

Úvod do numerické matematiky.

Přednáška pro posluchače informatiky.

Zimní resp. Letní semestr 2/2

Ivo Marek, Petr Mayer a Bohuslav Sekerka

1 Úvodní poznámky.

Vymezení problematiky vystihuje následující charakteristika.

Numerická (výpočtová) matematika je realizace matematických modelů na výpočetní technice.

Numerická matematika je tedy součástí matematiky. Je to však součást mající oproti teoretické matematice jeden významný nedostatek. Ten spočívá v tom, že čistě teoretické matematické modely podstatným způsobem operují s reálnými čísly zatímco numerická matematika je závislá na výpočetní technice a tudíž se musí bez tohoto prostředku obejít. Typickým projevem této skutečnosti jsou zaokrouhlovací chyby, což vede ve svých důsledcích k *numerické nestabilitě*.

Zdroj numerické nestability však nemusí být výhradně chyby ze zaokrouhlení. Numerická nestabilita může být zapřičiněna též vlastnostmi matematických modelů samých, např. realizace matematického špatně podmíněného modelu na počítači je přirozeným zdrojem numerické nestability. Zde použitý pojem *špatně podmíněný model* je charakterizován nespojitou závislostí výstupní informace na informacích vstupních.

Příklad 1.1 Budě

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

kde a_{jk} jsou reálná čísla.

Symbolom $h(A)$ označme hodnotu matice A . Snadno si uvědomíme, že h není spojitou funkcí prvků a_{jk} .

Abychom si tuto okolnost uvědomili, stačí položiti

$$A_0 = \begin{bmatrix} a_{11} & 0 \\ 0 & 0 \end{bmatrix}, \quad a_{11} \neq 0,$$

takže na jedné straně

$$h(A_0) = 1,$$

avšak, na straně druhé, pro libovolné reálné $b \neq 0$

$$h(A(b)) = 2,$$

kde

$$A(b) = \begin{bmatrix} a_{11} & 0 \\ 0 & b \end{bmatrix}.$$

Předchozí příklad ukazuje jednak nespojitost hodnosti matice jakožto funkce jejích prvků jednak další pozoruhodnou vlastnost hodnosti. Ta je obsahem následující věty.

Věta 1.1 *Hodnost matice je zdola polospojitou funkcí jejích prvků.*

Definice 1.1 Skalární reálná funkce f se nazývá *zdola polospojitou v bodě s_0* , jestliže pro každé $\varepsilon > 0$ množina

$$\{s : f(s) > f(s_0)\}$$

je otevřená, t. j. existuje $\delta > 0$ takové, že

$$f(s) - f(s_0) \geq \varepsilon$$

pro $s \in (s_0 - \delta, s_0 + \delta)$.

Dalším příkladem uvedeme numerickou nestabilitu poplatnou chybám ze zaokrouhlení.

Příklad 1.2 Určeme hodnoty integrálů

$$(1.1) \quad I_n = \frac{1}{e} \int_0^1 t^n e^t dt$$

pro $n = 1, 2, \dots$

Užitím metody "per partes" zjistíme, že

$$(1.2) \quad I_n = 1 - nI_{n-1}$$

při čemž

$$(1.3) \quad I_0 = \frac{e - 1}{e}.$$

Zřejmě platí

$$(1.4) \quad 0 \leq I_n \leq I_{n-1} \leq 1$$

jakož i

$$(1.5) \quad \lim_{n \rightarrow \infty} I_n = 0.$$

Výpočet na počítači se realizuje nikoliv v termínech zavedených veličin I_n avšak v termínech veličin počítaných \tilde{I}_n .

Vzhledem k zaokrouhlovacím chybám zjistíme, že existují indexy n_1 a n_2 tak, že

$$(1.6) \quad \tilde{I}_{n_1} < 0$$

a

$$(1.7) \quad \tilde{I}_{n_2} > 1.$$

Vidíme tedy, že vztahy (1.4) neplatí pro veličiny \tilde{I}_n .

Důvodem pro platnost vztahů (1.6) a (1.7) je skutečnost vyjádřená rovnostmi

$$\tilde{I}_n = I_n + \delta_n, \quad n = 1, 2, \dots$$

z nichž plyne, že

$$\delta_n = n\delta_{n-1},$$

odkud je patrné, že chyba se na každém kroku výpočtu násobí indexem kroku. Toto zesílení chyby je velikosti faktoriálu, což značí, že ”destrukce” se explicitně projeví po p krocích, kde p značí počet cifer používaných v průběhu výpočtu.

2 Diferenční rovnice

Připomeňme si pojem *obyčejné diferenciální rovnice*.

Buď $f = f(s, y, z)$ skalární reálná funkce, při čemž

$$(2.1) \quad s \in I = (a, b), \quad y \in S_1, \quad z \in S_2,$$

kde

$$S_1 \subset R, \quad S_2 \subset R,$$

při standartním značení množiny reálných čísel symbolem R .

Úlohu nalézti funkci $x = x(s)$, $s \in I$ takovou, že jednak

$$(2.2) \quad x(s) \in S_1, \quad x'(s) \in S_2$$

a jednak

$$(2.3) \quad f(s, x(s), x'(s)) = 0,$$

se nazývá *obyčejnou diferenciální rovnici 1. řádu*.

Příklad 2.1 Položme

$$f(s, y, z) = z - \alpha y.$$

Snadno se přesvědčíme, že

$$x(s) = C \exp\{\alpha s\}$$

je pro $s \in [0, T]$, kde $T > 0$, řešením diferenciální rovnice

$$x' = \alpha x$$

splňující počáteční podmínu

$$x(0) = C.$$

Příklad 2.2 Nechť $I = (-\infty, +\infty)$, $S_0 = S_1 = [-1, 1]$ a

$$f(s, y, z) = z^2 - (1 - y^2).$$

Pak máme co do činění s diferenciální rovnicí

$$x'^2 = 1 - x^2,$$

jejíž řešení jsou dána formulí

$$x(s) = \sin(s - a),$$

kde $a \in R$ je libovolné.

Dále si připomeňme pojem diferenciální rovnice N-tého řádu.

Budě

$$f = f(s, y_0, y_1, \dots, y_N)$$

skalárni reálná funkce definovaná pro $s \in I \subset R$ a y_0, \dots, y_N ležící v podmnožinách reálných čísel S_0, \dots, S_N . Problém nalézt funkci $x = x(s)$ definovanou pro $s \in I$, mající N derivací na I a splňující relace

$$(2.4) \quad x^{(k)}(s) \in S_k \quad k = 0, 1, \dots, N$$

a

$$(2.5) \quad f(s, x(s), x'(s), \dots, x^{(N)}(s)) = 0$$

pro všechna $s \in I$, se nazývá *obyčejnou diferenciální rovnici řádu N*; symbolicky píšeme

$$f(s, x(s), x'(s), \dots, x^{(N)}) = 0.$$

Příklad 2.3 Nechť $N = 2$ a $I = S_0 = S_2 = R$. Nechť dále

$$f(s, y_0, y_1, y_2) = y_0 + y_2.$$

Potom každá funkce x tvaru

$$x(s) = A \cos s + B \sin s,$$

kde A a B jsou konstanty, je řešením odpovídající diferenciální rovnice

$$x'' + x = 0.$$

Nyní přistoupíme k úlohám podobným diferenciálním rovnicím. Analogie bude zřejmá z definice. Budeme se zabývat diferenčními rovnicemi.

Uveďme napřed pojem *diferenční rovnice 1. řádu*.

Budiž dána soustava

$$f = \{f_n(y, z)\},$$

kde $f_n = f_n(y, z)$, $n = 1, 2, \dots$, jsou skalárni reálné funkce definované na množině $I \subset Z_\infty$, při čemž Z_∞ značí množinu celých čísel a $y, z \in S \subset R$. Úloha nalézt posloupnost $\{x_n\} \subset R$, $n \in I$, splňující

$$(2.6) \quad x_n \in S, \quad x_{n-1} \in S$$

a

$$(2.7) \quad f_n(x_n, x_{n-1}) = 0$$

se nazývá *diferenční rovnice 1. řádu*. Každá posloupnost $\{x_n\}$ splňující (2.6) a (2.7) se nazývá *řešením diferenční rovnice* (2.7).

Příklad 2.4 Buď I množina celých čísel. Potom

$$f_n(y, z) = y - z - 1$$

vede k diferenční rovnici

$$x_n - x_{n-1} = 1,$$

Její řešení mají tvar

$$x_n = n + c,$$

kde $c \in R$ je libovolné.

Podobně

Příklad 2.5 Nechť I je množina nezáporných celých čísel a

$$f_n(y, z) = y - z - n,$$

takže naše diferenční rovnice má tvar

$$x_n - x_{n-1} = n.$$

Jejím řešením je posloupnost $\{x_n\}$, kde

$$x_n = \frac{n(n-1)}{2} .$$

Příklad 2.6 Buď I množina všech celých čísel a

$$f_n(y, z) = y - qz,$$

Není obtížné ukázat, že řešení splňující $x_0 = 1$ odpovídající rovnice

$$x_n = qx_{n-1}$$

má tvar

$$x_n = q^n.$$

Nyní přistoupíme k vyšetřování *diferenčních rovnic* řádu $N \geq 1$.

Buď

$$f = \{f_n(y_0, y_1, \dots, y_N)\}$$

posloupnost funkcí definovaných na množině $I \subset Z_\infty$, přičemž $y_j \in S_j \subset R$, $j = 0, 1, \dots, N$.

Úlohu nalézt posloupnost $\{x_n\}$, $n \in I$ a splňující následují požadavky

$$(2.8) \quad x_n \in S_n, \quad x_{n-1} \in S_{n-1}, \quad \dots x_{n-N} \in S_{n-N}$$

a

$$(2.9) \quad f_n(x_n, x_{n-1}, \dots, x_{n-N}) = 0.$$

Posloupnost $\{x_n\}$ splňující (2.8) a (2.9) se nazývá *řešení diferenční rovnice* (2.9).

Příklad 2.7 Nechť $I = Z_\infty$ a

$$f_n(y_0, y_1, y_2) = y_0 - 2\cos\phi y_1 + y_2,$$

kde $\phi \in R$. Řešení takto vzniklé diferenční rovnice

$$x_n - 2\cos\phi x_{n-1} + x_{n-2} = 0$$

je dáno formulí

$$x_n = \cos n\phi.$$

Podobně jako v teorii diferenciálních rovnic se pro některé třídy úloh jednoznačné řešitelnosti dociluje vhodnou volbou počátečních podmínek.

Důležitým případem diferenčních rovnic jsou *diferenční rovnice lineární*. V tom případě

$$(2.10) \quad f_n(y_0, y_1, \dots, y_N) = a_{0n}y_0 + a_{1n}y_1 + \dots + A_{Nn}y_N + b_n,$$

kde $a_{jn}, b_n \in R$, $j = 0, 1, \dots, N$, $n = 0, 1, \dots$

Příklad 2.8 Všechny diferenční rovnice uvedené v příkladech 2.3 – 2.7 jsou lineární.

Podobně diferenční rovnice

$$x_n + 5nx_{n-1} + n^2x_{n-2} = 2$$

je lineární, zatímco

$$x_{n-1} - 2x_{n-1}^2 = 0$$

je nelineární.

Analogicky jako v teorii diferenciálních rovnic, obecné řešení diferenčních rovnic umíme plně charakterizovat pro případ diferenčních rovnic lineárních.

Zabývejme se úlohou (analytického) sestrojení řešení diferenčních rovnic 1. řádu. Tedy, řešme rovnici

$$(2.11) \quad x_n = a_n x_{n-1} + b_n, \quad n \in I \subset Z_\infty^+, \quad a_n \neq 0,$$

kde Z_∞^+ značí množinu všech nezáporných celých čísel.

Vyšetřujme nejprve *homogenní* diferenční rovnici

$$(2.12) \quad x_n = a_n x_{n-1}.$$

Snadno zjistíme, že řešení splňující podmínsku

$$(2.13) \quad x_0 = c$$

je dáno výrazem

$$(2.14) \quad x_n = c\pi_n,$$

kde

$$(2.15) \quad \pi_0 = 1, \quad \pi_n = \prod_{k=1}^n a_k, \quad n = 1, 2, \dots$$

K určení řešení *nehomogenní* rovnice (2.12) splňující $x_0 = c$, použijeme metody známé z teorie lineárních diferenciálních rovnic pod názvem *metoda variace konstanty*.

Položme

$$(2.16) \quad x_n = c_n \pi_n,$$

kde $\{c_n\}$ podléhá určení.

Z rovností

$$x_0 = c_0 \pi_0 = c_0,$$

vidíme, že

$$c_0 = c.$$

Po dosazení (2.16) do (2.11) zjistíme, že

$$c_n \pi_n = a_n c_{n-1} \pi_{n-1} + b_n = c_{n-1} \pi_n + b_n.$$

Z předpokladu $a_n \neq 0$ plyne, že též $\pi_n \neq 0$. Předchozí vztahy implikují rovnost

$$c_n = c_{n-1} + \frac{b_n}{\pi_n}$$

a dále pak rovnosti

$$c_n = c_0 + \sum_{k=1}^n (c_k - c_{k-1}) = c + \sum_{k=1}^n \frac{b_k}{\pi_k}.$$

Výsledek shrneme ve tvaru věty.

Věta 2.1 Nechť $a_n \neq 0$, $n = 1, 2, \dots$. Řešení diferenční rovnice

$$x_n = a_n x_{n-1} + b_n,$$

splňující $x_0 = c$, je dáno výrazem

$$(2.17) \quad x_n = \pi_n \left(c + \sum_{k=1}^n \frac{b_k}{\pi_k} \right), \quad n = 0, 1, \dots$$

při čemž

$$\pi_0 = 1, \quad \pi_n = \prod_{k=1}^n a_k.$$

Příklad 2.9 Buď I množina nezáporných čísel,

$$f_n(y_0, y_1, y_2) = y_0 - y_1 - y_2.$$

Odpovídající diferenční rovnice

$$(2.18) \quad x_n - x_{n-1} - x_{n-2} = 0$$

definuje *Fibonacciova čísla* ($x_0 = 0$, $x_1 = 1$).

Podobně jako v problematice diferenciálních rovnic je jedním z možných způsobů řešení lineárních rovnic řádu N jejich převod na soustavy lineárních rovnic 1. řádu.

Ukážeme si zmíněný postup na příkladě rovnice (2.18).

Nechť

$$X_n = \begin{pmatrix} x_n \\ x_{n-1} \end{pmatrix}, \quad X_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad n = 1, 2, \dots$$

a

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

Snadno nahlédneme, že rovnici (2.18) lze vyjádřiti ve tvaru

$$(2.19) \quad X_n = AX_{n-1}, \quad X_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Podobně jako v případě lineární diferenční rovnice 1. řádu, řešení rovnice (2.19) má tvar

$$X_n = A^n X_0.$$

Odtud plyne platnost zajímavého vztahu

$$\lim_{n \rightarrow \infty} \frac{x_n}{x_{n-1}} = \frac{1 + \sqrt{5}}{2}.$$

Vidíme tedy, že limitní poměr veličin x_n a x_{n-1} je dán číslem charakterizujícím t. zv. *zlatý řez*. Protože rovnice (2.18) má vskutku zajímavou biologickou interpretaci, můžeme ve výše uvedeném výsledku spatřovati úkaz estetiky projevující se v některých oblastech života na naší planetě. Model popisující jistý způsob rozmnожování králíků pochází ze 13. století; jeho autorem je právě Fibonacci (r. 1228).

Zcela jinou aplikací lineárních diferenčních rovnic je následující úloha.

Buď

$$p_n(x) = \sum_{k=0}^n b_k x^{n-k}, \quad b_0 \neq 0, \quad n = 0, 1, \dots, N$$

a položme si úlohu stanovit hodnotu

$$p_n(z),$$

kde z je daný bod na reálné ose.

Naším cílem je pokud možno minimalizovat při tom počet aritmetických operací ve snaze snížit riziko numerické nestability.

Algoritmus 2.1 Počítejme veličiny x_0, x_1, \dots, x_N rekurentně pomocí relací

$$(2.20) \quad x_0 = b_0, \quad x_n = zx_{n-1} + b_n, \quad n = 1, \dots, N.$$

Posloupnost $\{x_n\}$ daná v (2.20) je řešení diferenční rovnice (2.11) splňující podmínu $x_0 = b_0$, kde $a_n = z$ pro $n = 1, 2, \dots$. V tomto případě

$$\pi_n = z^n$$

a tudíž, na základě věty 2.1,

$$x_n = \sum_{k=0}^n b_k z^k$$

a, speciálně pro $n = N$,

$$x_N = \sum_{k=0}^N b_k z^k = p_N(z).$$

Platí tedy

Věta 2.2 Veličiny x_n tvořené pomocí algoritmu 2.1 jsou hodnoty polynomů p_n definovaných pomocí

$$p_n(x) = \sum_{k=0}^n b_k x^k, \quad n = 0, 1, \dots, N$$

v bodě $x = z$.

Algoritmus 2.1 se nazývá *Hornerovým schématem*.

Nyní si ukážeme jak lze uvedený algoritmus zobecnit pro výpočet hodnot derivací polynomu p_n v bodě $x = z$.

Opět je nutné zdůraznit, že přímá aplikace standardních pravidel derivování člen po členu není vhodná a algoritmy typu Hornerova jsou daleko stabilnější a jednodušší z hlediska teorie složitosti.

Vyjdeme ze známých vztahů

$$p(x) = p(z + h) = \sum_{k=0}^N c_k h^{N-k},$$

kde

$$(2.21) \quad c_{N-k} = \frac{1}{k!} p^{(k)}(z), \quad k = 0, 1, \dots, N.$$

$x_0^{(0)}$	$x_1^{(0)}$	$x_2^{(0)}$	$x_3^{(0)}$	<u>$x_4^{(0)}$</u>
		↓		
$x_0^{(1)}$	$x_1^{(1)}$	$x_2^{(1)}$	<u>$x_3^{(1)}$</u>	
$x_0^{(2)}$	$x_1^{(2)}$	<u>$x_2^{(2)}$</u>		
$x_0^{(3)}$	<u>$x_1^{(3)}$</u>			
<u>$x_0^{(4)}$</u>				

Figure 1: Obrázek č. 2.1

Algoritmus 2.2 Označme prvky tvořené pomocí algoritmu 2.1 symboly $\{x_n^{(0)}\}$. Pro $k = 1, 2, \dots, N$ tvořme posloupnosti $\{x_n^{(k)}\}$ rekurzivně pomocí diferenčních schémát

$$(2.22) \quad x_0^{(k)} = x_0^{(k-1)}, \quad x_n^{(k)} = zx_{n-1}^{(k)} + x_n^{(k-1)}, \quad n = 0, 1, \dots, N-k.$$

Věta 2.3 Označuje-li p_n polynom

$$p_n(x) = \sum_{k=0}^n b_k x^{n-k}, \quad n = 0, 1, \dots, N,$$

pak veličiny určované algoritmem 2.2 splňují

$$(2.23) \quad x_n^{(k)} = \frac{1}{k!} p_{n+k}^{(k)}(z), \quad n = 1, 2, \dots, N, \quad k = 0, 1, \dots, N-n.$$

Tudíž koeficienty c_{N-k} požadované ve formuli (2.21) jsou právě ty, jež se nalézají v dolní diagonále na obrázku 2.1 pro $N = 4$ v posloupnosti veličin tvořených algoritmem 2.2.

3 Iterační metody řešení nelineárních rovnic a jejich soustav

3.1 Postupné approximace

Začněme příkladem úlohy najít kořen funkce g , t. j. nalézt takový bod $\hat{x} \in I = [a, b]$, $-\infty < a \leq b < +\infty$, tak, aby

$$g(\hat{x}) = 0$$

za předpokladu, že g je daná spojitá funkce zobrazující $I = I_1$ do $I_2 \subset R$. Speciálně, g může být polynom. A již v tomto speciálním případě si můžeme uvědomit několik důležitých skutečností. Tak především, obecně nemusí v I existovat kořen funkce g . Ale i když takový kořen existuje, nemusí existovat analytické vyjádření tohoto kořene pomocí parametrů charakterizujících g (na př. koeficienty polynomu); takové formule jsou známy pro případ polynomů stupňů nepřevyšujících 4. A jak víme, cena takových formulí je sporná, protože takové formule jsou mnohdy prakticky nepoužitelné (t. zv. *casus irreducibilis*). Důsledky těchto skutečností jsou zřejmé: Naději, že se podaří nalézt kořen g mají spíše numerické než analytické metody.

Mezi numerickými metodami se z důvodů snadné algoritmické realizovatelnosti prosadily *metody iterační* a to navzdory jejich neuniversálnosti a pohříchu poměrně pomalé konvergenci. Právě uvedený nedostatek dal podnět k rozvoji metod urychlování konvergence obecných posloupností, jmenovitě pak takových generovaných iteračními procesy.

Základním iteračním postupem jsou *postupné approximace*. Jsou určeny k sestrojování approximací řešení rovnic typu

$$(3.1) \quad s = f(s),$$

kde f je daná funkce, tedy, $g(s) = s - f(s)$.

Algoritmus 3.1 Zvolme x_0 libovolně a tvořme posloupnost $\{x_n\}$ rekurzivně pomocí relace

$$(3.2) \quad x_n = f(x_{n-1}).$$

Všimněme si, že k tomu, abyhom mohli prováděti jednotlivé kroky dle (3.2) musí f kromě spojitosti míti ještě další vhodné vlastnosti. Budeme předpokládat, že obor hodnot funkce f patří do I_1 , tedy $f(I_1) \subset I_1$. V takovém případě spojitost f zaručuje neomezenou proveditelnost procesu (3.2).

Abyhom to ozřejmili, položme

$$g(s) = s - f(s), \quad a \leq s \leq b,$$

takže

$$g(a) \leq 0, \quad a \quad g(b) \geq 0.$$

Ze spojitosti f plyne spojitost g a odtud pak skutečnost, že g nabývá všech hodnot z intervalu $[g(a), g(b)] = I_3 \subset I_1$. Existuje tedy alespoň jeden bod \hat{s} takový,

$$g(\hat{s}) = 0.$$

To značí, že

$$\hat{s} = f(\hat{s})$$

a \hat{s} je námi hledaný kořen.

Obecně, množina M všech samodružných bodů funkce f , t. j.

$$M = \{s \in I_1 : s = f(s)\}$$

má mohutnost moh $M \geq 1$.

Jednoznačnosti se dociluje pomocí dalších předpokladů.

Definice 3.1 Předpokládejme, že f zobrazuje $I \subset R$ do R a že existuje konstanta L taková, že nerovnost

$$(3.3) \quad |f(s) - f(t)| \leq L|s - t|$$

platí pro libovolná $s, t \in I$. V takovém případě se f nazývá *lipschitzovskou* a konstanta L *Lipschitzovou konstantou* funkce f .

Všimněme si, že funkce lipschitzovská na I je nutně spojitá na I . Je-li f diferencovatelná na I a platí

$$\sup \{|f'(s)| : s \in I\} \leq L,$$

kde $L > 0$ je konstanta, pak je f na I lipschitzovská a L je Lipschitzovou konstantou funkce f .

Předpokládejme, že Lipschitzova konstanta L splňuje nerovnost

$$(3.4) \quad 0 < L < 1.$$

Snadno nahlédneme, že pro $s_1, s_2 \in M$ platí

$$|s_1 - s_2| = |f(s_1) - f(s_2)| \leq L|s_1 - s_2|$$

a tudíž

$$s_1 = s_2.$$

Můžeme přistoupit k formulaci základního výsledku o konvergenci postupných aproximací.

Věta 3.1 *Budě $I = [a, b]$ uzavřený konečný interval a nechť funkce f splňuje následující podmínky (i) – (iii).*

(i) $f(s) \in I$ pro $s \in I$.

(ii) f je lipschitzovská s Lipschitzovou konstantou $L < 1$.

Potom pro libovolné nulové přiblžení $x_0 \in I$ posloupnost $\{x_n\}$ definovaná pomocí algoritmu 3.1 konverguje k (jedinému) řešení rovnice $s = f(s)$.

Důkaz. Z předchozích úvah již víme, že existuje alespoň jedno řešení rovnice $s = f(s)$ v I a že, díky předpokladu o L , toto řešení je určeno jednoznačně. Stačí tedy dokázati konvergenci posloupnosti $\{x_n\}$.

Nechť $\hat{x} = f(\hat{x})$. Zřejmě

$$x_n - \hat{x} = f(x_{n-1}) - f(\hat{x}),$$

takže

$$|x_n - \hat{x}| \leq L|x_{n-1} - \hat{x}|$$

a tudíž

$$(3.5) \quad |x_n - \hat{x}| \leq L^n|x_0 - \hat{x}|.$$

Protože

$$\lim_{n \rightarrow \infty} L^n = 0,$$

obdržíme žádaný výsledek

$$\lim_{n \rightarrow \infty} |x_n - \hat{x}| = 0.$$

|||

Příklad 3.1 Pomocí postupných approximací hledejme řešení rovnice

$$s = e^{-s}.$$

Protože $e^{-s} > 0$ pro $s \in (-\infty, +\infty)$ a $s > 1 \geq e^{-s}$ pro $s > 1$, stačí omezit se na interval $[0, 1]$.

Snadno si uvědomíme, že obor hodnot funkce $f(s) = e^{-s}$ pro $s \in [0, 1]$ leží v intervalu $[e^{-1}, 1] \subset [0, 1]$ a že pro $s_1, s_2 \in [0, 1]$ platí

$$|f(s_1) - f(s_2)| = |f'(\tilde{s})||s_1 - s_2|,$$

při čemž

$$f'(\tilde{s}) = -e^{-\tilde{s}}, \quad \tilde{s} \in [s_1, s_2].$$

Protože

$$\max \{|f'(s)| : s \in [0, 1]\} = 1,$$

Lipschitzova konstanta pro f na $[0, 1]$ je rovna 1, což k platnosti věty 3.1 nestačí.

Avšak zvolíme-li za I interval $[\frac{1}{2}, \log 2]$, snadno zjistíme, že pro $\frac{1}{2} \leq s \leq \log 2$,

$$\frac{1}{2} = e^{-\log 2} \leq e^{-s} \leq e^{-\frac{1}{2}} < \log 2$$

a tudíž $[\frac{1}{2}, \log 2] = I \subset [0, 1]$. Navíc první derivace funkce f , kde $f(s) = e^{-s}$, je klesající, takže

$$\max \{|f'(s)| : s \in I\} = f'\left(\frac{1}{2}\right) \approx 0,606531 < 1.$$

Iterační proces (3.2) pro funkci $f(s) = e^{-s}$ je tedy konvergentní v I .

Odhad (3.5) nemá bezprostřední praktický význam, závisí totiž na znalosti hledaného řešení. V praxi jsme totiž vždy nuceni omezit se na určitý konečný počet iterací x_n v (3.2). Rádi bychom proto znali nějaký realistický odhad chyby po n -tém kroku.

Všimněme si, že pro libovolné $k \geq 1$ platí vztahy

$$\begin{aligned} |x_{k+1} - x_k| &= |f(x_k) - f(x_{k-1})| \\ &\leq L|x_k - x_{k-1}| \\ &\leq L^k|x_1 - x_0|. \end{aligned}$$

Budeme n zafixováno pevně a nechť $m = n + p > n$. Potom

$$x_m - x_n = \sum_{k=n}^{n+p-1} (x_{k+1} - x_k).$$

Po aplikaci trojúhelníkové nerovnosti obdržíme vztahy

$$\begin{aligned} |x_m - x_n| &\leq \sum_{k=n}^{n+p-1} L^k|x_1 - x_0| \\ &= L^n \sum_{k=0}^{p-1} L^k|x_1 - x_0| \\ &\leq L^n \sum_{k=0}^{\infty} L^k|x_1 - x_0| \\ &= L^n(1 - L)^{-1}|x_1 - x_0|. \end{aligned}$$

Protože pro $m \rightarrow \infty$ plyně, že $x_m \rightarrow \hat{x}$, tedy též $p \rightarrow \infty$, a platí

$$(3.6) \quad |x_n - \hat{x}| \leq \frac{L^n}{1-L} |x_1 - x_0|.$$

Důsledek 3.1 *Při splnění předpokladů věty 3.1 je odhad chyby po n krocích postupných approximací definovaných pomocí algoritmu 3.1 dán výrazem (3.6).*

Zkoumejme blíže chování posloupnosti $\{x_n\}$ za předpokladů věty 3.1 doplněných o předpoklad následující.

Nechť f je spojitě diferencovatelná a nechť v celém intervalu I je tato derivace nenulová.

To znamená, že f buď klesá nebo roste v I .

Je-li $x_0 \neq \hat{x} = f(\hat{x})$, pak $f(x_n) \neq \hat{x}$ pro libovolné $n = 1, 2, \dots$ Abychom to nahlédli, stačí si uvědomit, že kdyby

$$x_n = f(x_{n-1}) = f(x_n),$$

při čemž $x_n \neq x_{n-1}$, pak

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

kde ξ leží mezi x_{n-1} a x_n , t.j.

$$(x_{n-1} - \xi) (\xi - x_n) < 0.$$

Protože $x_{n-1} - x_n \neq 0$, musí být $f'(\xi) = 0$, ale to je ve sporu s předpokladem. Odtud vyplývá, že chyba

$$d_n = x_n - \hat{x}$$

není pro žádné $n \geq 0$ nulová.

Existuje limita

$$\lim_{n \rightarrow \infty} \frac{d_{n+1}}{d_n}?$$

A když ano, pak čemu je rovna?

Snadno zjistíme, že

$$\begin{aligned} d_{n+1} &= x_{n+1} - \hat{x} = f(x_n) - f(\hat{x}) \\ &= f(\hat{x} + d_n) - f(\hat{x}) = f'(\hat{x} + \theta_n d_n) d_n, \end{aligned}$$

kde

$$0 < |\theta_n| < 1.$$

Definujme ε_n pomocí relace

$$f'(\hat{x} + \theta_n d_n) = f'(\hat{x}) + \varepsilon_n.$$

Potom

$$(3.7) \quad d_{n+1} = (f'(\hat{x}) + \varepsilon_n) d_n$$

a tudíž, ježto $\varepsilon_n \rightarrow 0$ pro $n \rightarrow \infty$ plyne na základě spojitosti f' rovnost

$$(3.8) \quad \lim_{n \rightarrow \infty} \frac{d_{n+1}}{d_n} = f'(\hat{x}).$$

Lze této skutečnosti využít pro praktické účely? Veličinu \hat{x} a tudíž i $f'(\hat{x})$ neznáme.

Ukážeme, že explicitní znalost uvedených veličin není nutná k tomu, abychom rovnosti (3.8) efektivně využili k sestrojení algoritmu, který poskytne approximace konvergující k hledanému řešení \hat{x} rychleji než posloupnost $\{x_n\}$.

Předpokládejme na okamžik, že ve formuli (3.7) $\varepsilon_n = 0$ pro nějaké pevné $n \geq 1$. Nechť

$$f'(\hat{x}) = A,$$

takže

$$\begin{aligned} x_{n+1} - \hat{x} &= A(x_n - \hat{x}), \\ x_{n+2} - \hat{x} &= A(x_{n+1} - \hat{x}). \end{aligned}$$

Snadno zjistíme, že

$$A = \frac{x_{n+2} - x_{n+1}}{x_{n+1} - x_n},$$

takže

$$\begin{aligned} \hat{x} &= \frac{1}{1-A}(x_{n+1} - Ax_n) \\ &= x_n + \frac{1}{1-A}(x_{n+1} - x_n) = x_n - \frac{(x_{n+1} - x_n)^2}{x_{n+2} - 2x_{n+1} + x_n} \end{aligned}$$

Vidíme tedy, že za předpokladu, že $\varepsilon_n = 0$, pro nějaké $n \geq 1$, lze přesné řešení \hat{x} obdržeti použitím tří po sobě jdoucích iterací. To však je dost akademická situace.

Ačkoliv veličiny $\varepsilon_n \neq 0$, lze očekávat, že ve srovnání s modulem $f'(\hat{x})$ jsou malé.

Pro velká n tak posloupnost

$$(3.9) \quad y_n = x_n - \frac{(x_{n+1} - x_n)^2}{x_{n+2} - 2x_{n+1} + x_n}$$

skytá pro \hat{x} lepší approximaci než x_n .

Má tedy smysl vyšetřovat

Algoritmus 3.2 Bud $\{x_n\}$ libovolná posloupnost a nechť y_n je tvoreno pomocí (3.9).

