

Numerické řešení parciálních diferenciálních rovnic

Petr Knobloch

1 Úvod

V mnoha vědeckých a technických oblastech se setkáváme s matematickými modely sestávajícími z parciálních diferenciálních rovnic (PDR). Obvykle tyto modely nelze vyřešit analyticky a jejich řešení je potřeba aproximovat pomocí numerických metod. Za tím účelem bylo vyvinuto velké množství různých postupů. Cílem přednášky *Numerické řešení parciálních diferenciálních rovnic* je seznámit studenty s metodou konečných diferencí, k jejíž analýze postačují základní vědomosti z přednášek věnovaných matematické analýze. Budou vysvětleny hlavní myšlenky a prostředky analýzy metody konečných diferencí a studovány vlastnosti různých diskretizací pro základní příklady PDR (transportní rovnice, rovnice vedení tepla, rovnice difúze). Získané znalosti budou užitečné i pro získání představy o tom, s jakými problémy je nutno se vypořádávat při vývoji a aplikaci jiných numerických metod pro řešení PDR. Navíc je metoda konečných diferencí stále běžným prostředkem pro aproximaci časových derivací v PDR, zatímco diskretizace vzhledem k prostorovým proměnným je často získávána jinými přístupy (např. metodou konečných prvků). Jelikož však každá numerická metoda vede na soustavu algebraických rovnic, lze některé poznatky z této přednášky využít i při studiu vlastností diskretizací získaných jinými postupy.

1.1 Diferenční kvocienty

Metoda konečných diferencí je založena na tom, že se diferenciální rovnice uvažuje pouze v izolovaných bodech a derivace v těchto bodech se nahradí diferenčními kvocienty. Zmíněné body se získají jako uzly sítě pokrývající výpočetní oblast, a proto jiný název této numerické metody je metoda sítě.

Mějme funkci $v : \mathbb{R} \rightarrow \mathbb{R}$. Pak její první derivace v bodě $x \in \mathbb{R}$ je definována vztahem

$$(1) \quad v'(x) = \lim_{h \rightarrow 0} \frac{v(x+h) - v(x)}{h}.$$

Zvolíme-li $h > 0$, pak můžeme derivaci $v'(x)$ aproximovat vztahem

$$(2) \quad v'(x) \approx \frac{v(x+h) - v(x)}{h}.$$

Je-li funkce v dostatečně hladká, lze očekávat, že diferenční kvocient v (2) aproximuje $v'(x)$ tím lépe, čím je h menší. Zanedlouho uvidíme, že toto očekávání je správné.

Místo (1) můžeme též psát

$$(3) \quad v'(x) = \lim_{h \rightarrow 0} \frac{v(x) - v(x-h)}{h},$$

a tedy i (zprůměrujeme-li (1) a (3))

$$v'(x) = \lim_{h \rightarrow 0} \frac{v(x+h) - v(x-h)}{2h}.$$

Tomu odpovídají aproximace

$$v'(x) \approx \frac{v(x) - v(x-h)}{h}, \quad \text{resp.} \quad v'(x) \approx \frac{v(x+h) - v(x-h)}{2h},$$

popř.

$$v'(x) \approx \frac{v(x + \frac{h}{2}) - v(x - \frac{h}{2})}{h}.$$

Abychom zápis uvažovaných aproximací zpřehlednili, zavádíme pro diference následující označení (vždy předpokládáme, že $h > 0$):

dopředná diference:	$\Delta_+ v(x) = v(x+h) - v(x),$
zpětná diference:	$\Delta_- v(x) = v(x) - v(x-h),$
centrální diference s dvojnásobnou délkou intervalu:	$\Delta_0 v(x) = \frac{1}{2} [v(x+h) - v(x-h)],$
centrální diference:	$\delta v(x) = v(x + \frac{h}{2}) - v(x - \frac{h}{2}).$

Pak uvedené aproximace můžeme psát ve tvaru

$$v' \approx \frac{\Delta_+ v}{h}, \quad \text{resp.} \quad v' \approx \frac{\Delta_- v}{h}, \quad \text{resp.} \quad v' \approx \frac{\Delta_0 v}{h}, \quad \text{resp.} \quad v' \approx \frac{\delta v}{h}.$$

Je-li $v \in C^2(\mathbb{R})$, zjišťujeme pomocí Taylorova vzorce, že všechny tyto aproximace vedou k chybě $O(h)$ (tj. velikost chyby lze v každém bodě x omezit hodnotou Ch s konstantou C nezávislou na h). Je-li $v \in C^3(\mathbb{R})$, dostáváme navíc

$$\frac{\Delta_0 v}{h} = v' + O(h^2), \quad \frac{\delta v}{h} = v' + O(h^2).$$

Chceme-li aproximovat $v''(x)$, pak můžeme postupovat následovně:

$$v''(x) \approx \frac{\delta v'(x)}{h} \approx \frac{\delta^2 v(x)}{h^2}.$$

Platí

$$\delta^2 v(x) = v(x+h) - 2v(x) + v(x-h)$$

a tento výraz se nazývá centrální diference druhého řádu. Snadno ověříme, že

$$\delta^2 = \Delta_+ \Delta_- = \Delta_- \Delta_+.$$

Použitím Taylorova vzorce dostaneme pro $v \in C^4(\mathbb{R})$

$$\frac{\delta^2 v}{h^2} = v'' + O(h^2).$$

Podobně pro $v \in C^6(\mathbb{R})$ je

$$\frac{\delta^4 v}{h^4} = v^{(4)} + O(h^2).$$

Přitom

$$\delta^4 v(x) = v(x + 2h) - 4v(x + h) + 6v(x) - 4v(x - h) + v(x - 2h).$$

1.2 Příklady diskretizací Cauchyovy úlohy pro rovnici konvekce–difúze v jedné prostorové dimenzi

Ukažme si nyní, jak pomocí metody konečných diferencí lze získat diskretizaci Cauchyovy úlohy pro danou parciální diferenciální rovnici. Jako příklad budeme uvažovat rovnici konvekce–difúze v jedné prostorové dimenzi. Hledáme tedy funkci $u = u(x, t)$ závisující na prostorové proměnné $x \in \mathbb{R}$ a časové proměnné $t \in \mathbb{R}_0^+$ splňující parciální diferenciální rovnici

$$(4) \quad u_t = b u_{xx} - a u_x \quad \text{v } \mathbb{R} \times \mathbb{R}^+$$

s počáteční podmínkou

$$(5) \quad u(x, 0) = u^0(x) \quad \forall x \in \mathbb{R},$$

kde u^0 je daná funkce. Parametr b představuje koeficient difúze, a proto předpokládáme, že $b > 0$. Parametr a určuje rychlost šíření veličiny u vlivem konvekce (unášení). Pro jednoduchost předpokládáme, že a a b jsou konstanty.

Jak již bylo zmíněno výše, základem diskretizace metodou konečných diferencí je síť pokrývající výpočetní oblast, v jejíž uzlech aproximujeme parciální diferenciální rovnici rovnicemi algebraickými. Pro jednoduchost budeme nyní uvažovat rovnoměrnou síť s konstantním prostorovým krokem $h > 0$ a konstantním časovým krokem $\tau > 0$. Definujeme prostorové uzly $x_j = jh$, $j \in \mathbb{Z}$, a časové hladiny $t_n = n\tau$, $n \in \mathbb{N}_0$. Rovnoběžky se souřadnými osami protínající tyto osy v bodech x_j a t_n jsou takzvané síťové přímky, které vytvářejí síť pokrývající výpočetní oblast $\mathbb{R} \times \mathbb{R}_0^+$. Průsečíky síťových přímek (x_j, t_n) , kde $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$, nazveme uzly sítě. Řešení u Cauchyovy úlohy (4), (5) aproximujeme v uzlech (x_j, t_n) hodnotami U_j^n . Síťová funkce $\{U_j^n\}$ je pak přibližným řešením uvažované Cauchyovy úlohy. Pro přehlednost zavádíme též označení $u_j^n := u(x_j, t_n)$.

Přibližné řešení $\{U_j^n\}$ získáme tak, že parciální diferenciální rovnici v uzlech sítě nahradíme diferenčním schématem. Uvažujme libovolný uzel (x_j, t_n) s $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$. Pak rovnici (4) můžeme s využitím vztahů z předešlého oddílu aproximovat například

následovně:

$$\begin{aligned}
0 &= (u_t - b u_{xx} + a u_x)(x_j, t_n) \\
&\approx \frac{u(x_j, t_n + \tau) - u(x_j, t_n)}{\tau} - b \frac{u(x_j + h, t_n) - 2u(x_j, t_n) + u(x_j - h, t_n)}{h^2} \\
&\quad + a \frac{u(x_j + h, t_n) - u(x_j - h, t_n)}{2h} \\
&= \frac{u_j^{n+1} - u_j^n}{\tau} - b \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2} + a \frac{u_{j+1}^n - u_{j-1}^n}{2h} \\
&\approx \frac{U_j^{n+1} - U_j^n}{\tau} - b \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h}.
\end{aligned}$$

Použijeme-li symboly pro diference, můžeme uvedené aproximace zapsat v přehlednějším tvaru:

$$\begin{aligned}
0 = (u_t - b u_{xx} + a u_x)(x_j, t_n) &\approx \frac{\Delta_{+t} u_j^n}{\tau} - b \frac{\delta_x^2 u_j^n}{h^2} + a \frac{\Delta_{0x} u_j^n}{h} \\
&\approx \frac{\Delta_{+t} U_j^n}{\tau} - b \frac{\delta_x^2 U_j^n}{h^2} + a \frac{\Delta_{0x} U_j^n}{h}.
\end{aligned}$$

Nyní u symbolů pro diference vyznačujeme indexem x či t , vzhledem ke které proměnné je diference definována. Např.

$$\Delta_{+t} u(x, t) = u(x, t + \tau) - u(x, t), \quad \Delta_{+x} u(x, t) = u(x + h, t) - u(x, t).$$

Rovnici (4) tedy aproximujeme numerickým schématem

$$(6) \quad \frac{U_j^{n+1} - U_j^n}{\tau} = b \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} - a \frac{U_{j+1}^n - U_{j-1}^n}{2h} \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

které je nutno doplnit počáteční podmínkou

$$(7) \quad U_j^0 = u^0(x_j) \quad \forall j \in \mathbb{Z}.$$

Jedná se o příklad tzv. explicitního schématu, neboť hodnotu přibližného řešení v libovolném uzlu (x_j, t_n) s $j \in \mathbb{Z}$ a $n \in \mathbb{N}$ lze ze schématu vyjádřit pomocí hodnot přibližného řešení na předcházející časové vrstvě t_{n-1} . Přibližné řešení tedy existuje a je určeno jednoznačně. Označíme-li

$$(8) \quad \mu = \frac{\tau}{h^2} b, \quad \nu = \frac{\tau}{h} a,$$

můžeme schéma (6) psát ve tvaru

$$(9) \quad U_j^{n+1} = U_j^n + \mu (U_{j+1}^n - 2U_j^n + U_{j-1}^n) - \frac{\nu}{2} (U_{j+1}^n - U_{j-1}^n),$$

který je v některých případech vhodnější pro analýzu.

Jelikož existuje řada způsobů, jak pomocí diferenčních kvocientů aproximovat derivace, získáme různými kombinacemi diferenčních kvocientů mnoho rozličných numerických schémat pro aproximaci dané parciální diferenciální rovnice. Jak uvidíme později, takto získaná schémata mohou mít velmi odlišné vlastnosti. Ukažme si ještě dvě schémata, která lze definovat pro rovnici (4). Uvažujeme-li tuto rovnici v uzlu (x_j, t_{n+1}) s $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$, můžeme provést následující aproximace:

$$\begin{aligned} 0 = (u_t - b u_{xx} + a u_x)(x_j, t_{n+1}) &\approx \frac{\Delta_{-t} u_j^{n+1}}{\tau} - b \frac{\delta_x^2 u_j^{n+1}}{h^2} + a \frac{\Delta_{0x} u_j^{n+1}}{h} \\ &\approx \frac{\Delta_{-t} U_j^{n+1}}{\tau} - b \frac{\delta_x^2 U_j^{n+1}}{h^2} + a \frac{\Delta_{0x} U_j^{n+1}}{h}, \end{aligned}$$

což vede ke schématu

$$(10) \quad \frac{U_j^{n+1} - U_j^n}{\tau} = b \frac{U_{j+1}^{n+1} - 2U_j^{n+1} + U_{j-1}^{n+1}}{h^2} - a \frac{U_{j+1}^{n+1} - U_{j-1}^{n+1}}{2h} \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

k němuž opět musíme přidat počáteční podmínku (7). Nyní již v daném uzlu není možné jednoduchým způsobem vyjádřit hodnotu přibližného řešení pomocí hodnot z předchozí časové hladiny a hodnoty přibližného řešení na nové časové hladině je nutno získat vyřešením soustavy lineárních rovnic. Jedná se o příklad tzv. implicitního schématu.

Sečteme-li schémata (6) a (10), získáme po vydělení dvěma schéma

$$\frac{U_j^{n+1} - U_j^n}{\tau} = b \frac{\delta_x^2 U_j^{n+1} + \delta_x^2 U_j^n}{2h^2} - a \frac{U_{j+1}^{n+1} - U_{j-1}^{n+1} + U_{j+1}^n - U_{j-1}^n}{4h} \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0.$$

Toto schéma je opět implicitní a nazývá se schéma Crankovo–Nicolsonové. Jelikož hodnoty přibližného řešení, které toto schéma kombinuje, vykazují symetrické uspořádání nejen vzhledem k x_j , ale i vzhledem k $t_{n+1/2} := t_n + \tau/2$, aproximuje toto schéma parciální diferenciální rovnici (4) při $h, \tau \rightarrow 0$ s menší chybou než předchozí dvě schémata.

1.3 Chyba diskretizace

Chyba diskretizace je chyba, které se dopouštíme, když derivace v diferenciální rovnici nahradíme diferenčními kvocienty. Získáme ji dosazením přesného řešení do numerického schématu ve tvaru odpovídajícím nahrazení derivací diferenčními kvocienty a odečtením pravé strany od levé.

Jako příklad uvažujme schéma (6) pro numerické řešení rovnice (4). Pro definici chyby diskretizace použijeme tvar schématu s diferenčními kvocienty (tj. (6)), nikoli ekvivalentní tvar (9). Chyba diskretizace je pak dána vztahem

$$\varepsilon_j^n = \frac{u_j^{n+1} - u_j^n}{\tau} - b \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2} + a \frac{u_{j+1}^n - u_{j-1}^n}{2h},$$

kde $u_j^n = u(x_j, t_n)$. Tento vztah můžeme též psát ve tvaru

$$\varepsilon_j^n = \frac{\Delta_{+t} u_j^n}{\tau} - b \frac{\delta_x^2 u_j^n}{h^2} + a \frac{\Delta_{0x} u_j^n}{h}.$$

Můžeme také zavést chybu diskretizace $\varepsilon_{h,\tau}(x, t)$, která je definována v každém bodě (x, t) výpočetní oblasti, pro nějž $(x \pm h, t)$ a $(x, t + \tau)$ též leží ve výpočetní oblasti, vztahem

$$\varepsilon_{h,\tau} = \frac{\Delta_{+t} u}{\tau} - b \frac{\delta_x^2 u}{h^2} + a \frac{\Delta_{0x} u}{h}.$$

Zřejmě $\varepsilon_j^n = \varepsilon_{h,\tau}(x_j, t_n)$. Je-li u dvakrát spojitě derivovatelné podle t a čtyřikrát spojitě derivovatelné podle x , pak použitím Taylorova vzorce získáme

$$\varepsilon_{h,\tau}(x, t) = \frac{1}{2} u_{tt}(x, \eta) \tau - \frac{b}{12} u_{xxxx}(\xi, t) h^2 + \frac{a}{6} u_{xxx}(\zeta, t) h^2,$$

kde $\eta \in (t, t + \tau)$ a $\xi, \zeta \in (x - h, x + h)$. Říkáme, že schéma je prvního řádu přesnosti v čase a druhého řádu přesnosti v prostoru.

U každého schématu je základním požadavkem, aby chyba diskretizace konvergovala k nule pro $h, \tau \rightarrow 0$. Pak říkáme, že schéma je konsistentní s řešenou diferenciální rovnicí.

1.4 Chyba aproximace

Chyba aproximace je definována vztahem $e_j^n = U_j^n - u_j^n$, kde U_j^n je přibližné řešení a u_j^n je aproximované přesné řešení v uzlu (x_j, t_n) .

Všimněme si, že chyba aproximace je řešením příslušného schématu, v němž se na pravé straně objeví dodatečný člen $-\varepsilon_j^n$ nebo násobek této hodnoty (v závislosti na použitém tvaru schématu). Např. v případě schématu (6) máme

$$\frac{e_j^{n+1} - e_j^n}{\tau} = b \frac{e_{j+1}^n - 2e_j^n + e_{j-1}^n}{h^2} - a \frac{e_{j+1}^n - e_{j-1}^n}{2h} - \varepsilon_j^n,$$

resp.

$$e_j^{n+1} = e_j^n + \mu (e_{j+1}^n - 2e_j^n + e_{j-1}^n) - \frac{\nu}{2} (e_{j+1}^n - e_{j-1}^n) - \tau \varepsilon_j^n.$$

Později uvidíme, že na základě vztahů tohoto typu lze za určitých podmínek chybu aproximace odhadnout pomocí odhadu chyby diskretizace.

2 Von Neumannova analýza stability numerických schémat pro Cauchyovy úlohy

V teorii parciálních diferenciálních rovnic hraje důležitou roli Fourierova transformace. Pro funkci $u \in L^1(\mathbb{R})$ definujeme Fourierovu transformaci \hat{u} vztahem

$$(11) \quad \hat{u}(\xi) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u(x) e^{-ix\xi} dx, \quad \xi \in \mathbb{R}.$$

Za určitých předpokladů pak platí

$$(12) \quad u(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}(\xi) e^{ix\xi} d\xi$$

a je splněna Parsevalova rovnost

$$(13) \quad \|u\|_{L^2(\mathbb{R})} = \|\widehat{u}\|_{L^2(\mathbb{R})}.$$

Pravá strana vztahu (12) je inverzní Fourierova transformace. Vztah (12) vyjadřuje funkci u jako superpozici vln daných funkcemi $e^{ix\xi}$ s různými amplitudami $\widehat{u}(\xi)$. Funkce \widehat{u} představuje alternativní reprezentaci funkce u a může být komplexní, i když je funkce u reálná.

Podobně jako výše lze postupovat též v diskretním případě. Bud' $l \in \mathbb{R}$ a označme

$$\varphi_j(\xi) = \frac{1}{\sqrt{l}} e^{-i(2\pi j/l)\xi}, \quad \xi \in \mathbb{R}, \quad j \in \mathbb{Z}.$$

Pak pro libovolné reálné číslo a tvoří množina $\{\varphi_j\}_{j \in \mathbb{Z}}$ úplný ortonormální systém v prostoru $L^2(a, a+l)$. Je-li dána posloupnost $U = \{U_j\}_{j \in \mathbb{Z}}$ splňující $\sum_{j \in \mathbb{Z}} |U_j|^2 < \infty$, pak řada $\sum_{j \in \mathbb{Z}} U_j \varphi_j$ konverguje v prostoru $L^2(a, a+l)$ k funkci \widetilde{U} a platí $U_j = \int_a^{a+l} \widetilde{U} \overline{\varphi_j} d\xi$. Navíc je splněna Parsevalova rovnost $\sum_{j \in \mathbb{Z}} |U_j|^2 = \|\widetilde{U}\|_{L^2(a, a+l)}^2$. Předpokládejme nyní, že hodnoty U_j představují hodnoty síťové funkce v uzlech $x_j = jh$. Zvolme $a = -l/2$, $l = 2\pi/h$ a definujme funkci $\widehat{U} \in L^2(-\pi/h, \pi/h)$ vztahem

$$(14) \quad \widehat{U}(\xi) = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} h U_j e^{-ix_j \xi}$$

(tj. $\widehat{U} = \sqrt{h} \widetilde{U}$). Pak

$$(15) \quad U_j = \frac{1}{\sqrt{2\pi}} \int_{-\pi/h}^{\pi/h} \widehat{U}(\xi) e^{ix_j \xi} d\xi \quad \forall j \in \mathbb{Z}.$$

Vztahy (14) a (15) jsou diskretními analogiemi vztahů (11) a (12). Vztah (15) opět vyjadřuje U jako superpozici vln. Označíme-li

$$\|U\|_2 = \sqrt{\sum_{j=-\infty}^{\infty} h |U_j|^2}$$

diskretní analogii normy v prostoru $L^2(\mathbb{R})$, platí podobně jako v (13) Parsevalova rovnost

$$(16) \quad \|U\|_2 = \|\widehat{U}\|_{L^2(-\pi/h, \pi/h)}.$$

Uvažujme Cauchyovu úlohu najít funkci $u = u(x, t)$ definovanou pro $x \in \mathbb{R}$ a $t \geq 0$ a splňující

$$(17) \quad u_t + L u = 0 \quad \text{v } \mathbb{R} \times \mathbb{R}^+, \quad u(x, 0) = u^0(x) \quad \text{pro } x \in \mathbb{R},$$

kde u^0 je zadaná počáteční podmínka a L je lineární diferenciální operátor s konstantními koeficienty obsahující pouze derivace podle x . Obecné jednokrokové schéma pro úlohu (17) na stejnoměrné síti s uzly $x_j = jh$, $j \in \mathbb{Z}$, má tvar

$$(18) \quad \sum_{s=-M}^M \alpha_s U_{j+s}^{n+1} = \sum_{s=-M}^M \beta_s U_{j+s}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

$$(19) \quad U_j^0 = u^0(x_j) \quad \forall j \in \mathbb{Z},$$

kde přirozené číslo M závisí na způsobu aproximace derivací podle x v operátoru L . Předpokládáme, že přibližné řešení je uvedeným schématem určeno jednoznačně. Zapišeme-li $U^n = \{U_j^n\}_{j \in \mathbb{Z}}$ ve tvaru (15) pomocí Fourierovy transformace $\widehat{U}^n = \widehat{U}^n(\xi)$ definované vztahem (14), získáme dosazením do (18)

$$\int_{-\pi/h}^{\pi/h} e^{ix_j \xi} \left(\widehat{U}^{n+1}(\xi) \sum_{s=-M}^M \alpha_s e^{ish\xi} - \widehat{U}^n(\xi) \sum_{s=-M}^M \beta_s e^{ish\xi} \right) d\xi = 0 \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0.$$

Jelikož funkce $\{e^{ix_j \xi}\}_{j \in \mathbb{Z}}$ tvoří úplný ortogonální systém v prostoru $L^2(-\pi/h, \pi/h)$, platí

$$\widehat{U}^{n+1}(\xi) \sum_{s=-M}^M \alpha_s e^{ish\xi} = \widehat{U}^n(\xi) \sum_{s=-M}^M \beta_s e^{ish\xi} \quad \forall n \in \mathbb{N}_0.$$

Všimněme si, že tento vztah lze formálně získat tak, že diskrétní řešení U_j^n nahradíme v (18) výrazem $\widehat{U}^n(\xi) e^{ix_j \xi}$. Díky předpokladu o jednoznačné řešitelnosti schématu je součet na levé straně pro libovolné ξ nenulový, a můžeme proto definovat funkci

$$\lambda(\xi) = \frac{\sum_{s=-M}^M \beta_s e^{ish\xi}}{\sum_{s=-M}^M \alpha_s e^{ish\xi}}.$$

Pak obdržíme

$$(20) \quad \widehat{U}^{n+1}(\xi) = \lambda(\xi) \widehat{U}^n(\xi) \quad \forall n \in \mathbb{N}_0.$$

Tento vztah ukazuje, že provedení jednoho časového kroku schématu (18) je ekvivalentní přenásobení Fourierovy transformace diskrétního řešení amplifikačním faktorem $\lambda(\xi)$. Velikost amplifikačního faktoru $|\lambda(\xi)|$ představuje zesílení amplitudy $\widehat{U}^n(\xi)$ libovolné frekvence ξ při provedení jednoho časového kroku. Ze vztahu (20) získáme

$$(21) \quad \widehat{U}^n(\xi) = \lambda(\xi)^n \widehat{U}^0(\xi) \quad \forall n \in \mathbb{N}_0.$$

Vidíme, že všechny informace o schématu jsou obsaženy v amplifikačním faktoru. Mimo jiné lze z amplifikačního faktoru snadno získat informace o stabilitě a přesnosti příslušného

schématu. Fourierova transformace proto představuje standardní metodu pro studium vlastností diferenčních schémat.

Diferenční schémata jsou často pouze podmíněně stabilní, což znamená, že jsou stabilní, jen pokud prostorový krok h a časový krok τ splňují jistou podmínku. Množinu $\Lambda \subset \mathbb{R}^+ \times \mathbb{R}^+$ takovou, že pro libovolnou dvojici $(h, \tau) \in \Lambda$ je příslušná podmínka stability splněna, nazveme oblastí stability diferenčního schématu. Vždy budeme předpokládat, že množina Λ je omezená a že dvojice $(0, 0)$ je jejím hromadným bodem. Stabilitu schématu (18) lze definovat následujícím způsobem.

