ , napaka pa pada kot  $n^{-k}$ .

Ostale razlike med mrežnimi kvadraturnimi formulami in metodo Monte Carlo so:

- Mreža pri mrežnih kvadraturnih pravilih je pri velikem  $d$  prevelika. Če npr. uporabimo tenzorsko trapezno pravilo, potem potrebujemo vsaj  $2^d$  izračunov vrednosti funkcije.

- Monte Carlo metode so tudi v eni dimenziji primernejše za računanje Fourierovih koeficientov.
- Monte Carlo metoda nima težav z nezveznostmi, ki sicer zmanjšajo red kvadraturenih pravil.

Z računalnikom ne moremo generirati pravih naključnih števil, saj jih vedno generiramo z nekim procesom, ki ga lahko ponovimo. Tako v resnici delamo s psevdonaključnimi števili. Ena izmed preprostih metod za njihovo generiranje je *linearni kongruenčni model*, kjer računamo zaporedje

$$x_{r+1} = ax_r + c \pmod{m},$$

kjer sta  $a$  in  $c$  določeni pozitivni konstanti, izbiramo pa števila med 0 in  $m - 1$ . Pri tem  $a$ ,  $c$  in  $m$  izberemo tako, da je perioda (števila se najkasneje po  $m$  korakih začnejo ponavljati) čim večja.

Model ima še to težavo, da če generiramo naključne  $d$ -dimenzionalne vektorje, potem točke ležijo na  $(d - 1)$ -dimenzionalnih ravninah, ki jih je približno  $m^{1/d}$ .

### 7.13 Numerično integriranje v Matlabu

Na voljo so naslednje funkcije:

- `quad(f, a, b, tol)`: adaptivno Simpsonovo pravilo,
- `quadl(f, a, b, tol)`: adaptivno Gauss–Lobatto–Kronrodovo pravilo na 4+7 točkah,
- `dblquad(fun, a, b, c, d, tol)`: adaptivno Simpsonovo pravilo na pravokotniku.
- `trapz(x, y)`: sestavljeno trapezno pravilo.

Argumenti so:

- $f, fun$ : funkcija, podana v inline obliki, delovati mora tudi, če je argument vektor,
- $a, b, c, d$ : meje območij ( $[a, b]$  oziroma  $[a, b] \times [c, d]$ ),
- $tol$ : relativna natančnost, ko prenehamo rekurzivno deliti podinterval.
- $x, y$ : vektorja točk  $x_i$  in  $y_i = f(x_i)$ .

Primeri uporabe:

- `f=inline('cos(x.^2)'); quad(f,0,1)`
  - `f=inline('cos(x.^2)+sin(y.^3)'); dblquad(f,0,1,0,1)`
  - `x=linspace(0,pi,100); y=sin(x); trapz(x,y)`
-

## 7.14 Numerično seštevanje vrst

Če integral prevedemo na računanje vrste ali pa tudi kako drugače, se srečamo s problemom, ko je potrebno sešteti vrsto

$$S = \sum_{k=1}^{\infty} a_k,$$

za katero vemo, da je konvergentna. Iščemo torej limito zaporedja delnih vsot

$$S_n = \sum_{k=1}^n a_k.$$

V številnih primerih lahko z ustrežno transformacijo pohitrimo konvergenco delnih vsot in tako pridemo do dovolj točnega približka z računanjem relativno malo členov vrste.

Prva takšna transformacija je *Aitkenova  $\delta^2$ -transformacija*. Kadar je zaporedje delnih vsot  $S_n = \sum_{k=1}^n a_k$  počasi konvergentno in se obnaša podobno kot geometrijska vrsta, potem velikokrat dobimo hitrejše zaporedje z uporabo transformacije

$$T(S_n) = \frac{S_{n+1}S_{n-1} - S_n^2}{S_{n+1} - 2S_n + S_{n-1}}.$$

Transformacijo lahko ponovno uporabimo na novem zaporedju  $T(S_n)$ . Metoda se obnaša dobro pri konvergentnih alternirajočih zaporedjih.

Posplošitev Aitkenove transformacije je *Wynnova  $\epsilon$ -metoda*

$$\epsilon_{r+1}(S_n) = \epsilon_{r-1}(S_n) + \frac{1}{\epsilon_r(S_{n+1}) - \epsilon_r(S_n)},$$

kjer vzamemo  $\epsilon_0(S_n) = S_n$  in  $\epsilon_{-1}(S_n) = 0$ . Metoda se obnaša dobro pri konvergentnih alternirajočih zaporedjih.

Pri *Eulerjevi transformaciji* konvergentno alternirajočo vrsto

$$\sum_{k=0}^{\infty} (-1)^k a_k = a_0 - a_1 + a_2 - a_3 + \dots$$

spremenimo v hitreje konvergentno vrsto z isto vsoto

$$\sum_{k=0}^{\infty} \frac{(-1)^k \Delta^k a_0}{2^{k+1}},$$

kjer je

$$\Delta^k a_0 = \sum_{j=0}^k (-1)^j \binom{k}{j} a_{k-j}$$

$k$ -ta prema diferenca  $a_0$ .

Pomagamo si lahko tudi z Euler–Maclaurinovo formulo, ki smo jo uporabili pri razvoju Rombergove metode. Pri računanju določenega integrala vrednost le tega aproksimiramo s trapezno formulo. Če pa zamenjamo, lahko vrednost vsote za trapezno formulo aproksimiramo z integralom in preostalimi členi v obliki

$$\sum_{k=1}^N f(k) = \int_1^N f(x) dx + \frac{1}{2}(f(N+1) + f(1)) + \sum_{k=1}^m B_{2k} \frac{f^{(2k-1)}(N+1) - f^{(2k-1)}(1)}{(2k)!} + R_m,$$

kjer so  $B_m$  Bernoullijeva števila. Za napako velja

$$R_m = -N \cdot B_{2m+2} \frac{f^{(2m+2)}(\xi)}{(2m+2)!}.$$

Tako dobimo

$$\begin{aligned} \sum_{k=1}^{N-1} f(k) &= \int_1^N f(x) dx - \frac{1}{2}(f(0) - f(N)) + \frac{1}{12}(f'(N) - f'(0)) \\ &\quad - \frac{1}{720}(f^{(3)}(N) - f^{(3)}(0)) + \frac{1}{302240}(f^{(5)}(N) - f^{(5)}(0)) - \frac{1}{1209600}(f^{(7)}(N) - f^{(7)}(0)). \end{aligned}$$

Tako lahko npr. izpeljemo formulo za  $\sum_{k=1}^n k^4$ . Dobimo

$$\sum_{k=1}^n k^4 = \frac{1}{5}n^5 + \frac{1}{2}n^4 + \frac{1}{12}(4n^3) - \frac{1}{720}(24n),$$

višji odvodi pa so enaki 0. Ostane

$$\sum_{k=1}^n k^4 = \frac{1}{30}n(n+1)(2n+1)(3n^2+2n-1).$$

## Dodatna literatura

Numerično integriranje je obravnavano skoraj v vseh učbenikih numerične analize. V slovenščini lahko več o tem preberete v knjigi [17], osnovne informacije so tudi v knjigah [3] in [27], od tuje literature pa omenimo [6] in [16].

## Poglavje 8

# Diferencialne enačbe

### 8.1 Uvod

Rešujemo začetni problem prvega reda v obliki

$$\begin{aligned}y' &= f(x, y) \\ y(x_0) &= y_0,\end{aligned}$$

kjer je  $f$  dana dovolj gladka funkcija, zanima pa nas rešitev na intervalu  $[x_0, b]$ . Želimo, da je problem dobro definiran, kar pomeni, da rešitev obstaja in je enolična.

**Zgled 8.1** Začetni problem  $y' = 1 + y^2$ ,  $y(0) = 0$  ima analitično rešitev  $y(x) = \tan x$ . Rešitev desno od  $x = 0$  vedno bolj strmo narašča, v levo stran pa pada, in pri končni vrednosti  $x$  (pri  $\pm\pi/2$ ) pridemo do  $\pm\infty$ . Torej rešitev obstaja le na intervalu  $(-\pi/2, \pi/2)$ . Ta zgled kaže, da se lahko zgodi, da rešitev začetnega problema ni definirana na celotnem intervalu  $[x_0, b]$ .  $\square$

**Zgled 8.2** Začetni problem  $y' = y^{2/3}$ ,  $y(0) = 0$  ima dve rešitvi, saj tako trivialna rešitev  $y_1(x) = 0$  kot  $y_2(x) = x^3/27$  ustrežata tako diferencialni enačbi kot začetnemu pogoju. To je zgled začetnega problema, ki nima enolične rešitve.  $\square$

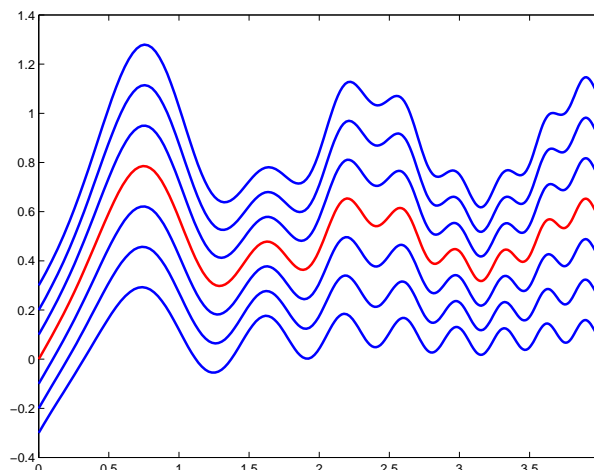
Numerična rešitev bo sestavljena iz zaporedja točk  $x_0 < x_1 < x_2 < \dots$  in pripadajočega zaporedja vrednosti  $y_0, y_1, y_2, \dots$ , kjer je  $y_n$  izračunani približek za rešitev začetnega problema v  $x_n$ , oziroma

$$y_n \approx y(x_n), \quad n = 0, 1, \dots,$$

kjer je  $y(x_n)$  točna vrednost rešitve v  $x_n$ . Pri tem sodobne metode delujejo na adaptivnem principu in avtomatično določajo velikosti korakov  $h_i = x_{i+1} - x_i$  tako, da napake približkov  $y_n$  ostanejo v vnaprej določenih mejah.

Metode za reševanje začetnega problema delimo na:

- *enokoračne metode*:  $y_{n+1}$  izračunamo iz  $y_n$ .
- *večkoračne metode*:  $y_{n+1}$  izračunamo iz  $y_n, y_{n-1}, \dots, y_{n-k+1}$ , kjer je  $k \geq 1$ .



Slika 8.1: Rešitev diferencialne enačbe  $y' = \cos(3x^2) + \sin(4x)y$  pri različnih začetnih pogojih  $y(0) = -0.3, -0.2, \dots, 0.3$ .

Metode ločimo še na:

- *eksplicitne metode*: imamo direktno formulo za  $y_{n+1}$ ,
- *implicitne metode*:  $y_{n+1}$  dobimo tako, da rešimo nelinearno enačbo.

Slika 8.1 prikazuje tipično obnašanje rešitev diferencialne enačbe pri različnih začetnih pogojih. Vsakemu začetnemu pogoju pripada svoja veja. Težava pri reševanju začetnega problema je, da je veja, ki ji želimo slediti, določena z vrednostjo v začetni točki. Ko pa med računanjem zaidemo s prave veje, nimamo več nobene informacije, ki bi nas vrnila nanjo, saj se metode obnašajo tako, kot da bi problem začeli reševati znova iz zadnje točke.

Sistem diferencialnih enačb prvega reda

$$\begin{aligned} y_1' &= f_1(x, y_1, y_2, \dots, y_k) \\ y_2' &= f_2(x, y_1, y_2, \dots, y_k) \\ &\vdots \\ y_k' &= f_k(x, y_1, y_2, \dots, y_k) \end{aligned}$$

z začetnimi pogoji

$$y_1(x_0) = y_{10}, y_2(x_0) = y_{20}, \dots, y_k(x_0) = y_{k0}$$

rešujemo tako kot eno enačbo, če ga zapišemo v vektorski obliki

$$Y' = F(x, Y),$$

kjer je  $Y = [y_1 \ y_2 \ \dots \ y_k]^T$ . Zaradi tega lahko vse metode, ki jih bomo spoznali v nadaljevanju, uporabljamo tako za eno samo enačbo kot tudi za sistem diferencialnih enačb.

## 8.2 Obstoje rešitve, občutljivost in stabilnost

V tem razdelku bomo na kratko ponovili teorijo s področja diferencialnih enačb, ki jo bomo potrebovali v nadaljevanju. Zanima nas, kdaj imamo zagotovljen obstoj enolične rešitve, za samo numerično reševanje pa sta pomembni tudi občutljivost problema in stabilnost rešitve.