Zavedeme označení

$$\Delta x_n = x_{n+1} - x_n, \quad n = 0, 1, \dots$$

a

$$\Delta^{k+1} x_n = \Delta(\Delta^k x_n), \quad k \geq 1.$$

Na př

$$\begin{aligned} \Delta^2 x_n &= \Delta(x_{n+1} - x_n) \\ &= x_{n+2} - 2x_{n+1} + x_n \end{aligned}$$

Formuli (3.9) lze pak psát ve tvaru

$$(3.10) \quad y_n = x_n - \frac{(\Delta x_n)^2}{\Delta^2 x_n}.$$

Posloupnost $\{y_n\}$ tvořená pomocí algoritmu 3.2 definuje t. zv. *Aitkenův urychlovací Δ^2 -proces*.

Jest na místě zdůraznit, že Aitkenův proces je definován pro posloupnosti ne nutně generované pomocí algoritmu 3.1. Pro takto obecný proces platí poměrně silné tvrzení.

Věta 3.2 *Budě $\{x_n\}$ daná posloupnost taková, že*

$$\hat{x} = \lim_{n \rightarrow \infty} x_n.$$

Nechť

$$(3.11) \quad d_n = x_n - \hat{x} \neq 0 \quad n = 0, 1, \dots$$

a

$$(3.12) \quad d_{n+1} = (A + \varepsilon_n) d_n,$$

kde A je konstanta taková, že

$$|A| < 1$$

a

$$\varepsilon_n \rightarrow 0 \text{ pro } n \rightarrow \infty.$$

Potom posloupnost $\{y_n\}$ tvořená pomocí (3.9) v algoritmu 3.2 je dobře definovaná pro dostatečně velká n a navíc

$$(3.13) \quad \lim_{n \rightarrow \infty} \frac{y_n - \hat{x}}{x_n - \hat{x}} = 0,$$

což značí, že posloupnost $\{y_n\}$ konverguje k \hat{x} rychleji než posloupnost $\{x_n\}$.

Důkaz. Použitím formule (3.12) obdržíme

$$d_{n+2} = (A + \varepsilon_{n+1}) (A + \varepsilon_n) d_n,$$

takže

$$\begin{aligned} \Delta^2 x_n &= x_{n+2} - 2x_{n+1} + x_n \\ &= d_{n+2} - 2d_{n+1} + d_n \\ &= [(A-1)^2 + \varepsilon'_n] d_n, \end{aligned}$$

kde

$$\varepsilon'_n = A[\varepsilon_n + \varepsilon_{n+1}] - 2\varepsilon_n + \varepsilon_n \varepsilon_{n+1}.$$

Z podmínky $\varepsilon_n \rightarrow 0$ pro $n \rightarrow \infty$ plyne, že

$$(3.14) \quad \varepsilon'_n \rightarrow 0 \text{ pro } n \rightarrow \infty.$$

Odtud dostáváme,

$$(A-1)^2 + \varepsilon'_n \neq 0$$

pro dostatečně velká n , řekněme pro $n > n_0$. Tudíž

$$\Delta^2 x_n \neq 0 \text{ pro } n > n_0.$$

Dále pak

$$\Delta x_n = \Delta d_n = (A - 1 + \varepsilon_n) d_n$$

a tudíž

$$\begin{aligned} y_n - \hat{x} &= d_n - \frac{(\Delta x_n)^2}{\Delta^2 x_n} \\ &= d_n - \frac{(A - 1 + \varepsilon_n)^2 d_n}{(A - 1)^2 + \varepsilon'_n}. \end{aligned}$$

Z $\varepsilon_n \rightarrow 0$ pro $n \rightarrow \infty$ a (3.14) plyne, že

$$\frac{y_n - \hat{x}}{d_n} = \frac{\varepsilon'_n - 2(A - 1) \varepsilon_n - \varepsilon_n^2}{(A - 1)^2 + \varepsilon'_n} \rightarrow 0,$$

což jsme měli dokázati. |||

Důsledek 3.2 *Předpokládejme, že f splňuje podmínky věty 3.1 a navíc, že f má spojitou 1. derivaci na I a $f'(s) \neq 0$ pro $s \in I$. Pro $x_0 \neq \hat{x}$ posloupnost $\{x_n\}$ vyhovuje předpokladům věty 3.2, takže algoritmus 3.2 poskytuje urychlení konvergence.*

Důkaz. Zřejmě stačí ověřit platnost vztahu (3.12). To však je vztah (3.7) odvozený za obecnějších předpokladů. Zbývá tudíž ověřit požadavek $|A| < 1$. Platnost tohoto vztahu však plyne odtud, že proces postupných approximací je konvergentní. Tímto konstatováním můžeme důkaz ukončit. |||

Znovu připomeňme širší aplikabilitu urychlovacího Aitkenova Δ^2 -procesu i mimo oblast metody postupných approximací.

3.2 Kvadratická konvergence a Newtonova metoda

V předchozím odstavci jsme předpokládali, že $f'(s) \neq 0$ pro $s \in I$, což zaručuje platnost předpokladu (3.12) a obdrželi jsme tak charakteristiku rychlosti konvergence danou pomocí (3.12), t. zv. *lineární konvergenci*. Dá se říci, že v takovém případě počet platných cifer je lineární funkcí indexu iterací.

Učiňme nyní předpoklad

$$(3.15) \quad f'(\hat{x}) = 0.$$

Tento předpoklad je poměrně velmi silný a zabezpečuje splnění některých předpokladů věty 3.1.

Věta 3.3 *Budě $I \subset (-\infty, +\infty)$, (nikoliv nutně ohraničený interval) a nechť f je definovaná na I a přitom platí*

- (i) *f a f' jsou spojité na I .*
- (ii) *Rovnice $s = f(s)$ má řešení $\hat{x} \in I$ takové, že platí (3.15).*

Potom existuje $\delta > 0$ takové, že algoritmus 3.1 poskytuje posloupnost $\{x_n\}$ konvergující k \hat{x} pro všechna $|x_0 - \hat{x}| \leq \delta$.

Důkaz. Označme

$$I_\delta = [\hat{x} - \delta, \hat{x} + \delta].$$

Pro dostatečně malá $\delta_1 > 0$, $I_{\delta_1} \subset I$. Potom existuje $\delta_0 > 0$ a $L > 0$ takové, že

$$|f'(s)| \leq L \text{ pro } s \in I_{\delta_0}.$$

Nechť $\delta = \min\{\delta_0, \delta_1\}$.

Podle věty o střední hodnotě, protože $f'(\hat{x}) = 0$ a tedy $0 \leq L < 1$,

$$|f(s) - \hat{x}| = |f(s) - f(\hat{x})| \leq |f'(\xi)||s - \hat{x}| \leq L|s - \hat{x}| < \delta.$$

Tudíž obor hodnot funkce f je obsažen v I_δ a algoritmus 3.1 vede ke konvergenci podle věty 3.1. |||

Zesilme naše požadavky na hladkost funkce f předpokladem existence spojité druhé derivace f'' na I_δ a její nenulovosti.

Podobně jako v předchozích úvahách předpoklad $x_0 \neq \hat{x}$ implikuje platnost vztahů $x_n \neq \hat{x}$, $n = 1, 2, \dots$. Postupné aproximace nemohou tedy poskytnout přesné řešení v konečném počtu kroků.

Užitím Taylorova rozvoje se zbytkem, obdržíme vztahy

$$\begin{aligned} d_{n+1} &= x_{n+1} - \hat{x} \\ &= f(x_n) - f(\hat{x}) \\ &= f'(\hat{x})d_n + \frac{1}{2}f''(\hat{x} + \theta_n d_n)d_n^2, \end{aligned}$$

kde $0 < |\theta_n| < 1$.

Z našich předpokladů plyne, že

$$(3.16) \quad d_{n+1} = \frac{1}{2}f''(\hat{x} + \theta_n d_n)d_n^2.$$

Protože $d_n \neq 0$, $n = 1, 2, \dots$ a $d_n \rightarrow 0$ pro $n \rightarrow \infty$, odvodíme, že

$$(3.17) \quad \lim_{n \rightarrow \infty} \frac{d_{n+1}}{d_n^2} = \frac{1}{2}f''(\hat{x}).$$

To je pozoruhodný vztah, podle něhož chyba po $(1+n)$ -tém kroku je úměrná čtverci chyby po n -tém kroku. Tato skutečnost se označuje jako *kvadratická konvergence*. Počet platných cifer se zdvojuje po každém iteračním kroku.

Příklad 3.2 Buď $a > 0$ a nechť $f(s) = (1/2)[s + (a/s)]$ pro $s > 0$. Rovnice $s = f(s)$ má řešení \sqrt{a} . Zřejmě $f'(\sqrt{a}) = 0$ a $f''(s) > 0$ pro $s > 0$ a $f''(\sqrt{a}) = \frac{1}{\sqrt{a}}$. Tudíž

$$(3.18) \quad x_{n+1} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right)$$

konverguje k \sqrt{a} kvadraticky pro x_0 dostatečně blízká k \sqrt{a} . Později ukážeme, že tento proces konverguje pro libovolné výchozí přiblížení $x_0 > 0$.

Nyní přistoupíme k snad nejpoužívanější metodě přibližného řešení nelineárních rovnic a lze říci, že nejen rovnic ale i nelineárních úloh - k *metodě Newtonově*.

V předchozích úvahách jsme dospěli ke kvadratické konvergenci za výrazně akademického předpokladu nulovosti první derivace v přesném řešení. Ani se nechce věřit, že tento akademismus má tak nádherné uplatnění, k jehož výkladu teď přistupujeme.

Budě F definovaná a dvakrát spojitě diferencovatelná na intervalu $I = [a, b] \subset R$ a nechť $F'(s) \neq 0$ pro $s \in I$. Dále nechť rovnice

$$(3.19) \quad F(s) = 0$$

má řešení $\hat{x} \in (a, b)$ (nikoliv nutně jediné). Dá se ukázat, že toto řešení lze nalézt pomocí postupných approximací s funkcí

$$f(s) = s + MF(s),$$

kde M je konstanta splňující určité podmínky, na př. $M = (\alpha + \beta)/2$, kde $I = [0, 1]$, $F(0) < 0 < F(1)$ a $0 < \alpha \leq F'(s) \leq \beta$. Víme též, že genericky obdržíme lineární konvergenci. Jak to tedy zařídit, abychom sestrojili algoritmus poskytující kvadratickou konvergenci?

Což připustit místo konstantního faktoru M vhodnou funkci, tedy snažit se sestrojit f ve tvaru

$$f(s) = s + h(s)F(s),$$

kde hodnoty funkce h se "snadno" počítají.

Snažme se tedy splnit tyto požadavky:

$$\hat{x} = f(\hat{x}) \text{ a } f'(\hat{x}) = 0,$$

To nás vede k relacím

$$f'(s) = 1 + h'(s)F(s) + h(s)F'(s)$$

a tudíž

$$h'(\hat{x})F(\hat{x}) + h(\hat{x})F'(\hat{x}) = -1$$

a ježto $F(\hat{x}) = 0$,

$$h(s) = -\frac{1}{F'(s)}.$$

Dospěli jsme tedy k následujícímu algoritmu.

Algoritmus 3.3 Volme x_0 a určeme posloupnost $\{x_n\}$ pomocí vztahů

$$(3.20) \quad x_{n+1} = x_n - \frac{1}{F'(x_n)}F(x_n), n = 0, 1, \dots$$

Algoritmus 3.3 nese název *Newtonova metoda*, či *metoda Newtonova - Raphsonova*.

Z předchozího výkladu již víme, že, ježto $f(s) = (1/F'(s))F(s)$,

$$f'(\hat{x}) = 0,$$

takže případná konvergence posloupnosti $\{x_n\}$ je kvadratická.

Jednotlivé kroky Newtonovy metody mají velice názornou interpretaci.

V bodě x_n approximuje graf funkce F pomocí tečny ke grafu F v tomto bodě. Průsečík této tečny s osou x určuje další approximaci x_{n+1} .

Umíme si snadno představit grafy funkce F , pro něž Newtonův proces diverguje. (To je velice hezké cvičení.)

Podobně jako pro všechny iterační metody typu postupných approximací, slabina Newtonovy metody spočívá mimo jiné též v silné závislosti na volbě počátečního přiblížení. Ukážeme si jak lze zabezpečit globální konvergenci; požadujeme vyšší hladkost F - totiž znaménkovou stálost druhé derivace F'' .

Věta 3.4 *Předpokládejme, že funkce F je definovaná a dvakrát spojitě differencovatelná v $I = [a, b]$ a nechť splňuje tyto požadavky*

- (i) $F(a)F(b) < 0$;
- (ii) $F'(s) \neq 0$ pro $s \in I$;
- (iii) *Budě $F''(s) \geq 0$ s $s \in I$ nebo $F''(s) \leq 0$, s $s \in I$;*
- (iv) *Nechť (α) $c = a$ jestliže*

$$|F'(a)| \leq |F'(b)| ,$$

- (β) $c = b$, jestliže

$$|F'(b)| \leq |F'(a)| .$$

Dále nechť

$$\left| \frac{F(c)}{F'(c)} \right| \leq b - a.$$

Potom Newtonova metoda konverguje pro libovolné výchozí přiblížení x_0 k jedinému řešení rovnice $F(s) = 0$.

Poznámky. Podmínka (i) zřejmě zaručuje existenci kořene díky spojitosti F . Z (ii) plyne jednoznačnost kořene v I . Podmínka (iii) říká, že F je buď konkávní nebo konvexní v I . Posléze (iv) zajistuje, aby tečna ke křivce $y = F(x)$ v koncovém bodě, kde $|F'(x)|$ je menší, protínala osu x v intervalu I .

Důkaz věty 3.4. Věta 3.4 zahrnuje tyto čtyři odlišné situace

- (a) $F(a) < 0$, $F(b) > 0$, $F''(s) \leq 0$ ($c = b$),
- (b) $F(a) > 0$, $F(b) < 0$, $F''(s) \geq 0$ ($c = b$),
- (c) $F(a) < 0$, $F(b) > 0$, $F''(s) \geq 0$ ($c = a$),
- (d) $F(a) > 0$, $F(b) < 0$, $F''(s) \leq 0$ ($c = a$).

Případy (b) a (d) se snadno redukují na (a) a (c) volbou $-F$ na místo F . (Tato úprava nemění dokonce ani prvky posloupnosti $\{x_n\}$.) Případ (c) se převede na (a) volbou $-x \rightarrow x$. (V tomto případě obdržíme místo $\{x_n\} \rightarrow \{-x_n\}$. Stačí se tedy zabývat případem (a). Máme pak co do činění s grafem rostoucí konkávní funkce F takové, že $F(a) < 0 < F(b)$.

Budě \hat{x} jediný kořen rovnice $F(s) = 0$. Napřed předpokládejme, že $a \leq x_0 \leq \hat{x}$. Potom, na základě předpokladu, že $F'(x_0) \geq 0$,

$$x_1 = x_0 - \frac{F(x_0)}{F'(x_0)} \geq x_0.$$

Indukcí lze ukázati, že

$$(3.21) \quad x_0 \leq x_1 \leq \dots \leq x_n \leq \dots \leq \hat{x}.$$

Indukční krok je realizován takto. Z věty o střední hodnotě,

$$-F(x_n) = F(\hat{x}) - F(x_n) = F'(\xi)(\hat{x} - x_n)$$

a přitom $x_n \leq \xi_n \leq \hat{x}$. Z podmínky $F''(s) \leq 0$ plyne, že F' je nerostoucí a tudíž $F'(\xi_n) \leq F'(x_n)$, takže

$$-F(x_n) \leq F'(x_n)(\hat{x} - x_n)$$

a tedy

$$x_{n+1} = x_n - \frac{F(x_n)}{F'(x_n)} \leq x_n + (\hat{x} - x_n) = \hat{x}.$$

Aplikací F k oběma stranám nerovnosti obdržíme (F je neklesající), že $F(x_{n+1}) \leq 0$ a tudíž

$$x_{n+2} = x_{n+1} - \frac{F(x_{n+1})}{F'(x_{n+1})} \geq x_{n+1}.$$

Tedy platí (3.21).

Na tomto místě aplikujeme jednu z fundamentálních vět matematické analýzy: *Každá ohrazená monotonní posloupnost reálných čísel je konvergentní.*

Označme tedy

$$\tilde{x} = \lim_{n \rightarrow \infty} x_n.$$

Z (3.21) plyne, že $\tilde{x} \leq \hat{x}$. Na druhé straně z definice posloupnosti $\{x_n\}$ obdržíme při $n = \infty$ rovnost

$$\tilde{x} = \tilde{x} - \frac{F(\tilde{x})}{F'(\tilde{x})}$$

a tedy $F(\tilde{x}) = 0$. Z jednoznačnosti pak $\tilde{x} = \hat{x}$. Tím je dokázáno tvrzení věty 3.4 pro případ $x_0 \leq \hat{x}$.

Dále vyšetřujme případ $x_0 \geq \hat{x}$. Opět použitím věty o střední hodnotě

$$F(x_0) = F(x_0) - F(\hat{x}) = F'(\xi_0)(x_0 - \hat{x}),$$

kde $\hat{x} \leq \xi_0 \leq x_0$ a díky monotonii F' , $F(x_0) \geq F'(\xi_0)(x_0 - \hat{x})$. Snadno zjistíme, že

$$x_1 = x_0 - \frac{F(x_0)}{F'(x_0)} \leq x_0 - (x_0 - \hat{x}) = \hat{x}$$

a tudíž

$$F(x_0) \leq F(b) - (b - x_0)F'(b).$$

Na základě podmínky (iv) obdržíme, že platí

$$\begin{aligned} x_1 &= x_0 - \frac{F(x_0)}{F'(x_0)} \geq x_0 - \frac{F(x_0)}{F'(b)} \\ &\geq x_0 - \frac{F(b)}{F'(b)} + b - x_0 \end{aligned}$$

$$= x_0 - (b - a) + b - x_0 = a.$$

Tudíž $a \leq x_1 \leq \hat{x}$. Tím jsme se dostali do situace, již jsme vyšetřovali v předchozí části důkazu s tím rozdílem, že na místě x_0 vystupuje x_1 . Podle již dokázaného, $\{x_n\}$ konverguje k \hat{x} . Tím je důkaz věty 3.4 proveden. |||

Vraťme se ještě k příkladu 3.2. O procesu v něm vyšetřovaném již víme, že je konvergentní pro počáteční přiblžení dostatečně blízká k přesnému řešení, t. j. k \sqrt{c} , kde $c > 0$ je dané kladné číslo, jehož odmocninu hledáme.

Položme tedy

$$F(s) = s^2 - c$$

a

$$f(s) = s - \frac{F(s)}{F'(s)}.$$

Snadno zjistíme, že

$$F'(s) = 2s > 0,$$

a

$$F''(s) = 2 > 0$$

a tudíž nastává v tomto případě situace (c) popsaná v důkazu věty 3.4. V každém intervalu (a, b) , kde $0 < a < \sqrt{c} < b$, platí, že

$$\left| \frac{F(a)}{F'(a)} \right| = \frac{c - a^2}{2a} \leq b - a,$$

pro každé b splňující nerovnost

$$b \geq \frac{1}{2}(a + c/a).$$

Z věty 3.4 plyne konvergence posloupnosti

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{c}{x_n} \right) \rightarrow \sqrt{c}$$

pro libovolné $x_0 > 0$.

Zvlášť výhodné je použití metody Newtonovy k hledání nulových bodů polynomů.

V tom případě využijeme té skutečnosti, že umíme efektivně stanoviti jak hodnoty $p(z)$ tak $p'(z)$, kde p je daný mnohočlen a z je daná veličina.

Bud tedy

$$p(z) = \sum_{k=0}^N a_k z^{N-k},$$

pak pomocí vztahů

$$(3.22) \quad b_0 = a_0, \quad b_n = z b_{n-1} + a_n, \quad n = 1, 2, \dots, N;$$

a

$$(3.23) \quad c_0 = b_0, \quad c_n = z c_{n-1} + b_n, \quad n = 1, 2, \dots, N-1,$$

obdržíme rovnosti

$$b_N = p(z) \text{ a } c_{n-1} = p'(z).$$

Schématicky postupujeme jak znázorněno níže

$$\begin{array}{ccccccc} a_0 & a_1 & a_2 & \dots & a_{N-1} & a_N \\ \downarrow & \downarrow & \downarrow & \dots & \downarrow & \downarrow \\ b_0 \Rightarrow & b_1 \Rightarrow & b_2 \Rightarrow & \dots & b_{N-1} \Rightarrow & \underline{b_N} \\ \downarrow & \downarrow & \downarrow & & \downarrow & \\ c_0 \rightarrow & c_1 \rightarrow & c_2 \rightarrow & \dots & \underline{c_{N-1}} & \end{array}$$

při čemž \rightarrow znamená sčítání a \Rightarrow násobení faktorem z a sčítání.

Bud \checkmark nyní z kořen polynomu p . Pak

$$p(x) = (x - z)q(x),$$

kde q je polynom stupně $N - 1$,

$$q(x) = \sum_{k=0}^{N-1} b_k x^{N-k-1}, \quad b_N = 0.$$

Snadno zjistíme, že

$$a_0 = b_0$$

a

$$a_{n-1} = -b_n z + b_n, \quad n = 1, \dots, N - 1.$$

Tudíž

$$\sum_{k=0}^N a_k x^{N-k} = \sum_{k=0}^{N-1} b_k x^{N-k} - z \sum_{k=0}^{N-1} b_k x^{N-k-1},$$

při čemž

$$a_n = b_n - z b_{n-1},$$

neboli,

$$(3.24) \quad b_n = a_n + z b_{n-1}.$$

Vidíme, že jsme obdrželi vztahy (2.20) algoritmu 2.1.

Platí tedy

Věta 3.5 Je-li z kořen polynomu p , kde

$$p(x) = \sum_{k=0}^N a_k x^{N-k}, \quad a_0 \neq 0,$$

a koeficienty b_0, \dots, b_{N-1} se počítají podle formule (3.24), pak

$$q(x) = (x - z)^{-1} p(x) = \sum_{k=0}^{N-1} b_k x^{N-k}, \quad b_0 = a_0.$$

Mohlo by se zdát, že je výhodné hledat další kořeny polynomu p jakožto nulové body polynomu nižšího stupně q . Ukazuje se však, že taková procedura je numericky nestabilní a je tudíž vhodnější počítat všechny kořeny polynomu pomocí původního polynomu p a nikoliv modifikovaného polynomu q .

Cvičení 3.1 Analyzujte způsobem obdobným analýze Newtonovy metody metodu následující, metodu regula falsi, definovanou předpisem

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})F(x_n)}{F(x_n) - F(x_{n-1})},$$

kde F je funkce, jejíž kořeny hledáme, t. j. $F(x) = 0$.

Cvičení 3.2 Jaké vlastnosti má Newtonova metoda v případě, kdy $F'(x^*) = 0$, kde $F(x^*) = 0$.

Cvičení 3.3 Analyzujte modifikovanou Newtonovu (Whittakerovu - Robinsonovu) metodu danou předpisem

$$x_{n+1} = x_n - \frac{1}{m}F(x_n),$$

kde m je pevně zvolená konstanta, na př. $m = F'(x_0)$.

4 Iterační metody řešení soustav nelineárních rovnic

4.1 Věta o kontrakci

Předmětem našeho studia bude vyšetřování následující úlohy. Nalézt prvek $x^* \in R^N$ takový, že pro $x = x^*$ platí

$$(4.1) \quad \begin{cases} x_1 &= f_1(x_1, \dots, x_N), \\ x_2 &= f_2(x_1, \dots, x_N), \\ &\vdots \\ &\vdots \\ x_N &= f_N(x_1, \dots, x_N), \end{cases}$$

neboli, psáno vektorově,

$$(4.2) \quad \mathbf{x} = \mathbf{f}(\mathbf{x}),$$

kde $\mathbf{x} = (x_1, \dots, x_N)^T$ a $\mathbf{f} = (f_1, \dots, f_N)^T$.

Příklad 4.1 Nechť $N = 2$ a $f_1(x_1, x_2) = x_1^2 + x_2^2$, $f_2(x_1, x_2) = x_1^2 - x_2^2$. Potom má (4.1) tvar

$$x_1 = x_1^2 + x_2^2,$$

$$x_2 = x_1^2 - x_2^2.$$

Výraz $x_1^2 + x_2^2 - x_1 = 0$ je rovnicí kružnice se středem v bodě $(1/2, 0)$ a $x_1^2 - x_2^2 - x_2 = 0$ je rovnicí hyperboly se středem v bodě $(0, -1/2)$. Obě tyto křivky procházejí počátkem, což značí, že jedním z nulových bodů je $x_1 = 0, x_2 = 0$. Pomocí grafů pro $(x_1 - 1/2)^2 + x_2^2 = 1/4$ a $(x_2 + 1/2)^2 - x_1^2 = 1/4$ snadno nahlédneme, že v okolí bodu $\tilde{x}_1 = 0, 8, \tilde{x}_2 = 0, 4$ leží další (netriviální) nulový bod soustavy (4.1).

Při analýze důkazů konvergence metod studovaných v předchozí kapitole zjistíme, že významným důkazovým prostředkem byly vlastnosti absolutní hodnoty reálných čísel:

$$|x| = 0 \Leftrightarrow x = 0;$$

$$|cx| = |c||x|;$$

a

$$|x + y| \leq |x| + |y|$$

pro c, x a $y \in R$.

V souvislosti s prostory R^N dimenze $N \geq 2$, to vede k nutnosti zobecnit pojem absolutní hodnoty. Matematika na přelomu 19. a 20. století dospěla při tomto zobecňování k pojmu *normy* a *normovaného prostoru*.

Definice 4.1 Budě E lineární prostor nad tělesem reálných čísel. Funkce $\nu : E \rightarrow R_+^1$ se nazývá normou na E , jestliže platí následující relace:

- (i) $\nu(x) = 0 \Leftrightarrow x = 0$;
- (ii) $\nu(cx) = |c|\nu(x)$ $x \in E, c \in R^1$;
- (iii) $\nu(x + y) \leq \nu(x) + \nu(y)$, $x, y \in E$.

Existuje-li na E norma, tak se prostor E se nazývá normovaným.

Poznámka. Obvykle se norma na E označuje symbolem $\| \cdot \|$. Je-li zapotřebí zvýraznit, že norma je brána na prostoru E , pak se značí následovně $\| \cdot \|_E$.

Definice 4.2 Posloupnost $\{x_n\}$, $x_n \in \mathcal{E}$ se nazývá konvergentní k $x \in \mathcal{E}$ a x její limitou, jestliže platí

$$\lim_{n \rightarrow \infty} \|x_n - x\| = 0.$$

Definice 4.3 Posloupnost $\{x_n\}$, $x_n \in \mathcal{E}$, se nazývá cauchyovskou, jestliže pro každé $\varepsilon > 0$ existuje N_0 takové, že $\|x_{n+p} - x_n\| < \varepsilon$ pro $n \geq N_0$ a $p \geq 1$.

Normovaný prostor E se nazývá Banachovým prostorem je-li v něm každá cauchyovská posloupnost konvergentní.

Příklad 4.2 Nechť $\mathcal{E} = \mathcal{R}^N$ a

(a)

$$\|\mathbf{x}\|_p = \left(\sum_{k=1}^N |x_k|^p \right)^{1/p}, \quad 1 \leq p < +\infty,$$

(b)

$$\|\mathbf{x}\|_\infty = \max \{|x_k| : k = 1, \dots, N\}.$$

Všimněme si, že pro $N = 1$

$$\|\mathbf{x}\|_\infty = \|x\|_p = |x|,$$

takže absolutní hodnota v R^1 je norma.

Dále nechť $C = (c_{jk})$ je regulární matice typu $N \times N$. Potom výraz

$$\|x\|_C = \|Cx\|_{R^N}$$

je norma na R^N .

Speciálně nechť $\mathcal{E} = (R^N, \|\cdot\|_2)$ a $C = (c_{jk}), \det C \neq 0$.

Dále nechť

$$(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^N x_k y_k.$$

Položme $B = C^T C$ a

$$\|\mathbf{x}\|_{2,B}^2 = (B\mathbf{x}, \mathbf{x}).$$

Snadno se ukáže, že $\|\cdot\|_{2,B}$ je norma na R^N .

Vraťme se však k vyšetřování soustav nelineárních rovnic v R^N .

Metoda postupných approximací je definována podobně jako pro případ jedné rovnice.

Algoritmus 4.1 Zvolme prvek $\mathbf{x} \in \mathcal{E}$ a počítejme posloupnost vektorů $\{\mathbf{x}_n\}$ rekuzivně pomocí formule

$$(4.3) \quad \mathbf{x}_{n+1} = \mathbf{f}(\mathbf{x}_n), \quad n = 1, \dots,$$

Opět si klademe otázku, zda je proces (4.3) konvergentní a jaké vlastnosti má případná limita; zřejmě se nabízí řešení rovnice (4.2).

Odpovědi na položené otázky jsou obsahem následujícího tvrzení.

Věta 4.1 Nechť $\mathcal{I} = \prod_{j=1}^N \mathcal{I}_j$, kde $\mathcal{I}_j = [a_j, b_j]$, $j = 1, \dots, N$, a nechť funkce f_1, \dots, f_N vyhovují následujícím podmínkám:

- (i) f_1, \dots, f_N jsou definovány a jsou spojité na \mathcal{I} ;
- (ii) Pro každý prvek $\mathbf{x} \in \mathcal{I}$ též vektor $\mathbf{f}(\mathbf{x}) \in \mathcal{I}$;
- (iii) Existuje konstanta L , $0 < L < 1$, taková, že pro libovolná $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{I}$ platí

$$(4.4) \quad \|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|.$$

Potom platí následující výroky

- (a) Rovnice (4.2) má právě jedno řešení $\mathbf{x}^* \in \mathcal{I}$.
- (b) Pro libovolný výchozí prvek $\mathbf{x}_0 \in \mathcal{I}$ posloupnost $\{\mathbf{x}_n\}$ určená v algoritmu 4.1 je definována pro každé n a konverguje k \mathbf{x}^* .
- (c) Pro libovolné $n \geq 1$ platí odhad

$$(4.5) \quad \|\mathbf{x}_n - \mathbf{x}^*\| \leq \frac{L^n}{1-L} \|\mathbf{x}_1 - \mathbf{x}_0\|.$$

Definice 4.4 Konstanta L , požadovaná v (4.4) se nazývá Lipschitzovou konstantou (vzhledem k normě $\|\cdot\|$) vektorové funkce \mathbf{f} . Vztah (4.4) se nazývá Lipschitzovou podmírkou pro \mathbf{f} vzhledem k normě $\|\cdot\|$.

Poznámky. Podmínky kladené ve větě 4.1 na Lipschitzovu konstantu L napovídají interpretaci metody. Vzdálenost obrazů je při takovém zobrazení \mathbf{f} ostře majorizována vzdáleností vzorů, což dává takovému zobrazení jméno: *kontraktivní zobrazení*, čili *kontrakce*.

Je na místě podotknout, že tvrzení (a) je v případě $N \geq 2$ daleko hlubším výrokem než je tomu pro $N = 1$, kdy existenci řešení lze nahlédnout na př. vyšetřením grafu zkoumané funkce. Tvrzení (b) vypovídá o konvergenci algoritmu, zatímco v (c) je dán horní odhad

chyby.

Důkaz věty 4.1. Z požadavku (ii) plyne, že posloupnost $\{\mathbf{x}_n\}$ je definována pro každé $n \geq 1$ a $\mathbf{x}_0 \in \mathcal{I}$.

Na základě (iii) odvodíme snadno, že

$$(4.6) \quad \|\mathbf{x}_{n+1} - \mathbf{x}_n\| \leq L^n \|\mathbf{x}_1 - \mathbf{x}_0\|, \quad n = 0, 1, \dots$$

Fixujme n pevně a vezměme $m > n$, $m = n + p$, $p \geq 1$. Vidíme, že

$$\mathbf{x}_m - \mathbf{x}_n = \sum_{k=1}^p (\mathbf{x}_{n+k} - \mathbf{x}_{n+k-1}),$$

a tedy,

$$\|\mathbf{x}_m - \mathbf{x}_n\| \leq \sum_{k=1}^p \|\mathbf{x}_{n+k} - \mathbf{x}_{n+k-1}\|.$$

Použitím (4.6) obdržíme vztahy

$$(4.7) \quad \begin{aligned} \|\mathbf{x}_m - \mathbf{x}_n\| &\leq \sum_{k=1}^p L^{n+k-1} \|\mathbf{x}_1 - \mathbf{x}_0\| \\ &\leq \frac{L^n}{1-L} \|\mathbf{x}_1 - \mathbf{x}_0\|. \end{aligned}$$

Protože $\mathcal{E} = R^N$ je úplný (Banachův) prostor, plyne odtud konvergence posloupnosti $\{\mathbf{x}_n\}$.