Definice 1 Diferenční schéma (18) je stabilní v oblasti stability Λ , pokud pro každý pevný čas $T > 0$ existuje konstanta C_T taková, že pro libovolnou počáteční podmínku U^0 platí

$$\|U^n\|_2 \leq C_T \|U^0\|_2 \quad \forall n \in \mathbb{N}, (h, \tau) \in \Lambda, n\tau \leq T.$$

Věta 1 Diferenční schéma (18) je stabilní v oblasti stability Λ právě tehdy, když existuje konstanta K nezávislá na ξ , h a τ taková, že

$$(22) \quad |\lambda(\xi)| \leq 1 + K\tau \quad \forall \xi \in \mathbb{R}, (h, \tau) \in \Lambda.$$

Je-li funkce $\lambda(\xi/h)$ na Λ nezávislá na h a τ , pak lze podmínku (22) nahradit podmínkou

$$(23) \quad |\lambda(\xi)| \leq 1 \quad \forall \xi \in \mathbb{R}, (h, \tau) \in \Lambda.$$

Důkaz. Z Parsevalovy rovnosti (16) a vztahu (21) plyne

$$(24) \quad \|U^n\|_2 = \|\lambda^n \widehat{U}^0\|_{L^2(-\pi/h, \pi/h)}.$$

Poznamenejme, že λ^n zde značí n -tou mocninu λ . Platí-li (22), je pro $n \leq T/\tau$

$$\|U^n\|_2 \leq (1 + K\tau)^n \|\widehat{U}^0\|_{L^2(-\pi/h, \pi/h)} \leq (1 + K\tau)^{T/\tau} \|U^0\|_2 \leq e^{KT} \|U^0\|_2,$$

tj. schéma (18) je stabilní v Λ . Předpokládejme nyní, že (22) neplatí pro žádnou konstantu K . Zvolme libovolné číslo $C > 0$. Jelikož funkce λ závisí na ξ spojitě a periodicky s periodou $2\pi/h$, existuje $(h, \tau) \in \Lambda$ a interval $(\xi_1, \xi_2) \subset (-\pi/h, \pi/h)$ tak, že $|\lambda(\xi)| > 1 + C\tau$ $\forall \xi \in (\xi_1, \xi_2)$. Necht'

$$\widehat{U}^0(\xi) = \frac{1}{\sqrt{\xi_2 - \xi_1}} \quad \text{pro } \xi \in (\xi_1, \xi_2) \quad \text{a} \quad \widehat{U}^0(\xi) = 0 \quad \text{pro } \xi \notin (\xi_1, \xi_2).$$

Pak dle (16) a (24) je $\|U^0\|_2 = \|\widehat{U}^0\|_{L^2(-\pi/h, \pi/h)} = 1$ a

$$\|U^n\|_2 = \|\lambda^n \widehat{U}^0\|_{L^2(\xi_1, \xi_2)} > (1 + C\tau)^n \geq (1 + C\tau_{max})^{n\tau/\tau_{max}} \|U^0\|_2,$$

kde τ_{max} je takové, že $\Lambda \subset \mathbb{R}^+ \times (0, \tau_{max})$ (využili jsme, že funkce $(1 + C\tau)^{1/\tau}$ je klesající). Pro libovolné $T > \tau_{max}$ a $n \in \mathbb{N}$ splňující $T/2 \leq n\tau \leq T$ je

$$\|U^n\|_2 > (1 + C\tau_{max})^{T/(2\tau_{max})} \|U^0\|_2.$$

Schéma (18) tedy není stabilní v Λ , neboť C lze volit libovolně velké. Je-li funkce $\lambda(\xi/h)$ na Λ nezávislá na h a τ , pak jsou zřejmě podmínky (22) a (23) ekvivalentní. \square

Poznámka 1 Uvedená věta pochází od von Neumanna, a analýza diferenčních metod založená na Fourierově metodě se proto obvykle nazývá von Neumannova analýza. Nerovnost (22) se většinou nazývá von Neumannova podmínka.

Uvažujme rovnici (4), přičemž opět předpokládáme, že $b > 0$ a a jsou konstanty, a ukažme si provedení analýzy stability na příkladu explicitního schématu (6). Použitím ekvivalentního tvaru (9) zjistíme, že amplifikační faktor je dán vztahem

$$\lambda(\xi) = 1 - 4\mu \sin^2 \frac{\xi h}{2} - i\nu \sin(\xi h).$$

Pro nalezení nutné a postačující podmínky pro stabilitu v diskretní L^2 normě je tedy třeba uvažovat podmínku (22). Snadno vypočítáme, že

$$|\lambda(\xi)|^2 = (1 - 4\mu s^2)^2 + 4\nu^2 s^2 (1 - s^2), \quad \text{kde } s = \sin \frac{\xi h}{2}.$$

Při $s^2 = 1$ stačí požadovat, aby $\mu \leq \frac{1}{2}$, neboť pak $|\lambda(\xi)| \leq 1$. Jelikož $\nu^2 = a^2 \mu \tau / b$ a $4s^2(1 - s^2) \leq 1$, dostáváme při $\mu \leq \frac{1}{2}$ nerovnost

$$|\lambda(\xi)| \leq \left(1 + \frac{a^2}{2b} \tau\right)^{1/2} \leq 1 + \frac{a^2}{4b} \tau.$$

Pro $\mu \leq \frac{1}{2}$ je tedy von Neumannova podmínka splněna a schéma je stabilní v diskretní L^2 normě. Nicméně při $\mu \leq \frac{1}{2}$ může pro některé hodnoty ξ být $|\lambda(\xi)| > 1$, což vede k růstu příslušné amplitudy v diskretním řešení, zatímco v řešení diferenciální rovnice jsou všechny amplitudy tlumeny. Pokud tedy exponenciální růst v čase, který von Neumannova podmínka umožňuje, neodpovídá vlastnostem aproximované parciální diferenciální rovnice, je von Neumannova podmínka v praxi příliš slabá. Zavádíme proto následující silnější definici stability.

Definice 2 Diferenční schéma se nazývá *silně stabilní*, jestliže platí: pokud Fourierova transformace řešení aproximované parciální diferenciální rovnice splňuje

$$(25) \quad |\widehat{u}(\xi, t + \tau)| \leq e^{\alpha \tau} |\widehat{u}(\xi, t)| \quad \forall \xi \in \mathbb{R}$$

pro nějaké $\alpha \geq 0$, pak amplifikační faktory diferenčního schématu splňují

$$|\lambda(\xi)| \leq e^{\alpha \tau} \quad \forall \xi \in \mathbb{R}.$$

Pro rovnici (4) platí (25) s $\alpha = 0$, a tedy požadujeme, aby $|\lambda(\xi)| \leq 1 \quad \forall \xi \in \mathbb{R}$. Jelikož

$$|\lambda(\xi)|^2 = 1 - 4(2\mu - \nu^2)s^2 + 4(4\mu^2 - \nu^2)s^4,$$

dostáváme

$$|\lambda(\xi)| \leq 1 \quad \forall \xi \in \mathbb{R} \quad \Leftrightarrow \quad (4\mu^2 - \nu^2)s^2 \leq 2\mu - \nu^2 \quad \forall s \in [-1, 1] \quad \Leftrightarrow \quad \nu^2 \leq 2\mu \leq 1.$$

Kromě očekávané podmínky $\mu \leq \frac{1}{2}$ jsme tedy dostali ještě další podmínku, kterou lze zapsat ve tvaru $\tau \leq 2b/a^2$. To může být velmi vážné omezení, neboť v praxi je často $b \ll |a|$.

3 Numerické řešení transportní rovnice

V této části se budeme zabývat numerickým řešením Cauchyovy úlohy pro transportní rovnici v jedné prostorové dimenzi, tj. úlohy

$$(26) \quad u_t + a u_x = 0 \quad \text{v } \mathbb{R} \times \mathbb{R}^+,$$

$$(27) \quad u(x, 0) = u^0(x) \quad \forall x \in \mathbb{R}$$

pro neznámou funkci $u = u(x, t)$ závisující na prostorové proměnné $x \in \mathbb{R}$ a časové proměnné $t \in \mathbb{R}_0^+$. Funkce u^0 je daná počáteční podmínka a koeficient a je rychlost šíření veličiny u . Náš hlavní zájem bude patřit rovnici (26) s konstantním koeficientem a , avšak v některých případech budeme uvažovat též rovnici (26) s nenulovou pravou stranou či s nekonztantním koeficientem a .

Je známo, že řešení rovnice (26) jsou konstantní podél charakteristik, což jsou křivky $x = x(t)$ splňující $x'(t) = a(x(t), t)$. Je-li funkce a Lipschitzovská v x a spojitá v t , charakteristiky se neprotínají. Jejich sestavením tedy získáme řešení dané vztahem $u(x(t), t) = u^0(x(0))$. Je-li a konstantní, pak charakteristiky jsou určeny rovnicí $x - at = \text{konst.}$ Jsou to tedy rovnoběžné přímky v rovině (x, t) se směrnici $1/a$. Z toho plyne, že pro $u^0 \in C^1(\mathbb{R})$ a konstantní a má úloha (26), (27) právě jedno klasické řešení, které je dáno vztahem

$$(28) \quad u(x, t) = u^0(x - at).$$

Vidíme tedy, že řešení u v čase t je rovno počáteční podmínce u^0 posunuté o vzdálenost $|a|t$ (doprava pro $a > 0$, doleva pro $a < 0$). Parametr a se nazývá rychlost šíření podél charakteristik. Řešení u lze tak považovat za vlnu šířící se rychlostí a beze změny tvaru. Proto se rovnice (26) také nazývá *jednosměrná vlnová rovnice*.

Samozřejmě vyvstává otázka, proč se zabývat numerickým řešením úlohy (26), (27), když je její řešení známo. Důvodem je to, že u komplikovanějších úloh, v nichž hrají důležitou roli transportní mechanismy, již řešení zpravidla nejsme schopni v analytickém tvaru získat a musíme ho aproximovat pomocí numerických metod. Přitom navržením vhodné metody není vůbec triviální a řada potíží, s nimiž se setkáváme, se vyskytuje již při numerickém řešení rovnice (26). Abychom příslušné jevy lépe pochopili, je rozumné je studovat na co nejjednodušší úloze, a proto uvažujeme rovnici (26). U metod, které se ukáží jako nevhodné pro řešení rovnice (26), nelze samozřejmě očekávat, že by dávaly dobré výsledky pro komplikovanější úlohy.

3.1 Příklady numerických schémat

Opět budeme uvažovat rovnoměrnou síť s konstantním prostorovým krokem h a konstantním časovým krokem τ a definujeme uzly sítě (x_j, t_n) , kde $x_j = jh$, $t_n = n\tau$, $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$. Řešení u Cauchyovy úlohy (26), (27) opět aproximujeme v uzlech (x_j, t_n) hodnotami U_j^n . Rovnici (26) aproximujeme v libovolném uzlu (x_j, t_n) diferenčním schématem získaným nahrazením parciálních derivací diferenčními kvocienty. Takto lze získat mnoho rozličných numerických schémat, z nichž několik nyní uvedeme. Pro jejich přehlednější zápis budeme užívat parametr $\nu = a\tau/h$, který jsme zavedli již v části věnované von Neumannově analýze stability.

Pro numerické řešení rovnice (26) lze uvažovat například schémata:

$$(29) \quad \frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_{j+1}^n - U_j^n}{h} = 0, \quad \text{tj.} \quad U_j^{n+1} = (1 + \nu) U_j^n - \nu U_{j+1}^n,$$

$$(30) \quad \frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_j^n - U_{j-1}^n}{h} = 0, \quad \text{tj.} \quad U_j^{n+1} = \nu U_{j-1}^n + (1 - \nu) U_j^n,$$

$$(31) \quad \frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} = 0, \quad \text{tj.} \quad U_j^{n+1} = \frac{\nu}{2} U_{j-1}^n + U_j^n - \frac{\nu}{2} U_{j+1}^n,$$

$$(32) \quad \frac{U_j^{n+1} - U_j^{n-1}}{2\tau} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} = 0, \quad \text{tj.} \quad U_j^{n+1} = U_j^{n-1} + \nu U_{j-1}^n - \nu U_{j+1}^n,$$

$$(33) \quad \frac{U_j^{n+1} - \frac{1}{2}(U_{j+1}^n + U_{j-1}^n)}{\tau} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} = 0, \quad \text{tj.} \\ U_j^{n+1} = \frac{1}{2}(1 + \nu) U_{j-1}^n + \frac{1}{2}(1 - \nu) U_{j+1}^n.$$

Všechna uvedená schémata jsou explicitní. Schéma (32) je známo pod označením *leapfrog scheme*. K určení přibližného řešení na nové časové vrstvě využívá toto schéma hodnot přibližného řešení z předcházejících dvou časových vrstev. Hovoříme proto o *dvoukrokovém* schématu. Ostatní schémata jsou jednokroková. Schéma (33) se nazývá Laxovo–Friedrichsovo schéma. Všimněme si, že

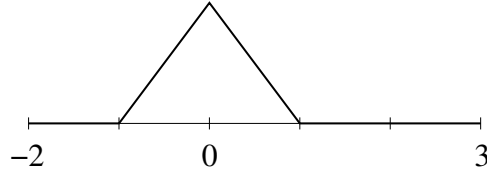
$$\frac{U_j^{n+1} - \frac{1}{2}(U_{j+1}^n + U_{j-1}^n)}{\tau} = \frac{U_j^{n+1} - U_j^n}{\tau} - \frac{h^2}{2\tau} \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2}.$$

Jelikož diferenční kvocient $(U_{j+1}^n - 2U_j^n + U_{j-1}^n)/h^2$ reprezentuje aproximaci $u_{xx}(x_j, t_n)$, lze tedy speciální aproximaci časové derivace použitou v (33) interpretovat jako přidání umělé difúze o velikosti $h^2/(2\tau)$ ke schématu (31). Laxovo–Friedrichsovo schéma tedy odpovídá diskretizaci rovnice

$$u_t - \frac{h^2}{2\tau} u_{xx} + a u_x = 0$$

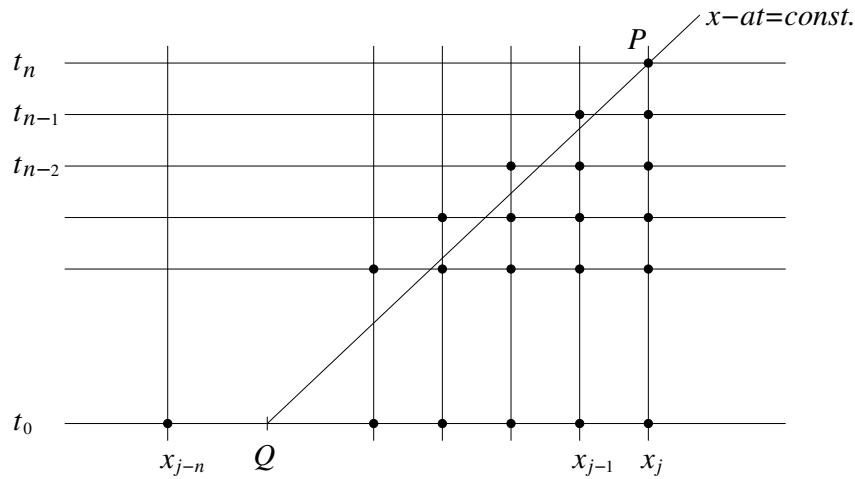
s použitím dopředné diference pro časovou derivaci a centrálních diferencí pro prostorové derivace.

Cvičení 1 Uvažujme rovnici $u_t + u_x = 0$ v oblasti $(-2, 3) \times \mathbb{R}^+$ s počáteční podmínkou u^0 znázorněnou v obr. 1. Definujme ekvidistanční prostorové uzly $-2 = x_0 < x_1 < \dots < x_J = 3$ s prostorovým krokem $h = 5/J$. Dále uvažujme konstatní časový krok $\tau > 0$ a časové hladiny $t_n = n\tau$, $n \in \mathbb{N}_0$. V čase t_0 je přibližné řešení určeno počáteční podmínkou. Na následujících časových hladinách pak můžeme hodnoty přibližného řešení ve vnitřních uzlech určit pomocí libovolného z výše uvedených schémat, přičemž položíme $U_0^n = 0$ a $U_J^n = U_{J-1}^n$ pro $n \in \mathbb{N}$ (pro $n = 1$ však nelze použít schéma (32)). Naprogramujte výše uvedená schémata a proveďte výpočty pro $J = 50$ a $\tau = 0.08$. U Laxova–Friedrichsova schématu otestujte, jaký vliv má zmenšení h na polovinu při nezměněném τ , zmenšení τ na polovinu při nezměněném h a současné zmenšení h i τ na polovinu, popř. na čtvrtinu.



Obrázek 1: Počáteční podmínka u^0 použitá ve cvičení 1.

Provedeme-li výpočty popsané ve cvičení 1, zjistíme, že v přibližných řešeních schémat (29) a (31) se již po několika málo časových krocích objeví velké oscilace a dostaneme zcela nepoužitelné výsledky. Na druhou stranu zbývající tři schémata vedou při $J = 50$ a $\tau = 0.08$ k přijatelným aproximacím řešení uvažované úlohy. Změníme-li však J či τ , objeví se v některých případech nestabilní chování i u těchto schémat. Naším cílem nyní bude vysvětlit, proč k uvedeným jevům dochází.



Obrázek 2: Oblast závislosti PDR a schématu (30).

Zvolme bod $P \in \mathbb{R} \times \mathbb{R}^+$ a uvažujme schéma (30) a takovou síť, že bod P je její uzel. Nechť např. $P = (x_j, t_n)$. Pak hodnota U_j^n je určena hodnotami U_{j-1}^{n-1} a U_j^{n-1} . Podobně tyto dvě hodnoty jsou určeny hodnotami U_{j-2}^{n-2} , U_{j-1}^{n-2} a U_j^{n-2} , viz obr. 2. Z toho plyne, že hodnota U_j^n závisí pouze na hodnotách počáteční podmínky u^0 v bodech $x_{j-n}, x_{j-n+1}, \dots, x_{j-1}, x_j$. Pokud bychom uvažovali v rovnici (26) nenulovou pravou stranu f , závisela by hodnota U_j^n též na hodnotách f v uzlech sítě ležících v trojúhelníku s vrcholy (x_j, t_n) , $(x_{j-n}, 0)$ a $(x_j, 0)$. Tento trojúhelník se nazývá *oblast závislosti* schématu v uzlu (x_j, t_n) .

Pro řešení u úlohy (26), (27) plyne ze vztahu (28), že $u(P) = u^0(Q)$, kde Q je průsečík přímky $t = 0$ a charakteristiky $x - at = \text{const.}$ procházející bodem P , viz obr. 2. V případě rovnice (26) s nenulovou pravou stranou f závisí hodnota $u(P)$ též na hodnotách funkce f v bodech úsečky PQ . Úsečka PQ je *oblastí závislosti* uvažované parciální diferenciální rovnice v bodě P . Nyní můžeme zformulovat následující podmínku.

Věta 2 (CFL podmínka; Courant, Friedrichs, Lewy (1928)) *Nutná podmínka konvergence diferenčního schématu je, aby oblast závislosti parciální diferenciální rovnice ležela uvnitř oblasti závislosti diferenčního schématu.*

Pro schéma (30) je CFL podmínka splněna, leží-li bod Q na přímce $t = 0$ mezi body x_{j-n} a x_j . To je právě tehdy, když $a \geq 0$ a $a\tau \leq h$. Uvažujme posloupnost sítí s $\tau/h = \text{const.}$ a $h \rightarrow 0$ a nechť P je uzlem všech těchto sítí. Pak oblast závislosti schématu (30) v uzlu P je pro všechny tyto sítě stejná. Pokud $a < 0$ nebo $a\tau > h$, pak pro všechny tyto sítě hodnota přibližného řešení U v uzlu P nezávisí na hodnotách počáteční podmínky u^0 v pevně daném okolí bodu Q (neměnném při $h \rightarrow 0$). Přibližná řešení v uzlu P tudíž obecně nemohou konvergovat k hodnotě $u(P)$. Z uvedeného je zřejmé, že platí též následující věta.

Věta 3 *Nutnou podmínkou konvergence explicitního schématu typu $U_j^{n+1} = \alpha U_{j-1}^n + \beta U_j^n + \gamma U_{j+1}^n$ pro rovnici (26) při $\tau/h = \text{const.}$ je, aby platilo $|a\tau/h| \leq 1$.*

Nerovnost $|a\tau/h| \leq 1$ z věty 3 se také často nazývá CFL podmínka.

Všimněme si, že pro schéma tvaru uvažovaného ve větě 3 se informace během jednoho časové kroku rozšíří o prostorový krok h . Numerická rychlost šíření informace je tudíž h/τ . Uvedená nerovnost tedy říká, že numerická rychlost šíření musí být větší nebo rovna rychlosti šíření odpovídající uvažované parciální diferenciální rovnici. Pokud diferenční schéma nemůže šířit řešení alespoň tak rychle, jako se šíří řešení parciální diferenciální rovnice, nemůže řešení schématu konvergovat k řešení parciální diferenciální rovnice.

Na CFL podmínku se lze dívat též jako na nutnou podmínku stability diferenčního schématu. Obecně není pro stabilitu postačující, ale její velkou předností je její jednoduchost. Umožňuje tak vyřadit řadu diferenčních schémat s nepatrnou námahou věnovanou jejich vyšetřování. Teprve schémata, která splňují CFL podmínku je vhodné vyšetřovat podrobněji použitím kritérií, která jsou pro jejich stabilitu postačující.

Vrátíme-li se ke schématům (29)–(33), pak vidíme, že ve všech případech je nutno splnit podmínku $|a\tau/h| \leq 1$. To vysvětluje, proč schéma (33) vede pro úlohu ze cvičení 1 při $h = 0.05$ a $\tau = 0.08$ k oscilacím. Je též zřejmé, že schéma (29) nelze pro úlohu ze cvičení 1 použít, neboť oblast závislosti parciální diferenciální rovnice neleží v oblasti závislosti schématu. Oscilace lze pozorovat též v případě schématu (31) při splnění CFL podmínky. Vyšetřeme proto stabilitu tohoto schématu pomocí Fourierovy metody (tj. provedme von Neumannovu analýzu stability).

Rigorózní postup určení amplifikačního faktoru schématu (31) spočívá ve vyjádření přibližného řešení pomocí jeho Fourierovy transformace a využití ortonormality funkcí $\{e^{i\xi j h}\}_{j \in \mathbb{Z}}$. Viděli jsme však, že formálně lze amplifikační faktor jednokrokového schématu získat nahrazením přibližného řešení síťovou funkcí $\lambda(\xi)^n e^{i\xi j h}$. Dosazením této funkce za U_j^n ve schématu (31) získáme

$$\lambda(\xi) = 1 + \frac{\nu}{2} (e^{-i\xi h} - e^{i\xi h}) = 1 - i\nu \sin(\xi h).$$

Pro libovolnou síť a každé ξ , pro které $\sin(\xi h) \neq 0$, je tedy $|\lambda(\xi)| > 1$. Pro každý proces zjemňování, při němž je ν pevné, je proto schéma nestabilní. Vidíme tedy, že CFL podmínka skutečně není pro stabilitu postačující.

Jak jsme viděli, při $a > 0$ lze k řešení rovnice (26) použít schéma (30). Bude-li $a < 0$, nelze toto schéma použít, neboť není splněna CFL podmínka, avšak můžeme použít schéma (29). To nás vede ke schématu

$$(34) \quad U_j^{n+1} = \begin{cases} (1 + \nu) U_j^n - \nu U_{j+1}^n, & \text{je-li } a < 0, \\ \nu U_{j-1}^n + (1 - \nu) U_j^n, & \text{je-li } a > 0. \end{cases}$$

Je-li a nekonstantní, pak a v (34) značí hodnotu $a(x_j, t_n)$. Všimněme si, že při $a > 0$ se informace šíří ve směru kladné poloosy x (tj. zleva doprava) a k diskretizaci v prostorovém bodě x_j je využita informace v bodě x_j a vlevo od něj. Na druhou stranu při $a < 0$ se informace šíří ve směru záporné poloosy x (tj. zprava doleva) a k diskretizaci v bodě x_j je použita informace v bodě x_j a vpravo od něj. Jedná se o tzv. *diskretizaci typu upwind*, kdy k diskretizaci v daném bodě využíváme informaci, která leží proti směru šíření (tj. která do daného bodu přichází).

Je-li $|a \tau / h| \leq 1$, splňuje schéma (34) CFL podmínku. Provedme též von Neumannovu analýzu stability (při konstantním a). Dosazením $U_j^n := \lambda(\xi)^n e^{i \xi j h}$ do schématu (34) dostaneme při $a < 0$

$$\lambda(\xi) = (1 + \nu) - \nu e^{i \xi h} = 1 + \nu - \nu \cos(\xi h) - i \nu \sin(\xi h)$$

a při $a > 0$

$$\lambda(\xi) = \nu e^{-i \xi h} + (1 - \nu) = 1 - \nu + \nu \cos(\xi h) - i \nu \sin(\xi h).$$

Tyto dva vztahy lze pro obě znaménka koeficientu a psát ve tvaru

$$(35) \quad \lambda(\xi) = 1 - |\nu| + |\nu| \cos(\xi h) - i \nu \sin(\xi h),$$

z čehož plyne, že

$$(36) \quad |\lambda(\xi)|^2 = 1 - 4 |\nu| (1 - |\nu|) \sin^2 \frac{\xi h}{2}.$$

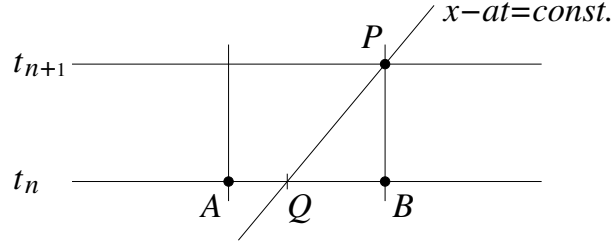
Zjišťujeme tedy, že $|\lambda(\xi)| \leq 1$ pro všechna $\xi \in \mathbb{R}$ právě tehdy, když $|\nu| \leq 1$. V tomto případě tedy dává CFL podmínka správné meze stability.