Funkcija  $f(x, y)$ , kjer je  $f : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^k$ , na območju  $D = [a, b] \times \Omega \subset \mathbb{R}^{k+1}$  zadošča Lipschitzovemu pogoju na  $y$  s konstanto  $L$ , če za poljubna  $(x, y_1), (x, y_2) \in D$  velja

$$\|f(x, y_1) - f(x, y_2)\| \leq L\|y_1 - y_2\|.$$

Zadostni pogoj, da  $f$  zadošča Lipschitzovemu pogoju, je, da je  $f$  zvezno odvedljiva po  $y$ , saj potem lahko vzamemo

$$L = \max_{(x,y) \in D} \|Jf_y(x, y)\|,$$

kjer je  $Jf_y$   $k \times k$  Jacobijeva matrika  $f$  glede na  $y$ :

$$Jf_y(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial y_1}(x, y) & \cdots & \frac{\partial f_1}{\partial y_k}(x, y) \\ \vdots & & \vdots \\ \frac{\partial f_k}{\partial y_1}(x, y) & \cdots & \frac{\partial f_k}{\partial y_k}(x, y) \end{bmatrix}.$$

V primeru navadne diferencialne enačbe, kjer je  $k = 1$ , je  $Jf_y(x, y)$  kar  $f_y(x, y)$ .

Če  $f$  ustreza Lipschitzovemu pogoju, potem za vsak  $(x_0, y_0) \in D$  obstaja podinterval  $[a, b]$ , ki vsebuje  $x_0$ , na katerem rešitev začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$  obstaja in je enolična.

### 8.2.1 Občutljivost rešitve začetnega problema

Pri občutljivosti nas zanima, kako se spremeni rešitev, če se spremenijo začetni podatki. Tu bomo ločili dva primera.

V prvem imamo motnjo samo v začetnem pogoju, kar ustreza temu, da rešujemo isto diferencialno enačbo, a smo na drugi veji. To se dogaja med numeričnim reševanjem začetnega problema, saj se v vsakem koraku zaradi lokalne napake premaknemo na drugo vejo. Denimo, da  $f$  ustreza Lipschitzovemu pogoju. Točna rešitev začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$  naj bo  $y(x)$ . Če je  $\tilde{y}(x)$  rešitev začetnega problema z zmotenim začetnim pogojem  $y' = f(x, y)$ ,  $y(x_0) = \tilde{y}_0$ , potem za poljuben  $x \geq x_0$  velja

$$\|\tilde{y}(x) - y(x)\| \leq e^{L(x-x_0)} \|\tilde{y}_0 - y_0\|.$$

V drugem primeru obravnavamo tako motnjo v začetnem pogoju kot tudi motnjo v diferencialni enačbi. Tu nas zanima, kaj se zgodi, če npr. zapleteno diferencialno enačbo nadomestimo z lažje izračunljivo aproksimacijo. Če namesto začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$ , rešujemo bližnji problem  $\hat{y}' = \hat{f}(x, \hat{y})$ ,  $\hat{f}(x_0) = \hat{y}_0$ , potem velja

$$\|\hat{y}(x) - y(x)\| \leq e^{L(x-x_0)} \|\hat{y}_0 - y_0\| + \frac{e^{L(x-x_0)} - 1}{L} \|\hat{f} - f\|,$$

kjer je  $\|\hat{f} - f\| = \max_{(x,y) \in D} \|\hat{f}(x, y) - f(x, y)\|$ ,



## 8.2.2 Stabilnost rešitve začetnega problema

Ponavadi smo stabilnost omenjali v povezavi z numeričnimi metodami. Pri diferencialnih enačbah pa govorimo o stabilnosti rešitve problema, ki ni povezana z numeričnimi metodami.

**Definicija 8.1** Pravimo, da je rešitev začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$  stabilna, če za vsak  $\epsilon > 0$  obstaja tak  $\delta > 0$ , da za  $\tilde{y}(x)$ , ki je rešitev  $\tilde{y}' = f(x, \tilde{y})$ ,  $\tilde{y}(x_0) = \tilde{y}_0$ , kjer je  $\|y_0 - \tilde{y}_0\| \leq \delta$ , velja  $\|\tilde{y}(x) - y(x)\| \leq \epsilon$  za vse  $x \geq x_0$ .

Če ima začetni problem stabilno rešitev, lahko pričakujemo, da bo rešitev pri zmotenem začetnem pogoju ostala blizu točne rešitve.

Pravimo, da je stabilna rešitev *asimptotično stabilna*, če gre  $\|\tilde{y}(x) - y(x)\|$  proti 0, ko gre  $x$  proti neskončnosti. Pri asimptotični rešitvi lahko pričakujemo, da bo napaka zmotenega začetnega pogoja šla proti 0, ko gre  $x$  proti neskončnosti.

**Zgled 8.3** Rešitev diferencialne enačbe  $y'(x) = \lambda y$ ,  $y(0) = y_0$ , kjer je  $\lambda = a + ib \in \mathbb{C}$ , je

$$y(x) = y_0 e^{\lambda x} = y_0 e^{ax} (\cos(bx) + i \sin(bx)).$$

Stabilnost rešitve je odvisna od  $a$ . Rešitev je asimptotično stabilna pri  $\operatorname{Re}(\lambda) < 0$ , stabilna pri  $\operatorname{Re}(\lambda) = 0$ , in nestabilna pri  $\operatorname{Re}(\lambda) > 0$ .  $\square$

**Zgled 8.4** Imamo sistem  $y' = Ay$ ,  $y_0 = y_0$ , kjer je  $A$  matrika  $k \times k$ . Če se da  $A$  diagonalizirati in so njene lastne vrednosti  $\lambda_1, \dots, \lambda_k$  z lastnimi vektorji  $v_1, \dots, v_k$ , potem je rešitev

$$y(x) = \sum_{i=1}^k \alpha_i e^{\lambda_i x} v_i,$$

kjer so koeficienti  $\alpha_i$  določeni z

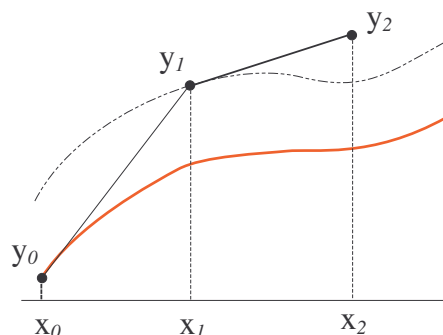
$$y_0 = \sum_{i=1}^k \alpha_i v_i.$$

Rešitev je asimptotično stabilna, če za vse lastne vrednosti velja  $\operatorname{Re}(\lambda_i) < 0$ ; stabilna, če velja  $\operatorname{Re}(\lambda_i) \leq 0$  za vse lastne vrednosti; in nestabilna, če obstaja lastna vrednost  $\lambda_i$ , kjer je  $\operatorname{Re}(\lambda_i) > 0$ .  $\square$

## 8.3 Enokoračne metode

### 8.3.1 Eulerjeva metoda

V prvem koraku reševanja začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$  se moramo iz točke  $(x_0, y_0)$  premakniti v točko  $(x_1, y_1)$ , kjer je  $y_1$  približek za vrednost  $y(x_1)$  in  $x_1 = x_0 + h$ . Poleg vrednosti  $(x_0, y_0)$  lahko v tej točki iz diferencialne enačbe izračunamo še smer tangente  $y'_0 = f(x_0, y_0)$ . Če je  $h$  dovolj majhen, ne bomo naredili velike napake, če bomo rešitev na intervalu  $[x_0, x_1]$  aproksimirali s premico in se premaknili v smeri tangente.



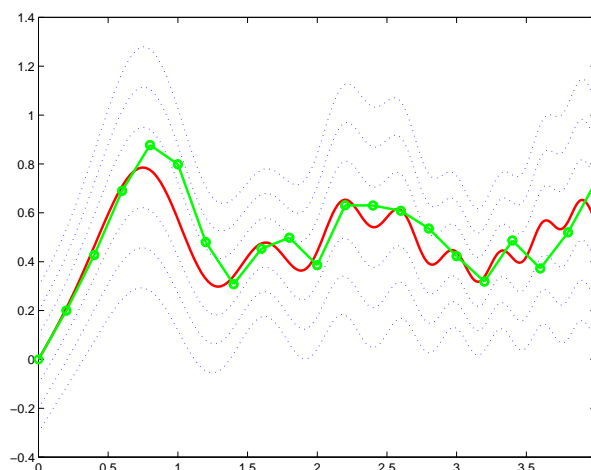
Slika 8.2: Eksplicitna Eulerjeva metoda.

Tako pridemo do najpreprostejše enokoračne metode, ki se imenuje *eksplicitna Eulerjeva metoda*. Nastavek za en korak je

$$\begin{aligned} y_{n+1} &= y_n + hf(x_n, y_n) \\ x_{n+1} &= x_n + h. \end{aligned}$$

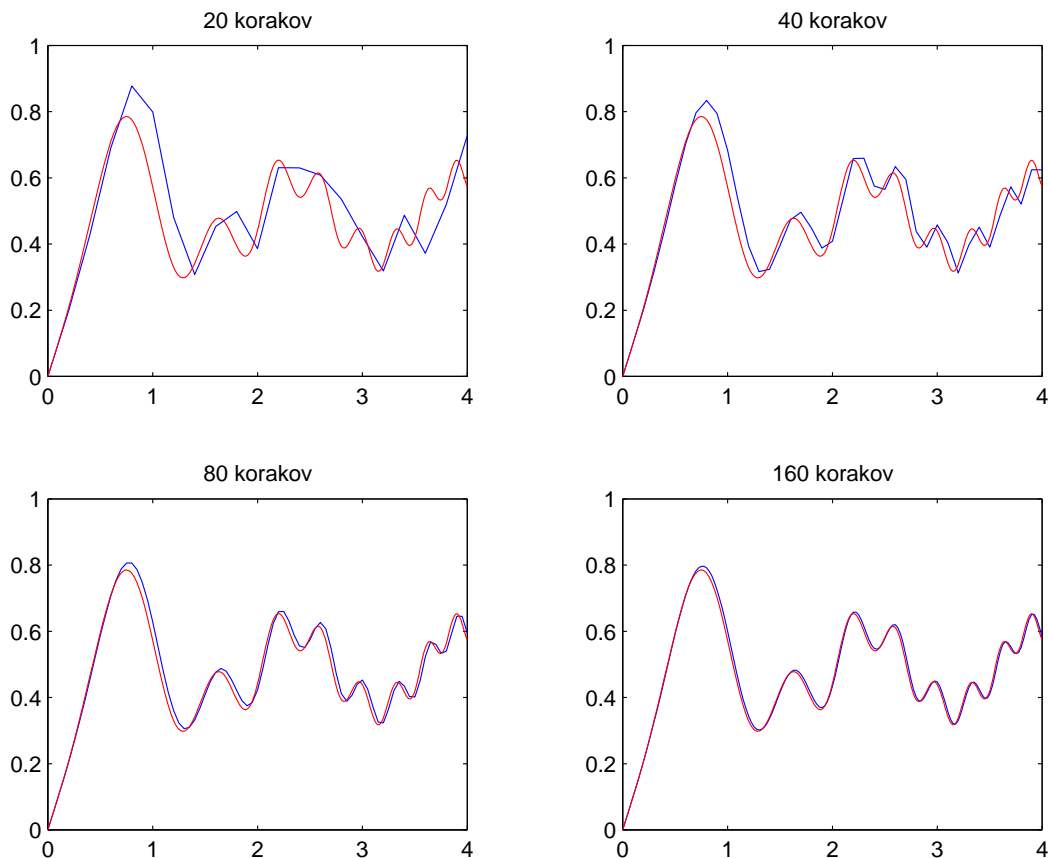
Metoda je prikazana na sliki 8.2. V vsaki točki se premaknemo v smeri tangente. Ponavadi uporabljamo fiksen  $h$ , lahko pa tudi velikost  $h$  v vsakem koraku zmanjšamo ali povečamo odvisno od obnašanja rešitve.

**Zgled 8.5** Oglejmo si delovanje eksplicitne Eulerjeve metode na primeru diferencialne enačbe  $y' = \cos(3x^2) + \sin(4x)y$ ,  $y(0) = 1$ . Denimo, da nas zanima rešitev na intervalu  $[0, 4]$ .



Slika 8.3: Delovanje eksplicitne Eulerjeve metode na enačbi  $y' = \cos(3x^2) + \sin(4x)y$ ,  $y(0) = 1$  na intervalu  $[0, 4]$  s premikom  $h = 0.2$ .

Če vzamemo  $h = 0.2$  in uporabimo eksplicitno Eulerjevo metodo, potem je graf numerično dobljene rešitve na sliki 8.3 predstavljen z zeleno barvo, medtem, ko je točna rešitev pobarvana rdeče. Vidimo, da je premik  $h = 0.2$  prevelik, da bi rešitev lahko sledila vseh lokalnim ekstremom, ki jih ima točna rešitev.