Nechť

$$\hat{\mathbf{x}} = \lim_{n \rightarrow \infty} \mathbf{x}_n,$$

což je totéž co

$$\lim_{n \rightarrow \infty} \|\mathbf{x}_n - \hat{\mathbf{x}}\| = 0.$$

Vzhledem ke spojitosti funkcí f_1, \dots, f_N obdržíme z (4.7)

$$\lim_{n \rightarrow \infty} \mathbf{f}(\mathbf{x}_n) = \mathbf{f}(\hat{\mathbf{x}})$$

a tudíž

$$\hat{\mathbf{x}} = \lim_{n \rightarrow \infty} \mathbf{x}_{n+1} = \lim_{n \rightarrow \infty} \mathbf{f}(\mathbf{x}_n) = \mathbf{f}(\hat{\mathbf{x}}),$$

a tak $\hat{\mathbf{x}} = \mathbf{x}^*$ je řešením soustavy rovnic (4.2).

Jednoznačnost se dokazuje zcela analogicky jako v případě $N = 1$.

Nechť \mathbf{x}_1^* a \mathbf{x}_2^* jsou dvě řešení soustavy (4.2). Potom, na základě (4.4)

$$\|\mathbf{x}_1^* - \mathbf{x}_2^*\| = \|\mathbf{f}(\mathbf{x}_1^*) - \mathbf{f}(\mathbf{x}_2^*)\| \leq L \|\mathbf{x}_1^* - \mathbf{x}_2^*\|,$$

což je možné pouze když $\mathbf{x}_1^* = \mathbf{x}_2^*$.

Platnost (4.5) plyne z (4.7) pro $p \rightarrow \infty$:

$$\|\mathbf{x}_n - \mathbf{x}^*\| = \lim_{n \rightarrow \infty} \|\mathbf{x}_n - \mathbf{x}_{n+p}\| \leq \frac{L^n}{1-L} \|\mathbf{x}_1 - \mathbf{x}_0\|.$$

Tím je důkaz věty 4.1 proveden. |||

4.2 Kvadratická konvergencie a Newtonova metoda pro soustavy

Zabýejme se otázkou explicitního vyjádření vektoru chyby

$$\mathbf{d}_n = \mathbf{x}_n - \mathbf{x}^*$$

při $n \rightarrow \infty$.

Předpokládejme, že \mathbf{f} splňuje předpoklady věty 4.1 a navíc, že funkce f_1, \dots, f_N mají v \mathcal{I} spojité parcální derivace až do rádu 2 včetně. Používáme značení

$$(f_j)_{x_k}(x_1, \dots, x_N) \equiv \frac{\partial f_j}{\partial x_k}(x_1, \dots, x_N), \quad j, k = 1, \dots, N$$

a

$$\begin{aligned} (f_j)_{x_k, x_l}(x_1, \dots, x_N) &\equiv \frac{\partial^2 f_j}{\partial x_k \partial x_l}(x_1, \dots, x_N) \\ &= \frac{\partial(f_j)_{x_k}}{\partial x_l}(x_1, \dots, x_N), \quad j, k, l = 1, \dots, N \end{aligned}$$

Aplikací Taylorovy věty pro funkce N proměnných zjistíme, že

$$(4.8) \quad \mathbf{d}_{n+1} = \mathbf{x}_{n+1} - \mathbf{x}^* = \mathbf{f}(\mathbf{x}_n) - \mathbf{f}(\mathbf{x}^*) = \mathbf{f}(\mathbf{x}^* + \mathbf{d}_n) - \mathbf{f}(\mathbf{x}_n) = \mathcal{J}(\mathbf{x}^*)\mathbf{d}_n + O(\|\mathbf{d}_n\|^2),$$

kde

$$(4.9) \quad \mathcal{J}(\mathbf{x}) = \begin{pmatrix} (f_1)_{x_1}(x_1, \dots, x_N) & \dots & (f_1)_{x_N}(x_1, \dots, x_N) \\ \dots & \dots & \dots \\ (f_N)_{x_1}(x_1, \dots, x_N) & \dots & (f_N)_{x_N}(x_1, \dots, x_N) \end{pmatrix}$$

a $O(\|\mathbf{d}_n^2\|)$ značí výrazy mající majorantu tvaru $c \|\mathbf{d}_n\|^2$, při čemž c je konstanta nezávislá jak na n tak na $\mathbf{x} \in \mathcal{I}$.

Definice 4.5 Matice definovaná v (4.9) se nazývá Jacobiovou maticí soustavy (4.1) resp. (4.2).

Všimněme si toho, že vztah (4.8) je N -dimenzionálním analogem vztahu (3.12).

Je-li $\mathcal{J}(\mathbf{x}^*)$ nenulová matice, pak (4.8) říká, že vektor chyby se v každém iteračním kroku násobí Jacobiovou maticí (4.9), což odpovídá *lineární konvergenci*. Je-li $\mathcal{J}(\mathbf{x}^*) = 0$, je z (4.8) patrné, že chyba po $1+n$ -tém kroku je úměrná čtverci chyby po n -tému kroku - v tom případě máme co do činění s *kvadratickou konvergencí*.

Opět můžeme formulovat analog věty 3.3 pro $N > 1$.

Věta 4.2 Nechť funkce f_1, \dots, f_N definované na \mathcal{I} splňují na \mathcal{I} tyto požadavky:

(i) Existují spojité derivace

$$\frac{\partial f_j}{\partial x_l}, \quad 1 \leq j, l \leq N;$$

(ii) Soustava (4.1) má uvnitř \mathcal{I} řešení \mathbf{x}^* takové, že $\mathcal{J}(\mathbf{x}^*) = 0$.

Potom existuje číslo $\delta > 0$ takové, že posloupnost definovaná pomocí algoritmu 4.1 konverguje k \mathbf{x}^* pro libovolný startovací vektor \mathbf{x}_0 , pro nějž

$$\|\mathbf{x}_0 - \mathbf{x}^*\| < \delta.$$

Tedy opět, jako pro $N = 1$, je-li výchozí přiblížení dostatečně blízké přesnému řešení, dochází ke konvergenci.

Globální konvergence je pro vícerozměrný případ problém principiálně komplikovanější než pro $N = 1$ (konvexita ev. konkavita řešení se jeví pro $N = 1$ dost akademicky). Podmínky zaručující globální konvergenci jsou dost složité a jejich analýza patří mezi speciální partie numerických metod.

Hledejme řešení soustavy

$$(4.10) \quad \mathbf{F}(\mathbf{x}) = 0,$$

kde $\mathbf{F} = (F_1, \dots, F_N)^T$.

Algoritmus 4.2 (*Newtonova metoda pro soustavy*) Volme $\mathbf{x}_0 \in \mathcal{I}$ a sestrojme posloupnost $\{\mathbf{x}_n\}$ podle schématu

$$(4.11) \quad \mathbf{x}_{n+1} = \mathbf{x}_n - \mathbf{h}_n,$$

kde $\mathbf{h}_n = (h_1, \dots, h_N)^T$ je řešením soustavy

$$(4.12) \quad \mathcal{J}(\mathbf{x}_n)\mathbf{h}_n = \mathbf{F}(\mathbf{x}_n)$$

a

$$\mathcal{J} = \begin{pmatrix} \frac{\partial F_1}{\partial x_1} & \dots & \frac{\partial F_1}{\partial x_N} \\ \vdots & \dots & \vdots \\ \frac{\partial F_N}{\partial x_1} & \dots & \frac{\partial F_N}{\partial x_N} \end{pmatrix}$$

Jak je patrné z (4.11) při provádění algoritmu 4.2 musíme při každém iteračním kroku řešiti soustavu lineárních algebraických rovnic (4.12). Jest proto přirozené požadovati

$$(4.13) \quad \det \mathcal{J}(\mathbf{x}_n) \neq 0, \quad n = 0, 1, \dots$$

Cvičení 4.1 Nalezněte matici $H = H(\mathbf{x})$ tak, aby užití postupných approximací na rovnici

$$\mathbf{x} = \mathbf{x} - H(\mathbf{x})\mathbf{F}(\mathbf{x})$$

poskytlo kvadratickou konvergenci. Rozhodněte, zda se vždy obdrží Newtonova metoda.

5 Řešení soustav lineárních algebraických rovnic

5.1 Obecné soustavy a řešení ve smyslu nejmenších čtverců

Budeme se zabývat řešením soustav typu

$$(5.1) \quad Ax = b,$$

kde

$$x \in R^N, b \in R^M,$$

$$A = (a_{jk}), 1 \leq j \leq M, 1 \leq k \leq N.$$

Poznámka. Je zřejmé, že lze vždy docílit situace, kdy $M = N$ a to rozšířením matice A o vhodný počet nulových sloupců pro případ $N < M$ či řádků pro případ $N > M$.

Řešitelnost soustavy v klasickém smyslu, t. j. kdy vztah (5.1) je chápán jako soustava rovností

$$\sum_{k=1}^N a_{jk} x_k = b_j, \quad j = 1, \dots, M,$$

nelze vyžadovat.

Jako příklad může sloužit častý případ, kdy soustava (5.1) reprezentuje nějaký model v němž některé z prvků a_{jk} a b_j se obdrží měřením. Obecně nelze vyloučit, aby odpovídající soustava pro dvě různá měření byla v jednom případě řešitelná a v druhém nikoliv.

V praktických aplikacích se ukazuje potřeba následujícího tvrzení.

Každá soustava typu (5.1) má právě jedno "řešení".

Takové tvrzení nemůže platit pro klasické řešení (podmínka o hodnostech). Aby požadované tvrzení platilo, musíme vhodným způsobem zobecnit pojem řešení. Proto zavedeme *zobecněné řešení*.

Buďte $(\cdot, \cdot)_M$ a $(\cdot, \cdot)_N$ skalární součiny na R^M a R^N .

Definujme

$$(5.2) \quad F(x, x) = (Ax - b, Ax - b)_M$$

a hledejme $\bar{x} \in R^N$ tak, aby

$$(5.3) \quad F(\bar{x}) = \min\{F(x) : x \in R^N\}.$$

Vektor \bar{x} splňující relaci (5.3) se nazývá *zobecněné řešení soustavy* (5.1).

Věta 5.1 Existuje alespoň jeden vektor $\bar{x} \in R^N$ takový, že platí (5.3).

Důkaz. Definujme *obor hodnot matice* A

$$\text{range}(A) = \{u \in \mathcal{R}^M : \exists x \in \mathcal{R}^N \text{ for which } u = Ax\}$$

a *nulovou množinu matice* neboli *jádro* A

$$\text{ker}(A) = \{x \in \mathcal{R}^N : x \in \mathcal{R} \text{ for which } Ax = 0\}.$$

Snadno ověříme, že platí relace

$$(5.4) \quad R^M = \text{range}A \oplus \text{ker}A^*$$

$$(5.5) \quad R^N = \text{range}A^* \oplus \text{ker}A,$$

kde A^* je definována pomocí vztahu $A^*v = y$, při čemž

$$(Ax, v)_M = (x, y)_N$$

pro libovolná $x \in R^N$ a $v \in R^M$.

Existence i jednoznačnost y je důsledkem Rieszovy věty o reprezentaci spojitého lineárního funkcionálu na Hilbertově prostoru $(R^N, (\cdot, \cdot)_N)$.

Na základě (5.4) a (5.5) obdržíme, že platí

$$x = A^*u + v = w + v, \quad Av = 0,$$

a

$$b = Ac + d, \quad A^*d = 0.$$

Odtud pak

$$F(x) = (Aw - Ac, Aw - Ac)_M + (d, d)_M.$$

Vzhledem k tomu, že d je pevný vektor, stačí tedy položiti

$$\bar{x} = c,$$

takže

$$(5.6) \quad F(\bar{x}) = (d, d)_M = \min\{F(x) : x \in R^N\}.$$

Tím je věta 5.1 dokázána.

Poznámka. Protože (5.6) je hodnota globálního minima kvadratické funkce F , je veličina (5.6) táz pro všechna eventuální zobecněná řešení x soustavy $Ax = b$. Tato hodnota se nazývá *cenou* uvedené úlohy najít zobecněné řešení, neboli cenou zobecněného řešení. Všimněme si té skutečnosti, že cena zobecněného řešení x je nulová právě když x je řešení klasické.

Nechť

$$\mathcal{M} = \{x \in R^N : F(x) = (d, d)_M\}.$$

Z věty 5.1 plyne, že $\mathcal{M} \neq \emptyset$ a ze spojitosti F pak, že \mathcal{M} je uzavřená. Díky subadditivitě normy snadno nahlédneme, že \mathcal{M} je konvexní. Existuje tedy jediný prvek $x^* \in \mathcal{M}$ takový, že

$$\|x^*\|_N = \min\{\|x\|_N : x \in \mathcal{M}\}.$$

Zobecněné řešení x^* se nazývá *normálním řešením* soustavy (5.1).

Věta 5.2 Existuje právě jedno normální řešení libovolné soustavy $Ax = b$.

Vidíme tedy, že námi požadované tvrzení uvedené na počátku této kapitoly, platí pro námi zavedené zobecněné řešení.

Poznámka. Zobecněná řešení a tedy též normální řešení jsou podstatně závislá na skalárních součinech $(\cdot, \cdot)_M$ a $(\cdot, \cdot)_N$.

Volíme-li speciálně

$$(5.7) \quad (u, v)_M = \sum_{j=1}^M (u)_j (v)_j,$$

kde $(u)_j$ značí j -tou komponentu vektoru $u \in R^M$, obdržíme funkci F ve tvaru

$$F(x) = \sum_{j=1}^M \left[\sum_{k=1}^N a_{jk}(x)_k - (b)_j \right]^2.$$

Odtud získalo zobecněné řešení takto zavedené název *řešení ve smyslu nejmenších čtverců*. Řešení ve smyslu nejmenších čtverců patří mezi nejrozšířenější. Toto zobecněné řešení bylo zavedeno a velmi zevrubně aplikováno již Gaussem.

Zobecněné inverzní operátory

Budě A matice typu $M \times N$.

Vyšetřujme soustavu relací

- | | |
|-----------------------|----------------------|
| (i) $AXA = A$, | (ii) $XAX = X$, |
| (iii) $(AX)^* = AX$, | (iv) $(XA)^* = XA$. |

Příklad.

Nechť $M = N$ a $\det A \neq 0$.

Položme $X = A^{-1}$.

Splnění vztahů (i) - (iv) s $X = A^{-1}$ je očividné.

Napřed si ukážeme, že matice X typu $N \times N$ splňující vztahy (i) - (iv) existuje nejvýše jedna.

Buďte tedy X_1 a X_2 dvě takové matice. Jest potom

$$\begin{aligned} X_1 &= X_1 AX_1 = X_1 (AX_1)^* = X_1 (AX_2 AX_1)^* = \\ &X_1 (AX_1)^* (AX_2)^* = X_1 AX_1 AX_2 = X_1 AX_2 = \\ &(X_1 A)^* X_2 AX_2 = (X_1 A)^* (X_2 A)^* X_2. \end{aligned}$$

Na druhé straně,

$$(X_2 AX_1 A)^* X_2 = X_1 AX_2 AX_2 = (X_1 A)^* (X_2 A)^* X_2,$$

takže

$$X_1 = (X_1 AX_2 A)^* X_2 = X_2 AX_2 = X_2.$$

Již víme, že v případě $M = N$ a $\det A \neq 0$ jediným řešením soustavy vztahů (i) - (iv) je inverzní matice A^{-1} .

Existuje řešení v obecném případě?

Ano!

Položme

$$T = A^*A.$$

Snadno nahlédneme, že matice T je typu $N \times N$, je symetrická a též je *pozitivně semidefinitní*, t. j.

$$(Tx, x)_N \geq 0.$$

Existují tedy projekce P_j , $j = 1, \dots, s$, $s \leq N$ takové, že platí vztahy

$$(5.8) \quad T = \sum_{k=1}^s \lambda_k P_k, \quad TP_0 = 0,$$

při čemž

$$(5.9) \quad P_j P_k = P_k P_j = \delta_{jk} P_k, \quad P_j^* = P_j, \quad \sum_{k=1}^s P_k = I - P_0, \quad j, k = 0, 1, \dots, N$$

a dále pak

$$(5.10) \quad (T - \lambda_j) P_j = 0, \quad \lambda_j > 0, \quad j = 1, \dots, N, \quad \lambda_0 = 0.$$

Položme

$$(5.11) \quad Z = \sum_{j=1}^N \lambda_j^{-1} P_j A^*.$$

Věta 5.3 Ke každé matici A typu $M \times N$ existuje právě jedna matice typu $N \times M$ taková, že pro ni platí vztahy (i) - (iv).

Definice 5.1 Matice, jejíž existenci a jednoznačnost zaručuje věta 5.3, se nazývá (Mooreovou - Penroseovou) pseudoinverzní maticí a označuje se symbolem A^+ .

Důkaz. Důkaz jednoznačnosti je elementární a lze jej přenechat čtenáři.

Na základě (5.11) a (5.8) snadno zjistíme, že

$$AZA = A \sum_{j=1}^N \lambda_j^{-1} P_j A^* A = A \sum_{j=1}^N \sum_{k=1}^N \lambda_j^{-1} \lambda_k P_j P_k.$$

Z tohoto vyjádření pomocí (5.9) odvodíme, že

$$(5.12) \quad AZA = A(I - P_0).$$

Budě $x \in \mathcal{R}^N$. Podle předpokladu (5.8) platí, že

$$0 = (TP_0x, P_0x)_N = (AP_0x, AP_0x)_N,$$

takže

$$AP_0x = 0.$$

Protože x je libovolný, plyne odtud, že

$$AP_0 = 0.$$

Dokázali jsme tak, že, díky (5.12), matice Z splňuje vztah (i).

Platnost zbývajících vztahů (ii) - (iv) lze dokázat analogicky, což přenecháme čtenáři a pokládáme tím větu 5.3 za dokázanou s tím, že

$$A^+ = Z,$$

při čemž Z je definována v (5.11). \square

Položme

$$(5.13) \quad x^+ = A^+b.$$

Dokážeme, že platí

Věta 5.4 Budět A matice typu $M \times N$ a $b \in \mathcal{R}^M$, $b = b_1 + b_2$, $b_1 \in \text{range}(A)$, $b_2 \in \ker(A)$.

Potom platí rovnost

$$(5.14) \quad x^+ = x^*.$$

Důkaz. Zřejmě

$$AA^+b = b_1$$

a proto

$$F(x^+) = (Ax^+ - b, Ax^+ - b)_M = (AA^+b - b, AA^+b - b)_M = (b_2, b_2)_M.$$

Stačí tedy ukázati, že $x^+ \in R(A^*)$. To však plyne odtud, že

$$P_0A^+ = \sum_{j=1}^s \frac{1}{\lambda_j} P_0P_jA^* = 0$$

a tedy též $P_0x^+ = 0$ a tudíž, $x^+ = (I - P_0)x^+$. Protože $A^+A = I - P_0$ a $(I - P_0)\text{range}(A^*) = \text{range}(A^*)$, je tvrzení dokázáno. \square

Vyšetřujme opět soustavu $Ax = b$, kde A je $M \times N$ matice a $b \in \mathcal{R}^N$ a $b = b_1 + b_2$, při čemž $b_1 \in \text{range}(A)$ t.j. $b_1 = Ac_1$, dále pak a $A^*b_2 = 0$.

Snadno zjistíme, že platí

$$A^*Ax = A^*b = A^*b_1,$$

kterážto soustava se nazývá *normální soustavou*. Tato soustava má vždy klasické řešení. To vyplývá odtud, že $(A^*A)^* = A^*A$ a pro každé řešení y homogenní soustavy $A^*Ay = 0$ platí vztahy

$$(y, A^*b)_N = (y, A^*Ac_1)_N = (A^*Ay, c_1)_N = 0.$$

Hledat řešení normální soustavy rovnic je numericky velmi náročné, protože číslo podmíněnosti

$$\kappa(A^*A) = \frac{s_{\max}}{s_{\min}},$$

kde

$$s_{\min} = \|A^* A\|_2^{1/2}$$

a s_{\min} je odmocnina nejmenší kladné vlastní hodnoty matice $A^* A$, je zpravidla velmi veliké; v případě čtvercové matice A je rovno čtverci čísla podmíněnosti matice A . Připomeňme, že číslo podmíněnosti čtvercové matice A je definováno jakožto výraz

$$\|A\| \|A^{-1}\|.$$

Normální řešení je tedy (klasické) řešení normální soustavy ležící v $\text{range}(A^+)$.

Dále platí

Lemma 5.1 *Každé zobecněné řešení $\bar{x} \in \mathcal{M}$ má tvar*

$$(5.15) \quad \bar{x} = x^+ + y,$$

kde

$$Ay = 0.$$

Důkaz. Budě \bar{x} další řešení soustavy $Ax = b$ a nechť

$$\bar{\bar{x}} = a_1 + a_2, \quad a_1 = A^* c_1, \quad Aa_2 = 0.$$

Potom platí, že

$$\bar{x} - \bar{\bar{x}} = x^+ - a_1 + y - a_2.$$

Protože $Aa_1 = b_1$, odvodíme, že platí vztahy

$$A(x^+ - a_1) = AA^+b_1 - b_1 = 0.$$

Jest tudíž

$$\|A^+b_1 - A^*c_1\|^2 = (A^*\tilde{b}_1 - A^+c_1)_N = (\tilde{b}_1 - c_1, A[(x^+ - a_1)_M] = 0,$$

kde $A^+b_1 = A^*\tilde{b}_1$ a tedy

$$\bar{\bar{x}} = \bar{x} = x^+.$$

Protože $\bar{x} - \bar{\bar{x}} \in \ker(A)$, tvrzení lemmatu je dokázáno. \square

5.2 Pomocné prostředky ke konstrukci speciálních reprezentací matic

Nejprve si zopakujme některé pojmy, jež se ukáží jako vhodné pro konstrukci a vyšetřování určitých algoritmů.

Předpokládáme, že všechny matice v tomto článku jsou obecně komplexní typu $N \times N$. Matice A je *hermiteovsky sdružená*, jestliže

$$A^H = A, \quad A^H = (a_{jk}^*)$$

při čemž

$$a_{jk}^* = \bar{a}_{kj}, \quad j, k = 1, \dots, N,$$

kde pro $\alpha = a + ib$, $a, b \in \mathcal{R}^1$, $i^2 = -1$, $\bar{\alpha} = a - ib$. Je-li matice A reálná, t. j. je-li

$$\bar{a}_{jk} = a_{jk}$$

nazýváme matici hermiteovskou *symetrickou*. Jest totiž

$$A^T = A, \quad A^T = (a_{jk}^T), \quad a_{jk}^T = a_{kj}, \quad j, k = 1, \dots, N.$$

Permutační matice $P = (p_{jk})$ splňuje vztahy

$$p_{jk} = \begin{cases} 0 \\ 1 \end{cases}, \quad \sum_{k=1}^N p_{jk} = \sum_{k=1}^N p_{kj} = 1,$$

a

$$\det P \neq 0.$$

Zřejmě

$$P^H = P^T = P^{-1}.$$

Projekční maticí či *projekcí* se nazývá matice splňující

$$P^2 = P.$$

Speciálně tedy, $I^2 = I$ a $0^2 = 0$.

Permutační matice, která realizuje záměnu l - té a k - té složky vektorůz \mathcal{R}^N , se nazývá *elementární permutační* maticí neboli *transpoziční* maticí.

Taková elementární permutační matice má tvar

$$k \rightarrow \begin{matrix} & \text{k} & & & \text{l} & & \\ & \downarrow & & & \downarrow & & \\ & 1 & \cdot & \cdot & \cdot & 0 & \cdot & \cdot & \cdot & 0 & \cdot & \cdot & \cdot & 0 \\ & \cdot & \cdot & & & \cdot & & & & & & & & \cdot \\ & \cdot & & \cdot & & & & & & & & & & \cdot \\ & \cdot & & & \cdot & & & & & & & & & \cdot \\ & 0 & & & 0 & \cdot & \cdot & \cdot & \cdot & 1 & & & 0 & \\ l \rightarrow & \cdot & & & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & & & & \cdot \\ & 0 & & & 1 & \cdot & \cdot & \cdot & \cdot & 0 & & & 0 & \\ & \cdot & & & & \cdot & & & & & & & & \cdot \\ & \cdot & & & & & \cdot & & & & & & & \cdot \\ & 0 & \cdot & \cdot & \cdot & 0 & \cdot & \cdot & \cdot & 0 & \cdot & \cdot & \cdot & 1 \end{matrix}$$

a při její aplikaci platí vztahy $y = Px$, při čemž

$$y_j = x_j, \quad j \neq k, l,$$

zatímco

$$y_k = x_l \text{ a } y_l = x_k.$$

Unitární matice U je definová na vztahy

$$UU^H = U^H U = I,$$

což pro reálnou matici V značí, že

$$VV^T = V^T V = I.$$

V takovém případě hovoříme o matici *ortogonální*.

Pro *involutorní* matici B platí

$$B^2 = I.$$

Tedy speciálně $I^2 = I$.

Buděte $y \in \mathcal{R}^N$ $v \in \mathcal{R}^M$. Definujeme matici

$$y \otimes v$$

kladouce

$$[y \otimes v] x = (x, v)_N y.$$

Jinými slovy,

$$y \otimes v = y v^T.$$

Matici E_k se nazývá *elementární dolní trojúhelníkovou indexu k*, jestliže

$$E_k = I_N - v \otimes e_k$$

a přitom

$$e_j^T v = 0, \quad j = 1, \dots, k.$$

To značí, že

$$E_k = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & -v_{k+1} & 1 \\ & & & & \ddots \\ & & & & & 1 \\ & & & & & & -v_N \end{pmatrix}.$$

Snadno prověříme, že $y \otimes v$ splňuje následující relaci

$$[y \otimes v]^2 x = (y, v)_N [y \otimes v] x.$$

takže, je-li $(y, v)_N = 0$,

$$[y \otimes v]^2 = 0$$

a tedy,

$$\sigma(y \otimes v) = \{0\}$$

neboli, $y \otimes v$ je *nilpotentní matici*.

V případě dolní elementární matice E_k vidíme, že $(e_k, v)_N = 0$, takže platí

Věta 5.5 Budě $E_k = I_N - v_k \otimes e_k$, při čemž

$$(v_k, e_j)_N = 0, \quad j = 1, \dots, k.$$

Potom existuje inverzní matice E_k^{-1} a platí

$$E_k^{-1} = I_N + v_k \otimes e_k.$$

Věta 5.6 Nechť $x^T = (\xi_1, \dots, \xi_N)$ a nechť $0 \neq \xi_k = e_k^T x$.

Potom existuje jednoznačně určená elementární dolní trojúhelníková matice indexu k E_k taková, že

$$E_k x = (\xi_1, \dots, \xi_k, 0, \dots, 0)^T.$$

Důkaz. Hledejme E_k ve tvaru

$$E_k = I_N - v_k \otimes e_k.$$

Z požadavku elementarity E_k plyne, že

$$(5.16) \quad \nu_1 = \dots = \nu_k = 0,$$

kde

$$v_k^T = (\nu_1, \dots, \nu_N).$$

Dále pak

$$E_k x = x - (x, e_k)_N v = x - \xi_k v$$

a protože žádáme, aby

$$(5.17) \quad \xi_j - \xi_k \nu_j = 0, \quad j = k+1, \dots, N,$$

z předpokladu $\xi_k \neq 0$ odvodíme vztahy

$$(5.18) \quad \nu_j = \frac{\xi_j}{\xi_k}, \quad j = k+1, \dots, N.$$

Tedy, E_k existuje a je vztahy (5.16), (5.17), (5.18) určena jednoznačně. |||

Householderovu matici zavedeme pomocí vztahů

$$H = I - 2ww^H, \quad w^H w = 1,$$

kde w je (obecně komplexní) sloupcový vektor a w^H odpovídající vektor řádkový.

Zřejmě Householderova matice je hermiteovsky sdružená $H^H = H$ a též unitární, neboť

$$H^H H = H^2 = (I - 2ww^H)^2 = I.$$

Všimněme si, že pro $x \in C^N$ a $y = Hx$ platí, že

$$y = x - 2w^H x w,$$

takže

$$y^H y = (y, y) = (Hx, Hx) = (H^H H x, x) = (x, x)$$

a

$$(x, y) = \overline{(y, x)}.$$

Lemma 5.2 K danému vektoru $v \neq 0$ existuje ortogonální matici Q taková, že

$$(5.19) \quad Qv = -\sigma \|v\| e_1,$$

kde

$$\sigma = \begin{cases} +1 & \text{pro } v_1 = (v, e_1) \geq 0 \\ -1 & \text{pro } v_1 < 0. \end{cases}$$

Důkaz. Nechť

$$u = v + \sigma \sqrt{(v, v)} e_1$$

a

$$Q = I - 2 \frac{uu^H}{u^H u}.$$

Zřejmě

$$Q^H = I - 2 \frac{uu^H}{u^H u} = Q$$

a

$$Q^H Q = Q^2 = I.$$

Dále pak

$$Qv = Qu - \sigma \|v\| Qe_1.$$

Avšak

$$Qu = -u$$

a

$$Qe_1 = e_1 - 2\sigma \|v\| u(u)_1 \frac{1}{(u, u)}.$$

Odtud plyne, že

$$Qv = -u - \sigma \|v\| e_1 + 2\sigma^2 \|v\|^2 \frac{(u)_1}{(u, u)} u,$$

a protože

$$(u)_1 = (v)_1 + \sigma \|v\|,$$

zjistíme, že

$$-u + 2\sigma \|v\| \frac{(v)_1 + \sigma \|v\|}{(u, u)} = \frac{1}{(u, u)} \left\{ -2\sigma^2 (v, v) - 2\sigma \|v\| (v)_1 + 2\sigma^2 \|v\|^2 + 2\sigma \|v\| (v)_1 \right\} = 0.$$

Tedy

$$Qv = -\sigma \|v\| e_1.$$

□

Matici $H = (h_{jk})$ taková, že

$$h_{jk} = 0 \text{ pro } j - k > 1 \text{ a } j, k = 1, \dots, N - 2,$$

se nazývá *horní Hessenbergova*.

Názorně je patrně rozmištění prvků na následujícím schématu

$$\begin{pmatrix} \times & \times & \times & \cdot & \cdot & \cdot & \times & \times \\ \times & \times & \times & \cdot & \cdot & \cdot & \times & \times \\ 0 & \times & \times & \cdot & \cdot & \cdot & \times & \times \\ 0 & 0 & \times & \cdot & \cdot & \cdot & \times & \times \\ \cdot & \cdot \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & \times & \times \end{pmatrix}.$$

Lemma 5.3 Ke každé matici A existuje unitární matice T , taková, že

$$A = THT^H,$$

kde matice H je horní Hessenbergova.

Čtvercová matice $A = (a_{jk})$ $a_{jk} \in \mathcal{R}$ $j, k = 1, \dots, N$, se nazývá slabě diagonálně dominantní, jestliže platí

$$(5.20) \quad \sum_{k=1, k \neq j}^N |a_{jk}| \leq |a_{jj}|, \quad j = 1, \dots, N,$$

při čemž existuje alespoň jeden index j_0 , $1 \leq j_0 \leq N$, tak, že

$$(5.21) \quad \sum_{k=1, k \neq j_0}^N |a_{j_0k}| < |a_{j_0j_0}|.$$

Nastane-li v (5.20) ostrá nerovnost pro všechny indexy $j = 1, \dots, N$, pak se matice A nazývá diagonálně dominantní.

Věta 5.7 (Geršgorinova) Budě $A = (a_{jk})$ matice $N \times N$, při čemž a_{jk} jsou komplexní čísla. Dále budě λ vlastní hodnota matice A . Potom existuje index j_0 , $1 \leq j_0 \leq N$ takový, že platí vztahy

$$(5.22) \quad |\lambda - a_{j_0j_0}| \leq \sum_{k=1, k \neq j_0}^N |a_{j_0k}|.$$

Důkaz. Budě $x \in \mathcal{R}^N \oplus i \mathcal{R}^N$ vlastní vektor odpovídající vlastní hodnotě λ . Tedy,

$$Ax = \lambda x, \quad x \neq 0.$$

Položme $|x_{j_0}| = \max\{|x_j| : j = 1, \dots, N\}$. Jest potom

$$\begin{aligned} |\lambda - a_{j_0}| &\leq \sum_{k=1, k \neq j_0}^N |a_{j_0k}| \frac{|x_k|}{|x_{j_0}|} \\ &\leq \sum_{k=1, k \neq j_0}^N |a_{j_0k}|. \end{aligned}$$

Je tak dokázána platnost relace (5.22) a tím i věta 5.7. |||

Důsledek 5.1 Každá diagonálně dominantní matice A je regulární.