Všimněme si, že schéma (34) lze přepsat do tvaru

$$\frac{U_j^{n+1} - U_j^n}{\tau} - \frac{|a| h}{2} \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} = 0.$$

Podobně jako Laxovo–Friedrichsovo schéma můžeme tedy i schéma (34) interpretovat jako přidání umělé difúze ke schématu (31). V tomto případě má umělá difúze velikost $|a|h/2$. Jelikož $|a|h/2 = |\nu| h^2 / (2\tau)$, je při $|\nu| < 1$ umělá difúze přidávaná schématem (34) vždy menší než umělá difúze odpovídající Laxovu–Friedrichsovu schématu. Při $|\nu| = 1$ jsou obě schémata stejná a dávají přesné řešení, neboť uzly sítě leží na charakteristikách.

Ukažme si ještě jiný způsob odvození schématu (34). Pro jednoduchost uvažujme případ $a > 0$ a nechť $a \tau \leq h$. Předpokládejme, že známe hodnoty U_A, U_B přibližného



Obrázek 3: Odvození schématu (34) pomocí charakteristik.

řešení v uzlech A , B a chceme určit hodnotu U_P přibližného řešení v uzlu P , viz obr. 3. Charakteristika procházející uzlem P protíná síťovou přímkou $t = t_n$ v bodě Q ležícím mezi body A a B . Řešení u rovnice (26) splňuje $u(P) = u(Q)$ a je tedy přirozené definovat hodnotu U_P jako aproximaci u v bodě Q získanou pomocí hodnot U_A a U_B . Použijeme-li lineární interpolaci, dostaneme vztah

$$U_P = \frac{|BQ|}{|AB|} U_A + \frac{|AQ|}{|AB|} U_B.$$

To je přesně schéma (30), tj. schéma (34) pro $a > 0$, neboť $|AB| = h$, $|BQ| = a\tau = \nu h$ a $|AQ| = |AB| - |BQ| = (1 - \nu)h$.

3.2 Fázová rychlost a disperze

Uvažujme úlohu (26), (27) s $a = const$. Definujme Fourierovu transformaci řešení u vztahem

$$\hat{u}(\xi, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u(x, t) e^{-i\xi x} dx, \quad \xi \in \mathbb{R}.$$

Pak

$$\hat{u}_t + i\xi a \hat{u} = 0 \quad \text{v } \mathbb{R} \times \mathbb{R}^+, \quad \hat{u}(\xi, 0) = \hat{u}^0(\xi) \quad \forall \xi \in \mathbb{R},$$

kde

$$\hat{u}^0(\xi) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u^0(x) e^{-i\xi x} dx, \quad \xi \in \mathbb{R}.$$

Je tedy

$$\hat{u}(\xi, t) = \hat{u}^0(\xi) e^{-i\xi a t}$$

a platí

$$u(x, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}(\xi, t) e^{i\xi x} d\xi = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}^0(\xi) e^{i\xi(x-at)} d\xi = u^0(x - at).$$

Z uvedeného plyne, že

$$(37) \quad \hat{u}(\xi, t + \tau) = \hat{u}(\xi, t) e^{-i\xi a \tau}.$$

Během jednoho časového kroku se tedy fáze Fourierovy transformace řešení úlohy (26), (27) změní o $-\xi a \tau$. Všimněme si též, že pro libovolné $\xi \in \mathbb{R}$ je $|\hat{u}(\xi, t)| = |\hat{u}^0(\xi)|$, tj. amplituda libovolného Fourierova členu je v čase konstantní.

V diskrétním případě je

$$U_j^n = \frac{1}{\sqrt{2\pi}} \int_{-\pi/h}^{\pi/h} \widehat{U}^n(\xi) e^{i\xi j h} d\xi, \quad \text{kde} \quad \widehat{U}^n(\xi) = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} h U_j^n e^{-i\xi j h}.$$

Všimněme si, že interval $[-\pi/h, \pi/h]$ obsahuje všechny frekvence, které lze na síti s krokem h reprezentovat, neboť pro libovolné $k \in \mathbb{Z}$ je

$$e^{i(\xi + \frac{2\pi}{h}k)jh} = e^{i\xi j h} e^{i2\pi k j} = e^{i\xi j h}.$$

Pro libovolné jedнокrokové diferenční schéma platí $\widehat{U}^{n+1}(\xi) = \lambda(\xi) \widehat{U}^n(\xi)$. Vzhledem ke vztahu (37) je tedy žádoucí, aby $\lambda(\xi)$ byla dobrá aproximace $e^{-i\xi a \tau}$. Abychom tato dvě čísla mohli lépe porovnat, zapíšeme amplifikační faktor ve tvaru $\lambda(\xi) = |\lambda(\xi)| e^{-i\xi \alpha(\xi) \tau}$. Hodnota $|\lambda(\xi)|$ vyjadřuje pokles amplitudy složky přibližného řešení o frekvenci ξ během jednoho časového kroku. Jelikož, jak jsme viděli, pro přesné řešení libovolný Fourierův člen není tlumený, požadujeme, aby $|\lambda(\xi)|$ bylo blízké hodnotě 1. Veličina $\alpha(\xi)$ se nazývá *fázová rychlost*. Je to rychlost, kterou diferenční schéma šíří vlny o frekvenci ξ . Pokud by bylo $\alpha(\xi) = a$ pro všechna ξ , pak by se vlny šířily správnou rychlostí, avšak s tím se obecně nesetkáme u žádného z diferenčních schémat. Rychlost $\alpha(\xi)$ je pouze aproximací a a je obecně různá pro různé hodnoty ξ , což se projevuje deformací tvaru řešení diferenčního schématu. Jev, kdy se vlny o různých frekvencích pohybují různými rychlostmi, se nazývá *disperze*. Rozdíl $a - \alpha(\xi)$ se nazývá *fázová chyba*.

Fázovou rychlost můžeme určit ze vztahu

$$\tan(\xi \alpha(\xi) \tau) = -\frac{\operatorname{Im} \lambda(\xi)}{\operatorname{Re} \lambda(\xi)}.$$

Je-li $|\lambda(\xi)| = 1$, platí $\sin(\xi \alpha(\xi) \tau) = -\operatorname{Im} \lambda(\xi)$. Fázovou chybu je vhodné vyšetřovat zejména pro malé hodnoty $|\xi h|$, neboť vlny s takovými frekvencemi lze na síti s krokem h dobře aproximovat. Přitom můžeme využít následující vztahy pro počítání s malými čísly:

$$\begin{aligned} \frac{1}{1+x} &= 1 - x + O(x^2), & \sin x &= x \left(1 - \frac{x^2}{6} + O(x^4) \right), & \cos x &= 1 - \frac{x^2}{2} + O(x^4), \\ \tan x &= x \left(1 + \frac{x^2}{3} + O(x^4) \right), & \tan^{-1} y &= y \left(1 - \frac{y^2}{3} + O(y^4) \right). \end{aligned}$$

Příklad 1 Vyšetřeme fázovou chybu pro schéma (34).

Řešení: Podle (35) a výše uvedených vztahů je

$$\begin{aligned} \tan(\xi \alpha(\xi) \tau) &= \frac{\nu \sin(\xi h)}{1 - |\nu| + |\nu| \cos(\xi h)} = \frac{\nu \xi h \left(1 - \frac{(\xi h)^2}{6} + O((\xi h)^4) \right)}{1 - |\nu| \frac{(\xi h)^2}{2} + O((\xi h)^4)} \\ &= \nu \xi h \left[1 - \left(\frac{1}{6} - \frac{|\nu|}{2} \right) (\xi h)^2 + O((\xi h)^4) \right], \end{aligned}$$

z čehož plyne (je $\nu \xi h = \xi a \tau$)

$$\begin{aligned}\alpha(\xi) &= a \left[1 - \left(\frac{1}{6} - \frac{|\nu|}{2} \right) (\xi h)^2 + O((\xi h)^4) \right] \left(1 - \frac{\nu^2}{3} (\xi h)^2 + O((\xi h)^4) \right) \\ &= a \left[1 - \frac{1}{6} (1 - |\nu|) (1 - 2|\nu|) (\xi h)^2 + O((\xi h)^4) \right].\end{aligned}$$

Vidíme, že fázová chyba je řádu $(\xi h)^2$ a znaménko závisí na ν . Pro $|\nu| \in (\frac{1}{2}, 1)$ dochází k předbíhání a pro $|\nu| \in (0, \frac{1}{2})$ ke zpoždování fáze. Pro $|\nu| = \frac{1}{2}$ dostaneme z (35)

$$\begin{aligned}\lambda(\xi) &= \frac{1}{2} [1 + \cos(\xi h) - i 2\nu \sin(\xi h)] = \frac{1}{2} [1 + \cos(2\nu \xi h) - i \sin(2\nu \xi h)] \\ &= \frac{1}{2} [1 + \cos^2(\nu \xi h) - \sin^2(\nu \xi h) - i 2 \sin(\nu \xi h) \cos(\nu \xi h)] \\ &= \cos(\nu \xi h) [\cos(\nu \xi h) - i \sin(\nu \xi h)] = \cos\left(\frac{\xi h}{2}\right) e^{-i\xi a \tau},\end{aligned}$$

a tudíž chyba fáze je nulová. Pro $|\nu| = 1$ víme, že schéma dává přesné řešení, což plyne i z toho, že podle (35) je

$$\lambda(\xi) = \cos(\xi h) - i\nu \sin(\xi h) = \cos(\nu \xi h) - i \sin(\nu \xi h) = e^{-i\xi a \tau}.$$

Ukažme si to např. pro $\nu = 1$. Pak $\widehat{U}^n(\xi) = \lambda(\xi)^n \widehat{U}^0(\xi) = e^{-i\xi n h} \widehat{U}^0(\xi)$, a tudíž

$$\begin{aligned}U_j^n &= \frac{1}{\sqrt{2\pi}} \int_{-\pi/h}^{\pi/h} \widehat{U}^0(\xi) e^{i\xi(j-n)h} d\xi = U_{j-n}^0 = u^0(x_{j-n}) = u^0(x_j - nh) \\ &= u^0(x_j - a t_n) = u(x_j, t_n).\end{aligned}$$

Ze vztahu (36) vidíme, že kromě případu $|\nu| = 1$ dochází u upwind schématu (34) vždy k chybě v amplitudě řádu $(\xi h)^2$ v jednom časovém kroku, což vede ke globální chybě řádu ξh (na daném časovém intervalu $[0, T]$). Pro většinu problémů je takovéto tlumení nepřijatelné.

V případě Laxova–Friedrichsova schématu (33) je

$$\lambda(\xi) = \cos(\xi h) - i\nu \sin(\xi h), \quad |\lambda(\xi)|^2 = 1 - (1 - \nu^2) \sin^2(\xi h).$$

Schéma je tedy stabilní pro $|\nu| \leq 1$, tj. opět pro ty hodnoty ν , pro něž je splněna CFL podmínka. Tlumení je stejného řádu jako u schématu (34), avšak pokles amplitudy je více než dvojnásobný (v jednom časovém kroku zhruba $(1 + |\nu|)(1 - |\nu|)(\xi h)^2$ namísto $|\nu|(1 - |\nu|)(\xi h)^2$ u schématu (34)). Analogicky jako výše dostaneme

$$\alpha(\xi) = a \left[1 + \frac{1}{3} (1 - \nu^2) (\xi h)^2 + O((\xi h)^4) \right].$$

Kromě případu $|\nu| = 1$, kdy Laxovo–Friedrichsovo schéma je totožné se schématem (34) a dává tedy přesné řešení, dochází podle uvedeného vztahu vždy k předbíhání fáze.

3.3 Diskrétní princip maxima pro Cauchyovy úlohy

Uvažujme Cauchyovu úlohu najít funkci $u = u(x, t)$ definovanou pro $x \in \mathbb{R}$ a $t \geq 0$ a splňující

$$u_t + Lu = 0 \quad \forall \mathbb{R} \times \mathbb{R}^+, \quad u(x, 0) = u^0(x) \quad \text{pro } x \in \mathbb{R},$$

kde u^0 je zadaná počáteční podmínka a L je lineární diferenciální operátor s konstantními koeficienty tvaru

$$L = \sum_{k=1}^m a_k \frac{\partial^k}{\partial x^k}.$$

Obecné jednokrokové schéma pro tuto úlohu na stejnoměrné síti s uzly (x_j, t_n) , kde $x_j = jh$, $t_n = n\tau$, $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$, má tvar

$$(38) \quad \sum_{s=-M}^M \alpha_s U_{j+s}^{n+1} = \sum_{s=-M}^M \beta_s U_{j+s}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

$$(39) \quad U_j^0 = u^0(x_j) \quad \forall j \in \mathbb{Z}.$$

Budeme předpokádat, že $\alpha_0 > 0$, $\alpha_s \leq 0 \quad \forall s \neq 0$, $\beta_s \geq 0 \quad \forall s$ a

$$(40) \quad \sum_{s=-M}^M \alpha_s = \sum_{s=-M}^M \beta_s > 0.$$

Poznámka 2 Je-li schéma (38) ve tvaru s diferenčními kvocienty, pak

$$\sum_{s=-M}^M \alpha_s u_{j+s}^{n+1} - \sum_{s=-M}^M \beta_s u_{j+s}^n \approx (u_t + Lu)(x_j, t_{n+\theta}),$$

kde $\theta \in [0, 1]$. Je rozumné požadovat, aby tato aproximace byla přesná, pokud u je lineární funkce (polynom prvního stupně v x a t). Volbou $u(x, t) = 1$ a $u(x, t) = t$ pak dostáváme podmínku (40).

Pro libovolnou síťovou funkci $\{U_j^n\}_{j \in \mathbb{Z}, n \in \mathbb{N}_0}$ zavedeme označení

$$U_{\min}^n = \inf_{j \in \mathbb{Z}} U_j^n, \quad U_{\max}^n = \sup_{j \in \mathbb{Z}} U_j^n, \quad n \in \mathbb{N}_0.$$

Věta 4 (Diskrétní princip maxima) *Nechť $\{U_j^n\}_{j \in \mathbb{Z}}$ je pro každé $n \in \mathbb{N}_0$ omezené a splňuje vztah (38). Pak*

$$(41) \quad U_{\min}^n \leq U_{\min}^{n+1} \leq U_{\max}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0.$$

Pokud místo (38) je pouze

$$(42) \quad \sum_{s=-M}^M \alpha_s U_{j+s}^{n+1} \leq \sum_{s=-M}^M \beta_s U_{j+s}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

platí

$$(43) \quad U_j^{n+1} \leq U_{\max}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0.$$

Pokud

$$(44) \quad \sum_{s=-M}^M \alpha_s U_{j+s}^{n+1} \geq \sum_{s=-M}^M \beta_s U_{j+s}^n \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0,$$

platí

$$(45) \quad U_{\min}^n \leq U_j^{n+1} \quad \forall j \in \mathbb{Z}, n \in \mathbb{N}_0.$$

Důkaz. Necht' platí (42). Pak pro libovolné $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$ je

$$\alpha_0 U_j^{n+1} \leq \sum_{\substack{s=-M \\ s \neq 0}}^M (-\alpha_s) U_{j+s}^{n+1} + \sum_{s=-M}^M \beta_s U_{j+s}^n \leq U_{\max}^{n+1} \sum_{\substack{s=-M \\ s \neq 0}}^M (-\alpha_s) + U_{\max}^n \sum_{s=-M}^M \beta_s$$

Označíme-li $A = \sum_{s=-M}^M \beta_s$, plyne z této nerovnosti díky (40)

$$\alpha_0 U_{\max}^{n+1} \leq (\alpha_0 - A) U_{\max}^{n+1} + A U_{\max}^n,$$

a tedy i (43). Je-li splněno (44), pak $-U_j^n$ splňuje (42), a tudíž pro všechna $j \in \mathbb{Z}$ a $n \in \mathbb{N}_0$ je $-U_j^{n+1} \leq \sup_{k \in \mathbb{Z}} (-U_k^n) = -U_{\min}^n$, tj. platí (45). Platí-li (38), platí dle předchozího (43) a (45), a tudíž i (41). \square

Diskrétní princip maxima (41) zaručuje, že se v přibližném řešení neobjeví narůstající oscilace pozorované u některých metod ve cvičení 1. Navíc lze díky diskrétnímu principu maxima odhadnout chybu aproximace pomocí odhadu chyby diskretizace. Ukažme si to na příkladu schématu (29). Pro toto schéma je chyba diskretizace dána vztahem

$$\varepsilon_j^n = \frac{u_j^{n+1} - u_j^n}{\tau} + a \frac{u_{j+1}^n - u_j^n}{h},$$

kde opět $u_j^n = u(x_j, t_n)$. Je-li u dvakrát spojitě derivovatelné, pak použitím Taylorova vzorce získáme

$$(46) \quad \varepsilon_j^n = \frac{1}{2} u_{tt}(x_j, \eta) \tau + \frac{a}{2} u_{xx}(\xi, t_n) h,$$

kde $\xi \in (x_j, x_j + h)$ a $\eta \in (t_n, t_n + \tau)$. Jedná se tedy o schéma prvního řádu přesnosti v prostoru i v čase. Chyba aproximace $e_j^n = U_j^n - u_j^n$ v případě schématu (29) splňuje

$$(47) \quad e_j^{n+1} = (1 + \nu) e_j^n - \nu e_{j+1}^n - \tau \varepsilon_j^n.$$

Schéma (29) splňuje předpoklady použité k důkazu diskrétního principu maxima právě tehdy, když $1 + \nu \geq 0$ a $-\nu \geq 0$, tj. $\nu \in [-1, 0]$. Tato podmínka zaručuje platnost CFL

podmínky i stabilitu ve smyslu von Neumannovy analýzy a její platnost budeme nadále předpokládat. Uvažujme rovnici (26) pouze na časovém intervalu $(0, T]$, kde $T > 0$ je daný pevně zvolený čas, tj. uvažujme úlohu

$$u_t + a u_x = 0 \quad \text{v } \mathbb{R} \times (0, T], \quad u(x, 0) = u^0(x) \quad \text{pro } x \in \mathbb{R}.$$

Schéma (29) tedy uvažujeme pro $n = 0, \dots, N_T - 1$, kde N_T je dolní celá část čísla T/τ . Pro $n = 0$ splňuje přibližné řešení vztah (39). Chyba aproximace pak splňuje rovnici (47) a $e_j^0 = 0$ pro $j \in \mathbb{Z}$. Předpokládejme, že druhé derivace u_{xx} a u_{tt} jsou v množině $\mathbb{R} \times [0, T]$ omezeny v absolutní hodnotě konstantou K a položme $D = K(\tau + ah)/2$. Pak podle (46) je $|\varepsilon_j^n| \leq D$ pro $j \in \mathbb{Z}$ a $n = 0, \dots, N_T - 1$. Z (47) tudíž dostáváme

$$\sup_{j \in \mathbb{Z}} |e_j^{n+1}| \leq |1 + \nu| \sup_{j \in \mathbb{Z}} |e_j^n| + |\nu| \sup_{j \in \mathbb{Z}} |e_j^n| + \tau D \quad \forall n = 0, \dots, N_T - 1.$$

Poněvadž $|1 + \nu| = 1 + \nu$ a $|\nu| = -\nu$, plyne odtud nerovnost

$$(48) \quad \sup_{j \in \mathbb{Z}} |e_j^{n+1}| \leq \tau D + \sup_{j \in \mathbb{Z}} |e_j^n| \quad \forall n = 0, \dots, N_T - 1.$$

Jelikož $e_j^0 = 0$ pro $j \in \mathbb{Z}$, dostáváme $\sup_{j \in \mathbb{Z}} |e_j^n| \leq n \tau D$ pro $n = 0, \dots, N_T$, a tudíž

$$|U_j^n - u_j^n| \leq \frac{1}{2} K T (\tau + ah) \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T.$$

To je požadovaný odhad chyby aproximace, který ukazuje, že chyba přibližného řešení konverguje stejnoměrně (dokonce lineárně v h a τ) do nuly pro $h, \tau \rightarrow 0$.

Z uvedeného odvození není zřejmé, jak odhad chyby aproximace souvisí s diskrétním principem maxima. Ukažme si proto nyní formální odvození vztahu (48), které platnost diskrétního principu maxima explicitně využívá. Nejprve definujme tzv. *srovnávací funkci* $\Phi_j^n = t_n D$ pro $j \in \mathbb{Z}$ a $n = 0, \dots, N_T$. Pak

$$\Phi_j^{n+1} = (1 + \nu) \Phi_j^n - \nu \Phi_{j+1}^n + \tau D \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T - 1.$$

S přihlédnutím k (47) zjišťujeme, že síťová funkce $V_j^n = e_j^n - \Phi_j^n$ splňuje

$$V_j^{n+1} \leq (1 + \nu) V_j^n - \nu V_{j+1}^n \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T - 1.$$

Podobně síťová funkce $W_j^n = -e_j^n - \Phi_j^n$ splňuje

$$W_j^{n+1} \leq (1 + \nu) W_j^n - \nu W_{j+1}^n \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T - 1.$$

Jelikož jsou pro $\nu \in [-1, 0]$ splněny předpoklady pro platnost diskrétního principu maxima formulované v části 3.3 a předchozí dvě nerovnosti odpovídají nerovnosti (42), platí podle věty 4 nerovnost (43), tj.

$$V_j^{n+1} \leq V_{\max}^n, \quad W_j^{n+1} \leq W_{\max}^n \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T - 1.$$

První z těchto nerovností implikuje, že pro libovolné $j \in \mathbb{Z}$ a $n = 0, \dots, N_T - 1$

$$\begin{aligned} e_j^{n+1} &\leq t_{n+1} D + V_{\max}^n = t_{n+1} D + \sup_{j \in \mathbb{Z}} (e_j^n - t_n D) = \tau D + \sup_{j \in \mathbb{Z}} e_j^n \\ &\leq \tau D + \sup_{j \in \mathbb{Z}} |e_j^n|. \end{aligned}$$

Podobně druhá nerovnost implikuje, že

$$\begin{aligned} -e_j^{n+1} &\leq t_{n+1} D + W_{\max}^n = t_{n+1} D + \sup_{j \in \mathbb{Z}} (-e_j^n - t_n D) = \tau D + \sup_{j \in \mathbb{Z}} (-e_j^n) \\ &\leq \tau D + \sup_{j \in \mathbb{Z}} |e_j^n|. \end{aligned}$$

Poslední dvě nerovnosti implikují, že

$$|e_j^{n+1}| \leq \tau D + \sup_{j \in \mathbb{Z}} |e_j^n| \quad \forall j \in \mathbb{Z}, \quad n = 0, \dots, N_T - 1,$$

z čehož plyne nerovnost (48).

Poznamenejme nakonec, že analýza založená na chybě diskretizace a diskrétním principu maxima má u úloh pro parciální diferenciální rovnice prvního řádu jen omezenou použitelnost, neboť příslušné derivace obvykle neexistují v celé výpočetní oblasti. Uvedený postup je tak užitečný pouze pro lokální analýzu.

3.4 Numerické okrajové podmínky

Dosud se naše teoretické úvahy týkaly pouze numerického řešení Cauchyových úloh. V praxi ovšem obvykle řešíme parciální diferenciální rovnice na omezených prostorových oblastech. Na jejich hranicích (či jejich částech) pak obvykle předepisujeme okrajové podmínky. Může se však stát, že použité numerické schéma vyžaduje okrajovou podmínku na části hranice, kde parciální diferenciální rovnice žádnou okrajovou podmínku nevyžaduje. V takovém případě je nutno předepsat tzv. *numerickou okrajovou podmínku*.

Uvažujme rovnici (26) na omezeném prostorovém intervalu (A, B) . Je-li $a > 0$ konstantní, pak charakteristiky vycházející z bodů ležících na úsečce $(A, B) \times \{0\}$ protínají polopřímku $\{B\} \times \mathbb{R}^+$, avšak nikoliv polopřímku $\{A\} \times \mathbb{R}^+$. Na polopřímce $\{A\} \times \mathbb{R}^+$ je proto nutno předepsat okrajovou podmínku, zatímco na polopřímce $\{B\} \times \mathbb{R}^+$ je řešení u určeno počáteční podmínkou a okrajovou podmínkou na $\{A\} \times \mathbb{R}^+$. Je-li $a < 0$ konstantní, pak je situace opačná: okrajovou podmínku je třeba předepsat na polopřímce $\{B\} \times \mathbb{R}^+$, zatímco na polopřímce $\{A\} \times \mathbb{R}^+$ je řešení u určeno počáteční podmínkou a okrajovou podmínkou na $\{B\} \times \mathbb{R}^+$. Pokud a není konstantní, může např. každý bod polopřímek $\{A\} \times \mathbb{R}^+$ a $\{B\} \times \mathbb{R}^+$ ležet na nějaké charakteristice vycházející z bodu na $(A, B) \times \{0\}$, a tudíž žádnou okrajovou podmínku nelze předepsat. Nebo naopak žádná charakteristika zmíněné polopřímky neprotíná a pak na obou polopřímkách je potřeba předepsat okrajové podmínky.