Slika 8.4: Delovanje eksplcitne Eulerjeve metode na enačbi  $y' = \cos(3x^2) + \sin(4x)y$ ,  $y(0) = 1$  na intervalu  $[0, 4]$  s premiki  $h = 0.2, 0.1, 0.05, 0.025$ .

Če zmanjšamo premik  $h$  se poveča čas računanja, a kot kaže slika 8.4, se z manjšanjem  $h$  povečuje natančnost dobljene rešitve. Ko gre  $h$  proti 0, numerična rešitev konvergira proti točni rešitvi.  $\square$

Poleg eksplcitne poznamo tudi *implicitno Eulerjevo metodo*, ki je podana z

$$\begin{aligned} y_{n+1} &= y_n + hf(x_n, y_n) \\ x_{n+1} &= x_n + h. \end{aligned}$$

Pri implicitni metodi rešitev na intervalu  $[x_n, x_{n+1}]$  še vedno aproksimiramo s premico, le da tokrat za smerni koeficient vzamemo vrednost odvoda v točki  $(x_{n+1}, y_{n+1})$ . Ker  $y_{n+1}$  nastopa na obeh straneh enačbe, je potrebno pri implicitni metodi v vsakem koraku rešiti nelinearni sistem za  $y_{n+1}$ .

### 8.3.2 Taylorjeva vrsta

Pri eksplcitni Eulerjevi metodi smo rešitev aproksimirali s tangento. Če bi imeli na voljo še višje odvode, potem bi lahko uporabili razvoj v Taylorjevo vrsto in rešitev aproksimirali s polinomom stopnje višje od ena.

Za razvoj v Taylorjevo vrsto potrebujemo višje odvode  $y$ . Dobimo jih tako, da odvajamo diferencialno enačbo. Če predpostavimo, da je  $f$  dovoljkrat zvezno odvedljiva, potem z odvajanjem enačbe  $y' = f(x, y)$  dobimo

$$\begin{aligned}y' &= f \\y'' &= f_x + f_y y' = f_x + f_y f \\y''' &= f_{xx} + 2f_{xy} f + f_{yy} f^2 + f_y (f_x + f_y f) \\&\vdots\end{aligned}$$

Preko razvoja v Taylorjevo vrsto lahko potem izračunamo

$$y(x_1) = y(x_0 + h) = y_0 + hy'_0 + \frac{h^2}{2}y''_0 + \dots$$

**Zgled 8.6** Z metodo razvoja v Taylorjevo vrsto bomo naredili en korak za začetni problem  $y' = xy + 1$ ,  $y(0) = 0$ . Vzeli bomo  $h = 0.2$  in uporabili prve štiri člene razvoja v Taylorjevo vrsto.

Z odvajanjem diferencialne enačbe dobimo

$$\begin{aligned}y' = xy + 1 &\implies y'(0) = 1 \\y'' = xy' + y &\implies y''(0) = 0 \\y''' = xy'' + 2y' &\implies y'''(0) = 2\end{aligned}$$

$$y(h) = h + \frac{1}{3}h^3 + \dots$$

Pri  $h = 0.2$  tako dobimo

$$y_1 = 0.2 + 0.00267 = 0.20267.$$

Če iščemo  $y(0.4)$ , nadaljujemo iz točke  $(0.2, 0.20267)$  in ponovimo postopek. □

Očitno bo napaka izračunane rešitve manjša, če uporabimo več členov razvoja v Taylorjevo vrsto. Na osnovi tega definiramo lokalno napako metode.

**Definicija 8.2** Pravimo, da je ima metoda lokalno napako reda  $k$ , če se pri točni vrednosti  $y_n = y(x_n)$  izračunani  $y_{n+1}$  ujema z razvojem  $y(x_n + h)$  v Taylorjevo vrsto okrog  $x_n$  do vključno členu  $h^k$ .

Lokalna napaka predstavlja razliko, ki nastopi v enem samem koraku. Pove nam, koliko se naslednji približek razlikuje od vrednosti na veji, kjer smo pričeli z računanjem, torej točki  $(x_n, y_n)$ .

Metoda iz zadnjega zgleda ima red 3, sicer pa z razvijanjem v Taylorjevo vrsto lahko dobimo metodo poljubnega reda. Eulerjeva metoda (tako eksplisitna kot implicitna) ima red 1.

### 8.3.3 Runge-Kutta metode

Runge-Kutta metode so še vedno enokoračne. Da določimo čim natančnejši premik, izračunamo vrednost funkcije  $f$  v  $m$  točkah. Tako dobimo  $m$  premikov v smereh tangent, na koncu pa sestavimo premik iz uteženega povprečja dobljenih premikov.

Najprej izračunamo  $m$  koeficientov

$$k_i = hf(x_n + \alpha_i h, y_n + \sum_{j=1}^i \beta_{ij} k_j), \quad i = 1, \dots, m,$$

nato pa je naslednji približek enak

$$y_{n+1} = y_n + \sum_{i=1}^m \gamma_i k_i.$$

Pri tem je  $m$  stopnja Runge-Kutta metode, ki je ne smemo zamenjevati z redom metode. Konstante  $\alpha_i$ ,  $\beta_{ij}$  in  $\gamma_i$  določimo tako, da se  $y_{n+1}$  čim boljše ujema z razvojem  $y(x_n + h)$  v Taylorjevo vrsto. Izkaže se, da mora veljati  $\alpha_i = \sum_{j=1}^i \beta_{ij}$  in  $\sum_{i=1}^m \gamma_i = 1$ .

V primeru, ko je  $\beta_{ii} = 0$  za  $i = 1, \dots, m$ , je metoda eksplicitna, sicer pa implicitna.

**Zgled 8.7** Dvostopenjska eksplicitna Runge-Kutta metoda ima obliko

$$\begin{aligned} k_1 &= hf(x_n, y_n) \\ k_2 &= hf(x_n + \alpha h, y_n + \beta k_1) \\ y_{n+1} &= y_n + \gamma_1 k_1 + \gamma_2 k_2. \end{aligned}$$

S pomočjo razvoja funkcije  $f$  v Taylorjevo vrsto kot funkcijo dveh spremenljivk dobimo razvoja

$$\begin{aligned} k_1 &= hf \\ k_2 &= hf + \alpha h^2 f_x + \beta h k_1 f_y + \mathcal{O}(h^3). \end{aligned}$$

Od tod sledi, da za naslednji približek  $y_{n+1}$  vzamemo

$$y_{n+1} = y_n + (\gamma_1 + \gamma_2)hf + \gamma_2 \alpha h^2 f_x + \gamma_2 \beta h^2 f f_y + \mathcal{O}(h^3),$$

kar primerjamo z

$$y(x_n + h) = y(x_n) + hf + \frac{1}{2}h^2(f_x + f f_y) + \mathcal{O}(h^3).$$

Sledi

$$\begin{aligned} \gamma_1 + \gamma_2 &= 1 \\ \alpha \gamma_2 &= \frac{1}{2} \\ \beta \gamma_2 &= \frac{1}{2}. \end{aligned}$$

Sistem ima več rešitev, saj za poljubni  $\gamma_2 \neq 0$  dobimo

$$\gamma_1 = 1 - \gamma_2, \quad \alpha = \frac{1}{2\gamma_2}, \quad \beta = \frac{1}{2\gamma_2}.$$

□

Primeri eksplicitne dvostopenjske Runge-Kutta metode drugega reda sta:

- Heunova metoda

$$\begin{aligned}k_1 &= hf(x_n, y_n) \\k_2 &= hf(x_n + h, y_n + k_1) \\y_{n+1} &= y_n + \frac{1}{2}(k_1 + k_2), \quad \text{lokalna napaka : } \mathcal{O}(h^3); \end{aligned}$$

- modificirana Eulerjeva metoda

$$\begin{aligned}k_1 &= hf(x_n, y_n) \\k_2 &= hf(x_n + \frac{1}{2}h, y_n + \frac{1}{2}k_1) \\y_{n+1} &= y_n + k_2, \quad \text{lokalna napaka : } \mathcal{O}(h^3). \end{aligned}$$

Zelo znana je tudi Runge-Kutta 4-stopenjska metoda reda 4, kjer vzamemo

$$\begin{aligned}k_1 &= hf(x_n, y_n) \\k_2 &= hf(x_n + \frac{1}{2}h, y_n + \frac{1}{2}k_1) \\k_3 &= hf(x_n + \frac{1}{2}h, y_n + \frac{1}{2}k_2) \\k_4 &= hf(x_n + h, y_n + k_3) \\y_{n+1} &= y_n + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4), \quad \text{lokalna napaka : } \mathcal{O}(h^5). \end{aligned}$$

**Zgled 8.8** Denimo, da s 4-stopenjsko Runge-Kutta metodo rešujemo  $y' = -y - 5e^x \sin(x)$ ,  $y(0) = 1$ , kjer vzamemo  $h = 0.1$ . Dobimo

$$\begin{aligned}k_1 &= hf(x_0, y_0) = 0.1 * f(0, 1) = -0.1 \\k_2 &= hf(x_0 + \frac{1}{2}h, y_0 + \frac{1}{2}k_1) = 0.1 * f(0.05, 0.95) = -0.12127 \\k_3 &= hf(x_0 + \frac{1}{2}h, y_0 + \frac{1}{2}k_2) = 0.1 * f(0.05, 0.93936) = -0.12021 \\k_4 &= hf(x_0 + h, y_0 + k_3) = 0.1 * f(0.1, 0.87979) = -0.14315 \\y_1 &= y_0 + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) = 0.87898. \quad \square \end{aligned}$$

Izkaže se, da je pri eksplicitnih Runge-Kutta metodah stopnje 5 ali več red vedno manjši od stopnje. Tako je npr. maksimalni red 5-stopenjske eksplicitne Runge-Kutta metoda enak 4, 6-stopenjske pa 5.

### 8.3.4 Adaptivna ocena koraka

Če imamo na voljo oceno lokalne napake, lahko  $h$  adaptivno prilagajamo. Denimo, da po 4-stopenjski Runge-Kutta metodi reda 4 iz  $y(x)$  izračunamo  $y(x + 2h)$  enkrat s korakom  $2h$ , drugič pa v dveh korakih s korakom  $h$ . Podobno kot pri Richardsonovi ekstrapolaciji dobimo

$$\begin{aligned}y(x + 2h) &= y^{(1)} + (2h)^5 \cdot C_1 + \mathcal{O}(h^6) \\y(x + 2h) &= y^{(2)} + 2(h)^5 \cdot C_2 + \mathcal{O}(h^6) \end{aligned}$$

in

$$\Delta = \frac{y^{(2)} - y^{(1)}}{15}$$

je ocena za lokalno napako približka  $y^{(1)}$ . Ko je ocena  $\Delta$  velika  $h$  zmanjšamo, če pa je ocena dovolj majhna lahko  $h$  povečamo. Za izračun  $\Delta$  in  $y^{(1)}$  potrebujemo 11 izračunov vrednosti funkcije  $f$  (4 izračune za vsako Runge-Kutta metodo, pri čemer se izračun  $f(x_n, y_n)$  ponovi dvakrat).

Boljša je *Runge-Kutta–Fehlbergova metoda*, kjer vzamemo 6 stopenjsko Runge-Kutta metodo

$$k_i = hf(x_n + \alpha_i h, y_n + \sum_{j=1}^{i-1} \beta_{ij} k_j), \quad i = 1, \dots, 6,$$

potem pa iz istih  $k_1, \dots, k_6$  sestavimo metodo reda 4

$$y_{n+1} = y_n + \sum_{i=1}^6 \gamma_i k_i$$

in kontrolno metodo reda 5

$$y_{n+1}^* = y_n + \sum_{i=1}^6 \gamma_i^* k_i.$$

Ocena za napako  $y_{n+1}$  je potem kar  $y_{n+1} - y_{n+1}^* = \sum_{i=1}^6 (\gamma_i - \gamma_i^*) k_i$ .

Podrobne formule za Runge-Kutta–Fehlbergovo metodo so

$$\begin{aligned} k_1 &= hf(x_n, y_n) \\ k_2 &= hf(x_n + \frac{1}{4}h, y_n + \frac{1}{4}k_1) \\ k_3 &= hf(x_n + \frac{3}{8}h, y_n + \frac{3}{32}k_1 + \frac{9}{32}k_2) \\ k_4 &= hf(x_n + \frac{12}{13}h, y_n + \frac{1932}{2197}k_1 - \frac{7200}{2197}k_2 + \frac{7296}{2197}k_3) \\ k_5 &= hf(x_n + h, y_n + \frac{439}{216}k_1 - 8k_2 + \frac{3680}{513}k_3 - \frac{845}{4104}k_4) \\ k_6 &= hf(x_n + \frac{1}{2}h, y_n - \frac{8}{27}k_1 + 2k_2 - \frac{3544}{2565}k_3 + \frac{1859}{4104}k_4 - \frac{11}{40}k_5) \\ y_{n+1} &= y_n + \frac{25}{216}k_1 + \frac{1408}{2565}k_3 + \frac{2197}{4104}k_4 - \frac{1}{5}k_5, \quad \text{lokalna napaka : } \mathcal{O}(h^5) \\ y_{n+1}^* &= y_n + \frac{16}{135}k_1 + \frac{6656}{12825}k_3 + \frac{28561}{56430}k_4 - \frac{9}{50} + \frac{2}{55}k_6, \quad \text{lokalna napaka : } \mathcal{O}(h^6). \end{aligned}$$

Opazimo lahko, da za izračun  $y_{n+1}$  potrebujemo le koeficiente  $k_1, \dots, k_5$ , torej je vrednost  $y_{n+1}$  v bistvu dobljena s 5-stopenjsko metodo. Koeficient  $k_6$  potrebujemo le za izračun kontrolne vrednosti  $y_{n+1}^*$ .