Důkaz. Stačí si uvědomit, že čtvercová matice je regulární právě když 0 není její vlastní hodnotou. To, že 0 není vlastní hodnotou dané čtvercové matice však je bezprostředním důsledkem diagonální dominance. Kdyby totiž 0 byla vlastní hodnotou matice A , platily by vztahy (5.22) pro $\lambda = 0$,

$$|a_{j_0 j_0}| = |\lambda - a_{j_0 j_0}| \leq \sum_{k=1, k \neq j_0}^N |a_{j_0 k}| < |a_{j_0 j_0}|.$$

Tento spor dokazuje platnost tvrzení. |||

Věta 5.8 Budě $0 \neq A = (a_{jk})$, diagonálně dominantní a budě $E_1 = I_N - v_1 \otimes e_1$ elementární dolní trojúhelníková matice indexu 1, kde

$$v_1 = \begin{pmatrix} 0 \\ a_{21} / a_{11} \\ \vdots \\ \vdots \\ a_{N1} / a_{11} \end{pmatrix}.$$

Potom

$$A^{(1)} = E_1 A$$

je též diagonálně dominantní.

Důkaz. Vyšetřujme výraz a_{jk}^1 rovný dle definice,

$$a_{jk}^1 = a_{jk} - \frac{a_{j1}}{a_{11}} a_{1k}, \quad 1, \dots, N.$$

Jest tedy,

$$\begin{aligned} & |a_{jj}^1| - \sum_{k=2, k \neq j}^N |a_{jk}^1| \\ & \geq |a_{jj}| - \frac{1}{|a_{11}|} |a_{j1} a_{1j}| - \sum_{k=2, k \neq j}^N |a_{jk}| - \frac{1}{|a_{11}|} \sum_{k=2, k \neq j}^N |a_{j1} a_{1k}| \\ & = |a_{jj}| - \sum_{k=2, k \neq j}^N |a_{jk}| - \frac{1}{|a_{11}|} \sum_{k=1, k \neq j}^N |a_{j1} a_{1k}| \\ & \geq |a_{jj}| - \sum_{k=1, k \neq j}^N |a_{jk}| \\ & \geq 0, \quad j = 2, \dots, N. \end{aligned}$$

|||

Čtvercová matice $A = (a_{jk})$ nad tělesem komplexních čísel se nazývá *pozitivně definitní*, jestliže existuje konstanta $\kappa > 0$ taková, že

$$(Ax, x)_N \geq \kappa(x, x)_N \quad \forall x \in \mathcal{C}.$$

Jsou-li $a_{jk} \in \mathcal{R}$, pak navíc požadujeme, aby $A = A^T$.

5.3 Přímé metody řešení soustav lineárních algebraických rovnic

Pod přímými metodami řešení soustav lineárních algebraických rovnic rozumíme takové metody, které jsou založeny na algoritmech, jež poskytují přesné řešení po provedení jistého konečného počtu aritmetických operací.

Je zřejmé, že pro soustavy, jejichž matici jsou trojúhelníkové - horní či dolní - se nabízejí půrozené přímé metody. Jsou to t. zv. *zpětná* a *přímé dosazení*.

Vyšetřujme soustavy

$$(5.23) \quad Ly = c,$$

a

$$(5.24) \quad Ux = b,$$

kde

$$L = \begin{pmatrix} l_{11} & \cdot & \cdot & \cdot & 0 \\ l_{21} & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ l_{N1} & \cdot & \cdot & \cdot & l_{NN} \end{pmatrix}$$

a

$$U = \begin{pmatrix} u_{11} & u_{12} & \cdot & \cdot & \cdot & u_{1N} \\ 0 & u_{22} & \cdot & \cdot & \cdot & u_{2N} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & u_{NN} \end{pmatrix}.$$

Nechť

$$l_{jj} \neq 0, \quad j = 1, \dots, N$$

a

$$u_{jj} \neq 0, \quad j = 1, \dots, N.$$

Potom

$$(5.25) \quad y_1 = \frac{c_1}{l_{11}}$$

a

$$(5.26) \quad y_k = \frac{c_k}{l_{kk}} - \sum_{j=1}^{k-1} \frac{l_{kj}c_j}{l_{jj}l_{kk}}, \quad k = 2, \dots, N,$$

a podobně

$$(5.27) \quad x_N = \frac{b_N}{u_{NN}},$$

a

$$(5.28) \quad x_{N-k} = \frac{b_{N-k}}{u_{N-k,N-k}} - \sum_{j=k+1}^N \frac{u_{N-k,j}b_j}{u_{N-k,N-k}u_{jj}}, \quad k = 1, \dots, N-1.$$

Nejznámější a přitom nejdůležitější přímou metodou řešení soustav lineárních algebraických rovnic je *Gaussova eliminaci metoda*. Její základní myšlenka spočívá v transformaci původní soustavy s maticí A na soustavu s trojúhelníkovou maticí \tilde{A} . Matici \tilde{A} lze získat postupným

eliminováním, jež lze reprezentovat elementárními eliminačními maticemi popsanými v odstavci 5.2.

Buď tedy $A = (a_{jk})$ daná regulární matice s $a_{11} \neq 0$ a

$$E_1 = I - v_1 \otimes e_1$$

elementární dolní trojúhelníková matice indexu 1, kde

$$v_1 = \frac{1}{a_{11}} \begin{pmatrix} 0 \\ a_{12} \\ \vdots \\ \vdots \\ a_{N1} \end{pmatrix}.$$

Jest tedy

$$A^{(0)} = A$$

a

$$A^{(1)} = \begin{pmatrix} a_{11} & a_{12} & \dots & \dots & a_{1N} \\ 0 & a_{22} - \frac{a_{12}}{a_{11}}a_{21} & \dots & \dots & a_{2N} - \frac{a_{1N}}{a_{11}}a_{21} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & a_{N2} - \frac{a_{12}}{a_{11}}a_{N1} & \dots & \dots & a_{NN} - \frac{a_{1N}}{a_{11}}a_{N1} \end{pmatrix}.$$

Zvláštní pozornosti si zasluzuje skutečnost, že první řádek transformované matice A^1 je totožný s prvním řádkem původní matice A .

Obecně, buď $1 < k \leq N - 1$, $a_{kk}^{k-1} \neq 0$, kde

$$A^{(k-1)} = \begin{pmatrix} a_{11}^{(k-1)} & \dots & \dots & \dots & \dots & a_{1N}^{(k-1)} \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & a_{kk}^{(k-1)} & \dots & a_{(kN)} \\ \vdots & & & \vdots & & \vdots \\ 0 & \dots & \dots & a_{Nk}^{k-1} & \dots & a_{NN}^{(k-1)} \end{pmatrix},$$

při čemž řádky s indexy $1 \leq j \leq k$, jsou totožné pro všechny matice A^l , $l = j, \dots, k - 1$.

Nechť

$$E_k = I_N - v_k \otimes e_k,$$

kde

$$v_k = \begin{pmatrix} a_{k+1,k}^{(k-1)} / a_{kk}^{(k-1)} \\ \vdots \\ a_{Nk}^{(k-1)} / a_{kk}^{(k-1)} \end{pmatrix}.$$

Potom

$$A^{(k)} = E_k A^{(k-1)} = \begin{pmatrix} a_{11}^{(k)} & \dots & \dots & a_{1k}^{(k)} & a_{1k+1}^{(k)} & \dots & \dots & a_{1N}^{(k)} \\ \vdots & & & & & & & \\ 0 & & & a_{kk}^{(k)} & a_{kk+1}^{(k)} & \dots & \dots & a_{kN}^{(k)} \\ & & & & & & & \\ 0 & \dots & \dots & 0 & a_{Nk+1}^{(k)} & \dots & \dots & a_{NN}^{(k)} \end{pmatrix},$$

při čemž řádky s indexy $1 \leq j \leq k+1$ jsou totožné pro matice $A^{(l)}$, kde $l = j, \dots, k$.

Právě provedený indukční krok ukazuje, že ve "fyziologickém" případě, t. j. když

$$a_{kk}^{(k-1)} \neq 0,$$

výše uvedený proces vede po konečném počtu kroků, $k \leq N-1$, k situaci, kdy $A^{(N-1)}$ je horní trojúhelníková matice.

Hledané řešení původní soustavy je evidentně totožné s řešením soustavy

$$A^{(N-1)}x = b^{(N-1)},$$

kde

$$b^{(N-1)} = E_{N-1} \dots E_1 b.$$

Schematicky lze tuto skutečnost postihnout tak, že sestrojíme obdélníkovou matici typu $N \times N+1$

$$B = (A, b)$$

a provádíme postupně transformace s touto maticí, tedy konstruujeme postupně matice

$$B^{(k)} = E_k B^{(k-1)}, \quad k = 1, \dots,$$

při čemž

$$B^{(0)} = B.$$

Po $N-1$ krocích obdržíme matici B^{N-1} mající tvar

$$B^{(N-1)} = \begin{pmatrix} \times & \times & \cdot & \cdot & \cdot & \times & \cdot & \cdot & \cdot & \times \\ 0 & \times & \cdot & \cdot & \cdot & \times & \cdot & \cdot & \cdot & \times \\ \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & \times & \cdot & \cdot & \cdot & \times \end{pmatrix}.$$

Je vcelku zřejmé, jak uvedený postup zobecnit na případ současného řešení soustav s danou maticí soustavy a s více pravými stranami b_1, \dots, b_p , $p \geq 1$. V takovém případě

$$B_{(p)} = (A, b_1, \dots, b_p)$$

s

$$B_{(p)}^{(k)} = E_k B_{(p)}^{(k-1)},$$

kde

$$B_{(p)}^{(0)} B_{(p)}.$$

Již víme, že řešení soustav

$$Ax = B_j, \quad j = 1, \dots, p,$$

se převádí na řešení soustav

$$A^{(N-1)}x = b_j^{(N-1)}, \quad j = 1, \dots, p,$$

kde b_j je $N+j$ -tý sloupec matice $B_{(p)}^{(N-1)}$. Matice těchto soustav jsou však horní trojúhelníkové a nalezení hledaných řešení je záležitostí rutinní. To vše za předpokladu, že

$$a_{kk}^{(k-1)} \neq 0.$$

Tento předpoklad je rozhodující. Numerická analýza spojená s tímto předpoladem si vynutí posléze úpravu vlastního algoritmu řešení soustav lineárních algebraických rovnic - t. zv. *pivotní strategie*.

Ještě než se budeme numerickými aspekty eliminačních metod zabývat, uvedme složitost výpočtu dle výše uvedeného schematu za předpokladu, že proces je algoritmicky proveditelný, t.j.

$$a_{kk}^{(k-1)} \neq 0, \quad k = 1, \dots, N-1.$$

Transformace matice $B_{(p)}$ ($= B_{(p)}^{(0)}$) do lichoběžníkového tvaru *stojí* asymptoticky

$$\sum_{j=1}^N (N-j+p)(N-1) \sim \frac{1}{3}N^3 + \frac{1}{2}pN^2.$$

Je-li na daném počítači časová spotřeba na jednu aritmetickou operaci ν sekund, pak stanovení řešení dané soustavy s maticí typu $N \times N$ pro p pravých stran se realizuje v následujících časových mezích:

N	$\frac{1}{3}N^3$ (eliminace)	$\frac{1}{2}N^2$ (zpětná substituce)
10	$3, 3 \cdot 10^{-4} \nu$ sek	$5 \cdot 10^{-5} \nu$ sek
10^2	$3, 3 \cdot 10^{-1} \nu$ sek	$5 \cdot 10^{-2} \nu$ sek
10^3	$3, 3 \cdot 10^5 \nu$ sek	$5 \cdot 10^{-1} \nu$ sek
10^6	$3, 3 \cdot 10^{12} \nu$ sek	$5 \cdot 10^5 \nu$ sek

Poznámka 5.1 Jistě stojí za povšimnutí, že spotřeba času uvedená v posledním řádku činí řádově $\nu \cdot 10^9$ hodin! To ukazuje, že pro velké dimenze jsou metody eliminační na soudobých seriálních počítacích prakticky nepoužitelné.

Uvedená poznámka osvětuje nejrůznější snahy obejít komplikace spojené s časovou náročností eliminačních metod. Tendence jsou zřejmé - využít bud speciální struktury počítáče (víceprocesorové) a nebo speciální struktury dané třídy matic řešených soustav (na př. pozitivně definitní a pásové).

Pivotní strategie

Snadno nahlédneme, že předpoklad nenulovosti diagonálních prvků $a_{kk}^{(k-1)}$, $k = 1, \dots, N-2$, je podstatný a že algoritmus eliminace se zhroutí při jeho nesplnění. Nepříjemné je též, že totéž lze říci o numerické realizaci, t. j. eliminační algoritmy se zhroutí, pakliže uvedené prvky $a_{kk}^{(k-1)}$ nejsou dostatečně velké v absolutní hodnotě.

Demonstrujme takovou situaci na příkladě následující soustavy.

Nechť tedy

$$(5.29) \quad \begin{cases} x_1 + x_2 + x_3 = 1 \\ x_1 + x_2 + 2x_3 = 2 \\ x_1 + 2x_2 + 2x_3 = 1 \end{cases}$$

Matrice této soustavy je regulární a existuje tudíž jediné řešení

$$x_1 = -x_2 = x_3 = 1.$$

Avšak po provedení prvního kroku eliminace obdržíme

$$(5.30) \quad \begin{cases} & x_3 = 1 \\ x_2 + x_3 = 0, \end{cases}$$

a dále eliminovat nelze

$$(5.31) \quad a_{22}^{(1)} = 0.$$

V obecném případě kdy $\det A \neq 0$ a přitom

$$a_{kk}^{k-1} = 0,$$

musí existovat pro nějaké $j > k$ prvek

$$a_{kj}^{(k-1)} \neq 0,$$

jinak by k -tý řádek byl nulový, což vylučuje regularitu matice soustavy. Vhodnou záměnou řádků či sloupců soustavy lze posléze docílit toho, aby po takové úpravě diagonální prvek byl roven vhodně vybranému prvku k -tého řádku.

Tato úvaha má oprávnění i v situaci, kdy posuzujeme nenulovost z pohledu daného počítáče. Jinými slovy, permutaci proměnných či řádků je nutné provádět kdykoliv absolutní hodnota prvku $a_{kk}^{(k-1)}$ není dostatečně velká.

Dokumentujme takový případ soustavou

$$(5.32) \quad \left\{ \begin{array}{l} x_1 + 0,0001x_2 + 9,999x_3 = 1 \\ 0,0001x_2 + x_3 = 1 \\ 9,999x_3 = 10 \end{array} \right.$$

Matice této soustavy je regulární horní trojúhelníková. Provedeme-li zpětnou substituci užívající tří platných cifer v pohyblivé desetinné čárce, zjistíme, že

$$\tilde{x}_1 = 0, \tilde{x}_2 = 0, \tilde{x}_3 = 1$$

splňují požadované vztahy. Avšak přesné řešení (vzaté se třemi polatnými ciframi) je dáno výrazy

$$x_1 = 1, -x_2 = x_3 = 1.$$

Jak vysvětlit tento rozpor tedy numerickou nestabilitu? Odpověď je smutná, je to v podstatě věci: *Gaussova eliminační metoda je numericky nestabilní*.

V souvislosti s předchozím příkladem si všímněme toho, že současnou záměnou 2. a 3. řádku a 2. a 3. sloupce, t. j.

$$y_1 = x_1, y_2 = x_3, y_3 = x_2,$$

se obdrží soustava

$$\begin{aligned} y_1 + & y_2 + y_3 = 1 \\ 9,999y_2 & = 10 \\ y_2 + & 0,0001y_3 = 1. \end{aligned}$$

Zde již

$$y_1 = y_2 = -y_3 = 1,$$

takže po zpětné substituci $y \rightarrow x$ obdržíme přesné řešení v rámci tří platných cifer.

Částečná pivotace

Hledejme přirozené číslo t tak, aby

$$t = \min\{j : k \leq j \leq N\}$$

a přitom

$$|a_{tk}^{(k-1)}| = \max \left\{ |a_{jk}^{(k-1)}| : j = 1, \dots, N \right\}.$$

Úplná pivotace

Hledejme přirozená čísla t a s taková, aby

$$\begin{aligned} t &= \min \{j : k \leq j \leq N\} \\ s &= \min \{j : k \leq l \leq N\} \end{aligned}$$

a přitom

$$|a_{ts}^{(k-1)}| = \max \left\{ |a_{jl}^{(k-1)}| : k \leq j \leq N, k \leq l \leq N \right\}.$$

Důležité doporučení plynoucí z praktických zkušeností praví, že částečná pivotáz je postačující; úplná pivotáz jednak výpočet zpomaluje, jednak numerickou stabilitu podstatně nezvyšuje.

Formální výraz předchozích úvah je obsahem následujících vět.

Věta 5.9 (o úplné pivotaci) Budě A matici typu $M \times N$ a nechť $A \neq 0$.

Potom existují celá kladná čísla $t \geq 1$ a $s \leq \min(M-1, N)$, elementární permutační matice P_j, Q_l , $j = 1, \dots, M-1$, $l = 1, \dots, N$ a elementární dolní trojúhelníkové matice E_j indexu $N-j$, $j = 1, \dots, s$, tak, že platí vyjádření

$$\begin{aligned}\hat{A} &= E_s P_s \dots E_1 P_1 A Q_1 \dots Q_t \\ &= \begin{pmatrix} \hat{A}_{11} & \hat{A}_{12} \\ 0 & 0 \end{pmatrix},\end{aligned}$$

při čemž \hat{A}_{11} je horní trojúhelníková matica typu $r \times r$ s

$$\det \hat{A} \neq 0.$$

Dále pak r je hodnota matici A .

Věta 5.10 (o částečné pivotaci) Budě A matici typu $M \times N$ a nechť $s = \min(M-1, N)$.

Potom existují elementární permutační matice P_j a elementární dolní trojúhelníkové matice E_j indexu $M-j$, $j = 1, \dots, s$ takové, že platí vyjádření

$$\tilde{A} = E_s P_s \dots E_1 P_1 A,$$

při čemž \tilde{A} je dolní lichoběžníková, t. j.

$$\tilde{a}_{jk} = 0, \quad j > k.$$

Již víme, že při řešení soustav typu

$$Ax = b,$$

kde A je matici typu $N \times N$, Gaussovou eliminační metodou je nutné provádět pivotaci.

Tedy, budě podle schématu úplné pivotace

$$A = E_{N-1} P_{N-1} \dots E_1 P_1 A Q_1 Q_{N-1}$$

nebo užitím pivotace částečné

$$A = E_{N-1} P_{N-1} \dots E_1 P_1 A.$$

Obě tyto strategie lze bez obtíží přenést na případ simultanního řešení p soustav typu

$$Ax = b_j, \quad j = 1, \dots, p.$$

V takovém případě se potřebné transformace provádějí s maticí

$$B_p = (A, b_1, \dots, b_p).$$

Zde je zřejmě podstatný předpoklad o nezávislosti vektoru b_j na řešení soustav s indexy $l < j$. Nejsou však řídké případy, kdy

$$b_{j+1} = b_{j+1}(x_1, \dots, x_p)$$

a v tom případě naznačený postup selže.

Věta 5.11 *Budě A matici typu $N \times N$ a A_k nechť označují matici typu $k \times k$ tvořenou prvky ležícími v průniku řádků a sloupců s indexy $1, \dots, k$.*

Je - li

$$(5.33) \quad \det A_k \neq 0, \quad \forall k = 1, \dots, N-1,$$

pak existuje právě jedna dolní trojúhelníková matici $L = (l_{jk})$ typu $N \times N$ taková, že

$$l_{jj} = 1, \quad j = 1, \dots, N,$$

a právě jedna horní trojúhelníková matici $U = (u_{jk})$ typu $N \times N$ taková, že

$$A = L U.$$

Důkaz. Nechť $N = 1$. Potom

$$l_{11} = 1, \quad u_{11} = a_{11}$$

a tvrzení platí. Nechť tvrzení platí pro všechny matice řádu $l \leq k-1$. Vyšetřujme

$$A_k = \begin{pmatrix} A_{k-1} & a \\ c^T & a_{kk} \end{pmatrix},$$

za předpokladu, že

$$A_{k-1} = L_{k-1} U_{k-1}.$$

Položme

$$L_k = \begin{pmatrix} L_{k-1} & 0 \\ l^T & 1 \end{pmatrix},$$

a

$$U_k = \begin{pmatrix} U_{k-1} & u \\ 0^T & u_{kk} \end{pmatrix},$$

kde a a c jsou známé vektory a l je vektor typu $(k-1)$, jenž je zapotřebí určit podobně jako vektor u a číslo u_{kk} .

Abychom určili součin $L_k U_k$, vycházíme z relací

$$(5.34) \quad L_{k-1} U_{k-1} = A_{k-1},$$

$$(5.35) \quad L_{k-1} u = a,$$

$$(5.36) \quad l^T U_{k-1} = c^T$$

a

$$(5.37) \quad l^T u + u_{kk} = a_{kk}.$$

Dle indukčního předpokladu jsou L_{k-1} a U_{k-1} určeny jednoznačně a navíc

$$\begin{aligned} \det U_{k-1} &= \det L_{k-1} \det U_{k-1} \\ &= \det A_{k-1} \neq 0, \end{aligned}$$

takže (5.35) je jednoznačně řešitelná vzhledem k u . Podobně je tomu s vektorem l v (5.36). Díky jednoznačnosti l a u je posléze jednoznačně určen i prvek u_{kk} pomocí (5.37). |||

Není - li předpoklad (5.33) splněn, věta 5.11 neplatí; na př.

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

Předpokládejme, že naopak,

$$A = L U,$$

kde

$$L = \begin{pmatrix} l_{11} & 0 \\ l_{21} & l_{22} \end{pmatrix}, \quad U = \begin{pmatrix} u_{11} & u_{12} \\ 0 & u_{22} \end{pmatrix},$$

takže

$$L U = \begin{pmatrix} l_{11}u_{11} & l_{11}u_{12} \\ l_{21}u_{11} & l_{21}u_{12} + l_{22}u_{22} \end{pmatrix}.$$

odkud plyne, že budě

$$l_{11} = 0$$

nebo

$$u_{11} = 0,$$

ale potom budě

$$L = \begin{pmatrix} 0 & 0 \\ l_{21} & l_{22} \end{pmatrix}$$

a nebo

$$U = \begin{pmatrix} 0 & u_{12} \\ 0 & 0 \end{pmatrix},$$

tedy spor.

Ukážeme nyní, že existence LU rozkladu pro A je zaručena vždy, pakliže lze Gaussův algoritmus realizovat bez pivotace.

Předpokládejme tedy, že A je taková, že Gaussova eliminace projde bez pivotace, to značí, že

$$A^{(N-1)} = E_{N-1} \dots E_1 A$$

a tudíž

$$A = E_1^{-1} \dots E_{N-1}^{-1} A^{(N-1)},$$

takže

$$A = L U,$$

kde

$$L = E_1^{-1} \dots E_{N-1}^{-1}$$

a

$$U = A^{(N-1)} = \begin{pmatrix} a_{11}^{(N-1)} & \cdot & \cdot & \cdot & a_{1N}^{(N-1)} \\ \cdot & \cdot & & & \cdot \\ \cdot & & \cdot & & \cdot \\ 0 & \cdot & \cdot & \cdot & a_{NN}^{N-1} \end{pmatrix}.$$

Protože zřejmě $l_{jj} = 1, j = 1, \dots, N$, výše uvedené tvrzení platí.

Obecně však se Gaussova eliminace, byť pro symetrické matice bez permutace řádků či sloupců provádět nedá. To ukazuje příklad

$$\begin{pmatrix} 0 & 1 \\ 1 & \epsilon \end{pmatrix}.$$

Pro pozitivně definitní symetrické matice však lze Gaussovou eliminaci provádět bez pivottace.

Pozitivně definitní matice jsou totiž charakterizovány skutečností, že platí

Věta 5.12 (Sylvesterovo kriterium) *Matrice $A = (a_{jk})$ typu $N \times N$ je pozitivně definitní právě když*

$$(5.38) \quad 0 < \det A_k, \quad k = 1, \dots, N,$$

kde A_k jsou tvořeny prvky ležícími v průseku prvních k řádků a sloupců.

Důkaz. Nutnost je triviální. Je-li A pozitivně definitní, potom, podle předpokladu platí nerovnosti

$$\exists \kappa > 0 : \kappa(x, x)_N \leq (Ax, x)_N,$$

takže volbou

$$x = \begin{pmatrix} \tilde{x} \\ 0 \end{pmatrix}, \quad \tilde{x} \in \mathcal{R}^k,$$

obdržíme, že platí

$$(A_k \tilde{x}, \tilde{x})_k \geq \kappa(\tilde{x}, \tilde{x})_k.$$

K důkazu postačitelnosti nechť naopak platí (5.38). Předpokládejme, že existuje vektor $y \in \mathcal{R}^N$ takový, že

$$(5.39) \quad (Ay, y)_N \leq 0,$$

Je tedy zaručena existence vlastní hodnoty $\lambda \leq 0$ a tudiž,

$$(5.40) \quad 0 = \det(\lambda I - A) = \sum_{k=0}^N a_{N-k} \lambda^k,$$

kde

$$\text{sign}(a_{N-k}) = (-1)^k, \quad k = 1, \dots, N.$$

Vidíme tedy, že rovnost (5.40) je s $\lambda \leq 0$ vyloučena.

Tím je důkaz věty 5.12 proveden. |||

Pozitivně definitní matice

Budě $A = (a_{jk})$ symetrická pozitivně definitní matice typu $N \times N$.

Ukážeme si, že provedení Gaussovy eliminační transformace zachovává symetrii těch hlavních submatic, které ještě nejsou horní trojúhelníkové. Tedy, je-li

$$A^{(k)} = \begin{pmatrix} a_{11}^{(k)} & \dots & \dots & a_{1k}^{(k)} & a_{1k+1}^{(k)} & \dots & \dots & a_{1N}^{(k)} \\ \vdots & & & \vdots & & & & \vdots \\ 0 & \dots & \dots & a_{kk}^{(k)} & a_{kk+1}^{(k)} & \dots & \dots & a_{kN}^{(k)} \\ 0 & \dots & \dots & 0 & a_{k+1k+1}^{(k)} & \dots & \dots & a_{k+1N}^{(k)} \\ \vdots & & & \vdots & & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & a_{Nk+1}^{(k)} & \dots & \dots & a_{NN}^{(k)} \end{pmatrix},$$

potom tvrdíme, že

$$(5.41) \quad a_{jl}^{(k)} = a_{lj}^{(k)}, \quad j, l = k+1, \dots, N.$$

Nechť tedy

$$(5.42) \quad a_{jl}^{(k-1)} = a_{lj}^{(k-1)}, \quad j, l = k, \dots, N.$$

Protože

$$a_{jl}^{(k)} = a_{jl}^{(k-1)} - \frac{a_{kl}^{(k-1)}}{a_{kk}^{(k-1)}} a_{jk}^{(k-1)},$$

plyne z (5.42), že

$$a_{jl}^k = a_{lj}^{(k-1)} - \frac{a_{lk}^{(k-1)}}{a_{kk}^{(k-1)}} a_{kj}^{(k-1)} = a_{lj}^k.$$

Platí tedy (5.41). Indukce dává obecný výsledek.

Poznámka 5.2 *Tato skutečnost zřejmě snižuje počet operací nutných k provedení eliminační transformace zhruba na polovinu.*

Ze Sylvesterova kriteria (věta 5.12) plyne, že

$$a_{kk} > 0, \quad k = 1, \dots, N$$

a dále

Věta 5.13 Je-li A symetrická pozitivně definitní matice, potom

$$(5.43) \quad |a_{jk}|^2 \leq a_{jj} a_{kk}, \quad j, k = 1, \dots, N.$$

Tudíž maximální prvek matice A leží na její diagonále.

Důkaz. Buděj P permutační matice. Potom

$$(PAP^T)^T = PA^TP^T = PAP^T$$

je symetrická.

Pro pevnou dvojici indexů j a k existuje elementární permutační matice P tak, že PAP^T je symetrická a pozitivně definitní matice. Na základě Sylvesterova kriteria potom platí relace

$$0 < \det \begin{pmatrix} a_{jj} & a_{jk} \\ a_{kj} & a_{kk} \end{pmatrix} = a_{jj}a_{kk} - a_{jk}^2$$

a tudíž (5.43).

Kdyby

$$a_{jk} = \max \{a_{lt} : l, t = 1, \dots, N\}, \quad j \neq k,$$

pak by na jedné straně, dle předpokladu,

$$a_{jk} > a_{kk}, \quad a_{jk} > a_{jj},$$

tedy

$$a_{jk}^2 > A_{jj}a_{kk}$$

a současně

$$a_{jk} \leq \sqrt{a_{jj}} \sqrt{a_{kk}},$$

spor. \square

Vidíme tedy, že Gaussovu eliminaci lze pro pozitivně definitní matice provádět vždy a to bez pivotace. Hrozí jen, že se při provádění transformací vytvoří dělením malými čísly prvky velké. Již víme, že ty však musí ležet na diagonále. Dále pak z vyjádření

$$a_{jj}^{(k)} = a_{jj}^{(k-1)} - \frac{a_{kj}^{(k-1)}}{a_{kk}^{(k-1)}} a_{jk}^{(k-1)}$$

odvodíme, že

$$a_{jj}^{(k)} \leq a_{jj}^{(k-1)} + \frac{|a_{kj}^{(k-1)}|}{a_{kk}^{(k-1)}} |a_{jk}^{(k-1)}| \leq 2a_{jj}^{(k-1)}.$$

Vidíme, tedy, že růst prvků transformovaných matic je nezávislý na velikosti prvků transformovaných matic. To nás vede k závěru, že Gaussova eliminace bez pivotace je pro pozitivně definitní matice bezpečná.

Analog věty o LU -rozkladu (Věta 5.11, str. 51) lze pro pozitivně definitní matice formulovat takto:

Věta 5.14 Budě A symetrická pozitivně definitní matici.

Potom existuje právě jedna horní trojúhelníková matici R taková, že platí

$$A = R^T R.$$

Důkaz. Nechť

$$A = L U,$$

kde $l_{jj} = 1$, $j = 1, \dots, N$. Tento rozklad existuje na základě věty 5.11 a věty 5.12.

Na základě tohoto rozkladu odvodíme, že

$$a_{11} = u_{11} > 0$$

a

$$u_{jj} = \frac{\det A_j}{\det A_{j-1}}, \quad j = 1, \dots, N.$$

Položme

$$D = \text{diag}\{u_{11}, \dots, u_{NN}\}.$$

Potom

$$A = L D D^{-1} U = L D \tilde{U}, \quad \tilde{U} = D^{-1} U,$$

při čemž

$$\text{diag}L = I = \text{diag}\tilde{U}.$$

Ze symetrie A plyne, že

$$A^T = A = \tilde{U}^T D L^T,$$

neboli, na základě jednoznačnosti LU -rozkladu,

$$(5.44) \quad \tilde{U}^T = L, \quad D L^T = U.$$

Položme dále

$$R = D^{-\frac{1}{2}} U,$$

kde

$$D^{-\frac{1}{2}} = \text{diag} \left\{ d_{11}^{-\frac{1}{2}}, \dots, d_{NN}^{-\frac{1}{2}} \right\}.$$

Potom

$$R^T R = U^T D^{-\frac{1}{2}} D^{-\frac{1}{2}} U = \tilde{U}^T U$$

a dle (5.44) tedy

$$R^T R = L U = A.$$

□

Třídiagonální matice

Matici $A = (a_{jk})$ se nazývá pásová, jestliže existují indexy $p \geq 1$ a $q \geq 1$ takové, že

$$a_{jk} = 0$$

pro

$$j > k + p \text{ a } j - q < k.$$

Speciálně pro $p = q = 1$ se takové pásové matice nazývají *třídiagonální*.

Je zřejmé, že Gaussova eliminace pásovost zachovává.