Výpočetní oblast $[A, B] \times \mathbb{R}_0^+$ pokryjme rovnoměrnou sítí s prostorovým krokem h a časovým krokem τ . Předpokládáme, že $h = |B - A|/J$, kde $J \in \mathbb{N}$. Pak síť obsahuje uzly

(x_j, t_n) , kde $x_j = A + jh$, $t_n = n\tau$, $j = 0, \dots, J$, $n \in \mathbb{N}_0$. V uzlech (x_j, t_n) s $j = 1, \dots, J-1$ a $n \in \mathbb{N}_0$ uvažujme Laxovo–Friedrichsovo schéma (33). Toto schéma vyžaduje hodnoty U_0^n a U_J^n , bez ohledu na to, zda rovnice (26) umožňuje tyto hodnoty předepsat. Je-li např. $a > 0$ konstatní, pak, jak jsme již uvedli, rovnice (26) neumožňuje předepsat okrajovou podmínku podél polopřímky $\{B\} \times \mathbb{R}^+$, a je tedy nutno předepsat numerickou okrajovou podmínku v uzlech (x_J, t_n) . Příklady numerických okrajových podmínek v uzlech (x_J, t_n) jsou následující vztahy:

$$(49) \quad U_J^n = U_{J-1}^n,$$

$$(50) \quad U_J^n = 2U_{J-1}^n - U_{J-2}^n,$$

$$(51) \quad U_J^n = U_{J-1}^{n-1},$$

$$(52) \quad U_J^n = 2U_{J-1}^{n-1} - U_{J-2}^{n-2}.$$

Vztahy (49) a (50) jsou jednoduché extrapolace přibližného řešení z vnitřních uzlů na hranici. Vztahy (51) a (52) se někdy nazývají kvazicharakteristické extrapolace, neboť využívají hodnot přibližného řešení v uzlech ležících na přímkách procházejících uzly (x_{J-1}, t_{n-1}) a (x_J, t_n) , jejichž směr přibližně odpovídá směru charakteristik.

Alternativou k výše uvedenému postupu je zavést fiktivní uzly (x_{-1}, t_n) a (x_{J+1}, t_n) , do nichž extrapolujeme přibližné řešení z původních uzlů sítě, a ve všech původních uzlech sítě (tj. i hraničních) použít schéma (33). Použijeme-li k extrapolaci vztah (50), pak získáme

$$\frac{U_J^{n+1} - \frac{1}{2}(U_{J+1}^n + U_{J-1}^n)}{\tau} + a \frac{U_{J+1}^n - U_{J-1}^n}{2h} = 0, \quad U_{J+1}^n = 2U_J^n - U_{J-1}^n,$$

z čehož plyne

$$\frac{U_J^{n+1} - U_J^n}{\tau} + a \frac{U_J^n - U_{J-1}^n}{h} = 0.$$

Vidíme tedy, že extrapolace do fiktivních uzlů je v tomto případě ekvivalentní aproximaci diferenciální rovnice v hraničních uzlech pomocí jednostranných diferencí, což je další typ numerické okrajové podmínky. Často je však jednodušší použít některý z extrapoláčnických vztahů (49)–(52) než zavádět fiktivní uzly nebo jednostranné diference.

Je důležité zdůraznit, že některé numerické okrajové podmínky mohou pro dané schéma způsobit nestabilní chování, které se projeví narůstajícími oscilacemi v přibližném řešení. Přitom numerická okrajová podmínka, kterou lze použít s daným numerickým schématem, může být s jiným schématem nestabilní a naopak. Pro každou zvolenou kombinaci schématu a numerické okrajové podmínky je proto potřeba stabilitu vyšetřovat zvlášť.

3.5 Laxovo–Wendroffovo schéma

Z jednokrokových schémat, která jsme dosud analyzovali, jsou pro řešení rovnice (26) použitelná (při libovolném znaménku a) pouze schémata (33) a (34). Obě tato schémata jsou pouze prvního řádu přesnosti (v případě schématu (33) za podmínky $h = O(\tau)$). Naším cílem je nyní odvodit schéma druhého řádu přesnosti v prostoru i v čase. Pro

jednoduchost budeme opět uvažovat konstantní rychlost a . Bude však poučné schéma odvodit pro nenulovou pravou stranu $f = f(x, t)$. Řešíme tedy rovnici

$$(53) \quad u_t + a u_x = f \quad \text{v } \mathbb{R} \times \mathbb{R}^+$$

s počáteční podmínkou (27).

Předpokládejme, že $u \in C^4(\mathbb{R} \times \mathbb{R}_0^+)$. Pak pro $(x, t) \in \mathbb{R} \times \mathbb{R}_0^+$ je

$$(54) \quad u(x, t + \tau) = u(x, t) + u_t(x, t) \tau + \frac{1}{2} u_{tt}(x, t) \tau^2 + O(\tau^3).$$

Derivace u_t a u_{tt} získáme z rovnice (53). V případě u_t je to triviální. Abychom získali u_{tt} , zderivujeme rovnici (53) podle t . Tím vznikne člen u_{xt} , který získáme z rovnice (53) zderivováním podle x . Máme tedy

$$u_{tt} = f_t - a u_{xt} = f_t - a(f_x - a u_{xx}) = a^2 u_{xx} - a f_x + f_t.$$

Dosazením do (54) obdržíme

$$u(x, t + \tau) = u - a \tau u_x + \tau f + \frac{a^2 \tau^2}{2} u_{xx} - \frac{a \tau^2}{2} f_x + \frac{\tau^2}{2} f_t + O(\tau^3),$$

kde členy na pravé straně jsou uvažovány v bodě (x, t) . Jelikož

$$u_x = \frac{\Delta_{0x} u}{h} + O(h^2), \quad u_{xx} = \frac{\delta_x^2 u}{h^2} + O(h^2),$$

dostáváme

$$\Delta_{+t} u = -\nu \Delta_{0x} u + O(\tau h^2) + \frac{\nu^2}{2} \delta_x^2 u + \frac{\tau}{2} [f + f(x, t + \tau)] - \frac{a \tau^2}{2h} \Delta_{0x} f + O(\tau^3).$$

To nás vede ke schématu

$$(55) \quad \frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} - \frac{a^2 \tau}{2} \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} \\ = \frac{1}{2} (f_j^{n+1} + f_j^n) - \frac{a \tau}{4h} (f_{j+1}^n - f_{j-1}^n),$$

které můžeme též zapsat ve tvaru

$$U_j^{n+1} = U_j^n - \frac{\nu}{2} (U_{j+1}^n - U_{j-1}^n) + \frac{\nu^2}{2} (U_{j+1}^n - 2U_j^n + U_{j-1}^n) \\ + \frac{\tau}{2} (f_j^{n+1} + f_j^n) - \frac{\nu \tau}{4} (f_{j+1}^n - f_{j-1}^n).$$

Z odvození je zřejmé, že chyba diskretizace splňuje $\varepsilon_{h,\tau} = O(h^2 + \tau^2)$, tj. skutečně jsme získali schéma druhého řádu přesnosti v prostoru i v čase. Pro $f \equiv 0$ je Laxovo-Wendroffovo schéma ekvivalentní přidání umělé difúze o velikosti $\frac{1}{2} a^2 \tau$ k nestabilnímu

schématu (31). V tomto případě můžeme Laxovo–Wendroffovo schéma odvodit též pomocí charakteristik, podobně jako schéma (34), viz obr. 3. Přidáme-li v obr. 3 na přímce $t = t_n$ uzel C tak, že B je střed úsečky AC , a definujeme-li v bodě Q aproximaci u pomocí kvadratické interpolace hodnot přibližného řešení v uzlech A , B a C , získáme Laxovo–Wendroffovo schéma.

Von Neumannovu analýzu provádíme pro $f \equiv 0$. Získáme

$$\lambda(\xi) = 1 - 2\nu^2 \sin^2 \frac{\xi h}{2} - i\nu \sin(\xi h), \quad |\lambda(\xi)|^2 = 1 - 4\nu^2 (1 - \nu^2) \sin^4 \frac{\xi h}{2},$$

a schéma je tedy stabilní pro $|\nu| \leq 1$, přičemž celý tento interval je povolen CFL podmínkou. Fázová rychlost splňuje

$$\alpha(\xi) = a \left[1 - \frac{1}{6} (1 - \nu^2) (\xi h)^2 + O((\xi h)^4) \right]$$

a při $|\nu| < 1$ tedy dochází vždy ke zpoždování fáze (pro malé $|\xi h|$). Pro $|\nu| = 1$ je $\lambda(\xi) = e^{-i\xi a\tau}$, a schéma tudíž dává přesné řešení. Fázová chyba je stejného řádu jako u schémat (33) a (34), avšak tlumení je podstatně menší (v jednom časovém kroku je chyba v amplitudě řádu $(\xi h)^4$ místo $(\xi h)^2$). Pro nejvyšší frekvenci na síti ($\xi = \pi/h$) je $\lambda(\pi/h) = 1 - 2\nu^2$, a tedy $\alpha(\pi/h) = 0$, což znamená, že tato rychle oscilující složka řešení se nepohybuje pryč. Při $|\nu| < 1$ je však tato složka řešení tlumena: např. počáteční podmínce $U_j^0 = (-1)^j$ odpovídá řešení $U_j^n = (1 - 2\nu^2)^n (-1)^j$. Nedostatkem Laxova–Wendroffova schématu je, že na rozdíl od schémat (33) a (34) nesplňuje diskrétní princip maxima.

Cvičení 2 *Proveďte výpočet ze cvičení 1 pro Laxovo–Wendroffovo schéma a srovnajte výsledek s přibližným řešením získaným pomocí Laxova–Friedrichsova schématu (33).*

Cvičení 3 *Uvažujme rovnici $u_t + u_x = 0$ v oblasti $(-1, 1) \times \mathbb{R}^+$ s počáteční podmínkou $u^0(x) = \sin(2\pi x)$. Definujme ekvidistantní prostorové uzly $-1 = x_0 < x_1 < \dots < x_J = 1$ s prostorovým krokem $h = 2/J$. Dále definujme uzel $x_{J+1} = 1 + h$ ležící mimo interval $[-1, 1]$. Zvolme $J = 20$ a uvažujme konstantní časový krok $\tau = 0.09$. V čase $t = 0$ je přibližné řešení určeno počáteční podmínkou. Na následujících časových hladinách $t_n = n\tau$, $n \in \mathbb{N}$ počítejte v uzlech x_1, \dots, x_J přibližná řešení pomocí schémat (33), (34) a (55). Pak položte $U_0^n = U_J^n$ a $U_{J+1}^n = U_1^n$. Uvažujme tedy periodické okrajové podmínky. U zmíněných schémat porovnejte fázovou chybu a tlumení a srovnajte výsledky provedených numerických experimentů s teoretickými výsledky uvedenými výše.*

3.6 Schéma Crankovo–Nicolsonové

Schéma Crankovo–Nicolsonové jsme v úvodní části o metodě konečných diferencí definovali jako aritmetický průměr explicitního a implicitního schématu. Nyní si ukážeme jiný způsob, jak toto schéma odvodit. Naše motivace bude opět získat jednokrokové schéma pro rovnici (53), které je druhého řádu v prostoru i v čase.

Odvození bude založeno na tom, že na diferenci $u(x, t + \tau) - u(x, t)$ můžeme pohlížet jako na centrální diferenci vzhledem k času $t + \frac{\tau}{2}$. Je tedy

$$\frac{u(x, t + \tau) - u(x, t)}{\tau} = \frac{\delta_t u(x, t + \frac{\tau}{2})}{\tau} = u_t(x, t + \frac{\tau}{2}) + O(\tau^2).$$

Dále platí

$$\frac{u(x, t + \tau) + u(x, t)}{2} = u(x, t + \frac{\tau}{2}) + \frac{u(x, t + \tau) - 2u(x, t + \frac{\tau}{2}) + u(x, t)}{2}.$$

Čitatel zlomku na pravé straně představuje centrální diferenci druhého řádu vzhledem k t v bodě $(x, t + \frac{\tau}{2})$ s krokem $\frac{\tau}{2}$. Jelikož $\delta_t^2 u = O(\tau^2)$, dostáváme

$$\frac{u(x, t + \tau) + u(x, t)}{2} = u(x, t + \frac{\tau}{2}) + O(\tau^2).$$

Z toho plyne, že

$$\begin{aligned} u_x(x, t + \frac{\tau}{2}) &= \frac{1}{2} [u_x(x, t + \tau) + u_x(x, t)] + O(\tau^2) \\ &= \frac{1}{2} \left[\frac{u(x + h, t + \tau) - u(x - h, t + \tau)}{2h} + \frac{u(x + h, t) - u(x - h, t)}{2h} \right] + O(h^2 + \tau^2). \end{aligned}$$

To ukazuje, že v bodě $(x, t + \frac{\tau}{2})$ může být vhodná rovnice (53) aproximovat schématem

$$\frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_{j+1}^{n+1} - U_{j-1}^{n+1} + U_{j+1}^n - U_{j-1}^n}{4h} = \frac{f_j^{n+1} + f_j^n}{2},$$

což lze psát též ve tvaru

$$-\frac{\nu}{4} U_{j-1}^{n+1} + U_j^{n+1} + \frac{\nu}{4} U_{j+1}^{n+1} = \frac{\nu}{4} U_{j-1}^n + U_j^n - \frac{\nu}{4} U_{j+1}^n + \frac{\tau}{2} (f_j^{n+1} + f_j^n).$$

Nyní již není možné hodnotu přibližného řešení v daném uzlu na nové časové vrstvě jednoduchým způsobem vyjádřit pomocí hodnot přibližného řešení na předcházející časové vrstvě a jedná se tedy o implicitní schéma. Hodnoty přibližného řešení na nové časové vrstvě jsou navzájem svázány a získají se vyřešením soustavy lineárních rovnic. Jelikož hodnota přibližného řešení v daném uzlu na nové časové vrstvě závisí na všech hodnotách přibližného řešení z předcházející časové vrstvy, nevede CFL podmínka k žádné podmínce na parametr ν .

Von Neumannova analýza (pro $f \equiv 0$) dává amplifikační faktor

$$\lambda(\xi) = \frac{1 - i \frac{\nu}{2} \sin(\xi h)}{1 + i \frac{\nu}{2} \sin(\xi h)}.$$

Je tedy $|\lambda(\xi)| = 1$ pro všechna $\xi \in \mathbb{R}$, a tudíž schéma je nepodmíněně stabilní. Nedochozí však k žádnému tlumení, takže drobné poruchy v řešení zůstávají trvale přítomny. Speciálně pro vlnu s nejvyšší frekvencí π/h opět platí, že fázová rychlost je nulová, avšak na rozdíl od Laxova–Wendroffova schématu se její amplituda v čase nemění: počáteční podmínce $U_j^0 = (-1)^j$ odpovídá řešení $U_j^n = (-1)^j$. Stejně jako u Laxova–Wendroffova schématu neplatí ani pro schéma Crankovo–Nicolsonové diskrétní princip maxima.

3.7 Analýza leapfrog scheme

Názvem *leapfrog scheme* je označováno dvoukrokové schéma (32). Je zřejmé, že se jedná o schéma druhého řádu přesnosti v prostoru i v čase, které nesplňuje diskrétní princip maxima. CFL podmínka je splněna, jestliže $|\nu| \leq 1$, což budeme nadále předpokládat.

Zabývejme se nyní stabilitou leapfrog scheme. Stejně jako v případě jednokrokových schémat vyjádříme přibližné řešení U_j^n v integrálním tvaru pomocí jeho Fourierovy transformace $\widehat{U}^n(\xi)$ a dosadíme do schématu. Tím získáme vztah

$$\int_{-\pi/h}^{\pi/h} e^{i\xi j h} \left(\widehat{U}^{n+1}(\xi) + i 2\nu \sin(\xi h) \widehat{U}^n(\xi) - \widehat{U}^{n-1}(\xi) \right) d\xi = 0 \quad \forall j \in \mathbb{Z}, n \in \mathbb{N},$$

z něhož díky ortogonalitě funkcí $\{e^{i\xi j h}\}_{j \in \mathbb{Z}}$ plyne

$$(56) \quad \widehat{U}^{n+1}(\xi) + i 2\nu \sin(\xi h) \widehat{U}^n(\xi) - \widehat{U}^{n-1}(\xi) = 0 \quad \forall n \in \mathbb{N}, \xi \in [-\pi/h, \pi/h].$$

Pro každé $\xi \in [-\pi/h, \pi/h]$ tedy máme soustavu lineárních diferenčních rovnic s konstantními koeficienty. Charakteristická rovnice této soustavy diferenčních rovnic je

$$\lambda^2 + i 2\nu \sin(\xi h) \lambda - 1 = 0.$$

Její kořeny jsou

$$\lambda_{\pm}(\xi) = -i\nu \sin(\xi h) \pm \sqrt{1 - \nu^2 \sin^2(\xi h)}$$

a platí $|\lambda_{\pm}(\xi)| = 1$. Pokud $\lambda_+(\xi) = \lambda_-(\xi) =: \lambda(\xi)$, pak obecné řešení soustavy diferenčních rovnic (56) má tvar

$$\widehat{U}^n(\xi) = \alpha(\xi) \lambda(\xi)^n + \beta(\xi) n \lambda(\xi)^{n-1},$$

kde $\alpha(\xi)$ a $\beta(\xi)$ jsou konstanty určené hodnotami $\widehat{U}^0(\xi)$ a $\widehat{U}^1(\xi)$. Je-li $\beta(\xi) \neq 0$, bude $\widehat{U}^n(\xi)$ růst lineárně s n . Toto nestabilní chování je třeba vyloučit. Jelikož $\lambda_+(\xi) = \lambda_-(\xi)$ může nastat pouze při $|\nu| = 1$, budeme požadovat, aby $|\nu| < 1$. Pak $\lambda_+(\xi) \neq \lambda_-(\xi)$ a obecné řešení soustavy diferenčních rovnic (56) je

$$(57) \quad \widehat{U}^n(\xi) = \alpha_+(\xi) \lambda_+(\xi)^n + \alpha_-(\xi) \lambda_-(\xi)^n.$$

Snadno zjistíme, že schéma je v tomto případě stabilní. Nutná a postačující podmínka pro stabilitu leapfrog scheme tedy je $|\nu| < 1$.

Leapfrog scheme (32) nelze použít k určení přibližného řešení na časové vrstvě t_1 a schéma je tedy potřeba inicializovat pomocí vhodného jednokrokového schématu. Lze ukázat, že k inicializaci vícekrokových schémat lze použít libovolné jednokrokové schéma, které je konzistentní s řešenou parciální diferenciální rovnicí. Toto schéma může být i nestabilní, neboť aplikací schématu pouze v několika prvních krocích dojde jen k malému růstu řešení (malému díky konzistenci) a tento růst nebude při použití stabilního vícekrokového schématu dále zesilován. Jeli $\tau/h = \text{const.}$, pak inicializační schéma může mít o 1 nižší řád přesnosti než vícekrokové schéma, přičemž celková přesnost bude odpovídat vícekrokovému schématu.

Ze vztahu (57) plyne

$$\widehat{U}^n(\xi) = \widehat{U}^0(\xi) \lambda_+(\xi)^n + \gamma(\xi) \frac{\lambda_+(\xi)^n - \lambda_-(\xi)^n}{\lambda_+(\xi) - \lambda_-(\xi)},$$

kde $\gamma(\xi) = \widehat{U}^1(\xi) - \widehat{U}^0(\xi) \lambda_+(\xi)$. Použijeme-li k inicializaci schéma (31), je $\widehat{U}^1(\xi) = (1 - i\nu \sin(\xi h)) \widehat{U}^0(\xi)$. Použitím vztahu $\sqrt{1-x} = 1 - \frac{1}{2}x + O(x^2)$ získáme

$$\lambda_{\pm}(\xi) = -i\nu \sin(\xi h) \pm 1 \mp \frac{1}{2}\nu^2 \sin^2(\xi h) + O((\xi h)^4),$$

a tudíž

$$\gamma(\xi) = (1 - i\nu \sin(\xi h) - \lambda_+(\xi)) \widehat{U}^0(\xi) = \left(\frac{\nu^2}{2} \sin^2(\xi h) + O((\xi h)^4) \right) \widehat{U}^0(\xi).$$

Vidíme, že pro malé hodnoty $|\xi h|$ je $\gamma(\xi)$ malé, a schéma se tudíž chová jako jednokrokové schéma s amplifikačním faktorem λ_+ . Pro větší hodnoty $|\xi h|$ velikost $\gamma(\xi)$ být malá nemusí. Část řešení příslušná λ_- se nazývá parazitická složka řešení. Jelikož $\lambda_-(0) = -1$, tato složka řešení v průběhu času rychle osciluje. Později ukážeme, že se parazitická složka řešení pohybuje špatným směrem (pro $a > 0$ se pohybuje zprava doleva).

Cvičení 4 *Uvažujte ve cvičení 1 okrajovou podmínku $U_j^n = 0$, která není konzistentní s řešenou rovnicí, a proveďte výpočet pro leapfrog scheme. Měli byste pozorovat, že v důsledku zmíněné okrajové podmínky vznikne v bodě $x = 3$ parazitická složka řešení, která je silně oscilující a pohybuje se zprava doleva. Okrajová podmínka v bodě $x = -2$ přemění parazitickou složku na neparazitickou.*

Parazitická složka řešení se objeví při každém výpočtu s víceukrokovými schémata. Obvykle nezpůsobuje podstatné potíže, avšak v některých případech je ji potřeba redukovat nebo odstranit. Vliv parazitických složek lze redukovat pomocí disipace.

Definice 3 *Schéma je disipativní řádu $2r$, jestliže amplifikační faktory splňují podmínku*

$$|\lambda_j(\xi)| \leq 1 - C \left(\sin \frac{\xi h}{2} \right)^{2r} \quad \forall \xi \in \mathbb{R},$$

kde C je kladná konstanta nezávislá na h a τ .

Pro leapfrog scheme a schéma Crankovo–Nicolsonové je velikost amplifikačních faktorů rovna 1 pro všechny frekvence. Říkáme, že tato schémata jsou *ostře nedisipativní*. Pro Laxovo–Friedrichsovo schéma je $|\lambda(\pi/h)| = 1$, avšak $|\lambda(\xi)| < 1$ pro $0 < |\xi| < \pi/h$ a $|\nu| < 1$. Toto schéma je proto nedisipativní, ale ne ostře. Redukuje amplitudu u většiny frekvencí, nikoli však u nejvyšší frekvence na síti.

Všimněme si, že síťová funkce $U_j^n = (-1)^{n+j} \eta$ je pro libovolné $\eta \in \mathbb{R}$ řešením leapfrog scheme. Vidíme tedy, že počáteční poruchy se šíří, aniž by byly tlumeny. O přítomnosti této šachovnicové složky řešení rozhodují pouze počáteční a okrajové podmínky.

Vlastnosti nedisipativního schématu lze zlepšit přidáním disipace. Je však třeba dát pozor, abychom tím nezhoršili řád přesnosti schématu. Např. leapfrog scheme lze modifikovat takto:

$$\frac{U_j^{n+1} - U_j^{n-1}}{2\tau} + a \frac{U_{j+1}^n - U_{j-1}^n}{2h} + \frac{\varepsilon}{2\tau} \left(\frac{1}{2} \delta_x\right)^4 U_j^{n-1} = 0,$$

kde ε je vhodná kladná konstanta. Jelikož $\delta_x^4 u = O(h^4)$, schéma je nadále druhého řádu přesnosti, pokud $\tau \geq C h^2$. Tato podmínka nepředstavuje žádné podstatné omezení, neboť časový krok volíme obvykle srovnatelný s prostorovým krokem. Pro vyšetření stability přepíšeme schéma do tvaru

$$U_j^{n+1} - U_j^{n-1} + \nu U_{j+1}^n - \nu U_{j-1}^n + \varepsilon \left(\frac{1}{2} \delta_x\right)^4 U_j^{n-1} = 0$$

a přibližné řešení vyjádříme v integrálním tvaru pomocí jeho Fourierovy transformace $\widehat{U}^n(\xi)$. Oproti případu bez disipace se navíc objeví člen $\varepsilon \left(\frac{1}{2} \delta_x\right)^4 \widehat{U}^{n-1}(\xi) e^{i\xi j h}$. Jak víme, $\delta_x^2 e^{i\xi j h} = -4 \sin^2 \frac{\xi h}{2} e^{i\xi j h}$, a tudíž $\left(\frac{1}{2} \delta_x\right)^4 e^{i\xi j h} = \sin^4 \frac{\xi h}{2} e^{i\xi j h}$. Místo soustavy diferenciálních rovnic (56), proto nyní získáme

$$\widehat{U}^{n+1}(\xi) + i 2\nu \sin(\xi h) \widehat{U}^n(\xi) - \left(1 - \varepsilon \sin^4 \frac{\xi h}{2}\right) \widehat{U}^{n-1}(\xi) = 0.$$

Kořeny odpovídající charakteristické rovnice jsou

$$\lambda_{\pm}(\xi) = -i\nu \sin(\xi h) \pm \sqrt{1 - \varepsilon \sin^4 \frac{\xi h}{2} - \nu^2 \sin^2(\xi h)}.$$

Je-li

$$(58) \quad \nu^2 \sin^2(\xi h) + \varepsilon \sin^4 \frac{\xi h}{2} \leq 1,$$

je

$$|\lambda_{\pm}(\xi)| = \sqrt{1 - \varepsilon \sin^4 \frac{\xi h}{2}} \leq 1 - \frac{\varepsilon}{2} \sin^4 \frac{\xi h}{2},$$

a schéma je tedy stabilní a disipativní řádu 4 (k tomu je opět nutné, aby $|\nu| < 1$). Podmínku (58) můžeme přepsat do tvaru

$$4\nu^2 \sin^2 \frac{\xi h}{2} \left(1 - \sin^2 \frac{\xi h}{2}\right) + \varepsilon \sin^4 \frac{\xi h}{2} \leq 1.$$

Stačí tedy nalézt takové hodnoty ε , že

$$4\nu^2 s(1-s) + \varepsilon s^2 \leq 1 \quad \forall s \in [0, 1].$$

Je-li $\nu^2 \in (0, \frac{1}{2}]$, můžeme volit $\varepsilon \in (0, 1]$, neboť pak $4\nu^2 s(1-s) + \varepsilon s^2 \leq 2s(1-s) + s^2 = s(2-s) \leq 1$. Je-li $\nu^2 \in [\frac{1}{2}, 1)$, volíme $\varepsilon \in (0, 4\nu^2(1-\nu^2)]$. Pak $4\nu^2 s(1-s) + \varepsilon s^2 = 4\nu^2 s + (\varepsilon - 4\nu^2) s^2 \leq 4\nu^2 s(1-\nu^2 s) \leq 1$.