V naslednjem koraku vzamemo razmik  $qh$ , kjer je  $\epsilon$  željena natančnost in

$$q = \left( \frac{\epsilon h}{2|y_{n+1}^* - y_{n+1}|} \right)^{1/4}.$$

Potrebno je poudariti, da pri Runge-Kutta–Fehlbergovi metodi za naslednji približek vzamemo  $y_{n+1}$ , kljub temu, da imamo na voljo po vsej verjetnosti natančnejši približek  $y_{n+1}^*$ . To naredimo zaradi tega, ker iz razlike  $y_{n+1} - y_{n+1}^* = \mathcal{O}(h^5)$  lahko ocenimo vodilni člen lokalne napake  $y_{n+1}$ , tega pa ne moremo narediti za  $y_{n+1}^*$ . Tako ima Runge-Kutta–Fehlbergova metoda red 4.

## 8.4 Stabilnost in konvergenca enokoračnih metod

Numerične metode so primerne za uporabo le, če so stabilne in konvergentne. V tem razdelku bomo pokazali, kdaj lahko ti lastnosti predpostavimo za enokoračne metode. Navedli bomo le definicije in glavne izreke. Podrobnosti skupaj z dokazi izrekov lahko najdete v [17].

Vsako enokoračno metodo lahko zapišemo v obliki

$$y_{n+1} = y_n + h\phi(x_n, y_n, h),$$

kjer je  $\phi$  funkcija prirastka.

**Definicija 8.3** Enokoračna metoda za numerično reševanje začetnega problema je konsistentna, če velja

$$\lim_{h \rightarrow 0} \phi(x, y, h) = f(x, y).$$

**Zgled 8.9** Funkcija prirastka za modificirano Eulerjevo metodo

$$y_{n+1} = y_n + hf(x_n + \frac{h}{2}, y_n + \frac{1}{2}hf(x_n, y_n))$$

je  $\phi(x_n, y_n, h) = f(x_n + \frac{h}{2}, y_n + \frac{1}{2}hf(x_n, y_n))$ . Metoda je očitno konsistentna. □

Lokalno napako pri izračunu  $y_{n+1}$  lahko zapišemo kot

$$T(x_{n+1}) = y(x_{n+1}) - y(x_n) - h\phi(x_n, y_n, h).$$

Če velja  $T(x_{n+1}) = \mathcal{O}(h^{p+1})$ , potem je metoda reda  $p$ .

Numerična metoda za reševanje začetnega problema je stabilna, kadar majhne motnje ne povzročijo divergence numerične rešitve.

**Definicija 8.4** Denimo, da iščemo rešitev začetnega problema na intervalu  $[a, b]$ , ki ga razdelimo z ekvidistantnimi točkami  $x_i = a + ih$ , kjer je  $h = (b - a)/n$  in  $i = 0, \dots, n$ . Pravimo, da je enokoračna metoda stabilna, če za vsako diferencialno enačbo, ki zadošča Lipschitzovemu pogoju, obstajata taki konstanti  $h_0 > 0$  in  $K > 0$ , da za poljubni dve rešitvi  $y_r, \tilde{y}_r$ , ki ju dobimo z  $0 < h \leq h_0$ , velja  $\|y_r - \tilde{y}_r\| \leq K\|y_0 - \tilde{y}_0\|$  za  $x_r \in [a, b]$ .

Stabilnost numerične metode ni povezana s stabilnostjo rešitve diferencialne enačbe, ki smo jo definirali v podrazdelku 8.2.2.

**Izrek 8.5** Če funkcija prirastka  $\phi$  zadošča Lipschitzovemu pogoju na  $y$ , potem je metoda stabilna.

**Definicija 8.6** Enokoračna metoda je konvergentna, če za vsako diferencialno enačbo, ki zadošča Lipschitzovemu pogoju, za vsak  $r$  velja

$$\lim_{h \rightarrow 0} \|y_r - y(x_r)\| = 0.$$

Pri tem predpostavimo, da v limiti upoštevamo le tako majhne  $h$ , da je  $rh \leq |b - a|$ .

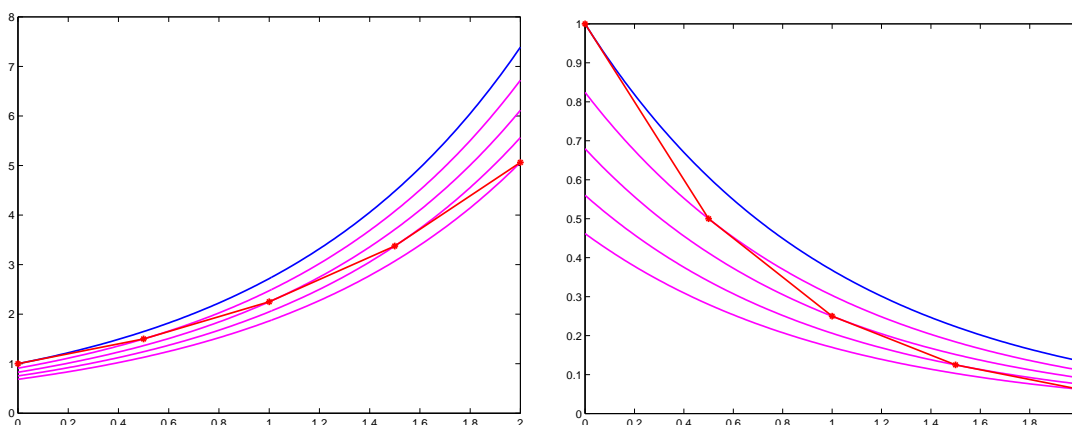


**Izrek 8.7** Če je funkcija prirastka  $\phi$  zvezna v  $x, y, h$  in zadošča Lipschitzovemu pogoju na  $y$ , potem je metoda konvergentna natanko tedaj, ko je metoda konsistentna.

**Izrek 8.8 (O globalni napaki)** Če funkcija prirastka  $\phi$  zadošča pogojem za konvergenco in za lokalno napako velja  $\|T(x)\| \leq Dh^{p+1}$ , potem za globalno napako velja

$$\|y_n - y(x_n)\| \leq Dh^p \frac{e^{L(x_n - x_0)} - 1}{L} + e^{L(x_n - x_0)} \|y_0 - y(x_0)\|.$$

Globalna napaka ni preprosto kar vsota lokalnih napak, v grobem pa velja, da ima metoda z lokalno napako reda  $k$  globalno napako reda  $k - 1$ . Primeri na sliki 8.5 prikazujeta dva primera povezave globalne in lokalne napake. V prvem primeru (levo) je globalna napaka večja od vsote absolutnih vrednosti lokalnih napak. V desnem primeru pa je globalna napaka manjša od posameznih lokalnih napak.



Slika 8.5: Povezava med globalno in lokalno napako.

**Zgled 8.10** Stabilnost diferencialne enačbe lahko razberemo iz obnašanja pri reševanju modelne enačbe  $y' = \lambda y$ ,  $y(0) = y_0$ . Ta diferencialna enačba ima rešitev  $y(x) = y_0 e^{\lambda x}$ . V primeru, ko je  $\text{Re}(\lambda) < 0$ , je rešitev asimptotično stabilna.

Če vzamemo eksplisitno Eulerjevo metodo, potem dobimo  $y_1 = (1 + \lambda h)y_0$  in

$$y_k = (1 + \lambda h)^k y_0.$$

Če naj bo pri  $\text{Re}(\lambda) < 0$  Eulerjeva metoda stabilna, mora veljati  $|1 + \lambda h| < 1$ , torej mora biti  $h$  omejen ( $h\lambda$  mora ležati v krogu z radijem 1 okrog  $-1$ ).

Če vzamemo implicitno Eulerjevo metodo, potem dobimo  $(1 - \lambda h)y_1 = y_0$  in

$$y_k = \left( \frac{1}{1 - \lambda h} \right)^k y_0.$$

Sedaj mora pri  $\text{Re}(\lambda) < 0$  veljati  $\left| \frac{1}{1 - \lambda h} \right| \leq 1$ , torej  $h$  ni omejen in za ta začetni problem je implicitna Eulerjeva metoda vedno stabilna.  $\square$

## 8.5 Večkoračne metode

Pri večkoračnih metodah poleg zadnjega uporabimo tudi še nekaj prejšnjih približkov. Splošni nastavek je

$$\sum_{i=0}^k \alpha_i y_{n+1-i} + h \sum_{i=0}^k \beta_i f_{n+1-i} = 0,$$

kjer je  $f_i = f(x_i, y_i)$ , privzamemo pa še  $\alpha_0 = 1$ . Če je  $\beta_0 = 0$ , je metoda eksplicitna, sicer pa implicitna. Metodo določata *rodovna polinoma*

$$\begin{aligned} \rho(z) &= \sum_{i=0}^k \alpha_{k-i} z^i, \\ \sigma(z) &= \sum_{i=0}^k \beta_{k-i} z^i. \end{aligned}$$

Ker na začetku potrebujemo več kot eno začetno vrednost, moramo te približke pridobiti s kakšno drugo metodo. Lahko uporabimo npr. Runge-Kutta metodo ali razvijanje v Taylorjevo vrsto, paziti pa moramo na to, da ne uporabimo metode, ki ima manjši red lokalne napake kot večkoračna metoda. Pri reševanju začetnega problema je namreč tako, da ko so napake enkrat prisotne, se tudi z natančnejšim računanjem v nadaljevanju ne moremo več približati točni rešitvi.

*Adamsove metode* dobimo tako, da enačbo  $y' = f(x, y)$  integriramo na  $[x_n, x_{n+1}]$

$$y_{n+1} - y_n = \int_{x_n}^{x_{n+1}} y' dx = \int_{x_n}^{x_{n+1}} f(x, y) dx,$$

$f$  pa nadomestimo z interpolacijskim polinomom na točkah:

a)  $x_n, x_{n-1}, \dots, x_{n-k+1}$ : eksplicitne *Adams–Bashworthove formule*

$$\begin{aligned} y_{n+1} &= y_n + hf_n, & \text{lokalna napaka } \mathcal{O}(h^2), \\ y_{n+1} &= y_n + \frac{h}{2}(3f_n - f_{n-1}), & \text{lokalna napaka } \mathcal{O}(h^3), \\ y_{n+1} &= y_n + \frac{h}{12}(23f_n - 16f_{n-1} + 5f_{n-2}), & \text{lokalna napaka } \mathcal{O}(h^4). \end{aligned}$$

b)  $x_{n+1}, x_n, \dots, x_{n-k+2}$ : implicitne *Adams–Moultonove formule*

$$\begin{aligned} y_{n+1} &= y_n + hf_{n+1}, & \text{lokalna napaka } \mathcal{O}(h^2), \\ y_{n+1} &= y_n + \frac{h}{2}(f_{n+1} + f_n), & \text{lokalna napaka } \mathcal{O}(h^3), \\ y_{n+1} &= y_n + \frac{h}{12}(5f_{n+1} + 8f_n - f_{n-1}), & \text{lokalna napaka } \mathcal{O}(h^4). \end{aligned}$$

Pri predpostavki, da je  $k = n$  in  $x_0 = 0$ , definiramo linearni funkcional

$$L(y) = \sum_{i=0}^k \left( \alpha_i y((k-i)h) + h\beta_i y'((k-i)h) \right) = \sum_{j=0}^k \left( \alpha_{k-j} y(jh) + h\beta_{k-j} y'(jh) \right).$$

Če  $y$  in  $y'$  razvijemo v Taylorjevo vrsto okrog 0, dobimo

$$L(y) = d_0 y(0) + d_1 h y'(0) + d_2 h^2 y''(0) + \dots$$

Metoda je reda  $p$ , če je  $L(y) = \mathcal{O}(h^{p+1})$  za vsako  $(p+1)$ -krat zvezno odvedljivo funkcijo  $y$ .

**Izrek 8.9** Za večkorlačno metodo je ekvivalentno:

- a)  $d_0 = d_1 = \dots = d_m = 0$ ,
- b)  $L(q) = 0$  za poljuben polinom  $q$  stopnje kvečjemu  $m$ ,
- c)  $L(y) = \mathcal{O}(h^{m+1})$  za vsako  $(m+1)$ -krat zvezno odvedljivo funkcijo  $y$ .