Obecně platí i zde, že pivotace je nezbytná i v případě symetrických matic. Výjimkami jsou matice diagonálně dominantní a matice pozitivně definitní.

Vyšetřujme případ třídiagonálních matic. Matici

$$A = \begin{pmatrix} a_1 & c_1 & 0 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ b_2 & a_2 & c_2 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ \cdot & \cdot \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & b_{N-1} & a_{N-1} & c_{N-1} \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & b_N & a_N \end{pmatrix}$$

se snažme vyjádřiti ve tvaru

$$A = \begin{pmatrix} 1 & 0 & 0 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ \beta_2 & 1 & 0 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ \cdot & \cdot \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & \beta_{N-1} & 1 & 0 \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & \beta_N & 1 \end{pmatrix} \times \begin{pmatrix} \alpha_1 & \gamma_1 & 0 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ 0 & \alpha_2 & \gamma_2 & \cdot & \cdot & \cdot & 0 & 0 & 0 \\ \cdot & \cdot \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & \alpha_{N-1} & \gamma_{N-1} \\ 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & 0 & a\alpha_N \end{pmatrix}$$

Odtud, za předpokladu, že

$$a_k \neq 0, \quad \forall k = 1, \dots, N,$$

$$b_k \neq 0, \quad \forall k = 1, \dots, N-1,$$

$$c_k \neq 0, \quad \forall k = 2, \dots, N,$$

vyplývají vztahy

$$\gamma_k = c_k \quad \forall k = 1, \dots, N-1,$$

$$\alpha_1 = a_1$$

$$\beta_k = \frac{1}{a_{k-1}} \quad \forall k = 2, \dots, N,$$

$$\alpha_k = a_k - \beta_k c_{k-1}, \quad \forall k = 2, \dots, N.$$

a z nich posléze algoritmus řešení. Řešení se poté obdrží obvyklým postupem, dopřednou eliminací a zpětnou substitucí.

Cvičení 5.1 Budě A třídiagonální pozitivně definitní matici. Dále budě b_j , $j = 1, \dots, p$, $p \geq 2$. Takto formulovanou úlohu lze řešit na př. těmito dvěma odlišnými postupy:

1. Napřed určíme inverzní matici A^{-1} a počítáme postupně $A^{-1}b_j$, $j = 1, \dots, p$.
2. Řešíme postupně soustavy

$$Ax = b_j, \quad j = 1, \dots, p.$$

Rozhodněte který z těchto postupů je výhodnější z hlediska složitosti algoritmů.

5.4 Metody řešení soustav lineárních algebraických rovnic založené na principu optimalizace

Metody uvedené v záhlaví tohoto odstavce jsou typické pro soustavy s pozitivně definitní maticí.

Budě tedy H pozitivně definitní matici typu $N \times N$. Vyšetřujme soustavu

$$(5.45) \quad H\mathbf{x} = \mathbf{b}, \quad \mathbf{b} \in \mathcal{R}^N.$$

Vyšetřujme funkcionál $f : \mathcal{R}^N \rightarrow \mathcal{R}^1$

$$(5.46) \quad f(\mathbf{x}) = (H\mathbf{x}, \mathbf{x}) - 2(\mathbf{b}, \mathbf{x}), \quad \mathbf{x} \in \mathcal{R}^N$$

a zkoumejme jeho vlastnosti.

Snadno se přesvědčíme, že

$$(5.47) \quad f(\mathbf{x}) = (H(\mathbf{x} - \hat{\mathbf{x}}), \mathbf{x} - \hat{\mathbf{x}}) - (\mathbf{b}, \hat{\mathbf{x}}), \quad \mathbf{x} \in \mathcal{R}^N,$$

kde

$$(5.48) \quad \hat{\mathbf{x}} = H^{-1}\mathbf{b},$$

takže díky nezápornosti výrazu (Hy, y) , $y \in \mathcal{R}^N$, existuje

$$(5.49) \quad \min \{f(\mathbf{x}) : \mathbf{x} \in \mathcal{R}^N\} = f(\hat{\mathbf{x}}) = -(\mathbf{b}, \hat{\mathbf{x}}),$$

kde $\hat{\mathbf{x}}$ splňuje (5.48).

Snadno nahlédneme, že úloha nalézt řešení soustavy (5.45) s pozitivně definitní maticí je ekvivalentní s úlohou nalézt (lokální) minimum funkcionálu (5.46). S podobnou situací jsme se již setkali při vyšetřování zobecněných řešení v odstaci 5.1.

Zabýejme se tedy otázkou sestrojování extrémů funkcionálů na \mathcal{R}^N .

Definice 5.2 Funkcionál $f : \mathcal{R}^N \rightarrow \mathcal{R}^1$ se nazývá kvadratický, jestliže je polynomem druhého stupně jakožto funkce N proměnných x_1, \dots, x_N .

Snadno zjistíme, že f je kvadratický funkcionál na \mathcal{R}^N právě když existuje symetrická matici H typu $N \times N$, vektor $\mathbf{x} \in \mathcal{R}^N$ a číslo $c \in \mathcal{R}^1$, takové, že platí

$$(5.50) \quad f(\mathbf{x}) = \frac{1}{2} (H\mathbf{x}, \mathbf{x}) + (\mathbf{b}, \mathbf{x}) + c, \quad \mathbf{x} \in \mathcal{R}^N.$$

Poznámka 5.3 Z definice (5.46) je patrné, že tímto vztahem definovaný funkcionál je kvadratický.

Symbolom $\mathcal{C}^{(m)}(\mathcal{S})$, $m \geq 0$, označujeme množinu funkcionálů, které mají spojité parciální derivace až do řádu m na množině \mathcal{S}), kde $\mathcal{S} \supset \{\mathbf{x} \in \mathcal{R}^N : \|\mathbf{x} - \mathbf{x}_0\| < \delta\}$, pro nějaké $\delta > 0$ a $\mathbf{x}_0 \in \mathcal{R}^N$.

Definice 5.3 Nechť $f \in \mathcal{C}^{(1)}(\mathcal{S})$. Vektor

$$\mathbf{g}(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} (\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_N} (\mathbf{x}) \end{pmatrix}$$

se nazývá gradientem funkcionálu f .

Jestliže $f \in \mathcal{C}^{(2)}(\mathcal{S})$, matice

$$H(\mathbf{x}) = \left(\frac{\partial^2 f}{\partial x_j \partial x_k} (\mathbf{x}) \right)$$

se nazývá Hessiánem funkcionálu f .

Je-li $f \in \mathcal{C}^{(1)}(\mathcal{S})$, pak

$$(5.51) \quad f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\mathbf{h}, \mathbf{g}(\mathbf{x})) + o(\|\mathbf{h}\|),$$

zatímco pro $f \in \mathcal{C}^{(2)}(\mathcal{S})$

$$(5.52) \quad f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\mathbf{h}, \mathbf{g}(\mathbf{x})) + \frac{1}{2} (H(\mathbf{x})\mathbf{h}, \mathbf{h}) + o(\|\mathbf{h}\|^2),$$

kde

$$o(\alpha) = F(\mathbf{x}; \alpha), \quad F : \mathcal{R}^N \times \mathcal{R}_+^1 \rightarrow \mathcal{R}^1,$$

při čemž

$$\lim_{\alpha \rightarrow 0} \frac{|F(\mathbf{x}; \alpha)|}{\alpha} = 0.$$

Poznámka 5.4 Všimněme si toho, že na základě (5.52) každý nelineární funkcionál, mající gradient a Hessián, se lokálně chová jako funkcionál kvadratický.

Definice 5.4 Prvek $\tilde{\mathbf{x}} \in \mathcal{S}$ se nazývá stacionárním prvkem funkcionálu $f \in \mathcal{C}^{(1)}(\mathcal{S})$, jestliže platí

$$\mathbf{g}(\tilde{\mathbf{x}}) = 0,$$

kde $\mathbf{g}(\mathbf{x})$ značí gradient f v $\tilde{\mathbf{x}}$.

Věta 5.15 Nechť $f \in C^{(1)}(\mathcal{S})$. Je-li $\hat{\mathbf{x}}$ prvkem lokálního minima funkcionálu f , pak f je stacionární v prvku \mathbf{x} .

Důkaz. V (5.51) položme $\mathbf{x} = \hat{\mathbf{x}}$ a $\mathbf{h} = -hg(\hat{\mathbf{x}})$, kde $h \in \mathcal{R}^1$ je reálná proměnná veličina z nějakého intervalu $[0, h_0]$.

Platí tedy

$$f(\hat{\mathbf{x}} + \mathbf{h}) = f(\hat{\mathbf{x}}) - h\|\mathbf{g}(\hat{\mathbf{x}})\|^2 + o(h).$$

Nechť f není stacionární v $\hat{\mathbf{x}}$. Potom $g(\hat{\mathbf{x}}) \neq 0$ a pro dostatečně malá h zjistíme, že platí

$$-h\|\mathbf{g}(\hat{\mathbf{x}})\|^2 + o(h) < 0,$$

neboli

$$f(\hat{\mathbf{x}} + \mathbf{h}) < f(\hat{\mathbf{x}})$$

a tudíž $\hat{\mathbf{x}}$ nemůže být prvkem lokálního minima. \square

Stacionarita funkcionálu v $\hat{\mathbf{x}}$ není však postačující podmínkou pro extrém v $\hat{\mathbf{x}}$.

Věta 5.16 Nechť $f \in C^{(2)}(\mathcal{S})$ a nechť f je stacionární v $\hat{\mathbf{x}} \in \mathcal{S}$. Je-li $H(\hat{\mathbf{x}})$ pozitivně definitní pak je $\hat{\mathbf{x}}$ prvkem ostrého lokálního minima funkcionálu f .

Důkaz. Z (5.52) a z podmínky $\mathbf{g}(\hat{\mathbf{x}}) = 0$ odvodíme, že

$$f(\hat{\mathbf{x}} + \mathbf{h}) = f(\hat{\mathbf{x}}) + \frac{1}{2}(H(\hat{\mathbf{x}})\mathbf{h}, \mathbf{h}) + o(\|\mathbf{h}\|^2).$$

Pozitivní definitnost $H(\hat{\mathbf{x}})$ poskytuje existenci $\gamma > 0$ takového, že

$$(H(\hat{\mathbf{x}})\mathbf{h}, \mathbf{h}) \geq \gamma(\mathbf{h}, \mathbf{h}).$$

Tudíž

$$(5.53) \quad f(\hat{\mathbf{x}} + \mathbf{h}) - f(\hat{\mathbf{x}}) \geq \frac{1}{2}\gamma(\mathbf{h}, \mathbf{h}) + o(\|\mathbf{h}\|^2)$$

a protože výraz v pravé části vztahu (5.53) je kladný pro dostatečně malá $\|\mathbf{h}\|^2$, $\hat{\mathbf{x}}$ je prvkem ostrého minima funkcionálu f . \square

Situace je podobná vyšetřujeme-li f z hlediska lokálních maxim. V tom případě totiž $-f$ má v odpovídajících prvcích lokální minima, takže je-li $\hat{\mathbf{x}}$ prvkem lokálního maxima, pak nutně f je stacionární v $\hat{\mathbf{x}}$ a je-li Hessián $H(\hat{\mathbf{x}})$ negativně definitní, $\hat{\mathbf{x}}$ je potom prvkem ostrého lokálního maxima.

Pro indefinitní Hessián nelze žádné závěry stran extrémů obecně zaručit. Není-li $\hat{\mathbf{x}}$ prvkem lokálního extrému, pak Hessián nemůže být definitní, tedy $H(\hat{\mathbf{x}})$ má ve svém spektru jak kladné tak záporné vlastní hodnoty. Dokonce, je-li Hessián $H(\hat{\mathbf{x}})$ semidefinitní, f může ale nemusí mít lokální extrém v $\hat{\mathbf{x}}$.

Dále si připomeňme klasickou definici z elementární matematické analýzy.

Definice 5.5 Nechť $f \in \mathcal{C}^{(2)}(\mathcal{S})$, $\mathbf{x} \in \mathcal{S} \subset \mathbb{R}^N$ s $\mathbf{y} \in \mathbb{R}^N$, $\|\mathbf{y}\| = 1$.

Výraz

$$(5.54) \quad f^{(m)}(\mathbf{x}; \mathbf{y}) = \frac{d^m f(\mathbf{x} + \tau \mathbf{y})}{d\tau^m} \Big|_{\tau=0}$$

se nazývá směrovou derivací funkcionálu f rádu m v prvku \mathbf{x} ve směru \mathbf{y} .

Výpočet směrových derivací se provádí použitím obvyklých pravidel derivování.

Tak pro $f \in \mathcal{C}^{(1)}(\mathcal{S})$, $\mathbf{x} \in \mathbb{R}^N$ a $\mathbf{y} \in \mathbb{R}^N$ s $\|\mathbf{y}\| = 1$,

$$(5.55) \quad f^{(1)}(\mathbf{x}; \mathbf{y}) = \sum_{j=1}^N \frac{\partial f}{\partial x_j}(x_1, \dots, x_N) y_j = (g(\mathbf{x}), \mathbf{y}).$$

Dále pak,

$$(5.56) \quad f^{(2)}(\mathbf{x}; \mathbf{y}) = (H(\mathbf{x})\mathbf{y}, \mathbf{y}).$$

Ze Schwarzovy nerovnosti odvodíme pro $f \in \mathcal{C}^{(1)}(\mathcal{S})$, že

$$\max \left\{ |f^{(1)}(\mathbf{x}; \mathbf{y})| : \mathbf{y} \in \mathbb{R}^N, \|\mathbf{y}\| = 1 \right\} = |f^{(1)}(\mathbf{x}; \hat{\mathbf{y}})| = (g(\mathbf{x}), \hat{\mathbf{y}}),$$

kde

$$\hat{\mathbf{y}} = \frac{g(\mathbf{x})}{\|g(\mathbf{x})\|}.$$

Z (5.55) je patrné, že $f^{(1)}(\mathbf{x}; \mathbf{y}) = 0$ pro všechny směry \mathbf{y} nastane právě když $g(\mathbf{x}) = 0$, tedy, je-li \mathbf{x} stacionárním prvkem funkcionálu f .

Z (5.56) odvodíme zase, že $f^{(2)}(\mathbf{x}; \mathbf{y}) > 0$ pro všechny směry $\mathbf{y} \neq 0$ právě když $H(\mathbf{x})$ je pozitivně definitní.

Předchozí věty lze posléze přeformulovat takto.

Věta 5.17 Nechť $f \in \mathcal{C}^{(1)}(\mathcal{S})$, a $\hat{\mathbf{x}} \in \mathcal{S}$ je prvkem lokálního minima pro f .

Potom $f^{(1)}(\hat{\mathbf{x}}; \mathbf{y}) = 0$ pro všechny směry $\mathbf{y} \in \mathbb{R}^N$.

Věta 5.18 Nechť $f \in \mathcal{C}^{(2)}(\mathcal{S})$ a $\hat{\mathbf{x}} \in \mathcal{S}$ je takový, že $f^{(1)}(\hat{\mathbf{x}}; \mathbf{y}) = 0$ pro všechny směry $\mathbf{y} \in \mathbb{R}^N$.

Potom $\hat{\mathbf{x}}$ je prvkem ostrého lokálního minima pro f v \mathcal{S} , jestliže $f^{(2)}(\hat{\mathbf{x}}; \mathbf{y}) > 0$ pro všechny směry $\mathbf{y} \in \mathbb{R}^N$ s $\|\mathbf{y}\| = 1$.

Poznámka 5.5 Věty 5.17 a 5.18 umožňují snadná zobecnění na případy, kdy vyšetřované funkcionály operují na obecnějších prostorech než jsou aritmetické vektorové prostory.

Definice 5.6 Předpokládejme, že funkcionál f je definován na $\mathcal{S} \subset \mathbb{R}^N$ a nechť $\eta \in \text{range } f = \{\alpha \in \mathbb{R}^1 : \exists \mathbf{x} \in \mathcal{S} \Rightarrow f(\mathbf{x}) = \alpha\}$.

Množina

$$\mathcal{L}_\alpha = \{\mathbf{x} \in \mathcal{S} : f(\mathbf{x}) = \alpha\}$$

se nazývá úrovňovou plochou funkcionálu f pro hodnotu α .

Situace s kvadratickými funkcionály je natolik jednoduchá, že lze používat názoru, který napovídá, že v okolí ostrého lokálního minima $\hat{\mathbf{x}}$ s $f(\hat{\mathbf{x}}) = \hat{\alpha}$ se úrovňové plochy \mathcal{L}_α s α blízkými k $\hat{\alpha}$ chovají jako nadelipsoidy (t. j. obecně jako elipsoidy ve vícerozměrných prostorech).

Obecně řečeno, numerické metody hledání ostrých lokálních extrémů se "chovají lépe", když úrovňové plochy v okolí extrémů jsou blízké plochám kulovým a "hůře, jsou-li tyto plochy výrazně deformovány od tvaru plochy kulové. Veličina, která výstihuje míru odchylky úrovňových ploch od tvaru sféry, je dána výrazem

$$(5.57) \quad \delta_\alpha = \inf \left\{ \frac{\sup \{ \| \mathbf{x} - \mathbf{y} \| : \mathbf{x} \in \mathcal{L}_\alpha \}}{\inf \{ \| \mathbf{x} - \mathbf{y} \| : \mathbf{x} \in \mathcal{L}_\alpha \}} : \mathbf{y} \in \mathcal{S} \right\},$$

při čemž

$$\mathcal{S}_\alpha = \text{Int} \mathcal{L}_\alpha = \left\{ \mathbf{x} \in \mathcal{L}_\alpha : \exists \rho > 0 \text{ tak, že } \forall \mathbf{h} \in \mathbb{R}^N \ \| \mathbf{h} \| < \rho \Rightarrow \mathbf{x} + \mathbf{h} \in \mathcal{L}_\alpha \right\}.$$

Zřejmě, $\delta_\alpha = 1$, $\alpha \rightarrow 0$, implikuje, že \mathcal{L}_α jsou kulové plochy.

Cvičení 5.2 Jednoduše určíme odchylku úrovňové plochy \mathcal{L}_α od plochy kulové pro případ kvadratického funkcionálu charakterizovaného pozitivně definitním Hessiánem H . Tehdy je totiž

$$\delta_\alpha = \frac{\lambda_{\max}}{\lambda_{\min}},$$

konstantní a přitom

$$\lambda_{\min} = \min \{ \lambda \in \sigma(H) \}, \quad \lambda_{\max} = \max \{ \lambda \in \sigma(H) \},$$

takže δ_α je rovno spektrálnímu číslu podmíněnosti matice H . V další části tohoto článku ještě uvidíme, jaký význam má tato skutečnost v analýze konvergence některých numerických metod. Dá se očekávat, že pro Hessiány H s velkými spektrálními čísly podmíněnosti budou odpovídající výpočtové procesy numericky nestabilní (odstrašující je případ, kdy \mathcal{L}_α je "doutník").

Iterační metody sestrojování prvků extrémů

Velké množství numerických metod hledání lokálního minima funkcionálů má iterační charakter a je zpravidla tvaru

$$(5.58) \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \tau_k \mathbf{d}^k,$$

v němž \mathbf{d}^k je *zkusmý směr* zatímco τ_k je parametr zajišťující lokální minimalizaci či alespoň redukci velikosti funkcionálu f na dané množině.

Iterační proces (5.58) vyžaduje tedy strategii pro určení dvou typů objektů: výběr \mathbf{d}^k a určení τ_k s ohledem na chování f podél přímky $\mathbf{x}^k + \tau \mathbf{d}^k$, $-\infty < \tau < +\infty$.

Definice 5.7 Předpokládejme, že pro funkcionál f a vektory \mathbf{x} a \mathbf{d} existuje $\tau_0 > 0$ takové, že

$$f(\mathbf{x} + \tau \mathbf{d}) < f(\mathbf{x}), \quad 0 < \tau < \tau_0.$$

Potom nazýváme směr \mathbf{d} směrem klesání funkcionálu f .

Věta 5.19 Nechť $f \in \mathcal{C}^{(1)}(\mathcal{R}^N)$ a nechť $[\mathbf{g}(\mathbf{x})]$ označuje (jako obvykle v tomto odstavci) gradient funkcionálu f v \mathbf{x} . Jestliže vektor \mathbf{d} splňuje relaci

$$(\mathbf{d}, \mathbf{g}(\mathbf{x})) < 0,$$

pak \mathbf{d} je směrem klesání funkcionálu f v \mathbf{x} .

Důkaz. Z (5.51) obdržíme, že

$$f(\mathbf{x} + \tau\mathbf{d}) = f(\mathbf{x}) + \tau (\mathbf{d}, \mathbf{g}(\mathbf{x})) + o(\tau),$$

a protože podle předpokladu $(\mathbf{d}, \mathbf{g}(\mathbf{x})) < 0$, musí platit nerovnost

$$\tau (\mathbf{d}, \mathbf{g}(\mathbf{x})) + o(\tau) < 0$$

pro dostatečně malá τ .

Poznámka 5.6 Je-li $\mathbf{g}(\mathbf{x}) = 0$, pak nelze stanovit, zda \mathbf{d} je či není směrem klesání f v \mathbf{x} bez další dodatečné informace.

Věta 5.20 Nechť $f \in \mathcal{C}^{(2)}(\mathcal{R}^N)$, nechť $(\mathbf{d}, \mathbf{g}(\mathbf{x})) = 0$ a $(H(\mathbf{x})\mathbf{d}, \mathbf{g}(\mathbf{x})) < 0$ pro nějaká \mathbf{x} a \mathbf{d} ($H(\mathbf{x})$ značí jako obvykle Heessián f).

Potom \mathbf{d} je směrem klesání f v \mathbf{x} .

Důkaz. Z (5.52) plyne

$$(5.59) \quad f(\mathbf{x} + \tau\mathbf{d}) = f(\mathbf{x}) + (1/2)\tau^2 (H(\mathbf{x})\mathbf{d}, \mathbf{d}) + o(\tau^2),$$

odkud tvrzení plyne. \square

Věta 5.21 Nechť $f \in \mathcal{C}^{(1)}(\mathcal{R}^N)$. Potom mezi všemi zkusmými směry \mathbf{d} funkcionálu f v \mathbf{x} tím směrem, v němž f klesá v okolí \mathbf{x} nejrychleji je $\mathbf{d} = -\mathbf{g}(\mathbf{x})$.

Důkaz. Naším úkolem je minimalizovat směrovou derivaci f v \mathbf{x} vzhledem ke všem zkusmým směrům \mathbf{d} . Z (5.55) plyne, že to znamená minimalizovat výraz $(\mathbf{d}, \mathbf{g}(\mathbf{x}))$ pro všechny směry \mathbf{d} , pro něž $\|\mathbf{d}\| = 1$. Protože $|(\mathbf{d}, \mathbf{g}(\mathbf{x}))| \leq |\mathbf{g}(\mathbf{x})|$, minimum se nabývá pro

$$\mathbf{d} = -\frac{\mathbf{g}(\mathbf{x})}{\|\mathbf{g}(\mathbf{x})\|}.$$

\square

Dále se zabývejme problémem určení τ jsou-li zadány \mathbf{d} a \mathbf{x} . Snažme se minimalizovat f podél přímky $\mathbf{y} = \mathbf{x} + \tau\mathbf{d}$, $-\infty < \tau < +\infty$. Pro obecné funkcionály stanovení parametru τ_k je poměrně složitá úloha a její řešení

vyžaduje další iterační proces. Pro případ kvadratického funkcionálu to však není žádný problém.

Nechť tedy

$$f(\mathbf{x}) = (1/2)(H\mathbf{x}, \mathbf{x}) + (\mathbf{b}, \mathbf{x}) + c.$$

Okamžitě lze ověřit, že

$$f(\mathbf{x} + \tau \mathbf{d}) - f(\mathbf{x}) = (1/2)\tau^2(H\mathbf{d}, \mathbf{d}) + \tau(\mathbf{d}, \mathbf{g}(\mathbf{x})).$$

Je-li H pozitivně definitní a $\mathbf{d} \neq 0$, potom $(H\mathbf{d}, \mathbf{d}) > 0$ a pro

$$\tau = \frac{(\mathbf{d}, \mathbf{g}(\mathbf{x}))}{(H\mathbf{d}, \mathbf{d})}$$

$f(\mathbf{x} + \tau \mathbf{d})$ nabývá minima vzhledem k $\tau \in (-\infty, +\infty)$.

Tedy, nezávisle na volbě \mathbf{d}^k v (5.58) je pro kvadratický funkcionál f

$$(5.60) \quad \tau = -\frac{(\mathbf{d}, \mathbf{g}(\mathbf{x}))}{(H\mathbf{d}, \mathbf{d})}$$

při pustná hodnota lokálního optimalizačního parametru.

Metoda největšího spádu

Na základě výsledku věty 5.21 je přirozené klást $\mathbf{d}^k = -\mathbf{g}(\mathbf{x}^k)$ v (5.58), takže

$$(5.61) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \tau_k \mathbf{g}(\mathbf{x}^k).$$

Tento vztah spolu se strategií určování τ_k definuje *metodu největšího spádu* minimalizace obecného funkcionálu f .

Pro případ kvadratického funkcionálu ve tvaru (5.59) poskytuje (5.60) evidentní způsob určení τ_k . Položme

$$\mathbf{g}^k = \mathbf{g}(\mathbf{x}^k), \quad \forall k = 0, 1, \dots$$

Pro (5.59) je metoda největšího spádu určena formulemi

$$(5.62) \quad \mathbf{g}^k = H\mathbf{x}^k - b,$$

$$(5.63) \quad \tau_k = \frac{(\mathbf{g}^k, \mathbf{g}^k)}{(H\mathbf{g}^k, \mathbf{g}^k)},$$

$$(5.64) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \tau_k \mathbf{g}^k,$$

kde $k = 0, 1, \dots$ a $\mathbf{g}^k = \mathbf{g}(\mathbf{x}^k)$.

Zřejmě (5.62) lze psát ve tvaru

$$(5.65) \quad \mathbf{g}^{k+1} = \mathbf{g}^k - \tau_k H\mathbf{g}^k,$$

a to umožňuje výpočet gradientů v (5.62) - (5.64) rekurentně.

Příslušný algoritmus se řídí formulemi

$$(5.66) \quad \tau_k = \frac{(\mathbf{g}^k, \mathbf{g}^k)}{(H\mathbf{g}^k, \mathbf{g}^k)},$$

$$(5.67) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \tau_k \mathbf{g}^k,$$

$$(5.68) \quad \mathbf{g}^{k+1} = \mathbf{g}^k - \tau_k H \mathbf{g}^k,$$

kde $k = 0, 1, \dots$ a $\mathbf{g}^0 = H\mathbf{x}^0 - b$.

Poznamenejme, že algoritmus určený formulemi (5.66) - (5.68) lze implementovat za použití jediného násobení maticí H na jednu iteraci na rozdíl od standardsního algoritmu (5.62) - (5.64), kdy násobení maticí H jsou zapotřebí dvě.

Analýza konvergence metody největšího spádu dává podnět k obecnější úvaze. Vzhledem k existenci obecně nekonečně mnoha norem jest možno klást si otázku, kterou normu zvolit k odhadu chyby metody tak, aby bylo obdrželi výsledek co možno nejpřirozenější, nejrealističtější, nevhodnější, nejjednodušší, nej- atd.

Na první pohled by se mohlo zdát, že výběr normy není nijak důležitý; vždyť všechny normy na \mathcal{R}^N jsou ekvivalentní. Při bližším zkoumání zjistíme, že pro některé normy je taková analýza dost obtížná. Rozhodujícím kriteriem pro volbu normy bude tedy co nejmenší analytická obtížnost vyšetřování konvergence. Tu se ukazuje, jako velice výhodná a vskutku přirozená t.zv. energetická norma.

Ovodme tedy větu o rychlosti konvergence metody největšího spádu v energetické normě.

Definice 5.8 Bud $\succsim H$ pozitivně definitní matici typu $N \times N$. Norma $\|\cdot\|_H$ definovaná pomocí skalárniho součinu

$$\phi_H(x, y) = (Hx, y),$$

kde $(x, y) = \sum_{k=1}^N (x)_k (y)_k$, se nazývá H -energetickou normou; tedy, $\|x\|_H^2 = (Hx, x)$.

Věta 5.22 Nechť $\{x^k\}$ je posloupnost získaná pomocí metody největšího spádu pro kvadratický funkcionál f . Nechť \hat{x} značí prvek globálního minima funkcionálu f .

Potom platí odhad

$$(5.69) \quad \|x^k - \hat{x}\|_H \leq \frac{(\kappa(H) - 1)^k}{(\kappa(H) + 1)^k} \|x^0 - \hat{x}\|_H.$$

Dále nechť pro libovolné $\epsilon > 0$ $p(\epsilon)$ značí nejmenší index takový, že $\|x^k - \hat{x}\|_H \leq \epsilon \|x^0 - \hat{x}\|_H$, pro každé $x^0 \in \mathcal{R}^N$, potom platí nerovnost

$$(5.70) \quad p(\epsilon) \leq \frac{1}{2} \kappa(H) \log \frac{1}{\epsilon}.$$

Důkaz. Nechť

$$E(\mathbf{x}^k - \hat{\mathbf{x}}) = (H(\mathbf{x} - \hat{\mathbf{x}}), \mathbf{x} - \hat{\mathbf{x}}) = (H^{-1}\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{x})).$$

Dále si uvědomme, že platí

$$(5.71) \quad E(\mathbf{x} - \hat{\mathbf{x}}) = (H^{-1}\mathbf{g}^k, \mathbf{g}^k).$$

Platnost (5.71) plyne odtud, že

$$\mathbf{f}(\mathbf{x}) - \mathbf{f}(\hat{\mathbf{x}}) = \frac{1}{2} (H(\mathbf{x} - \hat{\mathbf{x}}), \mathbf{x} - \hat{\mathbf{x}}),$$

$$E(\mathbf{x} - \hat{\mathbf{x}}) = 2\{\mathbf{f}(\mathbf{x}) - \mathbf{f}(\hat{\mathbf{x}})\}$$

a

$$\mathbf{g}(\mathbf{x}) = H(\mathbf{x} - \hat{\mathbf{x}}).$$

Z vyjádření

$$\mathbf{f}(\mathbf{x}) - \mathbf{f}(\hat{\mathbf{x}}) = \frac{1}{2} (H(\mathbf{x} - \hat{\mathbf{x}}), \mathbf{x} - \hat{\mathbf{x}})$$

totiž plyne, že

$$\mathbf{g}(\mathbf{x}) = \begin{pmatrix} \frac{\partial \mathbf{f}}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial \mathbf{f}}{\partial x_N}(\mathbf{x}) \end{pmatrix} = \frac{1}{2} H(\mathbf{x} - \hat{\mathbf{x}})$$

a tudíž

$$\mathbf{x} - \hat{\mathbf{x}} = H^{-1}\mathbf{g}(\mathbf{x}).$$

Položíme-li tedy $\mathbf{x} = \mathbf{x}^{k+1}$ a použijeme-li rekurentní formule pro prvky \mathbf{g}^k , obdržíme, že

$$\begin{aligned} E(\mathbf{x}^{k+1} - \hat{\mathbf{x}}) &= (H^{-1}\mathbf{g}^{k+1}, \mathbf{g}^{k+1}) \\ &= (H^{-1}(\mathbf{g}^k - \tau_k H \mathbf{g}^k), \mathbf{g}^k - \tau_k H \mathbf{g}^k) \\ &= (H^{-1}(I - \tau_k H) \mathbf{g}^k, (I - \tau_k H) \mathbf{g}^k) \\ &= (H^{-1}(I - \tau_k H)^2 \mathbf{g}^k, \mathbf{g}^k) \\ &= (H^{-1}\mathbf{g}^k, \mathbf{g}^k) - 2\tau_k (\mathbf{g}^k, \mathbf{g}^k) + \tau_k^2 (H\mathbf{g}^k, \mathbf{g}^k) \\ &= (H^{-1}\mathbf{g}^k, \mathbf{g}^k) - \frac{(\mathbf{g}^k, \mathbf{g}^k)^2}{(H\mathbf{g}^k, \mathbf{g}^k)^2} (H\mathbf{g}^k, \mathbf{g}^k). \end{aligned}$$

Výsledkem je vztah

$$E(\mathbf{x}^{k+1} - \hat{\mathbf{x}}) - E(\mathbf{x}^k - \hat{\mathbf{x}}) = \frac{(\mathbf{g}^k, \mathbf{g}^k)^2}{(H\mathbf{g}^k, \mathbf{g}^k)(H^{-1}\mathbf{g}^k, \mathbf{g}^k)} E(\mathbf{x}^k - \hat{\mathbf{x}}),$$

odkud pak

$$E(\mathbf{x}^{k+1} - \hat{\mathbf{x}}) = \left\{ 1 - \frac{(\mathbf{g}^k, \mathbf{g}^k)^2}{(H\mathbf{g}^k, \mathbf{g}^k)(H^{-1}\mathbf{g}^k, \mathbf{g}^k)} \right\} E(\mathbf{x}^k - \hat{\mathbf{x}}).$$

Aplikujeme-li nyní Bergstromovo lemma na pozitivně definitní matici H , mající vlastní čísla $\lambda_1, \dots, \lambda_N$, podle něhož

$$\begin{aligned} 1 &\leq \frac{(Hx, x)}{(x, x)} \frac{(x, H^{-1}x)}{(x, x)} \\ &\leq \frac{(\lambda_1 + \lambda_N)^2}{4\lambda_1 \lambda_N}, \quad \mathbf{x} \in \mathcal{R}^N, \end{aligned}$$

obdržíme odhad

$$E(\mathbf{x}^{k+1} - \hat{\mathbf{x}}) \leq \left[1 - \frac{4\lambda_1\lambda_N}{(\lambda_1 + \lambda_N)^2} \right] E(\mathbf{x}^k - \hat{\mathbf{x}}),$$

což vyjádřeno pomocí spektrálního čísla podmíněnosti $\kappa(H) = (\lambda_N/\lambda_1)$ poskytuje odhad

$$E(\mathbf{x}^{k+1} - \hat{\mathbf{x}}) \leq \left(\frac{\kappa(H) - 1}{\kappa(H) + 1} \right)^2 E(\mathbf{x}^k - \hat{\mathbf{x}}).$$

Vidíme tedy, že konvergenci metody největšího spádu lze popsát pomocí energetické normy vztahem

$$\|\mathbf{x}^{k+1} - \hat{\mathbf{x}}\|_H \leq \left(\frac{\kappa(H) - 1}{\kappa(H) + 1} \right) \|\mathbf{x}^k - \hat{\mathbf{x}}\|_H,$$

kterým je dokázána platnost prvního tvrzení dokazované věty.