Podobně jako výše můžeme modifikovat i metodu Crankovu–Nicolsonové pro řešení rovnice (53). Obdržíme schéma

$$\frac{U_j^{n+1} - U_j^n}{\tau} + a \frac{U_{j+1}^{n+1} - U_{j-1}^{n+1} + U_{j+1}^n - U_{j-1}^n}{4h} + \frac{\varepsilon}{\tau} \left(\frac{1}{2} \delta_x \right)^4 U_j^n = \frac{f_j^{n+1} + f_j^n}{2},$$

které je druhého řádu přesnosti (pokud $\tau \geq C h^2$) a disipativní řádu 4 pro malé hodnoty ε .

Cvičení 5 Zopakujte výpočet ze cvičení 4 pro leapfrog scheme s disipací (v prostorových uzlech x_1 a x_{J-1} uvažujte leapfrog scheme bez disipace).

Pro víceřadová schémata vždy existuje právě jeden amplifikační faktor $\lambda_0(\xi)$ takový, že $\lambda_0(0) = 1$. Tento amplifikační faktor použijeme k definici fázové rychlosti. Pro leapfrog scheme je to amplifikační faktor λ_+ . Platí

$$\sin(\xi \alpha(\xi) \tau) = \nu \sin(\xi h) = \nu \xi h \left(1 - \frac{(\xi h)^2}{6} + O((\xi h)^4) \right),$$

z čehož plyne

$$\begin{aligned} \alpha(\xi) &= a \left(1 - \frac{(\xi h)^2}{6} + O((\xi h)^4) \right) \left(1 + \frac{(\nu \xi h)^2}{6} + O((\xi h)^4) \right) \\ &= a \left(1 - \frac{1}{6} (1 - \nu^2) (\xi h)^2 + O((\xi h)^4) \right). \end{aligned}$$

Získali jsme tedy totéž jako pro Laxovo-Wendroffovo schéma a rovněž platí $\alpha(\pi/h) = 0$.

Pro parazitickou složku řešení leapfrog scheme odpovídající λ_- zavedeme fázovou rychlost vztahem $\lambda_-(\xi) = -|\lambda_-(\xi)| e^{-i\xi \alpha_-(\xi) \tau}$, neboť $\lambda_-(0) = -1$. Pak $\sin(\xi \alpha_-(\xi) \tau) = -\nu \sin(\xi h)$, a tudíž $\alpha_-(\xi) = -\alpha(\xi)$. Parazitická složka řešení se tedy pohybuje opačným směrem než je správný směr řešení. Navíc platí $\alpha_-(\pi/h) = 0$, a tudíž nejvyšší frekvence, které neobsahují přesnou informaci, se nešíří pryč.

3.8 Volba parametru ν

Vlastnosti schémat vyšetřovaných v předchozích odstavcích závisí na volbě parametru ν , tj. na poměru prostorového a časového kroku. V řadě případů jsme viděli, že pro $|\nu| = 1$ schémata dávají přesné řešení, což je ovšem dáno jednoduchostí uvažované modelové úlohy. Pro úlohy s nekonstantními koeficienty nebo komplikovanější problémy k tomu již nedochází. Nicméně naše teoretické úvahy ukazují, že je obecně vhodné volit $|\nu|$ blízko hranice stability, neboť tím dosáhneme malé disipace i disperze. Zajímá-li nás pouze určitá frekvence ξ_0 , měli bychom volit h tak, aby $|\xi_0 h| \ll \pi$ a vlna o frekvenci ξ_0 mohla být tudíž na použité síti dobře aproximována.

4 Numerické řešení rovnice vedení tepla

V této části se budeme zabývat numerickým řešením smíšené úlohy pro rovnici vedení tepla v jedné prostorové dimenzi. Hledáme funkci $u = u(x, t)$ definovanou pro $x \in [0, 1]$ a $t \geq 0$ takovou, že

$$(59) \quad u_t - u_{xx} = 0 \quad \text{v } (0, 1) \times \mathbb{R}^+,$$

$$(60) \quad u(0, t) = u(1, t) = 0 \quad \forall t > 0,$$

$$(61) \quad u(x, 0) = u^0(x) \quad \forall x \in [0, 1].$$

Budeme předpokládat, že úloha (59)–(61) má klasické řešení. To mimo jiné vyžaduje, aby počáteční a okrajové podmínky splňovaly podmínky kompatibility $u^0(0) = u^0(1) = 0$.

Úlohy uvedeného typu popisují šíření tepla v tenké izolované homogenní tyči konečné délky bez přítomnosti tepelných zdrojů. Úloha typu (59)–(61) též popisuje šíření tepla napříč nekonečnou deskou, přičemž x je souřadnice kolmá ke stěnám desky, na každé z nichž je předepsána konstantní teplota.

Poznámka 3 Transformací $v(x, t) = u(x, \kappa t) + \alpha(1 - x) + \beta x$ získáme úlohu s tepelnou difuzivitou κ a nehomogenními Dirichletovými okrajovými podmínkami:

$$v_t - \kappa v_{xx} = 0 \quad \text{v } (0, 1) \times \mathbb{R}^+,$$

$$v(0, t) = \alpha, \quad v(1, t) = \beta \quad \forall t > 0,$$

$$v(x, 0) = v^0(x) \equiv u^0(x) + \alpha(1 - x) + \beta x \quad \forall x \in [0, 1].$$

Uvažování $\kappa = 1$ a homogenních okrajových podmínek v (59)–(61) tedy nepředstavuje újmu na obecnosti. Podobně lze též jednoduchou transformací proměnné x přejít k úloze definované na libovolném zadaném prostorovém intervalu.

4.1 Řešení úlohy (59)–(61) Fourierovou metodou

Hledejme řešení úlohy (59)–(61) ve tvaru $u(x, t) = X(x)T(t)$. Dosazením do (59) a (60) snadno zjistíme, že

$$u(x, t) = a_m e^{-(m\pi)^2 t} \sin(m\pi x)$$

pro nějaké $m \in \mathbb{N}$ a $a_m \in \mathbb{R}$. Rovnice (59) a (60) splňuje zřejmě i libovolná lineární kombinace těchto funkcí. Nabízí se tedy hledat řešení u ve tvaru

$$(62) \quad u(x, t) = \sum_{m=1}^{\infty} a_m e^{-(m\pi)^2 t} \sin(m\pi x).$$

Z počáteční podmínky (61) plyne

$$(63) \quad u^0(x) = \sum_{m=1}^{\infty} a_m \sin(m\pi x),$$

a tudíž koeficienty a_m jsou Fourierovy koeficienty funkce u^0 při rozvoji do sinové řady. Platí

$$a_m = 2 \int_0^1 u^0(x) \sin(m\pi x) dx.$$

Je známo, že pro libovolné $u^0 \in L^2(0, 1)$ řada (63) konverguje v $L^2(0, 1)$. Součet řady (62) je pak nekonečně hladká funkce v $[0, 1] \times \mathbb{R}^+$, která řeší parciální diferenciální rovnici (59) a splňuje okrajovou podmínku (60). Pokud řada Fourierových koeficientů a_m konverguje absolutně, je součet řady (62) spojitá funkce na $[0, 1] \times \mathbb{R}_0^+$, která splňuje počáteční podmínku (61). K tomu stačí, aby funkce u^0 byla absolutně spojitá na $[0, 1]$, $(u^0)' \in L^2(0, 1)$ a $u^0(0) = u^0(1) = 0$ (konzistence s okrajovou podmínkou).

V praxi výsledek (62) umožňuje získat pouze numerickou aproximaci řešení u , neboť koeficienty a_m jsme obecně schopni určit pouze přibližně a navíc jsme schopni sečíst pouze konečně mnoho členů řady. Skutečným omezením uvedené metody však je, že ji nelze snadno zobecnit na komplikovanější úlohy. Je proto nutné hledat jiné způsoby výpočtu přibližného řešení úlohy (59)–(61) a jedním z možných postupů je aplikace metody konečných diferencí.

4.2 Explicitní schéma pro úlohu (59)–(61)

Podobně jako dříve zavedeme rovnoměrnou síť, kterou pokryjeme výpočetní oblast. Interval $[0, 1]$ nejprve rozdělíme na $J \in \mathbb{N}$ intervalů stejné délky $h = 1/J$, čímž vzniknou body $x_j = jh$, $j = 0, 1, 2, \dots, J$. Kromě prostorového kroku sítě h zavedeme též časový krok sítě $\tau > 0$ a definujeme časové hladiny $t_n = n\tau$, $n = 0, 1, 2, \dots$. Řešení u úlohy (59)–(61) aproximujeme v uzlech sítě (x_j, t_n) hodnotami U_j^n , kde $j = 0, 1, 2, \dots, J$ a $n = 0, 1, 2, \dots$. Opět položíme $u_j^n = u(x_j, t_n)$.

Uvažujeme-li rovnici (59) v uzlu (x_j, t_n) s $j \in \{1, \dots, J-1\}$ a $n \in \mathbb{N}_0$, můžeme provést následující aproximace:

$$0 = (u_t - u_{xx})(x_j, t_n) \approx \frac{\Delta_{+t} u_j^n}{\tau} - \frac{\delta_x^2 u_j^n}{h^2} \approx \frac{\Delta_{+t} U_j^n}{\tau} - \frac{\delta_x^2 U_j^n}{h^2}.$$

To vede k diferenčním rovnicím

$$(64) \quad U_j^{n+1} = U_j^n + \mu (U_{j+1}^n - 2U_j^n + U_{j-1}^n), \quad j = 1, 2, \dots, J-1,$$

kde

$$\mu = \frac{\tau}{h^2}.$$

Každou hodnotu na časové hladině t_{n+1} lze nezávisle spočítat z hodnot na časové hladině t_n a jedná se tedy o explicitní diferenční schéma.

K rovnicím (64) musíme ještě přidat okrajové a počáteční podmínky

$$(65) \quad U_0^n = U_J^n = 0, \quad n = 1, 2, \dots,$$

$$(66) \quad U_j^0 = u^0(x_j), \quad j = 0, 1, \dots, J.$$

Hodnoty přibližného řešení získáme tak, že nejprve definujeme hodnoty U_j^0 pomocí (66) a pak evolučně počítáme hodnoty přibližného řešení na následujících časových hladinách

z rovnic (64) a okrajových podmínek (65). Přibližné řešení U_j^n je tedy vztahy (64)–(66) jednoznačně určeno.

Cvičení 6 Uvažujte schéma (64)–(66) pro úlohu (59)–(61) s počáteční podmínkou

$$u^0(x) = \begin{cases} 2x & \text{pro } x \in [0, \frac{1}{2}], \\ 2 - 2x & \text{pro } x \in [\frac{1}{2}, 1]. \end{cases}$$

Zvolte $J = 20$ a proveďte výpočet pro $\tau = 0.0012$ a $\tau = 0.0013$.

Provedeme-li výpočty z cvičení 6, zjistíme, že pro $\tau = 0.0012$ získáme dobrou aproximaci řešení úlohy (59)–(61), zatímco pro $\tau = 0.0013$ se v přibližném řešení objeví oscilace, které se zvětšujícím se n rychle rostou. Jedná se o typický příklad stability či nestability numerického schématu. Jak uvidíme později, numerické výsledky zásadně závisejí na hodnotě μ , tj. na vztahu mezi prostorovým a časovým krokem.

Chyba diskretizace schématu (64) je dána vztahem

$$\varepsilon_j^n = \frac{\Delta_{+t} u_j^n}{\tau} - \frac{\delta_x^2 u_j^n}{h^2}$$

pro $j \in \{1, \dots, J-1\}$ a $n \in \mathbb{N}_0$. Je-li $x \in [h, 1-h]$ a $t \geq 0$, můžeme též položit

$$\varepsilon_{h,\tau}(x, t) = \frac{\Delta_{+t} u(x, t)}{\tau} - \frac{\delta_x^2 u(x, t)}{h^2}.$$

Pak $\varepsilon_j^n = \varepsilon_{h,\tau}(x_j, t_n)$. Použitím Taylorova vzorce získáme

$$\varepsilon_{h,\tau}(x, t) = \frac{1}{2} u_{tt}(x, \eta) \tau - \frac{1}{12} u_{xxxx}(\xi, t) h^2, \quad \xi \in (x-h, x+h), \eta \in (t, t+\tau).$$

Předpokládáme-li, že existují konstanty M_1 a M_2 takové, že $|u_{tt}| \leq M_1$, $|u_{xxxx}| \leq M_2$ na $[0, 1] \times [0, T]$, kde $T > 0$ je pevně zvolený čas, pak

$$(67) \quad |\varepsilon_{h,\tau}(x, t)| \leq \frac{1}{2} M_1 \tau + \frac{1}{12} M_2 h^2 = \frac{1}{2} \tau \left(M_1 + \frac{1}{6\mu} M_2 \right) \quad \forall x \in [h, 1-h] \times [0, T-\tau].$$

Vidíme tedy, že schéma je prvního řádu přesnosti v čase a druhého řádu přesnosti v prostoru. Pro pevný poměr μ (který je vhodné uvažovat pro zajištění stability, jak uvidíme později) se $\varepsilon_{h,\tau}$ chová jako $O(\tau)$ pro $\tau \rightarrow 0$ a v tomto smyslu je schéma prvního řádu přesnosti.

Poznámka 4 Jelikož $u_t = u_{xx}$, je $u_{tt} = u_{xxt} = (u_t)_{xx} = u_{xxxx}$, a tudíž

$$\varepsilon_{h,\tau}(x, t) = \frac{1}{2} \left(1 - \frac{1}{6\mu} \right) u_{xxxx}(x, t) \tau + O(\tau^2).$$

Pro $\mu = \frac{1}{6}$ je tedy schéma druhého řádu přesnosti. Jedná se ale o velmi speciální případ, se kterým se v obecnějších situacích nesetkáme.

4.3 Konvergence explicitního schématu

Věta 5 Uvažujme posloupnost $(h_i, \tau_i) \rightarrow (0, 0)$ pro $i \rightarrow \infty$ a předpokládejme, že $\mu_i \equiv \tau_i/h_i^2 \leq \frac{1}{2}$. Necht' $T > 0$ a $|u_{tt}| \leq M_1$, $|u_{xxxx}| \leq M_2$ v $[0, 1] \times [0, T]$. Pak pro libovolný bod $(x, t) \in [0, 1] \times [0, T]$ a libovolnou posloupnost $(j_i, n_i) \in \mathbb{N} \times \mathbb{N}$ takovou, že $j_i h_i \rightarrow x$, $n_i \tau_i \rightarrow t$, konvergují aproximace $U_{j_i}^{n_i}$ generované explicitním schématem (64)–(65) k řešení $u(x, t)$, přičemž tato konvergence je stejnoměrná v $[0, 1] \times [0, T]$.

Důkaz. Uvažujme libovolnou dvojici h, τ ($\tau < T$) a libovolný bod $(x_j, t_n) \in (0, 1) \times (0, T)$. Pro chybu aproximace $e_j^n = U_j^n - u_j^n$ platí

$$e_j^{n+1} = e_j^n + \mu [e_{j+1}^n - 2e_j^n + e_{j-1}^n] - \tau \varepsilon_j^n, \quad j = 1, \dots, J-1.$$

Definujeme-li $\|e^n\|_\infty = \max_{l=1, \dots, J-1} |e_l^n|$, pak (díky $e_0^n = e_J^n = 0$)

$$\|e^{n+1}\|_\infty \leq (|1 - 2\mu| + 2\mu) \|e^n\|_\infty + \tau \|\varepsilon^n\|_\infty \leq \|e^n\|_\infty + \tau \|\varepsilon^n\|_\infty.$$

Jelikož $e_l^0 = U_l^0 - u^0(x_l) = 0$, $l = 0, \dots, J$, dostáváme

$$\|e^n\|_\infty \leq \tau \sum_{k=0}^{n-1} \|\varepsilon^k\|_\infty \leq t_n \max_{k=0, \dots, n-1} \|\varepsilon^k\|_\infty.$$

Použitím (67) získáváme

$$|U_j^n - u(x_j, t_n)| \leq T \left(\frac{1}{2} M_1 \tau + \frac{1}{12} M_2 h^2 \right).$$

Tvrzení nyní plyne ze spojitosti u na $[0, 1] \times [0, T]$. □

4.4 Fourierova analýza chyby

Víme, že řešení úlohy (59)–(61) lze zapsat ve tvaru Fourierovy řady (62). Ukážeme, že v podobném tvaru lze zapsat i přibližné řešení splňující (64)–(66). Za tím účelem zapíšeme řady (62) a (63) pomocí komplexních exponencií:

$$(68) \quad u(x, t) = \sum_{m=-\infty}^{\infty} A_m e^{im\pi x - (m\pi)^2 t},$$

$$(69) \quad u^0(x) = \sum_{m=-\infty}^{\infty} A_m e^{im\pi x}.$$

Položíme-li $u^0(x) = -u^0(-x)$ pro $x \in [-1, 0)$, pak

$$A_m = \frac{1}{2} \int_{-1}^1 u^0(x) e^{-im\pi x} dx.$$

Snadno se ověří, že pro $m \in \mathbb{N}$ je $A_m = -A_{-m} = -i a_m/2$. Budeme předpokládat, že Fourierova řada (69) je absolutně konvergentní, tj.

$$(70) \quad \sum_{m=-\infty}^{\infty} |A_m| < \infty.$$

Jak víme, postačující podmínkou pro to je, aby funkce u^0 byla absolutně spojitá na $[0, 1]$, $(u^0)' \in L^2(0, 1)$ a $u^0(0) = u^0(1) = 0$.

Připomeňme, že funkce $e^{i m \pi x - (m \pi)^2 t}$ jsou řešeními rovnice (59). V uzlech sítě platí

$$e^{i m \pi x_j - (m \pi)^2 t_n} = e^{i k j h} \left[e^{-k^2 \tau} \right]^n,$$

kde jsem zavedli *vlňové číslo* $k = m \pi$. V analogii k tomu se můžeme ptát, kdy

$$U_j^n = e^{i k j h} \lambda^n$$

řeší diferenční rovnici (64). Dosazením do (64) získáme

$$e^{i k j h} \lambda^{n+1} = e^{i k j h} \lambda^n [1 + \mu (e^{i k h} - 2 + e^{-i k h})],$$

a tudíž

$$(71) \quad \lambda \equiv \lambda(k) = 1 - 2 \mu [1 - \cos(k h)] = 1 - 4 \mu \sin^2 \frac{k h}{2}.$$

Číslo $\lambda(k)$ se nazývá *amplifikační (zesilující) faktor* členu Fourierovy řady. Následující věta ukazuje, že přibližné řešení U_j^n můžeme zapsat ve tvaru podobném jako u .

Věta 6 *Nechť síťová funkce U_j^n splňuje (64)–(66). Pak*

$$(72) \quad U_j^n = \sum_{m=-\infty}^{\infty} A_m e^{i m \pi j h} [\lambda(m \pi)]^n, \quad j = 0, 1, \dots, J, \quad n \geq 0.$$

Důkaz. Jelikož amplifikační faktory jsou omezené a platí (70), konverguje řada (72) absolutně, a jelikož každý její člen řeší diferenční rovnici (64), řeší i součet řady (72) rovnici (64). Dále pro každé $m \in \mathbb{N}$ platí $A_m [\lambda(m \pi)]^n = -A_{-m} [\lambda(-m \pi)]^n$, a tudíž je pro $j = 0$ a $j = J$ součet řady 0. Jsou tedy splněny okrajové podmínky (65). Konečně pro $n = 0$ se řada (72) redukuje na řadu (69) s $x = j h$, a součet řady tudíž splňuje i počáteční podmínku (66). Součet řady tedy splňuje (64)–(66), a jelikož je řešení úlohy (64)–(66) určeno jednoznačně, platí (72). \square

Poznámka 5 Všimněme si, že z vlastností amplifikačního faktoru a koeficientů A_m plyne, že řadu (72) lze též přepsat do tvaru

$$(73) \quad U_j^n = \sum_{m=1}^{\infty} a_m \sin(m \pi j h) [\lambda(m \pi)]^n, \quad j = 0, 1, \dots, J, \quad n \geq 0.$$

Ze vztahu (68) plyne, že

$$u_j^n = \sum_{m=-\infty}^{\infty} A_m e^{im\pi jh} \left[e^{-(m\pi)^2 \tau} \right]^n.$$

Jelikož chceme, aby U_j^n dobře aproximovalo u_j^n , vidíme ze srovnání uvedené řady s řadou (72), že je potřeba, aby hodnoty $\lambda(k)$ byly dobrou aproximací hodnot $e^{-k^2 \tau}$, alespoň pro nízké hodnoty vlnového čísla k . To je předmětem následujícího lemmatu.

Lemma 1 *Platí*

$$(74) \quad |e^{-k^2 \tau} - \lambda(k)| \leq C(\mu) k^4 \tau^2 \quad \forall k, \tau > 0,$$

kde $C(\mu)$ závisí pouze na μ . Pokud $\mu = \frac{1}{6}$, pak

$$|e^{-k^2 \tau} - \lambda(k)| \leq \frac{4}{15} k^6 \tau^3 \quad \forall k, \tau > 0.$$

Důkaz. K důkazu využijeme Taylorův vzorec ve tvaru

$$f(x) = \sum_{k=0}^{l-1} \frac{1}{k!} f^{(k)}(0) x^k + \frac{1}{l!} f^{(l)}(\xi) x^l,$$

který platí pro libovolné $l \in \mathbb{N}$, $f \in C^l(\mathbb{R})$ a $x > 0$ s vhodným $\xi \in (0, x)$. Aplikujeme-li tento vzorec na $f(x) = e^{-x}$ a $f(x) = \cos x$, získáme

$$\begin{aligned} e^{-k^2 \tau} &= 1 - k^2 \tau + \frac{1}{2} e^{-\xi} k^4 \tau^2 \\ \cos(kh) &= 1 - \frac{1}{2} (kh)^2 + \frac{1}{24} (\cos \zeta) (kh)^4, \end{aligned}$$

kde $\xi \in (0, k^2 \tau)$ a $\zeta \in (0, kh)$. Podle (71) je tedy

$$\lambda(k) = 1 - 2\mu [1 - \cos(kh)] = 1 - k^2 \tau + \frac{1}{12} (\cos \zeta) k^4 h^2 \tau,$$

z čehož plyne

$$|e^{-k^2 \tau} - \lambda(k)| = \frac{1}{2} \left| k^4 \tau^2 e^{-\xi} - \frac{1}{6} k^4 h^2 \tau \cos \zeta \right| \leq \frac{1}{2} \left(1 + \frac{1}{6\mu} \right) k^4 \tau^2,$$

což dokazuje (74). Podobně získáme

$$|e^{-k^2 \tau} - \lambda(k)| = \frac{1}{360} |60 k^6 \tau^3 e^{-\tilde{\xi}} - k^6 h^4 \tau \cos \tilde{\zeta}| \leq \frac{4}{15} k^6 \tau^3 \quad \text{pro } \mu = \frac{1}{6},$$

kde opět $\tilde{\xi} \in (0, k^2 \tau)$ a $\tilde{\zeta} \in (0, kh)$. □

Uvedené lemma můžeme využít jako alternativní prostředek pro vyšetřování chyby diskretizace. Na základě vztahů (68) a (70) platí

$$\varepsilon_j^n = \frac{u_j^{n+1} - u_j^n}{\tau} - \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2}$$

$$\begin{aligned}
&= \sum_{m=-\infty}^{\infty} A_m e^{i m \pi x_j - (m \pi)^2 t_n} \left[\frac{e^{-(m \pi)^2 \tau} - 1}{\tau} - \frac{e^{i m \pi h} - 2 + e^{-i m \pi h}}{h^2} \right] \\
&= \sum_{m=-\infty}^{\infty} A_m e^{i m \pi x_j - (m \pi)^2 t_n} \frac{1}{\tau} \left[e^{-(m \pi)^2 \tau} - \lambda(m \pi) \right].
\end{aligned}$$

Použitím odhadu (74) dostáváme

$$|\varepsilon_j^n| \leq \sum_{m=-\infty}^{\infty} |A_m| \frac{1}{\tau} |e^{-(m \pi)^2 \tau} - \lambda(m \pi)| \leq \sum_{m=-\infty}^{\infty} |A_m| C(\mu) (m \pi)^4 \tau.$$

Abychom při konstantním μ získali stejnoměrný odhad chyby diskretizace prvního řádu přesnosti, je tedy nutné požadovat, aby $\pi^4 \sum_{m=-\infty}^{\infty} |A_m| m^4 < \infty$. Tato hodnota představuje odhad $|u_{xxxx}| = |u_{tt}|$ a dostáváme tedy odhad (67). Pro $h^2 = 6\tau$ lze pomocí lematu 1 získat přesnost druhého řádu (srv. pozn. 4).

Vyjádření řešení úlohy (59)–(61) ve tvaru řady (68) ukazuje, že pro libovolné vlnové číslo $k = m\pi$ zůstává amplituda příslušné složky řešení v čase omezená. Je přirozené tuto vlastnost požadovat i od přibližného řešení. Řekneme proto, že numerická metoda pro řešení úlohy (59)–(61) je stabilní, pokud existuje konstanta K nezávislá na k taková, že $|\lambda(k)^n| \leq K$ pro všechna k a n . To je ekvivalentní požadavku, aby $|\lambda(k)| \leq 1$ pro všechna k .