Red metode je  $p$ , če velja  $d_0 = d_1 = \dots = d_p = 0 \neq d_{p+1}$ .

Če v Taylorjevo vrsto razvijemo  $y$  in  $y'$ , dobimo

$$\begin{aligned} y(jh) &= \sum_{r=0}^{\infty} \frac{(jh)^r}{r!} y^{(r)}(0), \\ y'(jh) &= \sum_{r=0}^{\infty} \frac{(jh)^r}{r!} y^{(r+1)}(0). \end{aligned}$$

Sedaj dobimo naslednje formule za  $d_0, d_1, d_2, \dots$ :

$$\begin{aligned} d_0 &= \sum_{j=0}^k \alpha_{k-j}, \\ d_1 &= \sum_{j=0}^k (j\alpha_{k-j} + \beta_{k-j}), \\ d_2 &= \sum_{j=0}^k \left( \frac{j^2}{2} \alpha_{k-j} + j\beta_{k-j} \right), \\ &\vdots \\ d_q &= \sum_{j=0}^k \left( \frac{j^q}{q!} \alpha_{k-j} + \frac{j^{q-1}}{(q-1)!} \beta_{k-j} \right). \end{aligned}$$

**Zgled 8.11** Določi red večkorlačne metode  $y_n - y_{n-2} = \frac{h}{3}(f_n + 4f_{n-1} + f_{n-2})$ .

Koeficienti so  $\alpha_0 = -1$ ,  $\alpha_1 = 0$ ,  $\alpha_2 = 1$ ,  $\beta_0 = \frac{1}{3}$ ,  $\beta_1 = \frac{4}{3}$ ,  $\beta_2 = \frac{1}{3}$ . Po zgornjih formulah lahko izračunamo  $d_0 = d_1 = \dots = d_4 = 0$  in  $d_5 = -\frac{1}{90}$ , kar pomeni, da je red metode enak 4.  $\square$

Ponavadi pri večkorlačnih metodah uporabljamo *prediktor-korektor metode*, kjer za izračun prediktora  $y_{n+1}^P$  vzamemo eksplisitno metodo, za korektor pa vzamemo implicitno metodo, katere enačbo rešujemo iterativno tako, da na desno stran vstavimo približek  $y_{n+1}^P$ , na levi pa dobimo  $y_{n+1}^K$ . Tako se izognemo reševanju nelinearne enačbe oziroma je to kar en korak navadne iteracije za reševanje nelinearne enačbe. Postopek lahko ponovimo tako da dobljeni približek vstavimo na desno stran in izračunamo nov približek. Teorija nam zagotavlja, da bo ta postopek skonvergirala, če je  $h$  dovolj majhen.

Primer so Milneove metode, kjer  $y' = f(x, y)$  integriramo na  $[x_{n-k}, x_{n+1}]$ :

$$y_{n+1} - y_{n-k} = \int_{x_{n-k}}^{x_{n+1}} f(x, y) dx,$$

integral pa računamo preko Newton–Cotesovih formul. Pri odprti dobimo eksplicitno, pri zaprti pa implicitno metodo. Primer je Milneov par

$$\begin{aligned} y_{n+1} &= y_{n-3} + \frac{4h}{3}(2f_n - f_{n-1} + 2f_{n-2}), & \text{napaka } \mathcal{O}(h^5) & : \text{ prediktor} \\ y_{n+1} &= y_{n-1} + \frac{h}{3}(f_{n+1} + 4f_n + f_{n-1}), & \text{napaka } \mathcal{O}(h^5) & : \text{ korektor} \end{aligned}$$

**Zgled 8.12** Diferencialno enačbo  $y' = -y - 5e^{-x} \sin(5x)$  z začetnim pogojem  $y(0) = 1$  rešujemo z Milneovim parom prediktor-korektor

$$\begin{aligned} y_{n+1}^{(P)} &= y_{n-3} + \frac{4h}{3}(2f_n - f_{n-1} + 2f_{n-2}), \\ y_{n+1}^{(K)} &= y_{n-1} + \frac{h}{3}(f(x_{n+1}, y_{n+1}^{(P)}) + 4f_n + f_{n-1}). \end{aligned}$$

Denimo, da smo pri  $h = 0.1$  z drugimi metodami že prišli do približkov  $y_1 = 0.79407$ ,  $y_2 = 0.44235$ ,  $y_3 = 0.05239$ , sedaj pa bi z Milneovim parom radi izračunali  $y_4$ .

$$y_4^{(P)} = 1 + \frac{0.4}{3}(2 \cdot (-3.7472) + 3.8870 + 2 \cdot (-2.9631)) = -0.27115$$

$$y_4^{(K)} = 0.44235 + \frac{0.1}{3}(-2.7765 + 4 \cdot (-3.7472) - 3.8870) = -0.27939$$

Postopek lahko nadaljujemo,  $y_4^{(K)}$  vnesemo na desno stran in dobimo nov približek.

$$y_4^{(K2)} = 0.44235 + \frac{0.1}{3}(-2.7682 + 4 \cdot (-3.7472) - 3.8870) = -0.27912$$

$$y_4^{(K3)} = 0.44235 + \frac{0.1}{3}(-2.7685 + 4 \cdot (-3.7472) - 3.8870) = -0.27913$$

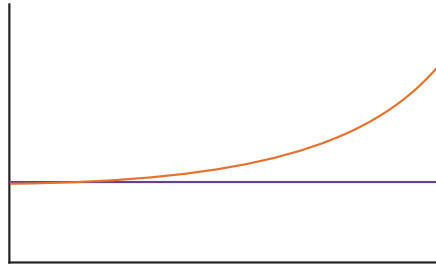
Sedaj nadaljujemo s točko  $y_4 = -0.27913$  in gremo na računanje  $y_5$ . □

## 8.6 Inherentna in inducirana nestabilnost

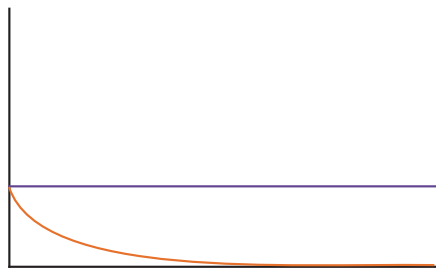
Inherentna nestabilnost je odvisna od problema in neodvisna od numerične metode.

Zgledi so:

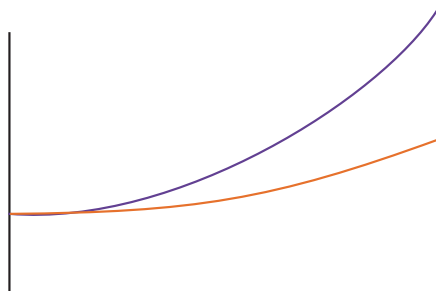
- a)  $y' = y - 1$ ,  $y(0) = 1$ , splošna rešitev je  $y(x) = Ce^x + 1$  in  $C = 0$ , numerično pa dobimo  $C \neq 0$ . Problem je inherentno absolutno in relativno nestabilen.



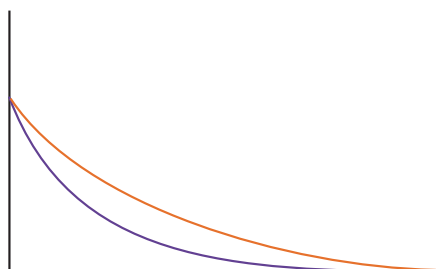
b)  $y' = -y + 1, y(0) = 1$ , splošna rešitev je  $y(x) = Ce^{-x} + 1$  in  $C = 0$ , numerično pa dobimo  $C \neq 0$ . Problem je inherentno absolutno in relativno stabilen.



c)  $y' = y + e^{2x}, y(0) = 1$ , splošna rešitev je  $y(x) = Ce^x + e^{2x}$  in  $C = 0$ , numerično pa dobimo  $C \neq 0$ . Problem je inherentno absolutno nestabilen in relativno stabilen.



d)  $y' = -y - e^{2x}, y(0) = 1$ , splošna rešitev je  $y(x) = Ce^{-x} + e^{-2x}$  in  $C = 0$ , numerično pa dobimo  $C \neq 0$ . Problem je inherentno absolutno stabilen in relativno nestabilen.



Pri relativni nestabilnosti včasih pomaga, če obrnemo interval:

a)  $y' = y - 1, y(x_0) = y_0$ . Vzamemo  $y(x_n) = \gamma$  in računamo nazaj,

b)  $y' = -y - e^{2x}$ ,  $y(x_0) = y_0$ . Vzamemo  $y(x_n) = \gamma$  in računamo nazaj.

Parameter  $\gamma$  moramo izbrati tako, da bo nazaj izračunana rešitev imela v točki  $x_0$  pravo začetno vrednost  $y_0$ . Za pravo izbiro  $\gamma$  lahko uporabimo katerokoli metodo za reševanje nelinearnih enačb, npr. bisekcijo ali sekantno metodo. Če z  $y(x; \gamma)$  označimo numerično rešitev, ki jo dobimo pri začetnem pogoju  $y(x_n) = \gamma$  in računanju nazaj, potem moramo rešiti enačbo  $y(x_0; \gamma) = y_0$ . Podoben postopek uporabimo pri streški metodi za reševanje robnega problema.

*Inducirana nestabilnost* je odvisna od izbire numerične metode. Pojavi se lahko le pri večkoračnih metodah. Za večkoračno metodo oblike

$$\sum_{i=0}^k \alpha_i y_{n-i} + h \sum_{i=0}^k \beta_i f_{n-i} = 0$$

pričakujemo, da je stabilna za enačbo  $y' = 0$  (temu pravimo *ničelna stabilnost*), katere rešitev je konstanta. Ničelna stabilnost pomeni, da tudi če začnemo z nekaj začetnimi vrednostmi, kjer je prisotna napaka, mora potem rešitev, dobljena z večkoračno metodo, ostati omejena.

V primeru  $y' = 0$  se večkoračna metoda spremeni v diferencialno enačbo

$$\sum_{i=0}^k \alpha_i y_{n-i} = 0.$$

Njen karakteristični polinom  $\rho(\zeta)$  ima  $k$  ničel  $\zeta_1, \dots, \zeta_k$ . Splošna rešitev ima pri enostavnih ničlah obliko

$$y_n = \sum_{j=1}^k A_j \zeta_j^n.$$

Za ničelno stabilnost mora tako veljati (Dahlquistov izrek)

- a)  $|\zeta_i| \leq 1$  za  $i = 1, \dots, k$ ,
- b) če je  $|\zeta_j| = 1$ , mora biti  $\zeta_j$  enostavna ničla.

Večkoračna metoda je konsistentna natanko tedaj, ko velja  $\rho(1) = 0$  in  $\rho'(1) + \sigma(1) = 0$ , ta dva pogoja pa sta ekvivalentna temu, da je red vsaj 1. O pogojih za konvergentno večkoračno metodo govori naslednji izrek. Tako kot pri enokoračnih metodah lahko dokaz in ostale podrobnosti najdete v [17].

**Izrek 8.10** *Večkoračna metoda je konvergentna natanko tedaj, ko je ničelno stabilna in konsistentna. ■*

## 8.7 Začetni problemi drugega reda

Rešujemo diferencialno enačbo drugega reda

$$\begin{aligned} y'' &= f(x, y, y') \\ y(x_0) &= y_0 \end{aligned}$$

$$y'(x_0) = y'_0.$$

To lahko prevedemo na sistem enačb prvega reda

$$\begin{aligned} y' &= p, & y(x_0) &= y_0, \\ p' &= f(x, y, p), & p(x_0) &= y'_0. \end{aligned}$$

**Zgled 8.13** Vzemimo začetni problem drugega reda

$$y'' = x + y^2$$

in  $y(0) = 1, y'(0) = 0$ . To prevedemo na sistem dveh enačb prvega reda

$$\begin{aligned} y' &= p, & y(0) &= 1, \\ p' &= x + y^2, & p(0) &= 0. \end{aligned} \quad \square$$

**Zgled 8.14** Problem dveh teles opisuje gibanje telesa zaradi sile težnosti telesa z veliko večjo maso. Če je težje telo v izhodišču, dobimo v kartezičnem koordinatnem sistemu enačbi

$$\begin{aligned} x''(t) &= -x(t)/r(t)^3 \\ y''(t) &= -y(t)/r(t)^3, \end{aligned}$$

kjer je  $r(t) = \sqrt{x(t)^2 + y(t)^2}$ .