Abychom odvodili platnost nerovnosti (5.70) uvědomme si, že platnost nerovnosti $\|\mathbf{x}^k - \hat{\mathbf{x}}\|_H \leq \epsilon \|\mathbf{x}^0 - \hat{\mathbf{x}}\|_H$ pro $k \geq \hat{k}$ je zaručena platností nerovnosti

$$\left(\frac{\kappa(H) - 1}{\kappa(H) + 1} \right)^{\hat{k}} \leq \epsilon$$

pro $k \geq \hat{k}$.

Hledejme tedy nejmenší k takové, aby

$$\frac{1}{\epsilon} \leq \left(\frac{\kappa(H) + 1}{\kappa(H) - 1} \right)^k.$$

Jest totiž

$$\log \frac{1}{\epsilon} \leq k \log(1 + \sigma),$$

kde

$$\sigma = \frac{2}{\kappa - 1}.$$

Protože

$$\frac{1}{\sigma} \log \frac{1}{\epsilon} \leq \frac{1}{2} (\kappa - 1) \log \frac{1}{\epsilon},$$

odvodíme snadno, že hledaný index $p(\epsilon)$ splňuje nerovnost

$$p(\epsilon) \leq \frac{1}{2} \kappa(H) \log \frac{1}{\epsilon} + 1.$$

Tím je dokázána platnost vztahu (5.70) a tedy též věty 5.22. \square

Předpodmínění.

Předpokládejme, že C je pozitivně definitní $N \times N$ matice faktorizovaná na tvar $C = EE^T$ a že

$$(5.72) \quad \mathbf{f}(\mathbf{x}) = \frac{1}{2}(H\mathbf{x}, \mathbf{x}) + (\mathbf{x}, \mathbf{b}) + \mathbf{c},$$

kde H je pozitivně definitní $N \times N$ matice. Zavedeme další funkcionál $\tilde{\mathbf{f}}$ předpisem

$$(5.73) \quad \tilde{\mathbf{f}}(\mathbf{y}) = \mathbf{f}(E^T \mathbf{y}) = \frac{1}{2}(\tilde{H}\mathbf{y}, \mathbf{y}) + (\mathbf{y}, \tilde{\mathbf{b}}) + \tilde{\mathbf{c}},$$

při čemž

$$(5.74) \quad \tilde{H} = E^{-1}H(E^T)^{-1}, \quad \tilde{\mathbf{b}} = E^{-1}\mathbf{b}, \quad \tilde{\mathbf{c}} = \mathbf{c}.$$

Zřejmě je \tilde{H} symetrická a navíc ze vztahu $(\tilde{H}\mathbf{y}, \mathbf{y}) = (H\mathbf{x}, \mathbf{x}) > 0$ pro všechna $\mathbf{x} = (E^T)^{-1}\mathbf{y}$, plyne, že \tilde{H} je pozitivně definitní.

Podobnostní transformace

$$(E^T)^{-1}\tilde{H}E^T = (E^T)^{-1}E^{-1}H = C^{-1}H$$

ukazuje, že matice \tilde{H} a $C^{-1}H$ mají totožná spektra. Označme vlastní hodnoty matice \tilde{H} symboly $\tilde{\lambda}_1, \dots, \tilde{\lambda}_N$. Číslo podmíněnosti

$$(5.75) \quad \kappa(\tilde{H}) = \frac{\tilde{\lambda}_N}{\tilde{\lambda}_1}$$

je určeno maticemi C a H presto, že \tilde{H} závisí ještě na rozkladu $C = EE^T$.

Aplikujme nyní metodu největšího spádu na funkcionál (5.73). Obdržíme tak tento postup

$$(5.76) \quad \tilde{\mathbf{g}}^k = \tilde{H}\mathbf{y}^k - \tilde{\mathbf{b}},$$

$$(5.77) \quad \tilde{\tau}_k = \frac{(\tilde{\mathbf{g}}^k, \tilde{\mathbf{g}}^k)}{(\tilde{H}\tilde{\mathbf{g}}^k, \tilde{\mathbf{g}}^k)},$$

$$(5.78) \quad \mathbf{y}^{k+1} = \mathbf{y}^k - \tilde{\tau}_k \tilde{\mathbf{g}}^k, \quad k = 0, 1, \dots$$

kde \mathbf{y}^0 lze volit libovolně.

Již víme, že platí

$$(5.79) \quad \lim_{k \rightarrow \infty} \mathbf{y}^k = \hat{\mathbf{y}} = \tilde{H}^{-1}\mathbf{b},$$

a že rychlosť konvergencie je určena číslem podmíněnosti $\kappa(\tilde{H})$.

Položme $\mathbf{x}^k = (E^T)^{-1}\mathbf{y}^k$ a $\mathbf{g}^k = H\mathbf{x}^k - \mathbf{b}$, $k = 0, 1, \dots$ Snadno se přesvědčíme, že platí vztahy

$$\tilde{\mathbf{g}}^k = (E^T)^{-1}\mathbf{g}^k, \quad \tilde{\tau}_k = \frac{(\mathbf{g}^k, \mathbf{h}^k)}{(H\mathbf{g}^k, \mathbf{g}^k)}, \quad \mathbf{y}^{k+1} = E^T(\mathbf{x}^k - \tilde{\tau}_k \mathbf{h}^k),$$

kde $\mathbf{h}^k = C^{-1}\mathbf{g}^k$.

Jelikož $\mathbf{y}^{k+1} = E^T \mathbf{x}^{k+1}$, výše uvedené vztahy lze realizovat též pomocí schématu

$$(5.80) \quad \mathbf{g}^k = H\mathbf{x}^k - \mathbf{b},$$

$$(5.81) \quad \mathbf{h}^k = C^{-1} \mathbf{g}^k,$$

$$(5.82) \quad \tau_k = \frac{(\mathbf{h}^k, \mathbf{g}^k)}{(H\mathbf{h}^k, \mathbf{g}^k)},$$

$$(5.83) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \tau_k \mathbf{h}^k, \quad k = 0, 1, \dots$$

Extrémní vektory pro \mathbf{f} a $\tilde{\mathbf{f}}$ jsou dány výrazy

$$\hat{\mathbf{x}} = H^{-1} \mathbf{b}, \quad \text{a} \quad \hat{\mathbf{y}} = \tilde{H}^{-1} \tilde{\mathbf{b}},$$

při čemž

$$\hat{\mathbf{y}} = E^T \hat{\mathbf{x}}.$$

Dále pak

$$\mathbf{y}^k - \hat{\mathbf{y}} = E^T (\mathbf{x}^k - \hat{\mathbf{x}})$$

a

$$\|\mathbf{y}^k - \hat{\mathbf{y}}\|_{\tilde{H}} = \|\mathbf{x}^k - \hat{\mathbf{x}}\|_{\tilde{H}}.$$

Na základě věty 5.22 však platí

$$\begin{aligned} \|\mathbf{y}^k - \hat{\mathbf{y}}\|_{\tilde{H}} &\leq \left(\frac{\kappa(\tilde{H})-1}{\kappa(\tilde{H})+1} \right)^k \|\mathbf{y}^k - \hat{\mathbf{y}}\|_{\tilde{H}} \\ &= \left(\frac{\kappa(H)-1}{\kappa(H)+1} \right)^k \|\mathbf{x}^k - \hat{\mathbf{x}}\|_H. \end{aligned}$$

Posloupnost $\{\mathbf{x}^k\}$, kde \mathbf{x}^0 lze volit libovolně, konstruovaná pomocí vztahů (5.80) - (5.83), konverguje k $\hat{\mathbf{x}}$ a její rychlosť konvergence je určena číslem $\kappa(\tilde{H})$. Lze-li najít matici C tak, aby $\kappa(\tilde{H}) < \kappa(H)$, pak odhad (5.70) napovídá, že rychlosť konvergence (5.80) - (5.83) je lepší než (5.66) - (5.68).

Matice C se nazývá *matice předpodmínění metody největšího spádu*. Metoda určená vzorci (5.76) - (5.79) se nazývá *transformovanou předpodmíněnou metodou největšího spádu*, zatímco metoda určená vzorci (5.80) - (5.83) se nazývá *netransformovanou předpodmíněnou metodou největšího spádu*. Názvy jsou odvozeny odtud, že extrém $\hat{\mathbf{y}} = E^T \hat{\mathbf{x}}$, pro metodu danou vzorci (5.76) - (5.79) je transformovaným extrémem pro (5.80) - (5.83), tedy pro $\hat{\mathbf{x}}$. Platí též, že $\mathbf{x}^k = (E^T)^{-1} \mathbf{y}^k$, a tuto transformaci lze ale není nutné provádět pro každé $k = 0, 1, \dots$

Výpočtový proces prováděný na základě výše uvedených algoritmů se budeme snažit ocenit z pozic jejich složitosti.

Potřebné gradienty lze stanovit rekurzivně. Jmenovitě pak (5.76) a (5.80) lze nahradit formulami

$$\tilde{\mathbf{g}}^k = \tilde{\mathbf{g}}^{k-1} - \tilde{\tau}_{k-1} \tilde{H} \tilde{\mathbf{g}}^{k-1}$$

a

$$\mathbf{g}^k = \tau_{k-1} H \mathbf{h}^{k-1}.$$

Vektory $\tilde{H} \tilde{\mathbf{g}}^{k-1}$ a $H \mathbf{h}^{k-1}$ jsou k dispozici z předchozího iteračního kroku.

Nás zajímají případy, v nichž rozměr matice H je značný a matice H sama je řídká. Matice E bývá zpravidla řídká trojúhelníková matice. Při realizaci transformované metody narázíme na skutečnost, že \tilde{H} nebývá řídká a tudíž se explicitnímu výpočtu snažíme vyhnout.

Snadno nahlédneme, že v rámci jednoho iteračního kroku je výpočet vektoru $\tilde{H}\tilde{\mathbf{g}}^k$ v (5.77) nejpracnější. Je tudíž záhadno provádět jeho výpočet nepřímo. Řešíme proto napřed soustavu

$$E^T \mathbf{z} = \mathbf{g}^k$$

vzhledem k \mathbf{z} , poté položíme

$$\mathbf{z}^* = H\mathbf{z}$$

a řešíme

$$E\mathbf{z}^{**} = \mathbf{z}^*$$

vzhledem k $\mathbf{z}^{**} = E^{-1}H(E^T)^{-1}\mathbf{g}^k = \tilde{H}^{-1}\mathbf{g}^k$. Je patrné, že abychom tento vektor určili, musíme řešit dvě soustavy s trojúhelníkovými maticemi a provézt jedno násobení maticí H . Dále pak se musí provézt transformace $\mathbf{x}^k = (E^T)^{-1}\mathbf{y}^k$ a to opět vyžaduje řešení jedné soustavy s trojúhelníkovou maticí.

Nejpodstatnější složkou výpočtu podle netransformované metody (5.80) - (5.83) je krok (5.81), kde se počítá $C^{-1}\mathbf{g}^k$ a krok (5.82) kde se stanovuje $H\mathbf{h}^k$. Vektor $C^{-1}\mathbf{g}^k$ se vypočítává tak, že se řeší

$$E\mathbf{z} = \mathbf{g}^k$$

vzhledem k \mathbf{z} a

$$E^T \mathbf{h}^k = \mathbf{z}$$

se řeší vzhledem k \mathbf{h}^k . Potom zřejmě $\mathbf{h}^k = C^{-1}\mathbf{g}^k$.

I v netransformované metodě se řeší dvě soustavy s trojúhelníkovými maticemi. Mohlo by se zdát, že netransformovaná metoda je výhodnější, protože poskytuje přímo hledaný extrémní vektor $\hat{\mathbf{x}}$. Existují však vhodné transformační matice, pro něž výpočtová náročnost není větší než pro metodu netransformovanou. Toto tvrzení se pokusíme ozřejmit pro případ předpodmíněné metody sdružených gradientů, tedy metody značně efektivnější než je metoda největšího spádu. Tvrzení samo je však pravdivé i pro metodu největšího spádu.

5.5 Metoda sdružených gradientů

V tomto článku ukážeme jak sestrojovat účinné approximace řešení soustavy lineárních algebraických rovnic

$$(5.84) \quad Ax = b, \quad x, b \in \mathcal{E} = \mathcal{R}^N,$$

kde A je pozitivně definitní $N \times N$ matice, pomocí extremalizace odpovídajícího funkcionálu

$$(5.85) \quad F(x) = (Ax, x) - (b, x) - (x, b).$$

Odvození provedeme simultánně jak pro základní verzi metody sdružených gradientů tak pro verzi předpodmíněnou.

Libovolně zvolme základní vektor approximace x_0 a položme $r_0 = b - Ax_0$. Dále definujme *Krylovovy prostory* (vzhledem k r_0 a A) kladouce

$$\mathcal{V}_k = \text{Lin} \left\{ r_0, Ar_0, \dots, A^{k-1}r_0 \right\},$$

při čemž symbol *Lin* znamená lineární obal vektorů $\{r_0, \dots, A^{k-1}r_0\}$.

Jakožto k -tou iteraci metody sdružených gradientů definujeme Galerkinovu approximaci \tilde{x}_k vektoru $x_k - x$ v prostoru \mathcal{V}_k , tedy

$$(5.86) \quad (A\tilde{x}_k, \phi) = (r_0, \phi), \quad \forall \phi \in \mathcal{V}_k, \quad k = 1, 2, \dots$$

Položme $x_k = \tilde{x}_k + x_0$. Jest potom (5.86) totéž co

$$(5.87) \quad (Ax_k, \phi) = (b, \phi) \quad \forall \phi \in \mathcal{V}_k$$

a $x_k - x_0$ patří do \mathcal{V}_k .

Dále odvodíme metodu sdružených gradientů s předpodmíněním. Budě B další symetrická pozitivně definitní matice typu $N \times N$. Jest tudíž soustava (5.84) ekvivalentní se soustavou

$$(5.88) \quad BAx = Bb.$$

Definujme dále nový skalární součin na \mathcal{E} formulí $[., .] = (B^{-1}., .)$. Snadno zjistíme, že operátor $\mathcal{A} = BA$ je symetrický a pozitivně definitní vzhledem k součinu $[., .]$, t.j. $[\mathcal{A}x, y] = x, \mathcal{A}y]$ a existuje $\gamma > 0$ takové, že $[\mathcal{A}x, x] \geq \gamma[x, x], x \in \mathcal{E}$.

Budě $y_0 \in \mathcal{E}$ libovolný startovní nástřel a $z_0 = Bb - \mathcal{A}y_0$. Definujme

$$\tilde{\mathcal{V}} = \text{Lin} \left\{ z_0, \mathcal{A}z_0, \dots, \mathcal{A}^{k-1}z_0 \right\}, \quad k = 1, 2, \dots$$

k -tý člen metody sdružených gradientů pro soustavu s maticí A a předpodmiňovačem B (PCG) je definována jako Galerkinova approximace \tilde{y}_k vektoru $x - y_0$ t.j. pro $k = 1, 2, \dots$

$$(5.89) \quad [\mathcal{A}\tilde{y}_k, \phi] = [z_0, \phi], \quad \forall \phi \in \tilde{\mathcal{V}}.$$

Pro $y_k = \tilde{y}_k + y_0$ je posléze (5.89) totéž co

$$(5.90) \quad (Ay_k, \phi) = (b, \phi)$$

$\forall \phi \in \tilde{\mathcal{V}}$ a $y_k - y_0$.

Porovnáme-li (5.86) a (5.89) zjistíme, že PCG je formálně podobné s CG dosadíme-li \mathcal{A} namísto A a $[., .]$ namísto $(., .)$. Za zmínu stojí též skutečnost, že PCG approximaci $y_k - y_0$ lze interpretovat jako standardní Galerkinovu approximaci vektoru $x - x_0$ v prostoru $\tilde{\mathcal{V}}_k$. Volíme-li $B = I$ pak $\tilde{\mathcal{V}}_k = \mathcal{V}_k$ a CG je tak speciální případ PCG.

Poznámka 5.1 Podle Hamiltonovy-Cayleyovy věty (blíže se o ní čtenář dozví v monografii I. Marek, K. Žitný: *Analytická teorie matic pro aplikované vědy*, Teubner Verlag, Díl I. z 1984 a díl II. z 1986) lze inverzní matici A^{-1} vyjádřiti jakožto polynom stupně nejvýše $N - 1$, kde N je dimenze prostoru \mathcal{E} . Protože $x - x_0 = A^{-1}b \in \mathcal{V}_N$, při exaktním počítání metoda CG docílí hledané řešení po nejvýše N krocích. Tato okolnost však ustupuje svou důležitostí do pozadí skutečnosti jiné, že totiž potřebnou přesnost approximace lze docílit za pomoci daleko menšího počtu kroků než je N . V praxi lze totiž docílit toho, že při vhodném předpodmiňovači počet kroků potřebný k zaručení žádané přesnosti je typická situace kdy, řekněme $k = 10$ či $k = 20$ zatímco dimenze $N = 10^6$ a ta může být i větší.

Analýza konvergence

Metody CG a PCG, vlastně však jen PCG, budeme analyzovat za použití energetické normy $[\mathcal{A}., .]^{1/2} = (A., .)^{1/2}$.

Z formule (5.89) plyne platnost vztahů

$$(5.91) \quad [\mathcal{A}(\tilde{y}_k - (x - y_0)), \phi] = 0, \quad \forall \phi \in \tilde{\mathcal{V}}_k.$$

Položíme-li tedy $e_k = \tilde{y}_k - (x - y_0) = y_k - x$ vyplynou nám následující vztahy

$$(5.92) \quad [\mathcal{A}e_k, e_k] = [\mathcal{A}e_k, y - (x - y_0)],$$

platné pro všechny prvky $y \in \tilde{\mathcal{V}}_k$. Nyní použijeme Schwarzovy nerovnosti pro skalární součin (5.92). Postupně obdržíme

$$\begin{aligned} [\mathcal{A}e_k, y - (x - y_0)] &= (\mathcal{A}e_k, y - (x - y_0)) \\ &= (A^{1/2}e_k, A^{1/2}(y - (x - y_0))) \\ &\leq (A^{1/2}e_k, A^{1/2}e_k)^{1/2}(A^{1/2}(y - (x - y_0)), A^{1/2}(y - (x - y_0)))^{1/2} \\ &= (Ae_k, e_k)^{1/2}(A(y - (x - y_0)), (y - (x - y_0)))^{1/2} \end{aligned}$$

takže též,

$$[\mathcal{A}e_k, e_k] \leq [\mathcal{A}e_k, e_k]^{1/2}[\mathcal{A}(y - (x - y_0)), (y - (x - y_0))]^{1/2}$$

a posléze

$$(5.93) \quad [\mathcal{A}e_k, e_k] \leq [\mathcal{A}(y - (x - y_0)), y - (x - y_0)]$$

pro jakýkoliv vektor $y \in \tilde{\mathcal{V}}_k$.

Protože $y \in \tilde{\mathcal{V}}_k$, lze zvolit libovolný polynom P_{k-1} stupně nejvyšše $k-1$ a položit

$$y = P_{k-1}(\mathcal{A})z_0 \equiv P_{k-1}(\mathcal{A})\mathcal{A}(x - y_0).$$

S touto volbou nabude formule (5.93) tvaru

$$(5.94) \quad [\mathcal{A}e_k, e_k] \leq [\mathcal{A}(P_{k-1}(\mathcal{A})\mathcal{A} - I)(x - y_0), (P_{k-1}(\mathcal{A})\mathcal{A} - I)(x - y_0)].$$

Nechť

$$Q_k(\mathcal{A}) = I - P_{k-1}(\mathcal{A})\mathcal{A}.$$

Potom

$$(5.95) \quad \begin{cases} [\mathcal{A}Q_k(\mathcal{A})(x - y_0), Q_k(\mathcal{A})(x - y_0)] \\ \leq [[Q_k(\mathcal{A})]]^2[\mathcal{A}(x - y_0), (x - y_0)], \end{cases}$$

kde $[[.]]$ značí normu indukovanou skalárním součinem $[., .]$. Tudíž (5.94) a (5.95) implikují platnost relace

$$(5.96) \quad (Ae_k, e_k) \leq [[Q_k(\mathcal{A})]]^2(\mathcal{A}(x - y_0), x - y_0)$$

pro každý polynom Q_k stupně k splňující podmínu $Q_k(0) = 1$. Protože \mathcal{A} je pozitivně definitní operátor, jeho spektrum je reálné a kladné. Budte $\underline{\lambda}$ a $\bar{\lambda}$ hranice intervalu, v němž leží spektrum $\sigma(\mathcal{A})$. Položme ještě

$$\mathcal{M} = I - \frac{2}{\underline{\lambda} + \bar{\lambda}} \mathcal{A}$$

a nechť pro každý polynom \tilde{Q}_k stupně k splňující $\tilde{Q}_k(1) = 1$ položme

$$Q_k(t) = \tilde{Q}_k\left(1 - \frac{2}{\underline{\lambda} + \bar{\lambda}} t\right).$$

Potom platí, že Q_k je polynom stupně k , splňuje $Q_k(0) = 1$ a

$$Q_k(\mathcal{A}) = \tilde{Q}_k(\mathcal{M}).$$

Tudíž,

$$(5.97) \quad [[Q_k(\mathcal{A})]] = [[\tilde{Q}_k(\mathcal{M})]].$$

Nyní pro takový polynom \tilde{Q}_k

$$(5.98) \quad [[\tilde{Q}_k(\mathcal{M})]] = \text{Max} \left\{ |\tilde{Q}_k(\lambda)| : \lambda \in \sigma(\mathcal{M}) \right\} = \text{Max} \left\{ |\tilde{Q}_k(\lambda)| : \lambda \in [-\rho, \rho] \right\},$$

kde

$$\rho = \frac{\bar{\lambda} - \underline{\lambda}}{\underline{\lambda} + \bar{\lambda}} < 1.$$

Z (5.96), (5.97) a (5.98) plyne platnost formule

$$(5.99) \quad (Ae_k, e_k) \leq \text{Max} \left\{ |\tilde{Q}_k(\lambda)| : \lambda \in [-\rho, \rho] \right\} (A(x - y_0), x - y_0)^{1/2}$$

pro jakýkoliv polynom \tilde{Q}_k splňující $\tilde{Q}_k(1) = 1$. Se znalostí jediné veličiny ρ je optimální volbou polynom určený pomocí Čebyshevových polynomů T_k , t.j.

$$\tilde{Q}_k(\lambda) = \frac{T_k(\lambda/\rho)}{T_k(1/\rho)}$$

Čebyševovy polynomy jsou definovány pomocí formulí

$$T_k(t) = \begin{cases} \cos(k \cos^{-1} t), & \text{je-li } |t| \leq 1 \\ \cosh(k \cosh^{-1} t), & \text{je-li } t \geq 1. \end{cases}$$

Platí tedy,

$$(5.100) \quad \text{Max} \left\{ |\tilde{Q}_k(\lambda)| : \lambda \in [-\rho, \rho] \right\} \leq \frac{1}{T_k(1/\rho)}.$$

Protože $\rho < 1$, můžeme klást $\sigma = \cosh^{-1}(1/\rho)$, takže

$$T_k(1/\rho) = \frac{e^{k\sigma} + e^{-k\sigma}}{2} = e^{k\sigma} \left[\frac{1 + e^{-2k\sigma}}{2} \right].$$

Dále pak,

$$(5.101) \quad \frac{1}{T_k(\frac{1}{\rho})} = e^{-k\sigma} \left[\frac{2}{1 + e^{-2k\sigma}} \right] \leq 2e^{-k\sigma}.$$

Avšak $\cosh \sigma = \frac{1}{\rho}$ a $\sinh \sigma = \sqrt{\frac{1}{\rho^2} - 1}$, takže $\cosh \sigma + \sinh \sigma = e^\sigma$ a tudíž, $\log \left[\frac{1}{\rho} + \sqrt{\frac{1}{\rho^2} - 1} \right]$ a tedy,

$$(5.102) \quad e^{-k\sigma} = \left[\frac{1}{\rho} + \sqrt{\frac{1}{\rho^2} - 1} \right]^{-k} = \left(\frac{\rho}{1 + \sqrt{1 - \rho^2}} \right)^k.$$

Do hry se dostává číslo podmíněnosti $\kappa(\mathcal{A}) = \frac{\bar{\lambda}}{\underline{\lambda}}$ takže jest potom

$$\rho = \frac{\kappa(\mathcal{A}) - 1}{\kappa(\mathcal{A}) + 1}.$$

Kombinujíce (5.101) a (5.102) zjistíme, že

$$(5.103) \quad \frac{1}{T_k(1/\rho)} \leq 2 \left(\frac{\kappa(\mathcal{A}) - 1}{\kappa(\mathcal{A}) + 1} \right)^k$$

a žádaný odhad má konečný tvar

$$(Ae_k, e_k)^{1/2} \leq 2 \left(\frac{\kappa(\mathcal{A}) - 1}{\kappa(\mathcal{A}) + 1} \right)^k (A(x - y_0), x - y_0)^{1/2}.$$

Za zmínku jistě stojí, že asymptoticky dává metoda sdružených gradientů lepší odhad než stacionární iterační proces s iterační maticí $\mathcal{A} = BA$. Jest totiž

$$\frac{\kappa(BA) - 1}{\kappa(BA) + 1} = \frac{\kappa(BA) - 1}{\kappa(BA) + 1 + 2\kappa^{1/2}(BA)} \leq \frac{\kappa(BA) - 1}{\kappa(BA) + 1}.$$

5.6 Zobecněné Bergstromovo lemma pro pozitivně definitní matice

Při vyšetřování konvergence některých metod řešení soustav lineárních algebraických rovnic bude užitečné následující lemma.

Lemma 5.4 *Budě H pozitivně definitní matici typu $N \times N$, $N \geq 1$. Nechť pro její spektrum $\sigma(H) = \{\lambda_j : j = 1, \dots, s\}$, $s \leq N$ platí vztahy*

$$\lambda_{\min} = \lambda_1 < \lambda_2 < \dots < \lambda_k < \dots < \lambda_s = \lambda_{\max}, \quad j, k = 1, \dots, s, \quad j \neq k.$$

Potom pro libovolný vektor $x \in \mathbb{R}^N$ platí nerovnosti

$$(5.104) \quad (x, x)^2 \leq (Hx, x) (H^{-1}x, x) \leq \frac{1}{\lambda_{\min} \lambda_{\max}} \left(\frac{\lambda_{\min} + \lambda_{\max}}{2} \right)^2 (x, x)^2.$$

Poznámka 5.7 Lemma 5.4 podává elegantní odhad dolní hranice čísla podmíněnosti matice H dané veličinou $V(x) = (Hx, x)$ ($H^{-1}x, x$) pomocí podílu čtverců aritmetického a geometrického průměru jejích extrémních vlastních hodnot.

Důkaz provedeme za použití následujících dvou lemmat.

Lemma 5.5 Necht $0 < \mu \leq \nu$,

$$f(x) = \mu \frac{1}{x} + \nu x, \quad 0 < a \leq x \leq 1.$$

Potom

$$(5.105) \quad f(x) \leq \max \left\{ \mu \frac{1}{a} + \nu a, \mu + \nu \right\}.$$

Důkaz. Snadno nahlédneme, že ze vztahů

$$f'(x) = -\frac{\mu}{x^2} + \nu = 0, \quad f''(x) = \frac{2\mu}{x^3} > 0,$$

plyne, že f je konvexní takže

$$\max\{f(x) : x \in [a, 1]\} = \max\{f(a), f(1)\}$$

a odtud tvrzení snadno plyne. \square

Necht $\zeta \in \mathcal{R}^N$, $\zeta^T = (\zeta_1, \dots, \zeta_N)$, $\zeta_j \geq 0, j = 1, \dots, s$ a

$$(5.106) \quad \sum_{j=1}^s \zeta_j = 1$$

Definujme dvě veličiny

$$(5.107) \quad \phi(\zeta) = \sum_{j=1}^s \zeta_j \lambda_j, \quad s = \text{card}\sigma(T),$$

$$(5.108) \quad \psi(\zeta) = \sum_{j=1}^s \zeta_j \lambda_j^{-1}.$$

Lemma 5.6 Pro veličiny definované v (5.107) a (5.108) platí následující tvrzení (a), (b), (c):

(a) $\lambda_{\min} \lambda_{\max} \psi(\zeta) \leq (\lambda_{\min} + \lambda_{\max}) - \phi(\zeta)$,

(b) Existuje $\lambda \in [\lambda_{\min}, \lambda_{\max}]$ takové, že $\phi(\zeta) = \lambda$.

(c)

$$\phi(\zeta) \psi(\zeta) \leq \left(\frac{\frac{1}{2} (\lambda_{\min} + \lambda_{\max})}{\sqrt{\lambda_{\min} \lambda_{\max}}} \right)^2.$$

Proof. Abychom ukázali platnosť tvrzení (a), vyšetřujme výrazy

$$\begin{aligned}\phi(\zeta) + \lambda_{\min} \lambda_{\max} \psi(\zeta) &= \lambda_{\min} \sum_{j=1}^s \zeta_j \frac{\lambda_j}{\lambda_{\min}} + \lambda_{\max} \sum_{j=1}^s \zeta_j \frac{\lambda_{\min}}{\lambda_j} \\ &= \sum_{j=1}^s \zeta_j \left[\lambda_{\min} \frac{\lambda_j}{\lambda_{\min}} + \lambda_{\max} \frac{\lambda_{\min}}{\lambda_j} \right].\end{aligned}$$

Aplikujeme-li Lemma 5.5 s

$$a = \frac{\lambda_{\min}}{\lambda_{\max}}, \quad \mu = \lambda_{\min}, \quad \nu = \lambda_{\max},$$

obdržíme, že

$$\phi(\zeta) + \lambda_{\min} \lambda_{\max} \psi(\zeta) \leq \sum_{j=1}^s \zeta_j (\lambda_{\min} + \lambda_{\max}) = \lambda_{\min} + \lambda_{\max}.$$

Tudíž, (a) platí.

Protože

$$\lambda_{\min} \leq \phi(\zeta) \leq \lambda_{\max},$$

platí i (b).