Věta 7 *Nechť pro dané hodnoty h a τ je $|\lambda(m\pi)| \leq 1$ pro všechna $m \in \mathbb{N}$. Pak řešení schématu (64)–(66) zůstává pro libovolnou počáteční podmínku u^0 splňující (70) omezené pro $n \rightarrow \infty$.*

Důkaz. Jelikož $|\lambda(m\pi)| \leq 1 \forall m \in \mathbb{N}$, je dle (72) a (70)

$$|U_j^n| \leq \sum_{m=-\infty}^{\infty} |A_m| |\lambda(m\pi)|^n \leq \sum_{m=-\infty}^{\infty} |A_m| < \infty.$$

□

Poznámka 6 *Nechť existuje $m_1 \in \mathbb{N}$ tak, že $|\lambda(m_1\pi)| > 1$. Buď $u^0(x) = \sin(m_1\pi x)$. Pak podle (73) je $U_j^n = \sin(m_1\pi jh) [\lambda(m_1\pi)]^n$, a tudíž $|U_j^n| \rightarrow \infty$ pro $n \rightarrow \infty$ a pro každé $j \in \{1, \dots, J-1\}$ takové, že $\sin(m_1\pi jh) \neq 0$.*

Nyní vidíme význam podmínky $\mu \leq \frac{1}{2}$. Je-li tato podmínka splněna, pak $|\lambda(m\pi)| \leq 1 \forall m \in \mathbb{N}$, a řešení tudíž zůstává omezené. Je-li $\mu > \frac{1}{2}$, pak pro některá $m \in \mathbb{N}$ je $\lambda(m\pi) < -1$ a velikost příslušných členů v (72) s postupujícím časem roste nade všechny meze. Teoreticky je možné zvolit počáteční podmínku tak, aby $A_m = 0$, kdykoli $\lambda(m\pi) < -1$. Avšak to je velmi speciální situace a vpraxi by vlivem zaokrouhlovacích chyb vznikly malé nenulové koeficienty u všech takovýchto členů a s postupujícím časem by tyto členy opět neomezeně rostly. Schéma (64) je tedy stabilní pouze při splnění podmínky $\mu \leq \frac{1}{2}$.

Fourierovu metodu můžeme použít též k důkazu konvergence. Její výhoda spočívá v tom, že nemusíme předpokládat dostatečnou hladkost řešení u a stejnoměrnou omezenost u_{xxxx} a u_{tt} . Naším jediným předpokladem o úloze (59)–(61) nyní bude absolutní konvergence řady (69) pro počáteční podmínku u^0 . Počáteční podmínka tedy nemusí být hladká. Budeme předpokládat, že μ je pevné a $\mu \leq \frac{1}{2}$. Chybu aproximace můžeme vyjádřit ve tvaru

$$e_j^n = U_j^n - u(x_j, t_n) = \sum_{m=-\infty}^{\infty} A_m e^{i m \pi x_j} \left\{ [\lambda(m \pi)]^n - [e^{-(m \pi)^2 \tau}]^n \right\}.$$

Při označení $\lambda_1 = \lambda(m \pi)$ a $\lambda_2 = e^{-(m \pi)^2 \tau}$ je $\lambda_1, \lambda_2 \in [-1, 1]$ a platí

$$\lambda_1^n - \lambda_2^n = \sum_{k=0}^{n-1} \lambda_1^{n-k} \lambda_2^k - \sum_{k=1}^n \lambda_1^{n-k} \lambda_2^k = (\lambda_1 - \lambda_2) \sum_{k=0}^{n-1} \lambda_1^{n-1-k} \lambda_2^k,$$

z čehož plyne

$$(75) \quad |\lambda_1^n - \lambda_2^n| \leq n |\lambda_1 - \lambda_2| \quad \forall \lambda_1, \lambda_2 \in [-1, 1].$$

Bud' $\varepsilon > 0$ libovolné a necht' $m_0 \in \mathbb{N}$ je takové, že

$$\sum_{|m| > m_0} |A_m| \leq \frac{\varepsilon}{4}.$$

Pak použitím (75) získáme

$$|e_j^n| \leq \frac{\varepsilon}{2} + \sum_{|m| \leq m_0} n |A_m| |\lambda(m \pi) - e^{-(m \pi)^2 \tau}|.$$

Z nerovnosti (74) plyne

$$|e_j^n| \leq \frac{\varepsilon}{2} + n \tau^2 C(\mu) \pi^4 \sum_{|m| \leq m_0} |A_m| m^4$$

a vidíme tedy, že pro τ dostatečně malé je $|e_j^n| \leq \varepsilon \forall (x_j, t_n) \in [0, 1] \times [0, T]$, neboť $n \tau \leq T$.

Při aplikaci Fourierovy metody jsme reprezentovali přibližné řešení pomocí nekonečné řady (72), neboť ji bylo možné snadno srovnávat s řadou pro přesné řešení. Jelikož však na uvažované prostorové síti s $J + 1$ uzly lze reprezentovat jen konečně mnoho různých frekvencí, lze přibližné řešení vyjádřit jako lineární kombinaci $2J$ po sobě jdoucích funkcí

$$(76) \quad e^{i m \pi j h} [\lambda(m \pi)]^n.$$

Snadno ověříme, že tyto funkce se skutečně nezmění, nahradíme-li m hodnotou $m + 2J$. Můžeme tedy např. uvažovat U_j^n jakožto lineární kombinaci funkcí (76) odpovídajících

$$m = -(J - 1), -(J - 2), \dots, -1, 0, 1, \dots, J.$$

4.5 Implicitní metoda pro úlohu (59)–(61)

Podmínka stability $\tau \leq h^2/2$ pro explicitní schéma (64) je velmi vážné omezení, které implikuje, že bude třeba velmi mnoho časových kroků. Navíc, budeme-li muset zmenšit h pro zvýšení přesnosti, velmi se zvýší celková výpočetní náročnost, neboť budeme muset též podstatně zmenšit τ . Ukážeme nyní, že zmíněných omezení se můžeme zbavit použitím zpětné diference pro diskretizaci časové derivace v (59), tj. použitím schématu

$$\frac{U_j^{n+1} - U_j^n}{\tau} = \frac{U_{j+1}^{n+1} - 2U_j^{n+1} + U_{j-1}^{n+1}}{h^2}.$$

Toto schéma lze přepsat do tvaru

$$(77) \quad -\mu U_{j-1}^{n+1} + (1 + 2\mu)U_j^{n+1} - \mu U_{j+1}^{n+1} = U_j^n, \quad j = 1, 2, \dots, J-1, \quad n \geq 0.$$

Na rozdíl od schématu (64) nyní danou hodnotu U_j^{n+1} nelze určit nezávisle na ostatních hodnotách na časové hladině t_{n+1} , nýbrž všechny tyto hodnoty je nutno vypočítat současně vyřešením soustavy $J-1$ lineárních rovnic pro $J-1$ neznámých. Jedná se tedy o implicitní metodu.

Stabilitu schématu (77) s okrajovými a počátečními podmínkami (65) a (66) můžeme vyšetřovat Fourierovou metodou analogicky jako pro schéma (64). Snadno zjistíme, že funkce $U_j^n = e^{ikjh} \lambda^n$ splňuje (77) právě tehdy, když

$$\lambda \equiv \lambda(k) = \frac{1}{1 + 4\mu \sin^2 \frac{kh}{2}}.$$

Tedy $\lambda(k) \in (0, 1]$ pro libovolné $\mu > 0$ a libovolné $k \in \mathbb{R}$, což znamená, že metoda je nepodmíněně stabilní.

4.6 Thomasův algoritmus

Soustava (77) je tridiagonální a můžeme ji zapsat ve tvaru

$$-a_j U_{j-1} + b_j U_j - c_j U_{j+1} = d_j, \quad j = 1, 2, \dots, J-1,$$

kde

$$U_0 = U_J = 0.$$

Budeme předpokládat, že

$$(78) \quad a_j > 0, \quad b_j > 0, \quad c_j > 0, \quad b_j > a_j + c_j,$$

což splňuje (77). Soustavu rovnic nejprve převedeme na soustavu s horní trojúhelníkovou maticí tvaru

$$(79) \quad U_j - e_j U_{j+1} = f_j, \quad j = 1, 2, \dots, J-1.$$

Máme-li v tomto tvaru k -tou rovnicí a chceme-li upravit $k+1$ -vou rovnicí, tj.

$$-a_{k+1} U_k + b_{k+1} U_{k+1} - c_{k+1} U_{k+2} = d_{k+1},$$

pak k této rovnici přičteme rovnici (79) s $j = k$ přenásobenou a_{k+1} , čímž získáme

$$(b_{k+1} - a_{k+1} e_k) U_{k+1} - c_{k+1} U_{k+2} = d_{k+1} + a_{k+1} f_k.$$

Algoritmus je tedy následující

```

e1 := c1/b1, f1 := d1/b1
for j = 2, ..., J - 2 do
  ej := cj/(bj - aj ej-1)
  fj := (dj + aj fj-1)/(bj - aj ej-1)
enddo
fJ-1 := (dJ-1 + aJ-1 fJ-2)/(bJ-1 - aJ-1 eJ-2)

```

Řešení U_j pak jednoduše určíme z rovnic (79).

Snadno lze ověřit, že $e_j \in (0, 1)$, $j = 1, 2, \dots, J - 2$. Algoritmus lze tedy vždy provést a je numericky stabilní (nevede ke vzrůstajícím chybám).

Uvedený algoritmus potřebuje pro vyřešení soustavy (77) na jeden uzel tři sčítání, tři násobení a dvě dělení, zatímco explicitní schéma (64) vyžaduje tři sčítání a dvě násobení (popř. čtyři sčítání a jedno násobení). Výpočetní náročnost implicitního schématu je tedy asi dvojnásobná oproti explicitnímu schématu. Důležité však je, že lze volit mnohem delší časové kroky (aniž by se zhoršila přesnost), neboť nyní není žádná podmínka stability omezující volbu τ . Proto je celková výpočetní náročnost implicitní metody pro dosažení zvoleného času T mnohem menší než u explicitní metody.

4.7 Dvukrokové metody

Výše uvažované metody pro řešení úlohy (59)–(61) jsou druhého řádu přesnosti v prostoru, avšak obecně pouze prvního řádu přesnosti v čase. Nabízí se proto uvažovat pro diskretizaci časové derivace místo jednostranné diference centrální diferenci, podobně jako u leapfrog scheme (32). To vede v nejjednodušším případě ke schématu

$$(80) \quad \frac{U_j^{n+1} - U_j^{n-1}}{2\tau} = \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2}, \quad j = 1, 2, \dots, J - 1, \quad n \geq 1.$$

Snadno zjistíme, že chyba diskretizace je v tomto případě druhého řádu v čase i v prostoru.

Uvedené schéma je příkladem dvukrokového schématu, neboť k určení hodnot U_j^{n+1} nestačí znát hodnoty U_j^n , ale potřebujeme též hodnoty U_j^{n-1} . Přibližné řešení tedy závisí nejen na počáteční podmínce, ale též na hodnotách v čase t_1 . Tyto hodnoty buď musíme předepsat, a nebo musíme stanovit postup, jak je určit. Obvykle se pro určení hodnot U_j^1 použije nějaké jednokrokové schéma. Metody, které zahrnují více než dvě časové hladiny, se souhrnně nazývají vícekroková schémata.

Hodnoty přibližného řešení v čase t_1 zapíšeme ve tvaru

$$(81) \quad U_j^1 = \sum_{m=-\infty}^{\infty} B_m e^{i m \pi j h}, \quad j = 0, 1, \dots, J,$$

a opět předpokládáme, že řada konverguje absolutně a že $B_m = -B_{-m} \forall m \in \mathbb{Z}$. Podobně jako při Fourierově analýze jednokrokových metod výše hledejme nejprve řešení diferenčního schématu ve tvaru se "separovanými proměnnými", tj. položme

$$U_j^n = e^{ikjh} \widehat{U}^n(k),$$

kde $k = m\pi$, $m \in \mathbb{Z}$. Dosazením do (80) získáme

$$(82) \quad \widehat{U}^{n+1}(k) + 2q(k)\widehat{U}^n(k) - \widehat{U}^{n-1}(k) = 0, \quad n \geq 1, \quad \text{kde } q(k) = \frac{4\tau}{h^2} \sin^2 \frac{kh}{2}.$$

Charakteristický polynom $\lambda^2 + 2q(k)\lambda - 1$ této soustavy diferenčních rovnic má dva různé reálné kořeny $\lambda_{\pm} = -q(k) \pm \sqrt{q(k)^2 + 1}$, a tudíž obecné řešení soustavy (82) je

$$(83) \quad \widehat{U}^n(k) = \alpha_+(k) \lambda_+(k)^n + \alpha_-(k) \lambda_-(k)^n,$$

kde koeficienty α_{\pm} jsou libovolná komplexní čísla. Tato čísla určíme tak, aby platilo $\widehat{U}^0(m\pi) = A_m$ a $\widehat{U}^1(m\pi) = B_m$. Z toho plyne

$$(84) \quad \alpha_+(m\pi) = \frac{B_m - A_m \lambda_-(m\pi)}{\lambda_+(m\pi) - \lambda_-(m\pi)}, \quad \alpha_-(m\pi) = \frac{A_m \lambda_+(m\pi) - B_m}{\lambda_+(m\pi) - \lambda_-(m\pi)}.$$

Řešení diskrétního problému (80), (66), (65), (81) je pak dáno řadou

$$(85) \quad U_j^n = \sum_{m=-\infty}^{\infty} e^{im\pi jh} \widehat{U}^n(m\pi), \quad j = 0, 1, \dots, J, \quad n \geq 0.$$

Bohužel složky řešení odpovídající kořenu λ_- jsou nestabilní, neboť $\lambda_-(k) < -1$, kdykoli $\sin(\frac{1}{2}kh) \neq 0$. Schéma (80) je tedy bezcenné, neboť je pro libovolnou volbu h a τ nestabilní.

Poznamenejme, že pokud charakteristický polynom soustavy diferenčních rovnic odpovídajících dvoukrokovému schématu má pro dané k jeden dvojnásobný kořen, tj. $\lambda_+(k) = \lambda_-(k) \equiv \lambda(k)$, má obecné řešení soustavy diferenčních rovnic tvar

$$\widehat{U}^n(k) = \alpha(k) \lambda(k)^n + \beta(k) n \lambda(k)^{n-1},$$

kde $\alpha(k), \beta(k) \in \mathbb{C}$. V tomto případě tedy je

$$(86) \quad \widehat{U}^n(m\pi) = A_m \lambda(m\pi)^n + [B_m - A_m \lambda(m\pi)] n \lambda(m\pi)^{n-1}.$$

Pokud je koeficient u druhého členu nenulový, je ke stabilitě nutné, aby $|\lambda(m\pi)| < 1$, neboť jinak \widehat{U}^n poroste lineárně v n .

Uvedená analýza schématu (80) samozřejmě neznamená, že každé dvoukrokové explicitní schéma je vždy nestabilní. Uvažujme např. schéma

$$(87) \quad \frac{U_j^{n+1} - U_j^{n-1}}{2\tau} = \frac{\delta_x^2 U_j^n + \delta_x^2 U_j^{n-1}}{2h^2}, \quad j = 1, 2, \dots, J-1, \quad n \geq 1.$$

V tomto případě získáme soustavu diferenčních rovnic

$$(88) \quad \widehat{U}^{n+1}(k) + q(k) \widehat{U}^n(k) + (q(k) - 1) \widehat{U}^{n-1}(k) = 0, \quad n \geq 1,$$

jejíž charakteristický polynom má kořeny $\lambda_+(k) = 1 - q(k)$ a $\lambda_-(k) = -1$. Je-li $q(k) \neq 2$, jsou oba kořeny různé a $\widehat{U}^n(k)$ je opět dáno vztahy (83) a (84). Zřejmě je $|\lambda_{\pm}(k)| \leq 1$ právě tehdy, když $q(k) \in [0, 2)$, z čehož plyne podmínka stability $\tau \leq h^2/2$. Je-li $q(k) = 2$, je $\lambda_+(k) = \lambda_-(k)$ a $\widehat{U}^n(k)$ je dáno vztahem (86) s $\lambda(k) = -1$. Příklad $q(m\pi) = 2$ za uvedené podmínky stability může nastat pouze tehdy, je-li $m = (2l + 1)J$, kde $l \in \mathbb{Z}$. Členy s $m = (2l + 1)J$ se však v řešení neprojeví, neboť řadu (85) můžeme díky vlastnosti $\widehat{U}^n(m\pi) = -\widehat{U}^n(-m\pi) \forall m \in \mathbb{Z}$ zapsat ve tvaru

$$U_j^n = 2i \sum_{m=1}^{\infty} \sin(m\pi j h) \widehat{U}^n(m\pi), \quad j = 0, 1, \dots, J, \quad n \geq 0.$$

Schéma (87) je tedy pro $\tau \leq h^2/2$ stabilní.

Stabilita metody však nemusí znamenat, že diskrétní řešení bude bez oscilací. Jak jsme viděli, platí v případě $q(k) \neq 2$

$$\widehat{U}^n(k) = \alpha_+(k) (1 - q(k))^n + \alpha_-(k) (-1)^n.$$

Označíme-li

$$V_j^n = \sum_{\substack{m=-\infty \\ q(m\pi) < 2}}^{\infty} e^{im\pi j h} \alpha_+(m\pi) (1 - q(m\pi))^n, \quad W_j^n = \sum_{\substack{m=-\infty \\ q(m\pi) < 2}}^{\infty} e^{im\pi j h} \alpha_-(m\pi),$$

můžeme psát

$$U_j^n = V_j^n + (-1)^n W_j^n, \quad j = 0, 1, \dots, J, \quad n \geq 0.$$

Síťová funkce W_j nezávisí na čase, a pokud je nenulová, bude pro velké n představovat dominantní složku řešení, neboť $V_j^n \rightarrow 0$ pro $n \rightarrow \infty$. Obvykle se projeví ve formě oscilací, jejichž velikost může být i větší než jsou hodnoty počáteční podmínky, jelikož velikost $\alpha_+(m\pi)$ a $\alpha_-(m\pi)$ může být podstatně větší než $|A_m|$. Ze stability metody však plyne, že tyto oscilace nebudou pro $n \rightarrow \infty$ narůstat. Navíc se jejich velikost zmenší, pokud zjemníme síť.

Časově nezávislá oscilující složka nebude v řešení schématu (87) přítomna, pokud řešení v čase t_1 určíme pomocí schématu (64). Podle (72) a (71) je pak totiž $B_m = A_m (1 - q(m\pi)) = A_m \lambda_+(m\pi)$, a tudíž $\alpha_-(m\pi) = 0$ pro každé $m \in \mathbb{Z}$.

4.8 Disipace

Z diskuse v předchozím odstavci plyne, že je žádoucí, aby diferenční schéma vedlo v každém časovém kroku k poklesu amplitud vysokofrekvenčních složek řešení, tj. aby velikost příslušných amplifikačních faktorů byla menší než 1. Jak již víme z části 3.7, tento pokles vysokofrekvenčních oscilací se nazývá disipace. Analogicky jako v definici 3 na str. 28

řekneme, že diferenční schéma je disipativní řádu $2r$ (má disipaci řádu $2r$), jestliže existuje kladná konstanta C nezávislá na h a τ taková, že každý amplifikační faktor $\lambda_j(k)$ splňuje

$$(89) \quad |\lambda_j(k)| \leq 1 - C \left(\sin \frac{kh}{2} \right)^{2r}$$

pro všechna vlnová čísla k . Poznamenejme, že pro každé k uvažujeme obecně více amplifikačních faktorů, abychom do našich úvah zahrnuli i víceukroková schémata. Pro ověření platnosti podmínky (89) může být užitečné si všimnout, že nerovnost (89) je ekvivalentní podmínce

$$|\lambda_j(k)|^2 \leq 1 - C' \left(\sin \frac{kh}{2} \right)^{2r}$$

s konstantou C' nezávislou na h a τ . Disipativnost schématu je v případě parabolických rovnic velmi přirozeným požadavkem, neboť pak dochází v průběhu času ke zhlazování diskrétního řešení, stejně jako je tomu u řešení aproximované diferenciální rovnice.

Použitím vztahu (71) snadno zjistíme, že explicitní schéma (64) je disipativní řádu 2 pro $\mu \in [\mu_0, \mu_1] \subset (0, \frac{1}{2})$, kde μ_0 a μ_1 jsou konstanty nezávislé na h a τ . Proto se obvykle nepoužívá volba $\mu = \frac{1}{2}$, při níž je schéma (64) ještě stabilní. Implicitní schéma (77) je disipativní řádu 2, je-li $\mu \geq \mu_0 > 0$, kde μ_0 je konstanta nezávislá na h a τ . Schéma (87) samozřejmě disipativní není.

4.9 θ -metoda (θ -schéma, metoda váženého průměru)

Uvažujeme-li vážený průměr explicitního schématu (64) a implicitního schématu (77), získáme šestibodové schéma

$$(90) \quad U_j^{n+1} - U_j^n = \mu [\theta \delta_x^2 U_j^{n+1} + (1 - \theta) \delta_x^2 U_j^n], \quad j = 1, 2, \dots, J - 1, \quad n \geq 0.$$

Budeme předpokládat, že $\theta \in [0, 1]$. Pro $\theta = 0$ získáváme explicitní schéma (64) a pro $\theta = 1$ plně implicitní schéma (77). Pro libovolné $\theta \in (0, 1]$ je k určení hodnot přibližného řešení v čase t_{n+1} nutno vyřešit tridiagonální soustavu lineárních rovnic

$$-\theta \mu U_{j-1}^{n+1} + (1 + 2\theta \mu) U_j^{n+1} - \theta \mu U_{j+1}^{n+1} = [1 + (1 - \theta) \mu \delta_x^2] U_j^n, \quad j = 1, 2, \dots, J - 1.$$

Koeficienty splňují nerovnosti (78), a můžeme tedy použít Thomasův algoritmus.

Stabilitu vyšetříme opět pomocí Fourierovy metody. Dosazením $U_j^n = e^{ikjh} \lambda^n$ do (90) získáme

$$\lambda = \frac{1 - 4(1 - \theta) \mu \sin^2 \frac{kh}{2}}{1 + 4\theta \mu \sin^2 \frac{kh}{2}}.$$

Zřejmě $\lambda \leq 1$. Nestabilita se může objevit pouze pokud $\lambda < -1$, což nastane právě tehdy, když

$$4(1 - 2\theta) \mu \sin^2 \frac{kh}{2} > 2.$$

Z toho plyne, že

$$(91) \quad \begin{cases} \text{Je-li } \theta \in [0, \frac{1}{2}), \text{ pak (90) je stabilní} & \iff \mu \leq \frac{1}{2(1-2\theta)}. \\ \text{Je-li } \theta \in [\frac{1}{2}, 1], \text{ pak (90) je stabilní} & \forall \mu > 0. \end{cases}$$

V prvním případě je tedy schéma podmíněně stabilní, v druhém nepodmíněně stabilní.

Chybu diskretizace schématu (90) je vhodné počítat v čase $t_{n+1/2} \equiv (n + \frac{1}{2})\tau$, tj. definujeme

$$\varepsilon_j^{n+1/2} = \varepsilon_{h,\tau}(x_j, t_{n+1/2}) = \frac{u(x_j, t_{n+1}) - u(x_j, t_n)}{\tau} - \frac{\theta \delta_x^2 u(x_j, t_{n+1}) + (1-\theta) \delta_x^2 u(x_j, t_n)}{h^2}.$$

Tedy

$$\begin{aligned} \varepsilon_{h,\tau}(x, t) &= \frac{\delta_t u(x, t)}{\tau} - \frac{\theta \delta_x^2 u(x, t + \frac{\tau}{2}) + (1-\theta) \delta_x^2 u(x, t - \frac{\tau}{2})}{h^2} \\ &= \frac{\delta_t u(x, t)}{\tau} - (\theta - \frac{1}{2}) \frac{\delta_t \delta_x^2 u(x, t)}{h^2} - \frac{\delta_x^2 u(x, t + \frac{\tau}{2}) + \delta_x^2 u(x, t - \frac{\tau}{2})}{2h^2}. \end{aligned}$$

Dosažením Taylorových rozvoju získáme

$$\begin{aligned} \varepsilon_{h,\tau}(x, t) &= [u_t + \frac{1}{24} u_{ttt} \tau^2 + \dots] - (\theta - \frac{1}{2}) [u_{xxt} \tau + \frac{1}{12} u_{xxxxt} h^2 \tau + \dots] \\ &\quad - [u_{xx} + \frac{1}{12} u_{xxxx} h^2 + \frac{2}{6!} u_{xxxxx} h^4 + \frac{1}{8} u_{xtt} \tau^2 + \dots]. \end{aligned}$$

Použitím (59) zjišťujeme, že obecně $\varepsilon_{h,\tau} = O(\tau + h^2)$, avšak pro $\theta = \frac{1}{2}$ je $\varepsilon_{h,\tau} = O(\tau^2 + h^2)$. Pro $\theta = \frac{1}{2}$ je tedy schéma (90) druhého řádu přesnosti v prostoru i v čase a nazývá se schéma Crankovo–Nicolsonové. Jelikož je nepodmíněně stabilní, můžeme uvažovat $h = O(\tau)$. Pak $\varepsilon_{h,\tau} = O(\tau^2)$ a jsme tedy schopni dosáhnout dobrou přesnost při malé výpočetní náročnosti. Při volbě $h = O(\tau)$ však schéma Crankovo–Nicolsonové není disipativní, což způsobuje, že při nehladké počáteční podmínce může být méně přesné než plně implicitní schéma (77), které je disipativní řádu 2.