To prevedemo na sistem

$$z(t) = \begin{bmatrix} x(t) \\ y(t) \\ x'(t) \\ y'(t) \end{bmatrix}, \quad z'(t) = \begin{bmatrix} z_3(t) \\ z_4(t) \\ -z_1(t)/r(t)^3 \\ -z_2(t)/r(t)^3 \end{bmatrix},$$

kjer je  $r(t) = \sqrt{z_1(t)^2 + z_2(t)^2}$ . □

V posebnem primeru, ko je  $y'' = f(x, y)$ , obstajajo posebne Runge-Kutta ali večkoračne metode, kot je npr. metoda Numerova

$$y_{n+1} - 2y_n + y_{n-1} = \frac{h^2}{12}(f_{n+1} + 10f_n + f_{n-1}) + \mathcal{O}(h^6).$$

## 8.8 Reševanje začetnih diferencialnih enačb v Matlabu

Na voljo imamo več metod. Izmed njih je verjetno najbolj pogosto uporabljena ode45. Kličemo jo v obliki

$$[x, y] = \text{ode45}(\text{fun}, \text{span}, y_0),$$

kjer je fun ime funkcije  $y_{odv} = \text{fun}(x, y)$ , ki iz argumentov  $x$  in  $y$ , kjer je zadnji v primeru sistema vektor, izračuna vrednosti odvoda iz diferencialne enačbe. Seveda je v primeru sistema diferencialnih enačb tudi odvod, ki ga vrne fun, v obliki vektorja. Argument span poda interval  $[a, b]$ , na katerem iščemo rešitev,  $y_0$  pa je začetni pogoj.

Zgled uporabe:

```

function yprime=fun(t,x)
yprime=-y-5*exp(-t)*sin(5*t);

tspan=[0 3]; yzero=1;
[x,y]=ode45('fun',tspan,y0)
plot(x,y,'*--');

```

Vse metode, ki jih uporablja Matlab, so adaptivne. Na voljo so naslednje metode:

- `ode45`: Temelji na Dormand–Princeovi metodi, ki je eksplicitna Runge-Kutta metoda reda 4 in 5 (podobno kot Fehlbergova metoda). To je metoda, ki najpogosteje daje dobre rezultate, zato je priporočljivo problem najprej poskusiti rešiti s to metodo.
- 
- `ode23`: Temelji na eksplicitnem Runge-Kutta pravilu reda 2 in 3. Kadar ne zahtevamo velike natančnosti in v primeru majhne togosti je lahko metoda bolj učinkovita kot pa `ode45`.
  - `ode113`: Gre za Adams–Bashworth–Moultonovo večkoračno prediktor-korektor metodo, kjer vedno naredimo en korak korektorja. Lahko je bolj učinkovita od `ode45` v primerih, ko zahtevamo visoko natančnost in pa kadar je funkcija  $f$  iz diferencialne enačbe še posebej zapletena in želimo imeti čim manj izračunov vrednosti  $f$ .

Prejšnje tri metode so namenjene za netoge sisteme. V primeru, ko ne dajejo dobrih rezultatov imamo lahko opravka s togim sistemom in takrat lahko uporabimo naslednje metode:

- `ode15s`: Večkoračna metoda, ki temelji na metodah za numerično odvajanje. Je sicer manj učinkovita kot `ode45`, a deluje na zmerno togih primerih, ko `ode45` odpove.
- `ode23s`: Uporablja modificirano Rosenbrockovo metodo reda 2. Kadar ne zahtevamo velike natančnosti je lahko bolj učinkovita kot `ode15s` in deluje na nekaterih zelo togih primerih, kjer `ode15s` odpove.
- `ode23t`: Varianta trapezne metode. Uporabna za zmerno toge primere kadar ne zahtevamo velike natančnosti.
- `ode23tb`: Varianta implicitne dvostopenjske Runge-Kutta metode, primerna za zelo toge sisteme.

Vse metode uporabljamo na enak način, zato je dovolj pogledati le uporabo `ode45`.

- Osnovna uporaba je `[x,y] = ode45('f',xab,y0)`, kjer metoda  $f$  določa diferencialno enačbo  $y' = f(x,y)$ ,  $xab = [a b]$  je interval, na katerem rešujemo enačbo in  $y_0$  je začetni pogoj  $y(a) = y_0$ . Kot rezultat dobimo vektorja izračunanih  $x$  in  $y$ .
- Dodatne opcije, s katerimi nastavimo npr. željeno natančnost, podamo v obliki `[x,y] = ode45('f',xab,y0,options)`. Pri tem moramo opcije `options` podati ali prej definirati s pomočjo funkcije `odeset`.

Tako npr. z `options=odeset('RelTol',1e-4,'AbsTol',1e-8)` povečamo relativno in absolutno natančnost (privzeti vrednosti sta  $1e-3$  in  $1e-6$ ).



- Če ima diferencialna enačba  $f$  poleg  $x$  in  $y$  še kakšne dodatne parametre, jih podamo v obliki  $[x, y] = \text{ode45}('f', x_{ab}, y_0, \text{options}, p_1, p_2, \dots)$ . Pri tem lahko kot `options` podamo prazen seznam `[]`, kadar ne želimo spreminjati nastavitvev.
- Z  $[x, y, s] = \text{ode45}('f', x_{ab}, y_0, \text{options})$  dobimo še vektor  $s$  s šestimi komponentami, ki nam povedo, kaj se je dogajalo med računanjem. Tako je  $s_1$  število uspešnih korakov,  $s_2$  število neuspešnih poskusov,  $s_3$  število izračunov  $f$ ,  $s_4$  število izračunov Jacobijeve matrike parcialnih odvodov  $f$  po  $y$ ,  $s_5$  število LU razcepov in  $s_6$  število reševanj linearnih sistemov.

Diferencialne enačbe, kjer poleg  $x$  in  $y$  nastopajo še dodatni parametri, uporabljamo v obliki  $[x, y] = \text{ode45}('f', x_{ab}, y_0, \text{options}, p_1, p_2, \dots)$ . Pri tem moramo paziti na to, da ima funkcija  $f$  po vrsti parametre  $x, y, flag, p_1, p_2, \dots$ . Parameter `flag` se uporablja za to, da diferencialna enačba po potrebi avtomatično poda začetne pogoje in interval, če tega ne uporabljamo pa moramo le pustiti prostor za nek parameter.

Tako npr. s kombinacijo

```
f=inline('x^2+y+p','x','y','flag','p')
[x1,y1]=ode45(f,[0 1],0,[],1)
[x2,y2]=ode45(f,[0 1],0,[],2)
```

rešimo diferencialni enačbi  $y' = x^2 + x + 1$  in  $y' = x^2 + y + 2$ .

Velikokrat rešujemo diferencialno enačbo  $y' = f(x, y)$ ,  $y(x_0) = y_0$ , kjer točen interval  $[x_0, x_n]$  ni znan, saj je izračun točke  $x_n$ , kjer pride do kakšnega dogodka, del same naloge. Tako lahko npr. rešujemo nalogo telesa, ki prosto pada z določene višine, zanima pa nas, v katerem trenutku zadane tla.

Matematično formulirano imamo dano funkcijo  $g(x, y)$  in iščemo rešitev začetnega problema  $y' = f(x, y)$ ,  $y(x_0) = y_0$  in točko  $x^*$ , kjer je  $g(x^*, y(x^*)) = 0$ . V Matlabu to lahko rešimo tako, da v nastavitvah določimo kontrolno funkcijo  $g$ , ki pove, kdaj se je dogodek zgodil.

```
function ydot = f(t,y)
ydot = [y(2); -1+y(2)^2];

function [gstop,isterminal,direction] = g(t,y)
gstop = y(1);      % Dogodek je, ko funkcija g vrne 0
isterminal = 1;    % Metoda se konča, ko pride do dogodka
direction = [];    % Do nič lahko pridemo z obeh strani

opts = odeset('events',@g);
y0 = [1; 0];
[t,y,tfinal] = ode45(@f,[0 Inf],y0,opts);
plot(t,y(:,1),'-',[0 tfinal],[1 0],'o')
```

## 8.9 Implicitno podane diferencialne enačbe

Diferencialna enačba je lahko namesto v eksplicitni obliki  $y' = f(x, y)$  podana v implicitni obliki

$$f(x, y, y') = 0. \quad (8.1)$$

Če se iz enačbe (8.1) ne da direktno izraziti  $y'$  kot funkcija  $x$  in  $y$ , potem lahko sistem numerično rešujemo tako, da v vsakem koraku pri danih vrednostih za  $y$  in  $x$  dobimo  $y'$  iz (8.1) z numeričnim reševanjem nelinearne enačbe.

Reševanje začetnega problema si lahko predstavljamo kot sledenje krivulji, ki se začne pri  $y(x_0) = y_0$ . V primeru implicitno podane diferencialne enačbe moramo poznati tudi vrednost  $y'(x_0)$ , saj ima enačba (8.1) lahko več kot le eno rešitev. Ko poznamo  $x_0, y_0$  in  $y'_0$  lahko naredimo npr. en korak Eulerjeve metode in pridemo do  $x_1 = x_0 + h, y_1 = y_0 + hy'_0$ . V novi točki moramo spet rešiti nelinearno enačbo za  $y'_1$ , pri čemer pa lahko, če  $h$  ni prevelik, za začetni približek vzamemo kar  $y'_0$ .

Podobno velja za nelinearne implicitne sisteme nelinearnih enačb. Princip je enak, le da moramo sedaj namesto nelinearne enačbe reševati nelinearne sisteme, da pridemo do vektorja odvodov.

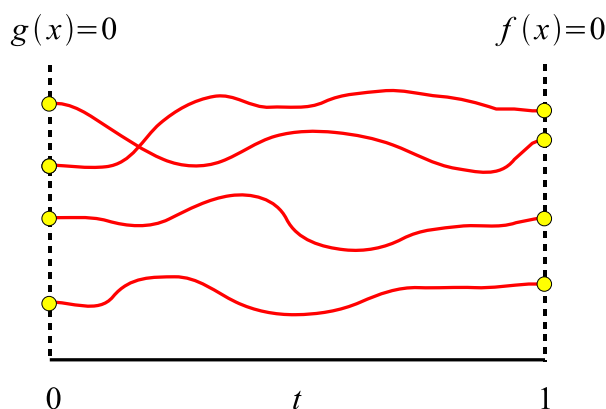
## 8.10 Metoda zveznega nadaljevanja

To je metoda za reševanje nelinearne enačbe  $f(x) = 0$ . Če je težko poiskati začetni približek (posebno, kadar je  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ), si lahko pomagamo z uvedbo dodatnega parametra. Ideja je, da izberemo sorodni enostavnejši problem  $g(x) = 0$ , ki ga znamo rešiti, potem pa z dodatnim parametrom  $t$ , ki teče od 0 do 1, enostavnejši problem zvezno prevedemo na  $f(x) = 0$ . Pri tem pričakujemo, da se bodo z zveznim spreminjanjem problema zvezno premikale tudi rešitve in bomo tako s sledenjem iz znanih rešitev enačbe  $g(x) = 0$  prišli do iskanih rešitev  $f(x) = 0$ .

Opazujemo npr. *konveksno homotopijo*

$$F(t, x) = t \cdot f(x) + (1 - t) \cdot g(x),$$

kjer je  $0 \leq t \leq 1$  in poznamo rešitve  $g(x) = 0$ . S sledenjem krivulji od  $t = 0$  do  $t = 1$  dobimo rešitve  $f(x) = 0$ .



Sledenje krivulje je povezano z reševanjem diferencialne enačbe. Z odvajanjem po  $t$  dobimo

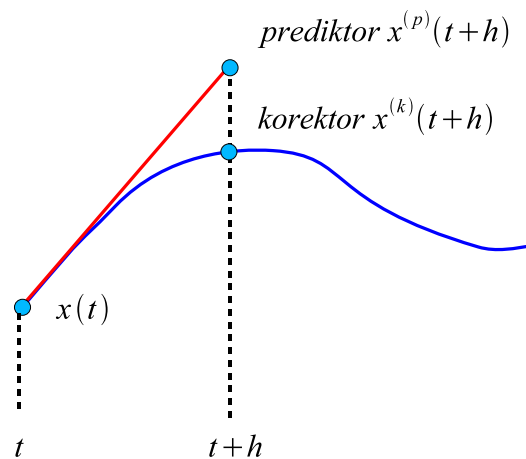
$$F_t(t, x) + F_x(t, x) \cdot \dot{x} = 0,$$

kar nam da navadno diferencialno enačbo

$$\dot{x} = -F_x(t, x)^{-1} \cdot F_t(t, x), \quad x(0) = x_0,$$

kjer je  $g(x_0) = 0$ .

Z metodami za reševanje navadnih diferencialnih enačb lahko sledimo rešitvi, poleg tega pa pri vsaki vrednosti  $t$  lahko upoštevamo še, da mora za  $x(t)$  veljati  $F(t, x(t)) = 0$ . Ponavadi uporabljamo kombinacijo prediktor–korektor, ko najprej uporabimo metodo za reševanje začetnega problema, potem pa iz dobljenega približka izračunamo novo točko na krivulji z kakšno metodo za reševanje nelinearnega sistema.