Konečně, na základě platnosti vztahů (a) a (b), a vědomi si toho, že pro libovolné reálné $\alpha > 0$ platí vztahy

$$\max \{(\alpha - \lambda)\lambda : 0 < \lambda_{\min} \leq \lambda \leq \lambda_{\max}\} = (1/4)\alpha^2$$

odvodíme, že

$$\begin{aligned}\phi(\zeta)\psi(\zeta) &\leq \frac{\lambda}{\lambda_{\min}\lambda_{\max}} (\lambda_{\min} + \lambda_{\max} - \lambda) \\ &\leq \frac{1}{\lambda_{\min}\lambda_{\max}} \left(\frac{1}{2} (\lambda_{\min} + \lambda_{\max}) \right)^2,\end{aligned}$$

Tím je důkaz lemmatu 5.5 proveden. \square

Důkaz lemmatu 5.4. Je li $s = 1$, pak $H = \lambda I$, a tvrzení platí.
Nechť tedy $s > 1$. Vyjděme ze spektrálního vyjádření

$$H = \sum_{j=1}^s \lambda_j P_j,$$

v němž

$$P_j^T = P_j = P_j^2,$$

$$P_j P_k = P_k P_j = \delta_{jk} P_j, \quad j, k = 1, \dots, s, \quad \sum_{j=1}^s P_j = I.$$

Pro $0 \neq x \in \mathcal{R}^N$ jest tedy,

$$V(x) = (Hx, x) (H^{-1}x, x) = \sum_j \lambda_j (P_j x, x) \sum_k \lambda_k^{-1} (P_k x, x)$$

a dále pak

$$(5.109) \quad V(x) = \sum_j (P_j x, x)^2 + \sum_j \sum_{k>j} \left(\frac{\lambda_j}{\lambda_k} + \frac{\lambda_k}{\lambda_j} \right) (P_j x, x) (P_k x, x),$$

takže

$$\begin{aligned} V(x) &\geq \sum_j (P_j x, x)^2 + 2 \sum_j \sum_{k>j} (P_j x, x) (P_k x, x) = \\ &\quad \left(\sum_j (P_j x, x) \right)^2 = (x, x)^2. \end{aligned}$$

Na druhé straně, položme

$$\zeta_j = \frac{(P_j x, x)}{(x, x)}, \quad j = 1, \dots, s$$

a aplikujme lemma 5.6. Tím však obdržíme platnost tvrzení dokazovaného v lemmatu 5.4.

□

Poznámka 5.8 Je zřejmé, že pro

$$H = \lambda I,$$

se nabude extrémní dolní (i horní) hranice

$$(Hx, x) (H^{-1}x, x) = (x, x)^2.$$

Podobně, pro

$$H = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix},$$

se nabývá extrémní horní hranice

$$\begin{aligned} (Hx, x) (H^{-1}x, x) &= x_1^4 + x_2^4 + \left(\frac{\lambda_1}{\lambda_2} + \frac{\lambda_2}{\lambda_1} \right) x_1^2 x_2^2 = \\ &= (x_1^2 + x_2^2)^2 \left(\frac{1}{2} (\lambda_1 + \lambda_2) \right)^2 \frac{1}{\lambda_1 \lambda_2}. \end{aligned}$$

Cvičení 5.3 Nalezněte všechny případy pozitivně definitních matic H , pro něž nastane některá z extrémních situací ve formuli (5.104).

Zobecněné Bergstromovo lemma pro operátory na nekonečně dimenzionálním Hilbertově prostoru

Předpokládejme, že \mathcal{E} je Hilbertův prostor a že $\mathbf{B}(\mathcal{E})$ značí prostor ohraničených lineárních operátorů zobrazujících \mathcal{E} do sebe.

Lemma 5.7 Budě $T \in \mathbf{B}(\mathcal{E})$ samoadjungovaný pozitivně definitní operátor.

Potom pro libovolný prvek $x \in \mathcal{E}$ platí vztahy

$$(5.110) \quad (x, x)^2 \leq (Hx, x) (H^{-1}x, x) \leq \frac{1}{\lambda_{\min} \lambda_{\max}} \left(\frac{\lambda_{\min} + \lambda_{\max}}{2} \right)^2 (x, x)^2,$$

při čemž λ_{\min} a λ_{\max} značí hranice spektra operátoru T , t. j. $\sigma(T) \subset [\lambda_{\min}, \lambda_{\max}]$ a $\lambda_{\min}, \lambda_{\max} \in \sigma(T)$.

Důkaz. Podobně jako formulace tak i důkaz lemmatu 5.7 je v podstatě totožný s důkazem lemmatu 5.4.

Budě $E = E(\lambda)$ spektrální míra operátoru T .

Vzhledem k předpokladu o nekonečnosti dimenze \mathcal{E} je nutné upravit důkaz části (a) v lemmatu 5.6, při čemž klademe

$$\phi(x) = \int_{\lambda_{\min}}^{\lambda_{\max}} \lambda (\mathrm{d}E(\lambda)x, x), \quad \psi(x) = \int_{\lambda_{\min}}^{\lambda_{\max}} \lambda^{-1} (\mathrm{d}E(\lambda)x, x).$$

Platí tedy,

$$\phi(x) + \lambda_{\min} \lambda_{\max} \psi(x) = \int_{\lambda_{\min}}^{\lambda_{\max}} \left[\lambda_{\min} \frac{\lambda}{\lambda_{\min}} + \lambda_{\max} \frac{\lambda_{\min}}{\lambda} \right] (\mathrm{d}E(\lambda)x, x)$$

a tedy též

$$(5.111) \quad \begin{cases} \phi(x) + \lambda_{\min} \lambda_{\max} \psi(x) \leq (\lambda_{\min} + \lambda_{\max}) \int_{\lambda_{\min}}^{\lambda_{\max}} (\mathrm{d}E(\lambda)x, x) \\ = (\lambda_{\min} + \lambda_{\max})(x, x). \end{cases}$$

Nerovnost ve vztahu (5.111) jsme odvodili pomocí lemmatu 5.5, kde jsme položili

$$a = \frac{\lambda_{\min}}{\lambda_{\max}}, \quad \mu = \lambda_{\min}, \quad \nu = \lambda_{\max}.$$

Dokázali jsme tak platnost tvrzení analogické k (a) v lemmatu 5.6. Ostatní části důkazu lemmatu 5.7 jsou shodné s odpovídajícími částmi důkazu lemmatu 5.6. Tímto konstatováním lze ukončit důkaz lemmatu 5.7. \square

6 Problémy vlastních hodnot

Jest na místě si uvědomit, že úloha nalézt vlastní číslo a příslušný vlastní vektor dané čtvercové matice je úloha *nelineární*. Zkušenost nám napovídá, že na rozdíl od lineárních úloh, pro které zpravidla jsme schopni nalézt finitní metody k nalezení přesných řešení, pro nelineární úlohy takové metody známy nejsou. Tak je tomu i v případě hledání vlastních čísel a vlastních vektorů. Dá se dokonce říci, že iterační metody jsou v tomto případě jedinými spolehlivými prostředky jak sestrojovat patřičné approximace hledaných vlastních prvků.

Základní metodou v problematice vlastních čísel je *mocninná metoda*.

6.1 Mocninná metoda

Konvergence mocninné metody je založena na jednoduchém pozorování, že totiž posloupnost $\{\lambda^k\}$, kde λ je komplexní číslo, je buď *divergentní* a to když $|\lambda| \geq 1$, $\lambda \neq 1$, *stacionární*, když $\lambda = 1$ anebo *konvergentní k nule*, když $|\lambda| < 1$.

Tuto skutečnost reflektuje posloupnost $\{T^k\}$, kde T je matice typu $N \times N$ s komplexními prvky t_{jk} , $j, k = 1, 2, \dots, N$.

Protože chování posloupnosti mocnin $\{T^k\}$ je pro vyšetřování algoritmů k nalezení vlastních hodnot matice T rozhodující, bude užitečné připomenout si některé důležité vlastnosti funkcí, jejichž argumentem je matice. K tomu nám poslouží známé tvrzení z lineární algebry

Tvrzení 6.1 *Budě $T = (t_{jk})$ matice typu $N \times N$ nad tělesem komplexních čísel. Budě*

$$\sigma(T) = \{\lambda_1, \dots, \lambda_N\}, \quad s \leq N,$$

spektrum matice T a nechť

$$r_1, \dots, r_s, \quad \sum_{j=1}^s r_j = N,$$

jsou násobnosti vlastních hodnot $\lambda_1, \dots, \lambda_s$ v charakteristickém mnohočenu

$$\chi_T(\lambda) = \det(\lambda I - T) = \prod_{j=1}^s (\lambda - \lambda_j)^{r_j}.$$

Dále nechť ϕ označuje minimální polynom matice T , t.j.

$$\phi(\lambda) = \prod_{j=1}^s (\lambda - \lambda_j)^{q_j},$$

při čemž $1 \leq q_j \leq r_j$.

Potom existuje $N \times N$ matice V , $\det V \neq 0$, taková, že

$$(6.1) \quad VTV^{-1} = \begin{pmatrix} \mathcal{J}_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \mathcal{J}_s \end{pmatrix}.$$

V tomto využádření je každý z bloků \mathcal{J}_{jr_j} sám o sobě blokově diagonální matici. Pro případ $q_j = 1$ jsou její diagonální bloky skalárni matice t.j. $\mathcal{J}_{jl} = \lambda_j I$, kde I je jednotková matica příslušného rozměru. Pro $q_j > 1$ mohou některé diagonální bloky být skalárni matice avšak alespoň jeden z bloků má tvar

$$J_{jl} = \begin{pmatrix} \lambda_j & 1 & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \\ & & & \lambda_j \end{pmatrix},$$

$$\dim J_{jl} = t_{j_l}$$

při čemž t_1, \dots, t_s jsou přirozená čísla, (tedy nikoliv nula) a platí rovnosti

$$\sum_{l=1}^{t_j} \dim J_{jl} = r_j,$$

$$\max\{\dim J_{jl} : j = 1, \dots, s; l = 1, \dots, t_j\} = q_j.$$

Posléze diagonální blok \mathcal{J}_j ve využádření (6.1) tvoří všechny bloky odpovídající vlastní hodnotě $\lambda_j, j = 1, \dots, s$.

Příklad 6.1 Nechť spektrum matici T tvoří vesměs navzájem různá vlastní čísla, tedy $s = N$. Potom $q_j = r_j = 1$ a všechny bloky $\mathcal{J}_{jl}, j = 1, \dots, N$, jsou nutně diagonální, takže

$$VT\bar{V}^{-1} = \text{diag}\{\lambda_1, \dots, \lambda_N\}.$$

Příklad 6.2 Nechť $T = \lambda I$, $N > 1$. Zřejmě $\sigma(T) = \{\lambda\}$. Tedy $s = 1 < N$ a $r_1 = N$, zatímco $q_1 = 1$.

Příklad 6.3 Nechť $N > 1$ a

$$T = \begin{pmatrix} \lambda & 1 & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \\ & & & \lambda \end{pmatrix}.$$

Zřejmě, $\sigma(T) = \{\lambda\}$, $s = 1$ a $r_1 = N$. Na rozdíl od příkladu 6.1, $q_1 = N$.

Definice 6.1 Matici T typu $N \times N$ má jednoduchou strukturu, jestliže $q_j = 1$ pro $j = 1, \dots, s$, kde s značí počet různých vlastních hodnot matice T .

Poznámka 6.1 Matice v příkladech 6.1 a 6.2 mají jednoduchou strukturu; matice z příkladu 6.3 jednoduchou strukturu nemá.

Zavedme následující označení.

$$B_{j1} = V^{-1} \begin{pmatrix} 0 & & & \\ & \ddots & & \\ & & I_{r_j} & \\ & & & \ddots \\ & & & & 0 \end{pmatrix} V,$$

kde I_{r_j} je jednotková matice rozměru $r_j = \sum_{l=1}^{t_j} \dim J_{jl}$. Vidíme, že

$$\sum_{j=1}^s B_{j1} = I,$$

kde $I = I_N$.

Dále nechť

$$B_{jk} = V^{-1} \begin{pmatrix} 0 & & & \\ & \ddots & & \\ & & \mathcal{J}_j - \lambda_j I_{r_j} & \\ & & & \ddots \\ & & & & 0 \end{pmatrix} V, j = 1, \dots, s, k = 1, 2, \dots$$

Snadno se přesvědčíme, že platí vztahy

$$B_{jk+1} = (T - \lambda_j I) B_{jk} = (T - \lambda_j I)^k B_{j1},$$

tedy,

$$B_{jk} = 0 \text{ pro } k > q_j.$$

Dále pak platí formule

$$B_{j1} B_{t1} = \delta_{jt} B_{j1}, \quad j, t = 1, \dots, s$$

a

$$B_{jk} B_{tl} = \delta_{jt} B_{jk+l-1};$$

speciálně

$$B_{j1}^2 = B_{j1}, \quad j = 1, \dots, s,$$

Podobně se lze přesvědčit, že

$$T = \sum_{j=1}^s (\lambda_j B_{j1} + B_{j2})$$

a tedy,

$$T^k = \sum_{j=1}^s (\lambda_j B_{j1} + B_{j2})^k$$

$$= \sum_{j=1}^s \sum_{t=1}^{q_j} \frac{f_k^{(t-1)}(\lambda_j)}{(t-1)!} B_{jt},$$

kde

$$f_k(\lambda) = \lambda^k,$$

takže

$$\begin{aligned} f'_k(\lambda) &= k\lambda^{k-1}, \\ f''_k(\lambda) &= k(k-1)\lambda^{k-2}, \\ &\dots \\ f_k^{(t)}(\lambda) &= k(k-1)\dots(k-t+1)\lambda^{k-t}. \end{aligned}$$

Typickou pro konvergenci mocninné metody je následující situace.

Nechť spektrum $\sigma(T) = \{\lambda_1, \dots, \lambda_s\}$ má strukturu danou vztahy

$$(6.2) \quad |\lambda_1| > |\lambda_j|, \quad j > 1, \quad q_1 = 1.$$

Potom

$$\left(\frac{1}{\lambda_1}T\right)^k = B_{11} + \sum_{j=2}^s \sum_{t=1}^{q_j} \frac{g_k^{(t-1)}(\lambda_j)}{(q_j-1)!} B_{jt},$$

kde

$$g_k(\lambda) = \left(\frac{\lambda}{\lambda_1}\right)^k.$$

Je okamžitě patrné, že

$$\lim_{k \rightarrow \infty} \left(\frac{1}{\lambda_1}T\right)^k = B_{11},$$

neboť

$$\lim_{k \rightarrow \infty} g_k^{(t)}(\lambda_j) = 0$$

pro $j = 2, \dots, s$ a $t = 0, 1, \dots, q_j - 1$, $j > 1$.

Věta 6.1 Je-li $q_1 = 1$, pak posloupnost

$$\left\{ \left(\frac{1}{\lambda_1}T\right)^k \right\}$$

konverguje k B_{11} s rychlosťí s jakou skalární posloupnost

$$\max \left\{ \left| \frac{\lambda_j}{\lambda_1} \right|^k k^{q_j-1} : j = 2, \dots, s \right\}$$

konverguje k nule.

Je-li $q_1 > 1$, pak platí

$$(6.3) \quad (q_1 - 1)! \lim_{k \rightarrow \infty} \frac{1}{\lambda_1^{k-q_1+1}} k^{-q_1+1} T^k = B_{1q_1}.$$

Předchozí teorie nabízí následující obecný algoritmus mocninné metody.

Algoritmus 6.1 Budě $\{x'_k\}$ posloupnost lineárních funkcionálů na C^N takových, že

$$(6.4) \quad x'_k(B_{11}x_0) \neq 0 \text{ a } \lim_{k \rightarrow \infty} x'_k = x'_\infty,$$

$$(6.5) \quad (B_{11})^T x'_\infty \neq 0.$$

Budě $\varepsilon > 0$ zadaná tolerance.

1. Zvolme x_0 .
2. Položme $k = 0$.
3. Sestrojme vektor x_k kladouče

$$x_k = Tx_k.$$

4. Určeme $\lambda^{(k)}$ pomocí předpisu

$$\lambda^{(k)} = \frac{x'_k(x_k)}{x'_k(x_{k-1})}$$

a

$$x_{k+1} = \frac{1}{\lambda^{(k)}} x_k.$$

5. Prověřme, zda

$$(6.6) \quad |\lambda^{(k)} - \lambda^{(k-1)}| \begin{cases} \geq \varepsilon & (\text{NE}), \\ < \varepsilon & (\text{ANO}) \end{cases}$$

5. Nastane - li v (6.6) případ NE, položme $k + 1 \rightarrow k$ a GOTO 3.

6. Nastane - li v (6.6) případ ANO, GOTO 7.

7. Položme $\lambda_1 = \lambda^{(k)}$ a STOP.

Potřebné vlastnosti algoritmu 6.1 jsou obsahem následující věty, již formulujeme ve zjednodušené verzi

Věta 6.2 Budě T matici typu $N \times N$, taková, že platí vztahy (6.2), (6.5) a $x'_k = x'$ pro $k = 0, 1, \dots$ takže (6.5) platí automaticky.

Potom platí rovnosti

$$(6.7) \quad \lim_{k \rightarrow \infty} x_k = x_\infty, \quad Tx_\infty = \lambda_1 x_\infty,$$

a

$$(6.8) \quad \lim_{k \rightarrow \infty} \lambda^{(k)} = \lambda_1.$$

Důkaz. Snadno ověříme, že

$$x_{k+1} = \prod_{l=0}^k \frac{1}{\lambda^{(l)}} T^{k+1} x_0$$

a

$$\begin{aligned} \prod_{l=0}^k \lambda^{(l)} &= \prod_{l=0}^k \frac{x'(Tx_l)}{x'_l(x)} \\ &= \frac{x'(T^{k+1}x_0)}{x'(x_0)} \end{aligned}$$

a tudíž

$$x_{k+1} = \frac{x'(x_0)}{x'(T^{k+1}x_0)} T^{k+1} x_0.$$

Na základě věty 6.1 odvodíme, že

$$\lim_{k \rightarrow \infty} x_k = x'(x_0) \lim_{k \rightarrow \infty} \frac{B_{1q_1}x_0 + Z_k x_0}{x'(B_{1q_1}x_0) + x'(Z_k)},$$

kde, vzhledem k (6.3)

$$\lim_{k \rightarrow \infty} Z_k x_0 = 0.$$

Posléze tedy,

$$\lim_{k \rightarrow \infty} x_k = \frac{1}{x'(B_{1q_1}x_0)} B_{1q_1}x_0.$$

Jako důsledek obdržíme vztahy

$$\begin{aligned} & \lim_{k \rightarrow \infty} \lambda^{(k)} = \\ &= \lim_{k \rightarrow \infty} \frac{x'(T^k x_0)}{x'(T^{k-1} x_0)} \\ &= \lambda_1 \lim_{k \rightarrow \infty} \frac{x'(k^{-q_1+1} (\frac{1}{\lambda_1} T)^k x_0)}{x'(k^{-q_1+1} (\frac{1}{\lambda_1} T)^{k-1} x_0)} = \lambda_1. \end{aligned}$$

Tím je důkaz věty 6.2 proveden. |||

6.2 Algoritmy LR a QR

Definice 6.2 Budě A matici typu $N \times N$ nad tělesem reálných čísel. **LR** rozkladem matici A nazveme dvojici matic typu $N \times N$ takovou, že platí:

(a) $L = (l_{jk})$ je dolní trojúhelníková s jednotkovou diagonálou, t.j.

$$l_{jk} = \begin{cases} 1 & \text{pro } j = k \\ 0 & \text{pro } j < k \end{cases} \quad j, k = 1, \dots, N,$$

(b) $R = (r_{jk})$ je horní trojúhelníková, t.j.

$$r_{jk} = 0 \text{ pro } j > k, \quad j, k = 1, \dots, N.$$

(c) $A = LR$.

Poznámka 6.2 Postačující podmínkou pro existenci LR rozkladu $N \times N$ matici A je plná regularita matici A , t.j. všechny hlavní minory matici A jsou regulární.

Algoritmus 6.2 Budě A matici typu $N \times N$ a $\varepsilon > 0$ daná tolerance.

(1) Položme $A = A_1$ a $k = 1$.

(2) Nechť

$$A_k = L_k R_k$$

je LR rozklad matici A_k .

(3) Dále položme

$$R_k L_k = A_{k+1}, \quad R_k = (r_{jl}^{(k)}).$$

(4) Posudíme, zda

$$\max \left\{ \left| r_{jl}^{(k)} \right| : j, l = 1, \dots, N, j > l \right\} \begin{cases} \geq \varepsilon & \text{NE} \\ < \varepsilon & \text{ANO.} \end{cases}$$

(5) Je - li v (4) NE, položme $k + 1 \rightarrow k$ a GOTO (2).

(6) Je - li v (4) ANO, položme

$$R_\infty = R_{k+1}$$

a GOTO (7).

(7) STOP.

Tvrzení 6.2 Nechť LR rozklady

$$(6.9) \quad A_k = L_k R_k$$

existují pro $k = 1, 2, \dots$

Potom platí

A_{k+1} a A_k jsou podobné matice; jmenovitě

$$(6.10) \quad A_{k+1} = L_k^{-1} A_k L_k,$$

$$(6.11) \quad A_{k+1} = (L_1 \dots L_k)^{-1} A_1 (L_1 \dots L_k).$$

Matice T_k , kde

$$T_k = L_1 \dots L_k$$

je dolní trojúhelníková a matice U_k ,

$$U_k = R_k \dots R_1$$

je horní trojúhelníková, při čemž

$$(6.12) \quad A^k = A_1^k = T_k U_k.$$

Důkaz. Vztahy (6.10) a (6.11) jsou zřejmé. Pomocí nich pak odvodíme, že

$$L_1 \dots L_k A_{k+1} = A_1 L_1 \dots L_k.$$

Dále pak,

$$\begin{aligned} T_k U_k &= (L_1 \dots L_{k-1}) (L_k R_k) (R_k \dots R_1) \\ &= (L_1 \dots L_{k-1}) A_k (R_{k-1} \dots R_1) \\ &= A_1 (L_1 \dots L_{k-1}) (R_{k-1} \dots R_1) \\ &= A_1 T_{k-1} U_{k-1} \\ &= \dots \dots \dots \\ &= A_1^k = A^k. \end{aligned}$$

|||

Matice T_k a U_k tvoří LR rozklady pro A^k , neboť zřejmě

$$\operatorname{diag} T_k = I.$$

Naším cílem je ukázat, že

$$(6.13) \quad \lim_{k \rightarrow \infty} L_k = I$$

a

$$(6.14) \quad \lim_{k \rightarrow \infty} A_k = \lim_{k \rightarrow \infty} R_k = \begin{pmatrix} \lambda_1 & \times & . & . & . & \times \\ 0 & \lambda_2 & . & . & . & \times \\ . & . & . & . & . & . \\ 0 & 0 & . & . & . & \lambda_N \end{pmatrix},$$

kde $\lambda_1, \dots, \lambda_N$ jsou vlastní čísla.

Věta 6.3 Nechť $A = A_1$ splňuje

1⁰. LR rozklady matic $R_k L_k$ existují pro $k = 1, 2, \dots$

2⁰. Nechť

$$(6.15) \quad |\lambda_1| > |\lambda_2| > \dots |\lambda_N| > 0.$$

3⁰. Předpokládejme, že pro matice X a $Y = X^{-1}$, kde $A = XDY$ je Jordanova forma matice A , existují LR rozklady

$$(6.16) \quad X = L_X R_X, \quad Y = L_Y R_Y, \quad \text{diag}L_X = \text{diag}L_Y = I.$$

Potom platí (6.13) a (6.14).

Důkaz. Podle předpokladu existuje D^{-1} takže

$$A^k = XD^kY = L_X R_X D^k L_Y R_Y = L_X R_X (D^k L_Y D^{-k}) D^k R_Y.$$

Avšak

$$D^k L_Y D^{-k} = (l_{jt}^{(k)})$$

je dolní trojúhelníková matice, při čemž

$$l_{jt}^{(k)} = \left(\frac{\lambda_j}{\lambda_t} \right)^k l_{jt}, \quad \text{kde } L_Y = (l_{jt} \text{ pro } j > t)$$

a

$$l_{jt}^{(k)} = \begin{cases} 1 & \text{pro } j = t \\ 0 & \text{pro } j < t \end{cases}$$

Protože

$$|\lambda_t| > |\lambda_j| \text{ pro } j > t,$$

odvodíme, že

$$\lim_{k \rightarrow \infty} l_{jt}^{(k)} = \lim_{k \rightarrow \infty} \left(\frac{\lambda_j}{\lambda_t} \right)^k = 0 \text{ pro } j > t.$$

Tudíž

$$D^{-k} L_Y D^k = I + E_k,$$

při čemž

$$\lim_{k \rightarrow \infty} E_k = 0.$$

Z předchozího plyne, že

$$\begin{aligned} A^k &= L_X R_X (I + E_k) D^k R_Y \\ &= L_X (I + R_X E_k R_X^{-1}) R_X D^k R_Y \\ &= L_X (I + F_k) R_X D^k R_Y, \end{aligned}$$

kde

$$F_k = R_X E_k R_X^{-1}.$$

Protože $E_k \rightarrow 0$, pro všechna dostatečně velká k existují LR rozklady

$$I + F_k = \tilde{L}_k \tilde{R}_k, \quad \text{diag} \tilde{L}_k = I, \quad \tilde{l}_{jt} = 0 \text{ pro } j < t.$$

Díky tomu, že $E_k \rightarrow 0$, též $F_k \rightarrow 0$, odvodíme snadno, že

$$\lim_{k \rightarrow \infty} \tilde{L}_k = I = \lim_{k \rightarrow \infty} \tilde{R}_k.$$

Posléze,

$$A^k = (L_X \tilde{L}_k) (\tilde{R}_k R_X D^k R_Y).$$

Z jednoznačnosti trojúhelníkových rozkladů plyne, že

$$T_k = L_1 \dots L_k = L_X \tilde{L}_k$$

a

$$U_k = R_k \dots R_1 = \tilde{R}_k R_X D^k R_Y.$$

Protože $\tilde{L}_k \rightarrow I$, $\tilde{R}_k \rightarrow I$, obdržíme, že platí vztahy

$$\begin{aligned} \lim_{k \rightarrow \infty} T_k &= L_X, \\ \lim_{k \rightarrow \infty} L_k &= \lim_{k \rightarrow \infty} T_{k-1}^{-1} T_k = I, \\ \lim_{k \rightarrow \infty} R_k &= \lim_{k \rightarrow \infty} U_k U_{k-1}^{-1} \\ &= \lim_{k \rightarrow \infty} \tilde{R}_k R_X D^k R_Y R_Y^{-1} D^{(-k+1)} R_X^{-1} \tilde{R}_k^{-1} \\ &= R_X D R_X^{-1}. \end{aligned}$$

Avšak

$$R_X D R_X^{-1} = \begin{pmatrix} \lambda_1 & \times & & \times & \times \\ 0 & \lambda_2 & & \times & \times \\ & & \ddots & & \\ & & & \ddots & \\ 0 & 0 & & 0 & \lambda_N \end{pmatrix}$$

a tím je důkaz věty 6.3 proveden. \square

7 Počáteční a okrajové problémy pro obyčejné diferenciální rovnice

7.1 Základní poznatky z teorie obyčejných diferenciálních rovnic

V této kapitole budeme vyšetřovati metod přibližného řešení rovnic typu

$$(7.1) \quad y' = f(x, y), \quad x \in [a, b],$$

$$(7.2) \quad y(a) = y_0,$$

kde $-\infty < a < b < +\infty$ a soustav takových rovnic.

$$(7.3) \quad \begin{cases} y'_1 = f_1(x, y_1, \dots, y_N) \\ \dots \\ y'_N = f_N(x, y_1, \dots, y_N) \end{cases}, \quad x \in [a, b]$$

a

$$(7.4) \quad y_j(0) = y_{j0}, \quad j = 1, \dots, N.$$

Poznamenejme, že problém (7.3) - (7.4) zahrnuje jako speciální úlohu

$$(7.5) \quad z^N = f(x, z, z', \dots, z^{N-1}), \quad x \in [a, b],$$

s podmínkami

$$(7.6) \quad z^{j-1}(a) = z_{j0}, \quad j = 1, \dots, N-1.$$

Stačí položit

$$z = y_1, \quad z' = y_2, \quad \dots, \quad z^{N-1} = y_N,$$

takže

$$y'_1 = y_2, \quad y'_2 = y_3, \quad \dots, \quad y'_{N-2} = y_{N-1}$$

a

$$y'_N = z^{(N)} = f(x, z, z', \dots, z^{N-1}) = f(x, y_1, \dots, y_N).$$

Kromě úloh typu (7.1) - (7.2) a (7.3) - (7.4) budeme vyšetřovat též úlohy typu

$$y' = f(x, y),$$

avšak místo počáteční podmínky (7.2) budeme požadovat, aby platila rovnost

$$(7.7) \quad r(y(a), y(b)) = 0,$$

kde r je vhodná funkce dvou proměnných.

Obecněji budeme vyšetřovat úlohy

$$y'_j = f_j(x, y_1, \dots, y_N), \quad j = 1, \dots, N,$$

s podmínkami

$$(7.8) \quad \mathbf{0} = \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)),$$

kde

$$(7.9) \quad \mathbf{r} = \mathbf{r}(\mathbf{u}, \mathbf{v}) = \begin{bmatrix} r_1(u_1, \dots, u_N; v_1, \dots, v_N) \\ \dots \\ r_N(u_1, \dots, u_N; v_1, \dots, v_N) \end{bmatrix}.$$

Úlohy typu (7.1) a (7.7), respektive (7.3) a (7.8) se nazývají *úlohami okrajovými* na rozdíl od úloh *počátečního* typu (7.1) a (7.2) respektive (7.3) a (7.4).

Numerické metody přibližného řešení výše uvedených úloh je založeno na vlastnostech řešení těchto úloh. Zformulujme některá tvrzení potřebná pro další vyšetřování.

Věta 7.1 Předpokládejme, že komponenty f_1, \dots, f_N vektoru \mathbf{f} jsou definovány a jsou spojité na oblasti

$$\mathcal{S} = \{(x, y) : a \leq x \leq b, \mathbf{y} \in R^N\}.$$

Dále nechť existuje lipschitzovská konstanta L tak; že

$$(7.10) \quad \|\mathbf{f}(x, \bar{\mathbf{y}}) - \mathbf{f}(x, \tilde{\mathbf{y}})\| \leq L\|\bar{\mathbf{y}} - \tilde{\mathbf{y}}\|,$$

platí pro všechna $x \in [a, b]$ a $\mathbf{y} \in R^N$ (Lipschitzova podmínka).

Potom ke každému $x_0 \in [a, b]$ a $\mathbf{y}^0 \in R^N$ existuje právě jedna vektor - funkce $\mathbf{y}(x)$ tak, že

1) \mathbf{y} je spojitá a spojitě diferencovatelná vektor - funkce v $[a, b]$;

2) $\mathbf{y}'(x) = \mathbf{f}(x, \mathbf{y}(x))$, $x \in [a, b]$;

3) $\mathbf{y}(x_0) = \mathbf{y}^0$.

Poznámka. Již víme, vze Lipschitzova podmínka je splněna, pakliže parciální derivace $\frac{\partial f_j}{\partial x}$, $j = 1, \dots, N$ existují a jsou spojité na \mathcal{S} .

V praxi se můžeme setkat převážně s případy, kdy f_j , $j = 1, \dots, N$, jsou spojité na \mathcal{S} , avšak derivace $\frac{\partial f_j}{\partial x}$ mohou být na \mathcal{S} neohraničené. V takovém případě má počáteční úloha (7.3) - (7.4) řešení, ale to nemusí být definováno na celé oblasti.

Příklad 7.1 Nechť $N = 1$ a

$$y' = y^2, \quad y(0) = 1, \quad 0 \leq x < +\infty.$$

Vidíme, že

$$y(x) = \frac{1}{1-x}$$

a to je ohraničená funkce jen pro $x < 1$.

Pro numerické počítání je důležité, že řešení závisí na počátečních datech spojitě. Tvrzení na toto téma je obsahem následující věty.

Věta 7.2 Za předpokladu věty 7.1 platí odhad

$$(7.11) \quad \|\mathbf{y}(x, \mathbf{a}^1) - \mathbf{y}(x, \mathbf{a}^2)\| \leq e^{L(x-a)}\|\mathbf{a}^1 - \mathbf{a}^2\|,$$

kde

$$(7.12) \quad \mathbf{y}'(x, \mathbf{a}^j) = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{a}^j, \quad j = 1, 2.$$

Významné pro teorii i praxi jsou lineární úlohy typu

$$(7.13) \quad \mathbf{Y}' = T(x)\mathbf{Y}, \quad \mathbf{Y}(a) = \mathbf{Y}^a,$$

přičemž $T = T(x)$ je matice typu $N \times N$. Růst řešení je charakterizován v následující větě.