Metodu druhého řádu přesnosti v čase lze získat též pro

$$\theta = \frac{1}{2} - \frac{h^2}{12\tau}, \quad \text{tj.} \quad \mu = \frac{1}{6(1-2\theta)}.$$

(Musí být $h^2 \leq 6\tau$, aby bylo $\theta \geq 0$.) Při této volbě je dle (91) schéma (90) stabilní a platí $\varepsilon_{h,\tau} = O(\tau^2 + h^4)$. Opět tedy můžeme používat velké časové kroky a metoda bude přitom pro hladké počáteční podmínky přesná a stabilní. Při volbě $h = O(\tau)$ však schéma opět není disipativní a pro malé h je blízké schématu Crankovo–Nicolsonové.

I když lze odvodit řadu dalších schémat pro řešení úlohy (59)–(61), nejpoužívanější je v praxi schéma (90). Nejlepší volba parametru θ však závisí na řešeném problému a často není jasné, které schéma je opravdu nejlepší.

4.10 Princip maxima a konvergence

Věta 8 θ -schéma (90) s $\theta \in [0, 1]$ a $\mu(1 - \theta) \leq \frac{1}{2}$ dává přibližné řešení $\{U_j^n\}$ splňující

$$U_{\min}^n \leq U_j^n \leq U_{\max}^n,$$

kde

$$\begin{aligned} U_{\min}^n &= \min\{U_0^m, 0 \leq m \leq n; U_j^0, 0 \leq j \leq J; U_J^m, 0 \leq m \leq n\}, \\ U_{\max}^n &= \max\{U_0^m, 0 \leq m \leq n; U_j^0, 0 \leq j \leq J; U_J^m, 0 \leq m \leq n\}. \end{aligned}$$

Důkaz. Schéma (90) zapíšeme ve tvaru

$$(92) \quad \begin{aligned} &(1 + 2\theta\mu)U_j^{n+1} \\ &= \theta\mu(U_{j-1}^{n+1} + U_{j+1}^{n+1}) + (1 - \theta)\mu(U_{j-1}^n + U_{j+1}^n) + [1 - 2(1 - \theta)\mu]U_j^n. \end{aligned}$$

Koeficienty na pravé straně jsou nezáporné a jejich součet je $(1 + 2\theta\mu)$ (koeficienty před dvojkleny počítáme dvakrát). Předpokládejme, že U nabývá svého maxima ve vnitřním bodě a necht' toto maximum je U_j^{n+1} . Pak hodnoty U na pravé straně vztahu (92) jsou menší nebo rovny U_j^{n+1} , a jelikož součet koeficientů je $(1 + 2\theta\mu)$, musí být $U = U_j^{n+1}$ v každém z pěti sousedních uzlů v (92), pokud příslušný koeficient je nenulový. Je-li tedy $\theta \neq 0$, dostáváme $U_j^{n+1} = U_0^{n+1} = U_J^{n+1}$. Je-li $\theta = 0$, můžeme zkonstruovat posloupnost bodů, až dosáhneme hranice. Tedy $U_j^{n+1} = U_{\max}^n$. Stejným způsobem lze postupovat pro minimum. \square

Věta 9 Uvažujme posloupnost $(h_i, \tau_i) \rightarrow (0, 0)$ pro $i \rightarrow \infty$ a necht' $\mu_i(1 - \theta) \leq \frac{1}{2}$. Necht' chyba diskretizace odpovídající schématu (90) konverguje k nule stejnoměrně v množině $[0, 1] \times [0, T]$. Necht' chyby v okrajových a počátečních podmínkách rovněž konvergují stejnoměrně k nule pro $i \rightarrow \infty$. Pak aproximace dané schématem (90) konvergují stejnoměrně v $[0, 1] \times [0, T]$ k řešení rovnice (59) s konzistentními okrajovými a počátečními podmínkami.

Důkaz. Dle definice chyby diskretizace je pro $e_j^n = U_j^n - u(x_j, t_n)$

$$(93) \quad \begin{aligned} (1 + 2\theta\mu)e_j^{n+1} &= \theta\mu(e_{j-1}^{n+1} + e_{j+1}^{n+1}) + (1 - \theta)\mu(e_{j-1}^n + e_{j+1}^n) \\ &+ [1 - 2(1 - \theta)\mu]e_j^n - \tau\varepsilon_j^{n+1/2}, \quad j = 1, 2, \dots, J-1, \quad n = 0, 1, 2, \dots \end{aligned}$$

Předpokládejme nejprve, že $e_j^0 = 0$, $j = 0, \dots, J$, $e_0^n = e_J^n = 0$, $n = 0, 1, 2, \dots$ a označme

$$\|e^n\|_\infty = \max_{j=0, \dots, J} |e_j^n|, \quad \|\varepsilon^{n+1/2}\|_\infty = \max_{j=1, \dots, J-1} |\varepsilon_j^{n+1/2}|.$$

Pak

$$(1 + 2\theta\mu)\|e^{n+1}\|_\infty \leq 2\theta\mu\|e^{n+1}\|_\infty + \|e^n\|_\infty + \tau\|\varepsilon^{n+1/2}\|_\infty,$$

a tudíž $\|e^{n+1}\|_\infty \leq \|e^n\|_\infty + \tau \|\varepsilon^{n+1/2}\|_\infty$, z čehož plyne

$$\|e^n\|_\infty \leq \tau \sum_{m=0}^{n-1} \|\varepsilon^{m+1/2}\|_\infty \leq n \tau \max_{m=0, \dots, n-1} \|\varepsilon^{m+1/2}\|_\infty \rightarrow 0 \quad \text{pro } i \rightarrow \infty.$$

Předpokládejme nyní, že chyby v okrajových a počátečních podmínkách jsou nenulové, tj.

$$(94) \quad e_j^0 = \eta_j^0, \quad j = 0, \dots, J, \quad e_0^n = \eta_0^n, \quad e_J^n = \eta_J^n, \quad n = 0, 1, 2, \dots$$

Pak $e_j^n = \bar{e}_j^n + \tilde{e}_j^n$, kde \bar{e}_j^n splňuje (93) s homogenními počátečními a okrajovými podmínkami a \tilde{e}_j^n splňuje (92) a (94). Pak $\|\bar{e}^n\|_\infty$ splňuje předchozí odhad a $\|\tilde{e}^n\|_\infty \leq \max\{|\eta_0^m|, 0 \leq m \leq n; |\eta_j^0|, 0 \leq j \leq J; |\eta_j^m|, 0 \leq m \leq n\}$ dle předchozí věty. \square

Podmínka pro platnost principu maxima je mnohem více omezující než podmínka stability plynoucí z Fourierovy analýzy. Například pro $\theta = \frac{1}{2}$ dostáváme $\mu \leq 1$.

Princip maxima představuje alternativní prostředek pro získání podmínek stability. Oproti Fourierově analýze má tu výhodu, že ho lze snadno aplikovat i na úlohy s nekonzstantními koeficienty. Avšak snadné je odvodit pouze postačující podmínky stability.

4.11 Obecnější okrajové podmínky¹

Nahraďme Dirichletovu okrajovou podmínku v bodě $x = 0$ okrajovou podmínkou

$$(95) \quad u_x(0, t) = \alpha(t) u(0, t) + g(t) \quad \forall t > 0,$$

kde $\alpha(t) \geq 0$.

Nejjednodušší aproximace okrajové podmínky (95) v čase $t = t_n$ je

$$\frac{U_1^n - U_0^n}{h} = \alpha^n U_0^n + g^n, \quad \alpha^n \equiv \alpha(t_n), \quad g^n \equiv g(t_n),$$

z čehož plyne

$$(96) \quad U_0^n = \beta^n U_1^n - \beta^n g^n h, \quad \text{kde} \quad \beta^n = \frac{1}{1 + \alpha^n h}.$$

Nyní můžeme definovat θ -schéma stejným způsobem jako výše. Soustava je opět tridiagonální a má nyní J rovnic. Rovnice (96) je první rovnicí této soustavy.

Z (96) plyne, že v prvním vnitřním bodě je

$$\delta_x^2 U_1^n = U_2^n - 2U_1^n + U_0^n = U_2^n - (2 - \beta^n) U_1^n - \beta^n g^n h.$$

Dostáváme tedy

$$U_1^{n+1} - U_1^n = \mu \theta \{U_2^{n+1} - (2 - \beta^{n+1}) U_1^{n+1} - \beta^{n+1} g^{n+1} h\} \\ + \mu (1 - \theta) \{U_2^n - (2 - \beta^n) U_1^n - \beta^n g^n h\},$$

¹Tato část nebyla odpřednášena a nebude zkoušena

z čehož plyne

$$[1 + \theta \mu (2 - \beta^{n+1})] U_1^{n+1} = [1 - (1 - \theta) \mu (2 - \beta^n)] U_1^n + \theta \mu U_2^{n+1} + (1 - \theta) \mu U_2^n - \mu h [\theta \beta^{n+1} g^{n+1} + (1 - \theta) \beta^n g^n].$$

Definujeme-li obvyklým způsobem chybu diskretizace $\varepsilon_1^{n+1/2}$, platí pro chybu aproximace $e_1^n = U_1^n - u(x_1, t_n)$

$$[1 + \theta \mu (2 - \beta^{n+1})] e_1^{n+1} = [1 - (1 - \theta) \mu (2 - \beta^n)] e_1^n + \theta \mu e_2^{n+1} + (1 - \theta) \mu e_2^n - \tau \varepsilon_1^{n+1/2}.$$

Jelikož se tato rovnice liší od rovnic v ostatních uzlech sítě, nemůžeme použít Fourierovu analýzu chyby. Lze však využít princip maxima, neboť pro $\mu (1 - \theta) \leq \frac{1}{2}$ jsou všechny koeficienty nezáporné a součet koeficientů napravo není větší než koeficient nalevo. Můžeme proto odhadnout chybu aproximace pomocí chyby diskretizace stejně jako výše.

Zbývá odhadnout $\varepsilon_1^{n+1/2}$. Uvažujme případ $\theta = 0$ (explicitní metoda). Pak

$$\begin{aligned} \frac{U_1^{n+1} - U_1^n}{\tau} &= \frac{\delta_x^2 U_1^n}{h^2} + \frac{1}{h^2} [-U_0^n + \beta^n U_1^n - \beta^n g^n h] \\ &= \frac{\delta_x^2 U_1^n}{h^2} + \frac{\beta^n}{h} \left[\frac{U_1^n - U_0^n}{h} - \alpha^n U_0^n - g^n \right]. \end{aligned}$$

Tedy (při označení $u_j^n = u(x_j, t_n)$)

$$\begin{aligned} \varepsilon_1^{n+1/2} &= \frac{u_1^{n+1} - u_1^n}{\tau} - \frac{\delta_x^2 u_1^n}{h^2} - \frac{\beta^n}{h} \left[\frac{\Delta_{-x} u_1^n}{h} - \alpha^n u_0^n - g^n \right] \\ &= [u_t + \frac{1}{2} u_{tt} \tau + \dots](x_1, t_n) - [u_{xx} + \frac{1}{12} u_{xxxx} h^2 + \dots](x_1, t_n) \\ &\quad - \frac{\beta^n}{h} [u_x + \frac{1}{2} u_{xx} h + \dots - \alpha u - g](x_0, t_n) \\ &\approx -\frac{1}{2} \beta^n u_{xx}(x_0, t_n). \end{aligned}$$

Chyba diskretizace tudíž nekonverguje k nule. Důkaz konvergence chyby aproximace lze sice zachránit, avšak diskrétní řešení je zatíženo poměrně velkou chybou.

Zkusme jiný postup. Zavedeme fiktivní hodnotu U_{-1}^n vně $[0, 1]$, takže okrajovou podmínku (95) můžeme aproximovat vztahem

$$(97) \quad \frac{U_1^n - U_{-1}^n}{2h} = \alpha^n U_0^n + g^n.$$

V bodě $x = 0$ aproximujeme rovnici jako ve vnitřních bodech a za U_{-1}^n dosadíme z (97). Pak

$$\begin{aligned} U_0^{n+1} - U_0^n &= \mu \theta [U_1^{n+1} - 2U_0^{n+1} + (U_1^{n+1} - 2h\alpha^{n+1}U_0^{n+1} - 2hg^{n+1})] \\ &\quad + \mu (1 - \theta) [U_1^n - 2U_0^n + (U_1^n - 2h\alpha^n U_0^n - 2hg^n)], \end{aligned}$$

z čehož plyne

$$\begin{aligned} [1 + 2\theta \mu (1 + \alpha^{n+1} h)] U_0^{n+1} &= [1 - 2(1 - \theta) \mu (1 + \alpha^n h)] U_0^n \\ &\quad + 2\theta \mu U_1^{n+1} + 2(1 - \theta) \mu U_1^n - 2\mu h [\theta g^{n+1} + (1 - \theta) g^n]. \end{aligned}$$

Pokud $\mu (1 - \theta) (1 + \alpha^n h) \leq \frac{1}{2}$, je možno chybu aproximace opět odhadnout pomocí chyby diskretizace postupem založeným na principu maxima. Chyba diskretizace je v tomto případě $\varepsilon_0^{n+1/2} = O(\tau + h)$.

4.12 Obecnější lineární rovnice²

Uvažujme nejdříve rovnici

$$u_t = b u_{xx} \quad \forall t > 0, \quad x \in (0, 1),$$

kde $b = b(x, t) > 0$. Explicitnímu schématu (64) pak odpovídá diskretizace

$$U_j^{n+1} = U_j^n + \frac{\tau}{h^2} b_j^n (U_{j+1}^n - 2U_j^n + U_{j-1}^n),$$

kde $b_j^n = b(x_j, t_n)$. Stejně jako dříve získáme

$$\varepsilon_{h,\tau}(x, t) = \frac{1}{2} u_{tt} \tau - \frac{1}{12} b(x, t) u_{xxxx} h^2 + \dots$$

Konvergenci lze dokázat stejným způsobem jako pro $b = 1$, ale podmínku stability je třeba nahradit podmínkou

$$\frac{\tau}{h^2} b(x, t) \leq \frac{1}{2}.$$

Odhad chyby pak je

$$|U_j^n - u(x_j, t_n)| \leq T \left(\frac{1}{2} M_1 \tau + \frac{1}{12} B M_2 h^2 \right),$$

kde $B \geq b(x, t) \forall (x, t) \in [0, 1] \times [0, T]$.

θ -schéma lze definovat různými způsoby. Jednou možností je uvažovat

$$U_j^{n+1} - U_j^n = \frac{\tau}{h^2} b^* [\theta \delta_x^2 U_j^{n+1} + (1 - \theta) \delta_x^2 U_j^n],$$

kde b^* je nějaká vhodná hodnota. Nabízí se položit $b^* = b_j^{n+1/2}$. Rozvoj chyby diskretizace je pak stejný jako dříve až na přenásobení faktorem b . Rovněž konvergenci lze dokázat jako dříve pomocí principu maxima, avšak potřebujeme, aby

$$\frac{\tau}{h^2} (1 - \theta) b(x, t) \leq \frac{1}{2}.$$

Je též možné položit $b^* = \frac{1}{2}(b_j^{n+1} + b_j^n)$, což nezhorší odhad chyby diskretizace, neboť $b^* = [b + \frac{1}{4} b_{tt} \tau^2 + \dots](x_j, t_{n+1/2})$.

Nejobecnější tvar lineární parabolické rovnice druhého řádu je

$$(98) \quad u_t = b u_{xx} - a u_x + c u + d \quad \forall t > 0, \quad x \in (0, 1),$$

kde $a = a(x, t)$, $b = b(x, t)$, $c = c(x, t)$, $d = d(x, t)$ jsou dané funkce, přičemž $b > 0$. Explicitní schéma je přirozené uvažovat ve tvaru

$$(99) \quad \frac{U_j^{n+1} - U_j^n}{\tau} = b_j^n \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} - a_j^n \frac{U_{j+1}^n - U_{j-1}^n}{2h} + c_j^n U_j^n + d_j^n.$$

²Tato část nebyla přednášena a nebude zkoušena

Označíme-li

$$\mu_j^n = \frac{\tau}{h^2} b_j^n, \quad \nu_j^n = \frac{\tau}{h} a_j^n,$$

zjistíme, že chyba aproximace splňuje

$$e_j^{n+1} = (1 - 2\mu_j^n + \tau c_j^n) e_j^n + (\mu_j^n - \frac{1}{2}\nu_j^n) e_{j+1}^n + (\mu_j^n + \frac{1}{2}\nu_j^n) e_{j-1}^n - \tau \varepsilon_j^n.$$

Abychom při odhadu chyby mohli postupovat jako dříve, musíme zajistit, že koeficienty jsou nezáporné a jejich součet není větší než 1. To vyžaduje

$$(100) \quad \frac{1}{2} |\nu_j^n| \leq \mu_j^n, \quad 2\mu_j^n - \tau c_j^n \leq 1, \quad c_j^n \leq 0.$$

Speciálně (dle první podmínky) musí být

$$h \frac{|a_j^n|}{2b_j^n} \leq 1$$

a toto omezení implikuje omezení τ prostřednictvím druhé podmínky:

$$\tau \leq \frac{h^2}{2b_j^n - h^2 c_j^n}.$$

V mnoha úlohách z praxe je $|a_j^n| \gg b_j^n$, což vyžaduje velmi malé prostorové a časové kroky.

Jednoduchý způsob, jak tento problém napravit, je použít aproximace

$$u_x(x_j, t_n) \approx \begin{cases} \frac{U_j^n - U_{j-1}^n}{h} & \text{je-li } a(x_j, t_n) \geq 0, \\ \frac{U_{j+1}^n - U_j^n}{h} & \text{je-li } a(x_j, t_n) < 0. \end{cases}$$

Funkci a můžeme interpretovat jako rychlost látky, v níž sledujeme rozložení veličiny u , ve směru kladné x -ové poloosy. K diskretizaci $u_x(x_j, t_n)$ tedy využíváme hodnoty u z té strany, odkud se do bodu x_j v čase t_n látka pohybuje. Hovoříme proto o diskretizaci typu *upwind*.

Předpokládejme pro jednoduchost, že $a(x, t) \geq 0$ a $c(x, t) = 0$. Explicitní schéma má pak tvar

$$\frac{U_j^{n+1} - U_j^n}{\tau} = b_j^n \frac{U_{j+1}^n - 2U_j^n + U_{j-1}^n}{h^2} - a_j^n \frac{U_j^n - U_{j-1}^n}{2h} + d_j^n,$$

což dává

$$e_j^{n+1} = (1 - 2\mu_j^n - \nu_j^n) e_j^n + \mu_j^n e_{j+1}^n + (\mu_j^n + \nu_j^n) e_{j-1}^n - \tau \varepsilon_j^n.$$

Aby všechny koeficienty na pravé straně byly nezáporné, potřebujeme nyní pouze podmínku $2\mu_j^n + \nu_j^n \leq 1$. Stabilita tedy nevyžaduje žádné omezení prostorového kroku h . Cenou za to je, že nyní máme pouze $\varepsilon_j^n = O(h + \tau)$.

Někdy se můžeme setkat s parabolickou rovnicí v samoadjungovaném tvaru

$$u_t = (p(x, t) u_x)_x \quad \forall t > 0, \quad x \in (0, 1),$$

kde $p > 0$. Rozderivováním můžeme tuto rovnici převést do tvaru (98), ale obvykle je výhodnější zkonstruovat diferenční aproximaci původního samoadjungovaného tvaru:

$$\begin{aligned} [(p u_x)_x](x_j, t_n) &\approx \frac{1}{h} [(p u_x)(x_{j+1/2}, t_n) - (p u_x)(x_{j-1/2}, t_n)] \\ &\approx \frac{1}{h^2} [p_{j+1/2}^n (u_{j+1}^n - u_j^n) - p_{j-1/2}^n (u_j^n - u_{j-1}^n)]. \end{aligned}$$

Explicitní diferenční schéma má tedy tvar

$$\frac{U_j^{n+1} - U_j^n}{\tau} = \frac{p_{j+1/2}^n (U_{j+1}^n - U_j^n) - p_{j-1/2}^n (U_j^n - U_{j-1}^n)}{h^2},$$

a tudíž

$$U_j^{n+1} = [1 - \mu (p_{j+1/2}^n + p_{j-1/2}^n)] U_j^n + \mu p_{j+1/2}^n U_{j+1}^n + \mu p_{j-1/2}^n U_{j-1}^n, \quad \mu = \frac{\tau}{h^2}.$$

Chybu aproximace můžeme tudíž vyšetřovat stejným způsobem jako dříve, budou-li všechny koeficienty nezáporné, což je splněno, pokud $\mu P \leq \frac{1}{2}$, kde P splňuje $p(x, t) \leq P$ v uvažované oblasti. Máme tedy omezení stejného typu jako dříve.

Zřejmým způsobem lze na výše uvažované obecnější rovnice zobecnit θ -schéma vedoucí k implicitnímu schématu.

5 Parabolické rovnice ve dvou prostorových dimenzích³

Nechť $\Omega \subset \mathbb{R}^2$ je omezená oblast. Hledáme funkci $u = u(x, y, t)$ definovanou pro $(x, y) \in \bar{\Omega}$ a $t \geq 0$ takovou, že

$$(101) \quad u_t = \Delta u \equiv u_{xx} + u_{yy} \quad \forall t > 0, \quad (x, y) \in \Omega,$$

$$(102) \quad u(x, y, t) = u^b(x, y, t) \quad \forall t > 0, \quad (x, y) \in \partial\Omega,$$

$$(103) \quad u(x, y, 0) = u^0(x, y) \quad \forall (x, y) \in \bar{\Omega},$$

kde u^b a u^0 jsou zadané funkce.

Předpokládejme nejdříve, že $\Omega = (0, X) \times (0, Y)$, zvolme $J_x, J_y \in \mathbb{N}$ a položme $h_x = X/J_x$, $h_y = Y/J_y$. Oblast Ω pokryjeme rovnoměrnou pravoúhloú sítí s krokem dělení h_x v x -ovém směru a h_y v y -ovém směru. Uzly prostorové sítě jsou $(x_r, y_s) = (r h_x, s h_y)$, kde $r = 0, 1, \dots, J_x$ a $s = 0, 1, \dots, J_y$. Přibližné řešení značíme

$$U_{r,s}^n \approx u(x_r, y_s, t_n), \quad r = 0, 1, \dots, J_x, \quad s = 0, 1, \dots, J_y, \quad n = 0, 1, 2, \dots$$

Nejjednodušší explicitní diferenční schéma je

$$\frac{U_{r,s}^{n+1} - U_{r,s}^n}{\tau} = \frac{\delta_x^2 U_{r,s}^n}{h_x^2} + \frac{\delta_y^2 U_{r,s}^n}{h_y^2}, \quad r = 1, \dots, J_x - 1, \quad s = 1, \dots, J_y - 1.$$

³Tato část nebyla odpřednášena a nebude zkoušena

Hodnota $U_{r,s}^{n+1}$ je určena hodnotami $U_{r,s}^n$, $U_{r+1,s}^n$, $U_{r-1,s}^n$, $U_{r,s+1}^n$, $U_{r,s-1}^n$; hovoříme o tzv. pětibodovém schématu.

Vlastnosti schématu lze analyzovat analogicky jako v jedné dimenzi. Chyba diskretizace je

$$\varepsilon(x, t) = \frac{1}{2} u_{tt} \tau - \frac{1}{12} [u_{xxxx} h_x^2 + u_{yyyy} h_y^2] + \dots$$

Za předpokladu omezenosti uvedených derivací a při

$$\mu_x + \mu_y \leq \frac{1}{2}, \quad \text{kde} \quad \mu_x = \frac{\tau}{h_x^2}, \quad \mu_y = \frac{\tau}{h_y^2},$$

dostaneme stejně jako v jedné dimenzi pro chybu aproximace

$$\|e^n\|_\infty \leq T \left(\frac{1}{2} M_1 \tau + \frac{1}{12} M_2^x h_x^2 + \frac{1}{12} M_2^y h_y^2 \right).$$

Lze rovněž aplikovat Fourierovu analýzu stability. Pro $U_{r,s}^n = \lambda^n e^{i[k_x r h_x + k_y s h_y]}$ dostaneme amplifikační faktor

$$\lambda = \lambda(\mathbf{k}) = 1 - 4 \left[\mu_x \sin^2 \frac{k_x h_x}{2} + \mu_y \sin^2 \frac{k_y h_y}{2} \right],$$

kde $\mathbf{k} = (k_x, k_y)$. Vidíme, že $|\lambda(\mathbf{k})| \leq 1 \forall \mathbf{k}$ právě tehdy, když $\mu_x + \mu_y \leq \frac{1}{2}$.

Zřejmým způsobem můžeme též rozšířit do dvou dimenzí θ -metodu. Speciálně metoda Crankova–Nicolsonové bude

$$\left(1 - \frac{1}{2} \mu_x \delta_x^2 - \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^{n+1} = \left(1 + \frac{1}{2} \mu_x \delta_x^2 + \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^n.$$

Soustava rovnic již není tridiagonální a její vyřešení je podstatně dražší než výpočet hodnot $U_{r,s}^{n+1}$ u explicitní metody. Hledáme proto jiné možnosti, jak získat nepodmíněně stabilní metodu.