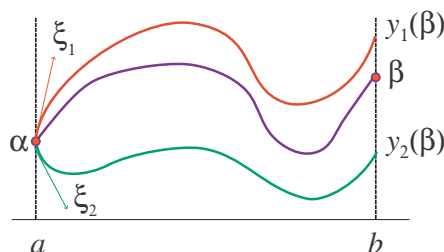


- prediktor: Iz točke  $(t, x(t))$  z eno izmed metod za reševanje začetnega problema (npr. z Eulerjevo metodo) izračunamo prediktor  $x^{(p)}(t+h)$ , ki je približek za rešitev v točki  $t+h$ .
- korektor: Z eno izmed metod za reševanje nelinearnega sistema (npr. z Newtonovo metodo) rešimo sistem  $F(t+h, x(t+h)) = 0$ , za začetni približek pa vzamemo  $(t+h, x^{(p)}(t+h))$ .

## 8.11 Robni problemi drugega reda

Pri robnih problemih imamo opravka z diferencialno enačbo višjega reda, kjer dodatni pogoji, ki določajo pravo vejo, niso podani v isti točki. Tako pri robnem problemu drugega reda rešujemo diferencialno enačbo drugega reda na intervalu  $[a, b]$  z robnima pogojema:

$$\begin{aligned} y'' &= f(x, y, y') \\ y(a) &= \alpha \\ y(b) &= \beta. \end{aligned}$$



Slika 8.6: Rešitev linearnega robnega problema lahko sestavimo kot linearno kombinacijo dveh začetnih problemov.

### 8.11.1 Linearni robni problem

Najpreprostejši je *linearni robni problem drugega reda*, ki ima obliko

$$\begin{aligned} -y''(x) + p(x)y'(x) + q(x)y(x) &= r(x), \\ y(a) = \alpha, \quad y(b) &= \beta. \end{aligned} \quad (8.2)$$

Numerični načini reševanja linearnega robnega problema so:

- a) *Kombinacija dveh začetnih problemov.* Izberemo različna  $\xi_1$  in  $\xi_2$  in rešimo začetna problema, ki ustrezata enačbi (8.2) in

$$\begin{aligned} y_1(a) &= \alpha, \quad y_1'(a) = \xi_1, \\ y_2(a) &= \alpha, \quad y_2'(a) = \xi_2. \end{aligned}$$

Za poljuben  $\lambda$  je  $y = \lambda y_1 + (1 - \lambda)y_2$  tudi rešitev enačbe (8.2) in ustreza pogoju  $y(a) = \alpha$ . Sedaj  $\lambda$  določimo tako, da bo  $y(b) = \beta$ , situacija je predstavljena na sliki 8.6. Veljati mora

$$\lambda y_1(b) + (1 - \lambda)y_2(b) = \beta,$$

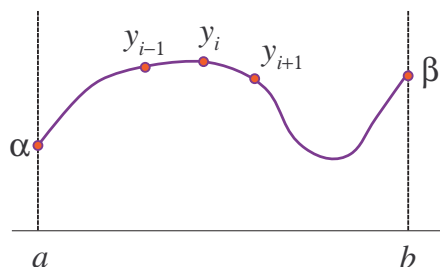
torej

$$\lambda = \frac{\beta - y_2(b)}{y_1(b) - y_2(b)}.$$

Dva začetna problema moramo reševati kot sistem, pa čeprav bi lahko enačbi rešili ločeno drugo od druge. Če jih rešujemo ločeno, se pri adaptivnih metodah lahko namreč zgodi, da rešitvi ne bosta tabelirani v istih točkah, potem pa ne moremo izračunati linearne kombinacije rešitev.

- b) *Diferenčna metoda.* Interval  $[a, b]$  ekvidistantno razdelimo na  $n + 1$  delov s točkami  $x_0 = a, x_1, \dots, x_n, x_{n+1} = b$ , kjer je  $x_i = x_0 + ih$  in  $h = (b - a)/(n + 1)$ . Odvode aproksimiramo s simetričnimi diferencami

$$\begin{aligned} y_i' &= \frac{y_{i+1} - y_{i-1}}{2h} + \mathcal{O}(h^2), \\ y_i'' &= \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} + \mathcal{O}(h^2). \end{aligned}$$



Slika 8.7: Približki pri diferenčni metodi.

Dobimo sistem enačb

$$\frac{-y_{i+1} + 2y_i - y_{i-1}}{h^2} + p_i \frac{y_{i+1} - y_{i-1}}{2h} + q_i y_i = r_i, \quad i = 1, \dots, n,$$

pri čemer je  $y_0 = \alpha$  in  $y_{n+1} = \beta$ . Sistem je linearen in tridiagonalen, zato ga lahko rešimo na preprost način. To nam tudi omogoča, da vzamemo velik  $n$ , kar izboljša natančnost rešitve.

Če v diferencialni enačbi ne nastopa  $y'$ , lahko uporabimo metodo Numerova, saj je potem napaka reda  $\mathcal{O}(h^6)$  namesto  $\mathcal{O}(h^2)$ . V primeru diferencialne enačbe

$$-y''(x) + q(x)y(x) = r(x)$$

tako za  $i = 1, \dots, n$  dobimo enačbo

$$y_{i+1} - 2y_i + y_{i-1} = \frac{h^2}{12} (q_{n+1}y_{n+1} - r_{n+1} + 10q_n y_n - r_n - q_{n-1}y_{n-1} + r_{n-1}).$$

**Zgled 8.15** Z diferenčno metodo in  $h = 1/2$  rešimo robni problem

$$\begin{aligned} (1 + x^2)y'' + 2xy' - x^2y &= 1 \\ y(-1) &= 0 \\ y(1) &= 0. \end{aligned}$$

Dobimo sistem enačb za  $y_1, y_2, y_3$ :

$$\begin{aligned} (1 + (-0.5)^2) \frac{y_2 - 2y_1 + 0}{(0.5)^2} + 2(-0.5) \frac{y_2 - 0}{2(0.5)} - (-0.5)^2 y_1 &= 1 \\ (1 + (0.0)^2) \frac{y_3 - 2y_2 + y_1}{(0.5)^2} + 2(0.0) \frac{y_3 - y_1}{2(0.5)} - (0.0)^2 y_2 &= 1 \\ (1 + (0.5)^2) \frac{0 - 2y_3 + y_2}{(0.5)^2} + 2(0.5) \frac{0 - y_3}{2(0.5)} - (0.5)^2 y_3 &= 1, \end{aligned}$$

ki ima rešitev

$$y_1 = -0.24, y_2 = -0.365, y_3 = -0.24. \quad \square$$

Diferenčno metodo lahko uporabljamo tudi pri robnih pogojih, v katerih nastopajo odvodi. Če imamo splošne robne pogoje oblike

$$\alpha_1 f(a) + \beta_1 f'(a) = 0, \quad \text{in} \quad \alpha_2 f(b) + \beta_2 f'(b) = 0,$$

si pomagamo s t.i. navideznimi točkami.

Npr., če imamo v robnem pogoju odvod v točki  $a$ , potem ta odvod aproksimiramo s simetrično diferenco

$$y'_0 = \frac{y_1 - y_{-1}}{2h}.$$

Nova neznanka  $y_{-1}$  je le navidezna. Ker mora tudi v točki  $a$  rešitev zadoščati diferencialni enačbi, dobimo še enačbo

$$\frac{-y_1 + 2y_0 - y_{-1}}{h^2} - p_0 \frac{y_1 - y_{-1}}{2h} + q_0 y_0 = r_0,$$

iz teh dveh enačb pa lahko izločimo  $y_{-1}$ .

Če podvojimo število točk, pričakujemo, da bomo dobili boljše približke. Ker pa za napako velja, da je reda  $h^2$ , lahko z uporabo ekstrapolacije, podobne kot pri Rombergovi metodi, dobimo še boljše približke v tistih točkah, ki nastopajo tako v grobi kot v fini delitvi.

### 8.11.2 Nelinearni robni problem

Pri *nelinearnem robnem problemu drugega reda* je  $f$  nelinearna funkcija  $y$  in  $y'$ . Zaradi nelinearnosti ne moremo uporabiti kombinacije dveh začetnih problemov. Uporabimo lahko diferenčno metodo, a dobimo nelinearni sistem z veliko neznankami.

Nelinearne enačbe imajo obliko

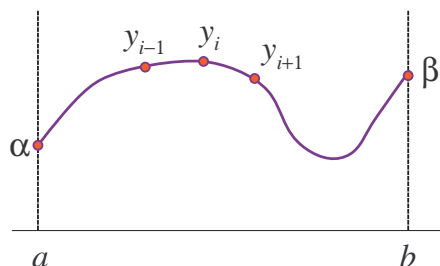
$$\frac{y_{i-1} - 2y_i + y_{i+1}}{h^2} = f\left(x, y_i, \frac{y_{i+1} - y_{i-1}}{2h}\right), \quad i = 1, \dots, n.$$

Če nelinearni sistem rešujemo z Newtonovo metodo, je Jacobijeva matrika tridiagonalna, tako da se da sistem reševati dokaj ekonomično. Za začetni približek lahko npr. vzamemo kar točke na daljci med  $(a, \alpha)$  in  $(b, \beta)$ , torej  $y_i = \alpha + i(\beta - \alpha)/n$ .

Na voljo imamo tudi *strelsko metodo*. Pri strelski metodi rešujemo začetni problem

$$\begin{aligned} y'' &= f(x, y, y'), \\ y(a) &= \alpha, \\ y'(a) &= \zeta. \end{aligned} \tag{8.3}$$

Rešitev je odvisna od parametra  $\zeta$ , ki ga je potrebno izbrati tako, da bo  $y(b; \zeta) = \beta$ .



Če definiramo  $E(\xi) := y(b; \xi) - \beta$ , potem iščemo tak  $\xi$ , da bo  $E(\xi) = 0$ . Uporabimo lahko katerokoli metodo za reševanje nelinearne enačbe. Če npr. želimo uporabiti tangентno metodo, potem potrebujemo tudi  $E'(\xi) = \frac{\partial y(b; \xi)}{\partial \xi}$ . Definiramo  $z := \frac{\partial y}{\partial \xi}$  in z odvajanjem

$$y'' = f(x, y, y'), \quad y(a) = \alpha, \quad y'(a) = \xi$$

dobimo

$$z'' = f_y(x, y, y')z + f_{y'}(x, y, y')z', \quad z(a) = 0, \quad z'(a) = 1. \quad (8.4)$$

Če rešujemo sistem (8.3) in (8.4), potem dobimo tudi vrednost odvoda funkcije  $E$ .

### 8.11.3 Prevedba na variacijski problem

Reševanje robnega problema

$$-\frac{d}{dx} \left( p(x) \frac{dy}{dx} \right) + q(x)y = f(x), \quad 0 \leq x \leq b,$$

z robnima pogoje  $y(0) = 0$  in  $y'(b) + \sigma y(b) = 0$ , kjer je  $p(x) \geq \delta > 0$ ,  $q(x) \geq 0$  in  $\sigma > 0$ , je ekvivalentno iskanju funkcije  $y \in C^2[0, b]$ , ki zadošča robnima pogoje in minimizira funkcijo

$$F(y) = \int_0^b (p(x)y'^2(x) + q(x)y^2(x) - 2f(x)y(x)) dx + p(b)\sigma y^2(b).$$

Ideja *Rayleigh–Ritzove metode* je, da namesto  $y \in C^2[0, b]$  iščemo približek oblike

$$y(x) = \sum_{i=1}^n c_i \phi_i(x),$$

kjer so  $\phi_1, \dots, \phi_n$  t.i. *bazne funkcije*.

V točki, kjer je dosežen minimum, morajo biti vsi parcialni odvodi  $F$  po  $c_i$  enaki 0. Tako dobimo linearni sistem  $Ac = b$  za koeficiente  $c = [c_1 \ \dots \ c_n]^T$ , kjer je

$$a_{ij} = \int_0^b (p(x)\phi_i'(x)\phi_j'(x) + q(x)\phi_i(x)\phi_j(x)) dx$$

in

$$b_i = \int_0^b \frac{\phi_i(x)}{f(x)} dx$$

za  $i, j = 1, \dots, n$ . Od izbire baznih funkcij je odvisno, ali bo sistem rešljiv in kako dober bo dobljen približek.

*Metoda končnih elementov* je poseben primer Rayleigh-Ritzove metoda, kjer za bazne funkcije izberemo t.i. *piramidne funkcije*

$$\phi_i(x) = \begin{cases} 0, & 0 \leq x \leq x_{i-1}, \\ (x - x_{i-1})/h_{i-1}, & x_{i-1} < x \leq x_i, \\ (x_{i+1} - x)/h_i, & x_i < x \leq x_{i+1}, \\ 0, & x_{i+1} < x \leq b, \end{cases}$$

kjer je  $0 = x_0 < x_1 < x_2 < \dots < x_n = b$  in  $x_i - x_{i-1} = h_{i-1}$ .