Věta 7.3 Nechť prvky matici $T(x)$ jsou spojité funkce na $[a, b]$ a nechť

$$(7.14) \quad \kappa(x) = \|T(x)\|.$$

Potom platí odhad

$$(7.15) \quad \|\mathbf{Y}(x) - I\| \leq \exp \left\{ \int_a^x \kappa(t) dt \right\} - 1 \text{ pro } x \geq a.$$

7.2 Racionální approximace exponenciální funkce

Jednou z řady ekvivalentních definic exponenciální funkce může posloužit formule daná Taylorovou řadou. Pro všechna komplexní čísla je platný rozvoj

$$(7.16) \quad \exp\{z\} \equiv e^z = \sum_{k=0}^{\infty} \frac{z^k}{k!}, \quad z \in \mathcal{C}^1.$$

Snadno si uvědomíme, že platí

$$(7.17) \quad |e^z| \leq 1 \text{ pro } z \in \mathcal{C}^1, \Re z \leq 0$$

a dále, že

$$(7.18) \quad \lim_{z \rightarrow \infty} e^{|z|} = 0.$$

Dále nechť $P = P(z)$ a $Q = Q(z)$ označují polynomy s komplexními koeficienty. Bude-li to vhodné, budeme indexy vyznačovat odpovídající stupeň. Tedy, na př.

$$P_0(z) = 1, \quad P_1(z) = 2z + 1/2, \quad Q_1(z) = 1 - z,$$

a tak podobně.

Budě $R = R(z)$ racionální funkce daná polynomy P a Q :

$$(7.19) \quad R(z) = \frac{P(zP)}{Q(z)}, \quad z \in \mathcal{C}^1.$$

V následujícím textu budeme předpokládat, že polynomy v (7.19) jsou nesoudělné, tedy, že neexistují polynomy \tilde{P}, \tilde{Q} , při čemž stupeň polynomu \tilde{Q} je nižší než stupeň polynomu Q , takové, že

$$R(z) = \frac{\tilde{P}(z)}{\tilde{Q}(z)}, \quad z \in \mathcal{C}^1.$$

Definice 7.1 Řekneme, že racionální funkce R je Padéovou approximací exponenciální funkce řádu p , pakliže platí vztahy

$$(7.20) \quad e^z - R(z) = g(z)z^{p+1} = o(z^{p+1}), \quad z \in \Omega_0,$$

kde Ω_0 je nějaké okolí bodu $0 \in \mathcal{C}^1$ a přitom

$$0 \neq \kappa = \sup \{ |g(z)| : z \in \Omega_0 \} < +\infty.$$

Příklady.

$$(7.21) \quad R_{1,0}(z) = \frac{P_1(z)}{Q_0(z)} = 1 + z, \quad P_1(z) = z, \quad Q_0(z) = 1,$$

$$(7.22) \quad R_{0,1}(z) = \frac{P_0(z)}{Q_1(z)} = \frac{P_0(z)}{Q_1(z)} = \frac{1}{1-z}, \quad P_0(z) = 1, \quad Q_1(z) = 1 - z,$$

$$(7.23) \quad R_{1,1}(z) = \frac{P_1(z)}{Q_1(z)} = \frac{2+z}{2-z}, \quad P_1(z) = 2+z, \quad Q_1(z) = 2-z.$$

Snadno zjistíme, že

$$R_{0,1}(z) = 1 + z + z^2 + \dots,$$

a

$$R_{1,1}(z) = 1 + z + \frac{1}{2}z^2 + \dots$$

a tedy,

$$(7.24) \quad |e^z - R_{1,0}(z)| = O(z^2), \quad p = 1,$$

$$(7.25) \quad |e^z - R_{0,1}(z)| = O(z^2), \quad p = 1,$$

$$(7.26) \quad |e^z - R_{1,1}(z)| = O(z^3), \quad p = 2.$$

Definice 7.2 *Racionální approximace R exponenciální funkce se nazývá A-přijatelnou, platí-li vztahy*

$$(7.27) \quad |R(z)| \leq 1 \text{ pro } \Re z \leq 0.$$

Příklady.

Snadno ověříme, že $R_{1,0}$ v (7.21) není A-přijatelná. Na druhé straně však, $R_{0,1}$ v (7.22) a $R_{1,1}$ v (7.23) A-přijatelné jsou a navíc pro ně platí

$$(7.28) \quad \lim_{z \rightarrow \infty} |R_{0,1}(z)| = 0$$

a

$$(7.29) \quad \lim_{z \rightarrow \infty} |R_{1,1}(z)| = 1.$$

Definice 7.3 *Racionální approximace $R_{1,0}, R_{0,1}$ a $R_{1,1}$ v (7.21)-(7.23) nesou názvy svých původních autorů. Approximace $R_{1,0}$ se tedy nazývá Eulerovou přímou approximací, $R_{0,1}$ Eulerovou zpětnou approximací a $R_{1,1}$ approximací Cranka-Nicholsonové.*

Test 7.1 Budě λ takové, že $\Re \lambda < 0$ je v absolutní hodnotě dostatečně velké číslo. Jest tedy číslo $|e^\lambda|$ velmi malé. Ptáme se, jak tuto skutečnost zprostředkovávají jednotlivé racionální approximace exponenciální funkce.

Protože v tomto případě,

$$\lim_{\Re \lambda \rightarrow -\infty} |e^\lambda| = \lim_{\Re \lambda \rightarrow -\infty} e^{\Re \lambda} = 0,$$

1⁰ Eulerova dopředná approximace dává

$$\begin{aligned} |e^\lambda - R_{1,0}(\lambda)| &= |e^\lambda - 1 - \lambda| \\ &\geq |\lambda| - 1 - e^{\Re \lambda}. \end{aligned}$$

Vidíme tedy, že tato approximace poskytuje libovolně velkou chybu.

Naproti tomu,

2⁰ pro zpětnou Eulerovu metodu zjistíme, že platí

$$\begin{aligned} |e^\lambda - R_{0,1}(\lambda)| &= \left| e^\lambda - \frac{1}{1-\lambda} \right| \\ &\leq \sum_{k=2}^{\infty} \frac{z^k}{k!} - \sum_{k=2}^{\infty} z^k \\ &= O(|z|^2) \end{aligned}$$

takže chyba approximace je pro $\Re \lambda < 0$ malá.

3⁰ Ještě uspokojivější se jeví situace pro approximaci Crankovu-Nicholsonové,

$$|e^\lambda - R_{1,1}(\lambda)| = \left| e^\lambda - \frac{2+\lambda}{2-\lambda} \right| = O(|z|^3).$$

7.3 Počáteční úlohy

Jako motivaci vyšetřujme skalární úlohu

$$(7.30) \quad y' = f(x, y), \quad y(x_0) = y_0, \quad x \in [0, 1], \quad y \in R^1.$$

Budě $h > 0$. Nahradme v (7.30) derivaci diferencí čímž obdržíme

$$\frac{y(x+h) - y(x)}{h} = f(x, y(x)),$$

neboli

$$(7.31) \quad y(x+h) = y(x) + h f(x, y(x)).$$

Obdrželi jsme tak algoritmus pro řešení.

Algoritmus 7.1 Budě M celé kladné. Zvolme

$$h = \frac{x - x_0}{M}$$

a počítejme rekursivně pomocí formulí

$$\eta_0 = y_0,$$

$$\eta_{k+1} = \eta_k + h f(x_k, \eta_k),$$

kde

$$x_{k+1} = x_k + h, \quad k = 0, 1, \dots, M - 1.$$

Algoritmus 7.1 popisuje t. zv. *Eulerovu polygonální metodu*. Všimněme si, že v Algoritmu 7.1 platí vztahy

$$\eta(x_0; h) = y_0,$$

a

$$\eta(x_0 + h; h) = \eta(x; h) + h f(x, \eta(x; h)).$$

To navádí na myšlenku definovati algortimus, v němž roli $f(x; \eta(x; h))$ převezme vhodná funkce $\Phi = \Phi(x, y; h, f)$.

Algoritmus 7.2 Volme $h > 0$ a nechť

$$(7.32) \quad \begin{cases} \eta_0 = y_0 \\ \eta_{k+1} = \eta_k + h \Phi(x_k, \eta_k; h; f) \text{ pro } k = 1, 2, \dots \end{cases}$$

kde

$$x_{k+1} = x_k + h.$$

Již víme, že pro Eulerovu polygonální metodu jest

$$\Phi(x, y, h, f) = f(x, y),$$

Φ tedy nezávisí na h .

Pro další účely označme pro libovolné $x \in [a, b]$ a $\mathbf{y} \in R^N$ symbolem $\mathbf{z} = \mathbf{z}(x)$ přesné řešení úlohy

$$(7.33) \quad \mathbf{z}'(t) = \mathbf{f}(x, \mathbf{z}(t)), \quad \mathbf{z}(x) = \mathbf{y}(x).$$

Definujme funkci Δ předpisem

$$(7.34) \quad \Delta(x, y; h; \mathbf{f}) = \begin{cases} \frac{\mathbf{z}(x+h) - \mathbf{y}}{h} & \text{pro } h > 0 \\ \mathbf{f}(x, \mathbf{y}) & \text{pro } h = 0. \end{cases}$$

V dalším textu nebudeme vyznačovati f v argumentu funkce Δ a Φ .

Pomocí funkce Δ definujeme t. zv. *lokální diskretizační chybu*, tedy veličinu

$$(7.35) \quad \tau(x, \mathbf{y}, h) = \Delta(x, \mathbf{y}; h) - \Phi(x, \mathbf{y}; h).$$

V algoritmu 7.2 zavedené metody se nazývají *jednokrokové*. To evidentně proto, že k získání hodnoty přiblížení v bodě x_k se vyžaduje znalost hodnoty v příslušném bodě předchozím x_{k-1} .

Jednokroková metoda se ukazuje jako vhodná, jestliže

$$(7.36) \quad \lim_{h \rightarrow 0} \tau(x, \mathbf{y}; h) = 0,$$

což, díky tomu, že

$$(7.37) \quad \lim_{h \rightarrow 0} \Delta(x, \mathbf{y}, h) = \mathbf{f}(x, \mathbf{y}),$$

vede k podmínce

$$(7.38) \quad \Phi(x, \mathbf{y}, 0) = \mathbf{f}(x, \mathbf{y}).$$

Říkáme, že jednokroková metoda je *konzistentní*, jestliže platí relace (7.38).

Poznámka. Eulerova polygonální metoda je zřejmě konzistentní.

Příklad 7.2 Předpokládejme, že $N = 1$ a f má vyšší hladkost než je pouhá spojitost, řekněme, že $\frac{\partial y}{\partial x}$ jsou spojité na $[a, b]$. Navíc nechť

$$z(x+h) = z(x) + hz'(x) + \frac{1}{2}h^2 z''(x) + \dots + \frac{h^p}{p!} z^{(p)}(x + \theta h), \quad 0 < \theta < 1.$$

Protože

$$\begin{aligned} z''(x) &= \frac{d}{dt} f(t, z(t))|_{t=x} \\ &= f_x(t, z(t))|_{t=x} + f_y(t, z(t))|_{t=x} \\ &= f_x(x, y) + f_y(x, y)f(x, y), \end{aligned}$$

a

$$z'''(x) = f_{xx}(x, y) + 2f_{xy}(x, y)[f(x, y)]^2 + f_y(x, y)z''(x),$$

zjistíme, že

$$(7.39) \quad \left\{ \begin{array}{l} \Delta(x, y; h) = z'(x) + \frac{h}{2}z''(x) + \dots + \frac{h^p}{p!} z^{(p)}(x + \theta h) \\ = f(x, y) + \frac{h}{2}[f_x(x, y) + f_y(x, y)f(x, y)] \\ + \dots \end{array} \right.$$

Na př. pro Eulerovu metodu

$$\begin{aligned} \tau(x, y; h) &= \Delta(x, y; h) - \Phi(x, y; h) \\ &= \frac{h}{2}[f_x(x, y) + f_y(x, y)f(x, y)] + O(h). \end{aligned}$$

Obecně, říkáme, že *jednokroková metoda je řádu p*, jestliže platí

$$(7.40) \quad \tau(x, \mathbf{y}; h) = O(h^p)$$

pro $x \in [a, b]$, $\mathbf{y} \in R^N$ a pro všechny $\mathbf{f} \in [C^{(p)}([a, b])]^N$.

Z předchozího příkladu odvodíme poměrně obecný způsob konstrukce metod vyšších řádů.

Na př.

$$(7.41) \quad \Phi(x, y; h) = f(x, y) + \frac{h}{2} [f_x(x, y) + f_y(x, y)f(x, y)]$$

definuje metodu 2. řádu.

Na tomto místě je nutné uvézt na pravou míru skutečnost, že metody vysokých řádů (> 2) jsou v praxi víceméně bezcenné, neúčinné, složité. To je způsobeno tím, že výpočet vyšších derivací je problém složitý v tom smyslu, že je *špatně podmíněný* (*non well posed*). Připomeňme, že to značí, že hodnota výsledku výpočet není spojitě závislá na datech.

Uvedeme některé příklady metod vyšších řádů.

Položme

$$(7.42) \quad \Phi(x, y; h) = a_1 f(x, y) + a_2 f(x + p_1 h, y + p_2 h y)),$$

kde a_1, a_2, p_1, p_2 jsou vhodné konstanty, které jest určit tak, aby platily vztahy určující řád metody.

Pro Φ z (7.42) platí

$$\Phi(x, y; h) = (a_1 + a_2)f(x, y) + a_1 h [p_1 f_x(x, y) + p_2 f_y(x, y)f(x, y)] + O(h^p),$$

odkud vyplývá, splnění podmínek

$$(7.43) \quad a_1 + a_2 = 1, \quad a_1 p_1 = 1/2, \quad a_2 p_2 = 1/2,$$

je postačující aby metoda byla 2. řádu.

Speciální volbou

$$a_1 = a_2 = \frac{1}{2}, \quad p_1 = p_2 = 1$$

obdržíme *metodu Heunova*

$$(7.44) \quad \Phi(x, y; h) = \frac{1}{2} [f(x, y) + f(x + h, y + h f(x, y))].$$

Je patrné, že tato metoda vyžaduje výpočet dvou funkčních hodnot na jeden krok metody a to je cenou, již je nutno zaplatit za vyšší řád.

Jinou metodu obdržíme volbou

$$a_1 = 0, \quad a_2 = 1, \quad p_1 = p_2 = 1/2.$$

Ta je dána *schématem Collatzovým*

$$(7.45) \quad \Phi(x, y; h) = f\left(x + \frac{h}{2}, y + \frac{h}{2}f(x, y)\right).$$

Toto schéma je opět řádu 2. a vyžaduje výpočet dvou funkčních hodnot na jeden krok metody.

Další metody vyšší řádů jsou spojeny se jmény *Runge* a *Kutta*.

Hledejme jenokrokovou metodu pomocí schématu

$$(7.46) \quad \Phi(x, y; h) = \frac{1}{6} [k_1 + 2k_2 + 2k_3 + k_4],$$

kde

$$\begin{aligned} k_1 &= f(x, y), \\ k_2 &= f\left(x + \frac{h}{2}, y + \frac{h}{2}k_1\right), \\ k_3 &= f\left(x + \frac{h}{2}, y + \frac{1}{2}hk_2\right), \\ k_4 &= f(x + h, y + hk_3). \end{aligned}$$

Dá se ukázati, že

$$\Phi(x, y; h) - \Delta(x, y; h) = O(h^4).$$

Znamená to, že tato, dá se říci, základní metoda *Runge - Kuttovy* třídy, je řádu 4.

Aplikujme některé jednokrokové metody na případ úlohy

$$y'(x) = f(x), \quad y(x_0) = y_0.$$

Přesné řešení této úlohy je dáno výrazem

$$y(x) = \int_{x_0}^x f(t)dt + y_0.$$

Snadno nahlédneme, že v problematice numerické kvadratury, což je vlastně námi vyšetřovaná úloha, Eulerova metoda vede na známou metodu obdélníkovou zatímco Heunova metoda na metodu lichoběžníkovou.

Runge Kuttova metoda poskytuje metodu approximace integrálu, jež je známá pod názvem metoda Simpsonova, kterážto metoda integruje polynomy až do řádu 3. stupně včetně.

Otzázkou konvergence jednokrokových metod zodpovídá následující věta. Soustředíme se na skalární případ s tím, že odpovídající věty mají takřka shodné formulace i důkazy.

Věta 7.4 Nechť $f \in C([a, b])$, $x_0 \in [a, b]$ a $y_0 \in R^1$. Budě $y = y(x)$ řešení počáteční úlohy

$$y' = f(x, y), \quad y(x_0) = y_0.$$

Budě Φ funkce spojitá na

$$\mathcal{G} = \{(x, y, h) : a \leq x \leq b, |y - y(x)| \leq \gamma, 0 \leq |h| \leq h_0, h_0 > 0, \gamma > 0\}$$

a nechť existují konstanty κ_1, κ_2 tak, že

$$|\Phi(x, y_1; h) - \Phi(x, y_2; h)| \leq \kappa_1 |y_1 - y_2|$$

a

$$|\tau(x, y(x); h)| = |\Delta(x, y(x); h) - \Phi(x, y(x); h)| \leq \kappa_2 |h|^p, \quad p > 0,$$

platí pro všechna $x \in [a, b]$, $|h| \leq h_0$.

Potom existuje \hat{h} , $0 < \hat{h} \leq h_0$ tak, že pro globální diskretizační chybu e , kde

$$e = e(x) = \eta(x; h) - y(x)$$

platí odhad

$$|e(x; h_k)| \leq |h_k|^p \kappa_2 \frac{\exp\{\kappa_1|x - x_0|\} - 1}{\kappa_1}$$

pro všechna $x \in [a, b]$ a všechna

$$h_k = \frac{x - x_0}{k}, \quad k = 1, 2, \dots,$$

pro něž $|h_k| \leq \hat{h}$. Pro $\gamma = \infty$ je $\hat{h} = h_0$.

Obecnější metody pro řešení počátečních úloh lze shrnout pod obecné schéma

$$(7.47) \quad \eta_{k+r} + a_{r-1}\eta_{k+r-1} + \dots + a_0\eta_k = hF(x_k, \eta_{k+r}, \eta_{k+r-1}, \dots, \eta_k; h; f).$$

Zde $F = F(x, u_1, \dots, u_{r+1}; h; f)$ je vhodná funkce a a_0, \dots, a_{r-1} jsou parametry metody. Je zřejmé, že k nastartování metod typu (7.47) je nutné znát r počátečních hodnot $\eta_0, \dots, \eta_{r-1}$.

Teorie vícekrokových a obecněji, metod typu (7.47), byť přirozeně složitější než ta, již jsme naznačili pro metody jednokrokové, je do šíře i do hloubky propracována a dovedena do realizace v podobě softwarových balíků.

7.4 Okrajové úlohy

Elementy matematické analýzy slabých řešení.

V tomto článku se omezíme na lineární okrajové úlohy s poznámkou, že nelineární úlohy lze v principu vyšetřovat a řešit podobným způsobem jako v problematice lineární; rozdíl spočívá pouze v tom, že výsledné diskrétní soustavy jsou nelineární.

Významným nástrojem vyšetřování okrajových úloh je pojem *slabé (variační) formulace* okrajových úloh.

Výklad budeme provádět na příkladě typické, byť akademické, úlohy

$$(7.48) \quad -y'' = f \text{ v } \Omega = [a, b],$$

$$(7.49) \quad \mu_1 y'(a) + \mu_0 y(a) = \alpha \text{ a } \nu_1 y'(b) + \nu_0 y(b) = \beta.$$

O funkci f předpokládáme, že patří do $L^2(a, b)$; prakticky však vždy je f po částech spojitá, t. j. f má v $[a, b]$ konečný počet bodů nespojnosti.

Formálně vyžaduje t. zv. *variační přístup* pojem zobecněné derivace.

Funkce $g \in L^2(a, b)$, t. j. g , pro níž

$$\int_a^b |g(t)|^2 dt < +\infty,$$

má v $[a, b]$ zobecněnou derivaci, pakliže platí, že existuje $h \in L^2(a, b)$ taková, že

$$(7.50) \quad \int_a^b g(t)v'(t)dt = - \int_a^b h(t)v(t)dt + g(b)v(b) - g(a)v(a).$$

Pak říkáme, že h je *zobecněnou derivací (1. řádu)* funkce g .

Dá se ukázati, že v našem případě, kdy vyšetřujeme oblast $\Omega \subset \mathbb{R}^1$, funkce g má v $[a, b]$ zobecněnou derivaci 1. řádu právě když (*klasická*) derivace g' existuje v $[a, b]$ s vyjímkou množiny $E \subset [a, b]$ mající Lebesgueovu míru $m(E) = 0$. Je známo, že spojitá funkce g má v $[a, b]$ zobecněnou derivaci h právě když g je *absolutně spojitá* v $[a, b]$. Množina funkcí mající zobecněné derivace 1. řádu integrovatelné se čtvercem na Ω se značí symbolem $W_2^{(1)}(\Omega)$ a nazývá se *Sobolevovým prostorem*. V tomto prostoru lze zavézti skalární součin

$$(7.51) \quad (u, v)_{W_2^1(\Omega)} = (u', v')_2 + (u, v)_2,$$

kde pro $u, v \in L^2(\Omega)$,

$$(u, v)_2 = \int_{\Omega} u(t)v(t)dt$$

a u' a v' značí zobecněné derivace funkcí u a v . Předpokládá se, že $u', v' \in L^2(\Omega)$.

Předpokládejme, že $w \in L^2(\Omega)$ je funkce splňující okrajové podmínky

$$(7.52) \quad \mu_0 w(a) + \mu_1 w'(a) = \alpha \text{ a } \nu_0 w(b) + \nu_1 w'(b) = \beta.$$

Dá se ukázati, že existuje hustá ve smyslu normy prostoru $L^2(\Omega)$ množina tvořená dostatečně hladkými funkcemi splňujícími (7.52); na př. množina všech polynomů splňujících (7.52).

Položme v (7.48)

$$y(s) = u(s) - w(s).$$

Jest tedy,

$$-y'' = w'' + f \text{ v } \Omega$$

a

$$\mu_1 y'(a) + \mu_0 y(a) = \mu_0 (u(a) - w(a)) + \mu_1 (u'(a) - w'(a)) = 0,$$

$$\nu_1 y'(b) + \nu_0 y(b) = \nu_1 (u'(b) - w'(b)) + \nu_0 (u(b) - w(b)) = 0.$$

Známe-li y , známe automaticky i hledané u .

K určení y máme tedy vztahy

$$(7.53) \quad -y'' = f + w'' \text{ v } \Omega$$

$$(7.54) \quad \mu_0 y(a) + \mu_1 y'(a) = \nu_0 y(b) + \nu_1 y'(b) = 0.$$

Protože w'' je známá funkce, lze (7.53) psát ve tvaru

$$(7.55) \quad -y'' = F \text{ v } \Omega,$$

při čemž F je daná funkce.

Všimněme si, že ve formulaci (7.48) - (7.49), neboli (7.54) - (7.55), vyžadujeme splnění těchto vztahů identicky v $[a, b]$. To však vyžaduje znalost F ve *všech bodech* oblasti Ω a obecněji, znalost všech dat v Ω . Toto je jednak ne příliš praktické a jednak dokonce ne vždy splnitelné. Ve vícerozměrných oblastech zadání okrajových podmínek na některých částech hranice může být nedefinovatelné. Na př. nechť $\Omega = [0, 1] \times [0, 1]$ a nechť $u(s, 0) = 0$ pro $s \in (0, 1)$, zatímco $u(0, y) = 1$ pro $y \in (0, 1)$. Otázkou je, jakou hodnotu může nabývat u v bodě $(1, 0)$, t. j., čemu je rovno $u(1, 0)$. Uvědomme si, že u má být spojitá na Ω .

V takových problémech se ukáže jako velice vhodná *slabá (variační)* formulace problému.

Označme symbolem V podprostor Sobolevova prostoru $W_2^1(\Omega)$ splňujících okrajové podmínky (7.54).

Vynásobme rovnici (7.55) funkcií $v \in V$ a integrujme podél Ω . Obdržíme vztahy

$$-\int_a^b y''(t)v(t)dt = \int_a^b F(t)v(t)dt, \quad v \in V,$$

takže požadované řešení lze hledat pomocí vztahů

$$(7.56) \quad \int_a^b y'(t)v'(t)dt = \int_a^b F(t)v(t)dt, \quad v \in V,$$

což plyne z (7.56) protože v splňuje tytéž okrajové podmínky (7.54) jako hledané řešení y (viz Cvičení 7.3).

Příklad 7.3 Buď \mathcal{L} množina všech funkcí y majících spojité derivace 2. řádu na $[a, b]$ a splňujících okrajové podmínky (7.54). Ukažte, že rovnosti

$$\int_a^b y''(t)v(t)dt = \int_a^b u(t)v''(t)dt$$

platí pro všechny prvky $y \in \mathcal{L}$, přičemž v je funkce mající na $[a, b]$ spojitou derivaci 2. řádu, právě když $v \in \mathcal{L}$.

Pro jednoduchost položme

$$a = 0 \text{ a } b = 1$$

a

$$\mu_1 = \nu_1 = 0.$$

V tomto případě se okrajové podmínky (7.54) redukují na

$$(7.57) \quad y(0) = y(1) = 0$$

a tedy též

$$(7.58) \quad v(0) = v(1) = 0.$$

Formulujme nyní naši okrajovou úlohu *variačně* neboli *slabě*, či přesněji řečeno, *ve slabém smyslu*.

Úloha (U). *Nalézt $y \in \tilde{W}_2^{(1)}(0, 1)$ tak, aby vztahy*

$$(7.59) \quad \int_0^1 y'(t)v'(t)dt = \int_0^1 F(t)v(t)dt$$

platily pro všechny prvky $v \in \tilde{W}_2^{(1)}(0, 1)$, kde $\tilde{W}_2^{(1)}(0, 1)$ značí uzávěr prostoru $C_0^\infty([0, 1])$ funkcí s kompaktním nosičem na $[0, 1]$ majících na $[0, 1]$ spojité derivace všech rádu. Uzávěr se míní ve smyslu metriky prostoru $W_2^{(1)}(0, 1)$.

Dá se ukázat, že úloha **U** má právě jedno řešení. Má-li okrajová úloha

$$y'' = F \text{ v } [0, 1]$$

s podmínkami (7.58) klasické řešení, pak toto klasické řešení je totožné s řešením úlohy **U**. Tato tvrzení jsou založena na platnosti slavného *Lax - Milgramova lemmatu*

Lemma 7.1 (Lax - Milgram) *Budě V Hilbertův prostor. Budě $B = B(u, v)$ bilineární forma, jež je spojitá a V -eliptická, t.j. existují konstanty $c > 0$ a $\tilde{c} > 0$ tak, že*

$$(7.60) \quad |B(u, v)| \leq c\|u\|\|v\|, \quad u, v \in V$$

a

$$(7.61) \quad |B(u, u)| \geq \tilde{c}\|u\|^2, \quad u \in V.$$

Potom pro každý prvek $f \in V$ existuje právě jeden prvek $u^* \in V$ takový, že

$$(7.62) \quad B(u^*, v) = (F, v), \quad \forall v \in V.$$

Platí přitom

$$(7.63) \quad \|u^*\| \leq \kappa\|F\|,$$

kde κ nezávisí na F .

V našem případě $V = \tilde{W}_2^{(1)}(0, 1)$ a

$$B(u, v) = \int_0^1 u'(t)v'(t)dt.$$

Protože

$$\|u\|^2 = \int_0^1 [\|u\|^2 + \|u'\|^2] dt,$$

platnost požadavku (7.61) je důsledkem *Friedrichsovy nerovnosti*

$$\int_\Omega u'(t)^2 dt \geq \tau \int_\Omega u^2 dt, \quad \forall u \in \tilde{W}_2^{(1)}(0, 1),$$

kde τ je kladná konstanta nezávislá na u ani u' .

Požadavek (7.60) lemmatu 7.1 je automaticky splněn díky tomu, že skalární součin v $\tilde{W}_2^{(1)}(0, 1)$ je dán výrazem (7.51).

Na základě lemmatu 7.1 tedy platí

Důsledek 7.1 Existuje právě jedno slabé řešení úlohy (U) .

Diskretizace.

Budě $V_h \subset V = \tilde{W}_2^1(0, 1)$, kde

$$\dim V_h = N < +\infty.$$

Nechť

$$\{\phi_1, \dots, \phi_N\}$$

tvoří basi prostoru V_h , tedy

$$V_h = \text{Span}\{\phi_1, \dots, \phi_N\}.$$

Tudíž,

$$(7.64) \quad u_h = \sum_{k=1}^N \nu_k \phi_k,$$

pro každý prvek $u_h \in V_h$.

Diskretizací úlohy (U) se rozumí úloha

(\mathbf{U}_h) Nalézt $u_h^* \in V_h$ takové, že platí

$$(7.65) \quad B(u_h^*, v_h) = (f, v_h), \quad \forall v_h \in V_h.$$

Věta 7.5 Splňuje-li forma $B = B(u, v)$ požadavky Laxova - Milgramova lemmatu, pak lze nalézt $h_0 > 0$ tak, že pro každé $0 < h \leq h_0$, existuje právě jedno řešení u_h^* úlohy (U_h) .

Navíc platí odhad

$$\|u^* - u_h^*\|_V \leq \kappa d_h(u_h^*) \|f\|,$$

kde κ nezávisí ani na h ani na f a

$$d_h(v) = \inf \{\|v - v_h\|_V : v_h \in V_h\}.$$

Příklad 7.4 Pro případ, kdy v problému (U_h)

$$h = \frac{1}{N}, \quad \text{a } x_k = kh, \quad k = 0, 1, \dots, N,$$

a basi approximačního podprostoru tvoří po částech lineární funkce definované formulami

$$\phi_k(x) = \begin{cases} 0 & \text{pro } 0 \leq x \leq x_{k-1}, \\ \frac{1}{h}(x - (k-1)h) & \text{pro } x_{k-1} \leq x_k, \\ 1 - \frac{1}{h}(x - x_k) & \text{pro } x_k \leq x \leq x_{k+1}, \\ 0 & \text{pro } x_{k+1} \leq x \leq 1, \end{cases}$$

lze odvodit, že

$$d_h(u) = h\|u\|_V.$$

Důsledkem je posléze odhad chyby v metodě konečných prvků pro approximaci úlohy (U) ve tvaru

$$\|u^* - u_h^*\|_{\tilde{W}_2^1(0,1)} \leq \kappa h \|f\|_{L^2(0,1)}.$$

Vyplyná odtud speciálně, že

$$\|u^* - u_h^*\|_{L^2(0,1)} \leq \hat{\kappa} h^2 \|f\|_{L^2(0,1)}.$$

Po dosazení vyjádření (7.64) do (7.65) obdržíme vztahy

$$(7.66) \quad \sum_{k=1}^N \nu_k B(\phi_k, \phi_j) = (f, \phi_j), \quad j = 1, \dots, N.$$

Vyplnění podmínek Laxova - Milgramova lemmatu zaručuje, že

$$\det(B(\phi_k, \phi_j)) \neq 0$$

a soustava (7.66) má tudíž právě jedno řešení $\nu^* = (\nu_1^*, \dots, \nu_N^*)^T$.

Matice

$$(B(\phi_k, \phi_j))$$

se díky její interpretaci v mechanice kontinua nazývá *maticí tuhosti* úlohy (U_h).

Je-li forma $B = B(u, v)$ symetrická a jsou-li splněny podmínky Laxova - Milgramova lemmatu, pak matice tuhosti úlohy (U_h) je pozitivně definitní, t.j.

$$(B(u_h, u_h)) \geq \tau(u_h, u_h) \quad \forall u_h, u_h \in V_h, \quad \tau > 0.$$

Výpočet prvků matice tuhosti spočívá ve stanovení veličin

$$\int_{a+(p-1)h}^{a+ph} \phi'_k(x) \phi'_j(x) dx, \quad j, k = 1, \dots, N, \quad p = 1, \dots, N-1,$$

při čemž $\Omega = (a, b)$, $-\infty < a < b < +\infty$.

Pro výše uvedený případ jednorozměrné oblasti s Ω a pro částech lineárních funkcí ϕ_k lze tyto integrály stanovit přesně. Obecně je nutné prvky matice tuhosti určovat přibližně pomocí vhodných kvadraturních vzorců.