Funkcija  $\phi_i$ , kjer je  $i = 1, \dots, n-1$ , ima nosilec na intervalu  $(x_{i-1}, x_{i+1})$ , izjema je le funkcija  $\phi_n$  ki ima neničelni del samo na intervalu  $(x_{n-1}, x_n]$ . Zaradi tega je matrika  $A$  sedaj tridiagonalna. Elementi sistema  $Ax = b$  so

$$\begin{aligned} a_{ii} &= \int_{x_{i-1}}^{x_i} \frac{1}{h_{i-1}^2} p(x) dx + \int_{x_i}^{x_{i+1}} \frac{1}{h_i^2} p(x) dx \\ &\quad + \int_{x_{i-1}}^{x_i} \frac{1}{h_{i-1}^2} (x - x_{i-1})^2 q(x) dx + \int_{x_i}^{x_{i+1}} \frac{1}{h_i^2} (x_{i+1} - x)^2 q(x) dx, \quad i = 1, \dots, n-1, \\ a_{i,i+1} &= \int_{x_i}^{x_{i+1}} \frac{-1}{h_i^2} p(x) dx + \int_{x_i}^{x_{i+1}} \frac{1}{h_i^2} (x_{i+1} - x)(x - x_i) q(x) dx, \quad i = 1, \dots, n-1, \\ a_{i-1,i} &= \int_{x_{i-1}}^{x_i} \frac{-1}{h_{i-1}^2} p(x) dx + \int_{x_i}^{x_{i+1}} \frac{1}{h_i^2} (x_{i+1} - x)(x - x_{i-1}) q(x) dx, \quad i = 2, \dots, n, \\ b_i &= \int_{x_{i-1}}^{x_i} \frac{1}{h_{i-1}} (x - x_{i-1}) f(x) dx + \int_{x_i}^{x_{i+1}} \frac{1}{h_i} (x_{i+1} - x) f(x) dx, \quad i = 1, \dots, n-1, \\ a_{nn} &= \int_{x_{n-1}}^{x_n} \frac{1}{h_{n-1}^2} p(x) dx + \int_{x_{n-1}}^{x_n} \frac{1}{h_{n-1}^2} (x - x_{n-1})^2 q(x) dx + p(b)\sigma, \\ b_n &= \int_{x_{n-1}}^{x_n} \frac{1}{h_{n-1}} (x - x_{n-1}) f(x) dx. \end{aligned}$$

## 8.12 Reševanje robnih problemov v Matlabu

V Matlabu imamo za reševanje robnih problemov na voljo metodo `bvp4c`, ki uporablja kolokacijsko metodo. S to metodo lahko rešujemo robne probleme, ki jih zapišemo v obliki

$$\frac{d}{dx} y(x) = f(x, y(x), p), \quad g(y(a), y(b), p) = 0.$$

Tu je  $y(x)$  vektor,  $f$  je podana funkcija  $x$  in  $y$ , ki opisuje diferencialno enačbo. Vektor  $p$  je neobvezen in opisuje neznane parametre, ki jih iščemo. Rešitev iščemo na intervalu  $[a, b]$ , s funkcijo  $g$  pa podamo robne pogoje. Metoda potrebuje tudi začetni približek za rešitev  $y(x)$ .

Oblika za klic `bvp4c` je `sol=bvp4c(@odefun,@bcfun,solinit,options)`, pri čemer:

- funkcija `odefun` poda diferencialno enačbo, oblika je `yodv = odefun(x,y)`;
- funkcija `bcfun` poda robne pogoje (vrne ostanek, ki mora biti pri izpolnjenih pogojih enak 0), oblika je `res = bcfun(ya,yb)`;
- `solinit` je struktura, s katero podamo začetni približek, strukturo zgradimo s pomožno funkcijo `bvpinit` kot `solinit = bvpinit(linspace(a,b,n),@initf)`, kjer je `initf` funkcija oblike `y = initf(x)`, s katero podamo robne pogoje;
- `options` je neobvezen argument, kjer lahko nastavimo delovanje metode, za nastavitve uporabimo metodo `bvpset`.

**Zgled 8.16** Presek oblike kapljice vode na ravni podlagi je podan z rešitvijo robnega problema

$$u''(x) + (1 - u(x))(1 + u'(x)^2)^{3/2} = 0, \quad u(-1) = 0, \quad u(1) = 0.$$



Za uporabo `bvp4c` to spremenimo v sistem enačb prvega reda:

$$\begin{aligned}y_1'(x) &= y_2(x), \\y_2'(x) &= (y_1(x) - 1)(1 + y_2(x)^2)^{3/2},\end{aligned}$$

Za začetni približek vzamemo  $y_1(x) = \sqrt{1 - x^2}$  in  $y_2(x) = -x / (0.1 + \sqrt{1 - x^2})$ .

Ustrezna koda za rešitev problema v Matlabu je:

```
function yodv = kaplja(x,y)
yodv = [ y(2); (y(1)-1)*((1+y(2)^2)^(3/2)) ];

function res = kapljarp(ya,yb)
res = [ ya(1); yb(1) ];

function y = kapljainit(x)
y = [ sqrt(1-x^2); -x/(0.1+sqrt(1+x^2)) ];

resinit = bvpinit(linspace(-1,1,20), @kapljainit);
res = bvp4c(@kaplja, @kapljarp, resinit);
plot(res.x, res.y(1,:))
```

□

**Zgled 8.17** Iščemo lastno vrednost  $\lambda$ , pri kateri ima robni problem

$$y''(x) + (\lambda - 10 \cos(2x))y(x) = 0, \quad y'(\pi) = 0, \quad y'(0) = 0$$

neničelno rešitev. Ker lahko vsako tako rešitev pomnožimo s poljubnim skalarjem, potrebujemo še en pogoj, npr.  $y(0) = 1$ . Podobno kot v prejšnjem zgledu enačbo spremenimo v sistem enačb prvega reda:

$$y_1'(x) = y_2(x), \quad y_2'(x) = (10 \cos(2x) - \lambda)y_1(x).$$

Za začetni približek vzamemo  $y(x) = \cos(4x)$  in  $\lambda = 15$ .

Ustrezna koda za rešitev problema v Matlabu je:

```
function yodv = mth(x,y,lambda)
yodv = [ y(2); (10*cos(2*x)- lambda)*y(1) ];

function res = mthrp(ya,yb,lambda)
res = [ ya(2); yb(2); ya(1)-1 ];

function y = mthinit(x)
y = [ cos(4*x); -4*sin(4*x) ];

resinit = bvpinit(linspace(0,pi,10), @mthinit, 15);
res = bvp4c(@mth, @mthrp, resinit);
fprintf('Lastna vrednost je %7.3f.\n', res.parameters)
plot(res.x, res.y(1,:))
```

□

**Zgled 8.18** Naslednje enačbe opisujejo pretok v dolgem navpičnem kanalu, kjer je tekočina vbrizgana skozi eno izmed stranic:

$$\begin{aligned} f''' - R[(f')^2 - ff''] + RA &= 0, \\ h'' + Rfh' + 1 &= 0, \\ \theta'' + Pf\theta' &= 0, \end{aligned}$$

pri čemer je  $R$  Raynoldsovo število in  $P = 0.7R$ . Ker parameter  $A$  spada med neznanke, imamo 8 robnih pogojev:

$$\begin{aligned} f(0) &= f'(0) = 0, & f(1) &= 1, & f'(1) &= 0, \\ h(0) &= h(1) = 0, \\ \theta(0) &= 0, & \theta(1) &= 1. \end{aligned}$$

Problem bi radi rešili pri  $R = 10000$ , vendar nimamo nobenih dobrih začetnih približkov. Zato uporabimo metodo zveznega nadaljevanja. Najprej problem rešimo za  $R = 100$ , potem to uporabimo kot začetni približek za  $R = 1000$ , tega pa za  $R = 10000$ .

Ustrezna koda za rešitev problema v Matlabu je:

```
R = 100; resinit = bvpinit(linspace(0,1,10), ones(7,1), 1);
res1 = bvp4c(@pretok, @pretokbc, resinit);
fprintf('Parameter A pri R=100 : %7.4f\n',res1.parameters)

R = 1000; res2 = bvp4c(@pretok, @pretokbc, res1);
fprintf('Parameter A pri R=1000 : %7.4f\n',res2.parameters)

R = 10000; res3 = bvp4c(@pretok, @pretokbc, res2);
fprintf('Parameter A pri R=10000 : %7.4f\n',res3.parameters)
plot(res1.x, res1.y(2,:), res2.x, res2.y(2,:), res3.x, res3.y(2,:));

function yodv = pretok(x,y,A)
    P = 0.7*R;
    yodv = [ y(2); y(3); R*(y(2)^2 - y(1)*y(3) - A);
            y(5); -R*y(1)*y(5) - 1; y(7); -P*y(1)*y(7) ];
end

function res = pretokbc(ya,yb,A)
    res = [ya(1); ya(2); yb(1) - 1; yb(2); ya(4); yb(4); ya(6); yb(6) - 1];
end
```

□

---

## Dodatna literatura

Več o metodi zveznega nadaljevanja lahko najdete v [1].

# Literatura

- [1] E. L. Allgower, K. Georg, *Numerical Continuation Methods: An Introduction*, Springer-Verlag, New York, 1990.
- [2] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, H. van der Vorst, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*, SIAM, Philadelphia, 1994.
- [3] Z. Bohte, *Numerične metode, 4. ponatis*, DMFA Slovenije, Ljubljana, 1991.
- [4] Z. Bohte, *Numerično reševanje nelinearnih enačb*, DMFA Slovenije, Ljubljana, 1993.
- [5] Z. Bohte, *Numerično reševanje sistemov linearnih enačb*, DMFA Slovenije, Ljubljana, 1994.
- [6] R. L. Burden, J. D. Faires: *Numerical Analysis*, Brooks/Cole, Pacific Grove, 1997.
- [7] B. N. Datta, *Numerical Linear Algebra and Applications*, Brooks/Cole, Pacific Grove, 1995.
- [8] J. W. Demmel, *Uporabna numerična linearna algebra*, DMFA Slovenije, Ljubljana, 2000.
- [9] G. H. Golub, C. F. van Loan, *Matrix Computations, 3rd edition*, The Johns Hopkins University Press, Baltimore, 1996.
- [10] G. H. Golub, F. Uhlig, *The QR algorithm: 50 years later its genesis by John Francis and Vera Kublanovskaya and subsequent developments*, IMA Journal of Numerical Analysis **29**, 2009, str. 467—485.
- [11] P. C. Hansen, *Regularization Tools. A Matlab Package for Analysis and Solution of Discrete Ill-Posed Problems*, 2008, <http://www2.imm.dtu.dk/~pch/Regutools/RTv4manual.pdf>
- [12] M. T. Heath, *Scientific Computing: An Introductory Survey, 2nd edition*, McGraw-Hill, 2002.
- [13] D. J. Higham, N. J. Higham, *Matlab Guide, Second Edition*, SIAM, Philadelphia, 2005.
- [14] N. J. Higham, *Accuracy and stability of numerical algorithms*, SIAM, Philadelphia, 1996.
- [15] C. T. Kelley, *Iterative Methods for Linear and Nonlinear Equations*, SIAM, Philadelphia, 1995.
- [16] D. Kincaid, W. Cheney, *Numerical Analysis*, Brooks/Cole, Pacific Grove, 1996.
- [17] J. Kozak, *Numerična analiza*, DMFA Slovenije, Ljubljana, 2008.
- [18] C.-Y. Lam, *Applied Numerical Methods For Partial Differential Equations*, Prentice Hall, Singapore, 1994.

- [19] K. W. Morton, D. F. Mayers, *Numerical Solution of Partial Differential Equations*, Cambridge University Press, Cambridge, 2002
- [20] T. Papler, *Stabilno računanje lastnih vektorjev simetrične tridiagonalne matrike v  $\mathcal{O}(n^2)$* , diplomsko delo, Ljubljana, 2005.
- [21] Y. Saad, *Iterative Methods for Sparse Linear Systems, Second Edition*, SIAM, Philadelphia, 2003.
- [22] G. D. Smith, *Numerical Solution of Partial Differential Equations: Finite Difference Methods, 3rd edition*, Clarendon Press, Oxford, 1985.
- [23] F. Tisseur, K. Meerbergen: *The quadratic eigenvalue problem*, SIAM Review **43** (2001), 235–286.
- [24] S. Van Huffel, J. Vandewalle, *The Total Least Square Problem. Computational Aspects and Analysis*, SIAM, Philadelphia, 1991.
- [25] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford University Press, Oxford, 1965.
- [26] T. J. Ypma, *Historical development of the Newton-Raphson method*, SIAM Review **37**, 1995, str. 531–555.
- [27] E. Zakrajšek, *Uvod v numerične metode, druga izdaja*, DMFA Slovenije, Ljubljana, 2000.