5.1 Metoda střídavých směrů

Uvažujme následující modifikaci schématu Crankova–Nicolsonové:

$$(104) \quad \left(1 - \frac{1}{2} \mu_x \delta_x^2\right) \left(1 - \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^{n+1} = \left(1 + \frac{1}{2} \mu_x \delta_x^2\right) \left(1 + \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^n.$$

Dodatečné členy ve schématu nezhorší chybu diskretizace, neboť

$$\begin{aligned} \frac{1}{4} \mu_x \mu_y \delta_x^2 \delta_y^2 \frac{u(x, y, t + \frac{1}{2} \tau) - u(x, y, t - \frac{1}{2} \tau)}{\tau} &\approx \frac{\tau^2}{4 h_x^2 h_y^2} \delta_x^2 \delta_y^2 u_t(x, y, t) \\ &\approx \frac{1}{4} u_{xyyt}(x, y, t) \tau^2, \end{aligned}$$

a tudíž $\varepsilon_{h_x, h_y, \tau} = O(h_x^2 + h_y^2 + \tau^2)$. Přednost schématu tkví v tom, že výpočet $U_{r,s}^{n+1}$ lze rozložit na řešení problémů s tridiagonálními maticemi. Zavedeme mezivýsledek $U_{r,s}^{n+1/2}$ a (104) zapíšeme v ekvivalentním tvaru

$$(105) \quad \left(1 - \frac{1}{2} \mu_x \delta_x^2\right) U_{r,s}^{n+1/2} = \left(1 + \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^n, \quad r = 1, \dots, J_x - 1, \quad s = 1, \dots, J_y - 1,$$

$$(106) \quad \left(1 - \frac{1}{2} \mu_y \delta_y^2\right) U_{r,s}^{n+1} = \left(1 + \frac{1}{2} \mu_x \delta_x^2\right) U_{r,s}^{n+1/2}, \quad r = 1, \dots, J_x - 1, \quad s = 1, \dots, J_y - 1.$$

Výpočet $U_{r,s}^{n+1/2}$ z (105) představuje vyřešení $J_y - 1$ soustav lineárních rovnic řádu $J_x - 1$. Podobně (106) sestává z $J_x - 1$ soustav lineárních rovnic řádu $J_y - 1$. Každá ze soustav má tridiagonální matici a provedení jednoho časového kroku je tak mnohem rychlejší než řešení soustavy lineárních rovnic odpovídající Crankově–Nicolsonové metodě. Výpočetní náročnost je asi trojnásobná ve srovnání s jedním krokem explicitního schématu. Dosazením Fourierova členu $\lambda^n e^{i[k_x r h_x + k_y s h_y]}$ do (104) získáme

$$\lambda(\mathbf{k}) = \frac{\left(1 - 2\mu_x \sin^2 \frac{k_x h_x}{2}\right) \left(1 - 2\mu_y \sin^2 \frac{k_y h_y}{2}\right)}{\left(1 + 2\mu_x \sin^2 \frac{k_x h_x}{2}\right) \left(1 + 2\mu_y \sin^2 \frac{k_y h_y}{2}\right)}$$

z čehož plyne, že schéma (104) je nepodmíněně stabilní.

6 Numerické řešení eliptických parciálních diferenciálních rovnic 2. řádu

6.1 Diskretizace Poissonovy rovnice

Uvažujme modelovou úlohu

$$(107) \quad -\Delta u = f \quad \text{v } \Omega := (0, 1)^2, \quad u = 0 \quad \text{na } \partial\Omega.$$

Množinu $\bar{\Omega}$ pokryjeme rovnoměrnou čtvercovou sítí s J intervaly v každém směru, čímž vzniknou uzly $(x_r, y_s) := (r h, s h)$, $r, s = 0, \dots, J$, kde $h = 1/J$. Označíme

$$\begin{aligned} N_\Omega &= \{(x_r, y_s); r, s \in \{1, \dots, J-1\}\}, \\ N_{\partial\Omega} &= \{(x_r, 0), (x_r, 1), (0, y_s), (1, y_s); r, s \in \{0, \dots, J\}\}. \end{aligned}$$

Pak $N_\Omega \subset \Omega$ a $N_{\partial\Omega} \subset \partial\Omega$, tj. N_Ω je množina vnitřních uzlů a $N_{\partial\Omega}$ je množina hraničních uzlů. Řešení u úlohy (107) budeme v uzlech $(x_r, y_s) \in N_\Omega \cup N_{\partial\Omega}$ aproximovat hodnotami $U_{r,s}$. Dále zavedeme označení $u_{r,s} := u(x_r, y_s)$ a $f_{r,s} := f(x_r, y_s)$.

K diskretizaci úlohy (107) použijeme v každém vnitřním uzlu $(x_r, y_s) \in N_\Omega$ aproximaci

$$\begin{aligned} (\Delta u)(x_r, y_s) &= u_{xx}(x_r, y_s) + u_{yy}(x_r, y_s) \approx \frac{\delta_x^2 u_{r,s}}{h^2} + \frac{\delta_y^2 u_{r,s}}{h^2} \\ &= \frac{u_{r+1,s} - 2u_{r,s} + u_{r-1,s}}{h^2} + \frac{u_{r,s+1} - 2u_{r,s} + u_{r,s-1}}{h^2}, \end{aligned}$$

což nás vede k následující definici přibližného řešení $\{U_{r,s}\}_{r,s=0}^J$:

$$(108) \quad \frac{4U_{r,s} - U_{r+1,s} - U_{r-1,s} - U_{r,s+1} - U_{r,s-1}}{h^2} = f_{r,s}, \quad r, s = 1, \dots, J-1,$$

$$(109) \quad U_{r,0} = U_{r,J} = U_{0,s} = U_{J,s} = 0, \quad r, s = 0, \dots, J.$$

K určení přibližného řešení je tedy třeba vyřešit soustavu $(J-1)^2$ lineárních algebraických rovnic.

6.2 Konvergence přibližných řešení

Nechť $U = \{U_{r,s}\}_{r,s=0}^J$ je libovolná síťová funkce a definujme operátor L_h vztahem

$$(L_h U)_{r,s} := \frac{4U_{r,s} - U_{r+1,s} - U_{r-1,s} - U_{r,s+1} - U_{r,s-1}}{h^2}, \quad r, s = 1, \dots, J-1.$$

Místo $(L_h U)_{r,s}$ budeme užívat jednodušší značení $L_h U_{r,s}$. Pak úlohu (108)–(109) můžeme ekvivalentně zapsat ve tvaru

$$(110) \quad L_h U_P = f_P \quad \forall P \in N_\Omega, \quad U_P = 0 \quad \forall P \in N_{\partial\Omega},$$

kde pro uzly sítě nyní používáme značení P místo (x_r, y_s) . Později ukážeme (viz důsledek 1 na str. 61), že operátor L_h splňuje diskrétní princip maxima

$$(111) \quad L_h U_P \leq 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \max_{P \in N_\Omega} U_P \leq \max_{Q \in N_{\partial\Omega}} U_Q.$$

Podobně jako u evolučních úloh můžeme princip maxima využít k odhadu chyby aproximace. Nejprve definujeme chybu diskretizace

$$\varepsilon_{r,s} := L_h u_{r,s} - f_{r,s}, \quad r, s = 1, \dots, J-1,$$

což můžeme ekvivalentně zapsat ve tvaru

$$(112) \quad \varepsilon_P := L_h u_P - f_P \quad \forall P \in N_\Omega.$$

Použitím Taylorova rozvoje získáme pro libovolné $P \in N_\Omega$

$$(113) \quad |\varepsilon_P| \leq \frac{1}{12} (M_1 + M_2) h^2, \quad \text{kde} \quad M_1 = \max_{\bar{\Omega}} |u_{xxxx}|, \quad M_2 = \max_{\bar{\Omega}} |u_{yyyy}|.$$

Pomocí (112), (110) a (107) zjišťujeme, že chyba aproximace $e_P := U_P - u_P$ splňuje

$$(114) \quad L_h e_P = -\varepsilon_P \quad \forall P \in N_\Omega, \quad e_P = 0 \quad \forall P \in N_{\partial\Omega}.$$

Vztah (114) nám umožňuje pomocí diskrétního principu maxima (111) a odhadu chyby diskretizace (113) odvodit odhad chyby aproximace. Za tím účelem je potřeba definovat vhodnou funkci splňující levou stranu implikace (111), k čemuž podobně jako u odvození odhadu chyby aproximace pro přibližné řešení transportní rovnice využijeme tzv. *srovnávací funkci* Φ . V tomto případě ji můžeme definovat vztahem

$$\Phi(x, y) := (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2.$$

Funkci Φ přiřadíme síťovou funkci $\{\Phi_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$, kde $\Phi_P = \Phi(P)$. Jelikož funkce Φ má nulové čtvrté derivace, plyne ze (113), že

$$L_h \Phi_P = (-\Delta\Phi)(P) = -4 \quad \forall P \in N_\Omega.$$

Označme

$$\Psi_P := e_P + \frac{1}{4} \frac{h^2}{12} (M_1 + M_2) \Phi_P.$$

Pak

$$L_h \Psi_P = L_h e_P - \frac{h^2}{12} (M_1 + M_2) = -\varepsilon_P - \frac{h^2}{12} (M_1 + M_2) \leq 0 \quad \forall P \in N_\Omega$$

a podle (111) je

$$\Psi_P \leq \frac{1}{4} \frac{h^2}{12} (M_1 + M_2) \max_{Q \in N_{\partial\Omega}} \Phi_Q = \frac{1}{8} \frac{h^2}{12} (M_1 + M_2) \quad \forall P \in N_\Omega.$$

Jelikož $e_P \leq \Psi_P$, dostáváme

$$U_P - u_P \leq \frac{1}{96} (M_1 + M_2) h^2.$$

Označíme-li $\Psi_P := -e_P + \frac{1}{4} \frac{h^2}{12} (M_1 + M_2) \Phi_P$, získáme stejný odhad pro $-(U_P - u_P)$. Je tedy

$$|U_{r,s} - u(x_r, y_s)| \leq \frac{1}{96} (M_1 + M_2) h^2 \quad \forall r, s \in \{0, \dots, J\}.$$

6.3 Obecnější rovnice difúze

Uvažujme úlohu

$$(115) \quad -\operatorname{div}(a \nabla u) = f \quad \text{v } \Omega, \quad \alpha_0 u + \alpha_1 \frac{\partial u}{\partial n} = g \quad \text{na } \partial\Omega,$$

kde Ω je omezená oblast v \mathbb{R}^2 , n je vnější jednotková normála k $\partial\Omega$, a je hladká funkce splňující $a(x, y) \geq a_0 > 0$ a α_0, α_1 jsou konstanty splňující

$$\alpha_0 \geq 0, \quad \alpha_1 \geq 0, \quad \alpha_0 + \alpha_1 > 0.$$

Rovnice (115) popisuje difúzi veličiny u v nehomogenním izotropním prostředí. Pokud je $a(x, y) = \varepsilon$, jedná se o difúzi v homogenním izotropním prostředí a rovnice (115) se redukuje na rovnici (107).

Oblast Ω pokryjeme pravidelnou obdélníkovou sítí s krokem h_x ve směru osy x a s krokem h_y ve směru osy y . Uzly sítě, které leží vedle sebe na některé ze sítových přímk nazýváme sousední. Každý uzel má tedy čtyři sousední uzly. Uzly ležící v Ω , jejichž všichni čtyři sousedé leží v $\bar{\Omega}$ nazýváme *regulární uzly*. Uzly ležící v Ω , jejichž některý soused neleží v $\bar{\Omega}$ nazýváme *neregulární uzly*. Pokud Ω je mnohoúhelník, jehož hrany jsou rovnoběžné se souřadnými osami, můžeme zkonstruovat síť tak, aby neobsahovala neregulární uzly. U oblastí s komplikovanější hranicí to však obecně nelze.

Pro diskretizaci úlohy (115) se nabízí levou stranu parciální diferenciální rovnice rozderivovat a derivace funkce u aproximovat diferenčními kvocienty uvažovanými dříve. Označíme-li $b = -a_x$, $c = -a_y$, můžeme rovnici (115) zapsat ve tvaru

$$(116) \quad -a \Delta u + b u_x + c u_y = f \quad \text{v } \Omega$$

a v regulárních uzlech (x_r, y_s) aproximovat použitím centrálních diferencí, čímž získáme

$$(117) \quad -a_{r,s} \left(\frac{\delta_x^2 U_{r,s}}{h_x^2} + \frac{\delta_y^2 U_{r,s}}{h_y^2} \right) + b_{r,s} \frac{\Delta_{0x} U_{r,s}}{h_x} + c_{r,s} \frac{\Delta_{0y} U_{r,s}}{h_y} = f_{r,s}.$$

Snadno lze ukázat, že chyba diskretizace je druhého řádu v h_x a h_y . Později uvidíme, že k platnosti diskrétního principu maxima je potřeba, aby koeficienty schématu (117) v uzlech sousedících s uzlem (x_r, y_s) byly nekladné. Avšak koeficient u $U_{r-1,s}$ je roven $-(2a_{r,s} + b_{r,s}h_x)/(2h_x^2)$ a koeficient u $U_{r+1,s}$ je $-(2a_{r,s} - b_{r,s}h_x)/(2h_x^2)$. Musí tedy být $|b_{r,s}|h_x \leq 2a_{r,s}$ a podobně $|c_{r,s}|h_y \leq 2a_{r,s}$. Z toho tedy plyne, že tam, kde a je malé, ale rychle se mění (a má tudíž velké první derivace), musíme použít dostatečně jemnou síť.

Je proto výhodnější diskretizovat přímo rovnici (115), která má podle definice divergence tvar

$$-(a u_x)_x - (a u_y)_y = f \quad \text{v } \Omega.$$

Ukažme si podrobně aproximaci výrazu $(a u_x)_x$ v regulárním uzlu (x_r, y_s) . Označíme-li $x_{r \pm \frac{1}{2}} := x_r \pm \frac{1}{2} h_x$, pak

$$\begin{aligned} (a u_x)_x(x_r, y_s) &\approx \frac{(a u_x)(x_{r+\frac{1}{2}}, y_s) - (a u_x)(x_{r-\frac{1}{2}}, y_s)}{h_x} \\ &\approx \frac{1}{h_x} \left(a_{r+\frac{1}{2},s} \frac{u_{r+1,s} - u_{r,s}}{h_x} - a_{r-\frac{1}{2},s} \frac{u_{r,s} - u_{r-1,s}}{h_x} \right) \\ &\approx \frac{1}{h_x^2} \left(a_{r+\frac{1}{2},s} (U_{r+1,s} - U_{r,s}) - a_{r-\frac{1}{2},s} (U_{r,s} - U_{r-1,s}) \right). \end{aligned}$$

Analogicky postupujeme u výrazu $(a u_y)_y$, čímž získáme pro aproximaci rovnice (115) schéma

$$(118) \quad \begin{aligned} &-\frac{1}{h_x^2} \left(a_{r+\frac{1}{2},s} (U_{r+1,s} - U_{r,s}) - a_{r-\frac{1}{2},s} (U_{r,s} - U_{r-1,s}) \right) \\ &-\frac{1}{h_y^2} \left(a_{r,s+\frac{1}{2}} (U_{r,s+1} - U_{r,s}) - a_{r,s-\frac{1}{2}} (U_{r,s} - U_{r,s-1}) \right) = f_{r,s}. \end{aligned}$$

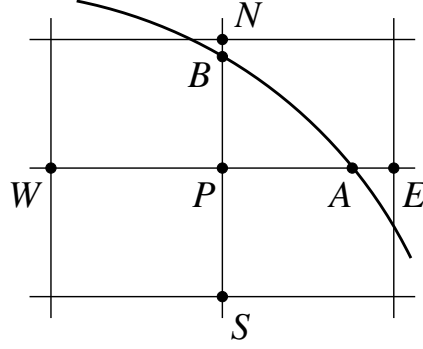
Toto schéma lze zapsat v kompaktním tvaru

$$-\left(\frac{\delta_x(a \delta_x U)}{h_x^2} + \frac{\delta_y(a \delta_y U)}{h_y^2} \right)_{r,s} = f_{r,s}.$$

Je vidět, že koeficienty schématu (118) jsou v uzlech sousedících s uzlem (x_r, y_s) nekladné bez jakéhokoli omezení na síť.

6.4 Diskretizace úlohy (115) v neregulárních uzlech

V této části si ukážeme, jakým způsobem lze diferenční schémata (117) a (118) modifikovat v neregulárních uzlech.



Obrázek 4: Neregulární uzel v blízkosti zakřivené hranice.

Uvažujme nejprve situaci znázorněnou v obr. 4. Uzel P leží v oblasti Ω v blízkosti zakřivené hranice, přičemž jeho sousední uzly S a W leží v $\bar{\Omega}$, zatímco zbývající dva sousední uzly E a N v $\bar{\Omega}$ neleží. Průsečík hranice s úsečkou PE je označen A a průsečík hranice s úsečkou PN je označen B . Předpokládejme, že na uvažované části hranice je předepsána Dirichletova okrajová podmínka (tj. v (115) je $\alpha_0 > 0$ a $\alpha_1 = 0$). V bodech A a B tudíž známe hodnoty řešení u a tedy i přibližného řešení U .

Nechť $|PA| = \theta h_x$, $\theta \in (0, 1)$. Pak hodnoty funkce u v bodech A a W můžeme vyjádřit pomocí následujících Taylorových rozvojų:

$$u_A = \sum_{k \geq 0} \frac{1}{k!} \frac{\partial^k u}{\partial x^k}(P) (\theta h_x)^k, \quad u_W = \sum_{k \geq 0} \frac{1}{k!} \frac{\partial^k u}{\partial x^k}(P) (-h_x)^k.$$

Tedy

$$\begin{aligned} u_A + \theta u_W &= (1 + \theta) u_P + \frac{1}{2} \theta (1 + \theta) u_{xx}(P) h_x^2 + O(\theta h_x^3), \\ u_A - \theta^2 u_W &= (1 - \theta^2) u_P + \theta (1 + \theta) u_x(P) h_x + O(\theta^2 h_x^3), \end{aligned}$$

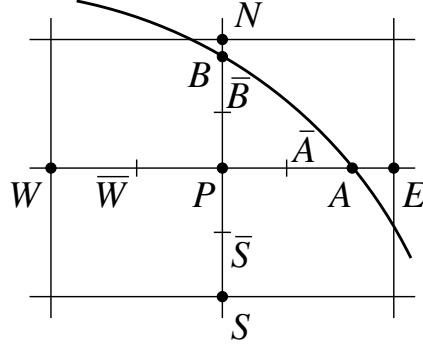
z čehož plyne

$$(119) \quad u_x(P) = \frac{u_A - \theta^2 u_W - (1 - \theta^2) u_P}{\theta (1 + \theta) h_x} + O(h_x^2),$$

$$(120) \quad u_{xx}(P) = \frac{u_A + \theta u_W - (1 + \theta) u_P}{\frac{1}{2} \theta (1 + \theta) h_x^2} + O(h_x).$$

Podobně můžeme postupovat pro body S , P , B a odvodit aproximace derivací $u_y(P)$ a $u_{yy}(P)$. Dosazením do (116) získáme modifikaci schématu (117) v neregulárním uzlu P . Pro platnost principu maxima musíme samozřejmě obecně opět požadovat omezení kroku sítě, jako tomu bylo v regulárních uzlech.

Aproximujeme-li přímo rovnici (115), můžeme i v neregulárním uzlu postupovat analogicky jako při odvození schématu (118). Označíme \bar{A} , \bar{B} , \bar{W} a \bar{S} středy úseček PA , PB ,



Obrázek 5: Další značení v okolí neregulárního uzlu.

PW a PS , viz obr. 5, a uvažujeme aproximace

$$\begin{aligned}
 (a u_x)_x(P) &\approx \frac{(a u_x)(\bar{A}) - (a u_x)(\bar{W})}{|\bar{A} - \bar{W}|} \\
 &\approx \frac{1}{\frac{1}{2}|\bar{A}\bar{W}|} \left(a_{\bar{A}} \frac{u_A - u_P}{|\bar{A}P|} - a_{\bar{W}} \frac{u_P - u_W}{|\bar{W}P|} \right) \\
 &\approx \frac{1}{\frac{1}{2}(1 + \theta) h_x} \left(a_{\bar{A}} \frac{U_A - U_P}{\theta h_x} - a_{\bar{W}} \frac{U_P - U_W}{h_x} \right).
 \end{aligned}$$

Analogicky lze aproximovat $(a u_y)_y$. Označíme-li $\phi := |PB|/h_y$, dostáváme v bodě P síťovou rovnici

$$\begin{aligned}
 (121) \quad &-\frac{1}{\frac{1}{2}\theta(1 + \theta)h_x^2} \left(a_{\bar{A}}(U_A - U_P) - \theta a_{\bar{W}}(U_P - U_W) \right) \\
 &-\frac{1}{\frac{1}{2}\phi(1 + \phi)h_y^2} \left(a_{\bar{B}}(U_B - U_P) - \phi a_{\bar{S}}(U_P - U_S) \right) = f_P.
 \end{aligned}$$

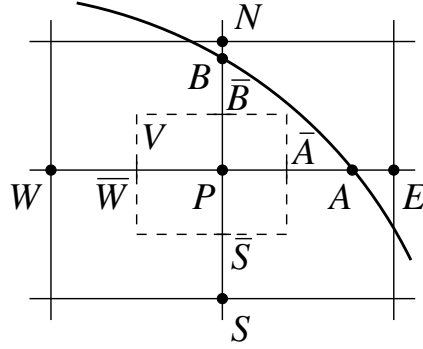
Pro $\theta = \phi = 1$ získáváme rovnici (118). Je-li $a = \text{const.}$, je síťová rovnice shodná se síťovou rovnicí získanou pomocí (120). Výhodou (121) však je, že v obecném případě koeficienty splňují podmínky požadované pro princip maxima, aniž by bylo nutno omezit velikost kroku sítě.

Rovnici (121) lze též odvodit pomocí integrálního tvaru rovnice (115). Nechť V je obdélník (tzv. *kontrolní objem*), jehož hrany jsou rovnoběžné se síťovými přímkami a procházejí body \bar{A} , \bar{B} , \bar{W} a \bar{S} , viz obr. 6. Integrací rovnice (115) přes kontrolní objem V a aplikací Gaussovy věty o integraci získáme

$$(122) \quad \int_V f \, dx = - \int_V \text{div}(a \nabla u) \, dx = - \int_{\partial V} a \frac{\partial u}{\partial n} \, d\sigma,$$

kde n je vnější jednotková normála k ∂V . Pro levou stranu identity (122) použijeme aproximaci

$$\int_V f \, dx \approx |V| f_P = \frac{1}{4} (1 + \theta)(1 + \phi) h_x h_y f_P.$$



Obrázek 6: Konstrukce obdélníku V .

Integrál na pravé straně identity (122) rozdělíme na integrály přes jednotlivé hrany obdélníku V a každý z nich aproximujeme zvlášť. Např. integrál přes hranu V procházející bodem \bar{S} můžeme aproximovat výrazem

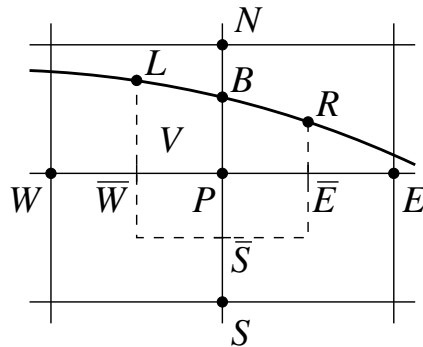
$$\frac{1}{2} (1 + \theta) h_x a_{\bar{S}} \frac{\partial u}{\partial n}(\bar{S}) \approx \frac{1}{2} (1 + \theta) h_x a_{\bar{S}} \frac{U_S - U_P}{h_y}.$$

Aproximujeme-li obdobně integrály přes ostatní hrany V , získáme síťovou rovnici (121).

Použití integrálního tvaru (122) pro odvození síťové rovnice v neregulárním uzlu je zejména výhodné, pokud je na části hranice v blízkosti tohoto uzlu předepsána okrajová podmínka obsahující derivaci, neboť pak je odvození vhodné diskretizace podstatně komplikovanější než v dirichletovském případě. Pro jednoduchost předpokládejme, že hranice oblasti Ω v blízkosti neregulárního uzlu P protíná pouze svislé síťové přímky (viz obr. 7) a že je na této části hranice předepsána Neumannova okrajová podmínka

$$\frac{\partial u}{\partial n} = g$$

(tj. v (115) je $\alpha_0 = 0$ a $\alpha_1 = 1$). Kontrolní objem V nyní definujeme tak, že část jeho hranice je tvořena částí LBR hranice oblasti Ω a zbývající část jeho hranice sestává z úseček kolmých k síťovým přímkám a protínajících úsečky PE , PW a PS v jejich



Obrázek 7: Konstrukce kontrolního objemu V v blízkosti neumannovské hranice.

středech \overline{E} , \overline{W} a \overline{S} . Označíme $|\overline{WL}| = \phi_1 h_y$, $|\overline{ER}| = \phi_2 h_y$ a $|LR| = \psi h_x$. V integrálním tvaru (122) aproximujeme levou stranu vztahem

$$\int_V f \, dx \approx \frac{1}{2} (1 + \phi_1 + \phi_2) h_x h_y f_P,$$

integrály přes rovné části ∂V aproximujeme analogicky jako výše a pro integrál přes oblouk LBR použijeme vztah

$$\int_{LBR} a \frac{\partial u}{\partial n} \, d\sigma \approx a_B g_B \psi h_x.$$

To vede k diferenciální rovnici

$$(123) \quad -\frac{2\phi_2 + 1}{(1 + \phi_1 + \phi_2) h_x^2} a_{\overline{E}} (U_E - U_P) - \frac{2\phi_1 + 1}{(1 + \phi_1 + \phi_2) h_x^2} a_{\overline{W}} (U_W - U_P) \\ - \frac{2}{(1 + \phi_1 + \phi_2) h_y^2} a_{\overline{S}} (U_S - U_P) = f_P + \frac{2\psi a_B g_B}{(1 + \phi_1 + \phi_2) h_y}.$$

Pokud hranice oblasti Ω protíná síťové přímky v blízkosti uzlu P tak jako na obr. 4, je základní myšlenka odvození síťové rovnice v uzlu P stejná jako výše, avšak konstrukce kontrolního objemu V a aproximace integrálu přes jeho hranici je komplikovanější.

6.5 Princip maxima pro diskretizace eliptických úloh

Předpokládejme, že jsme diskretizací úlohy (107) či jiné lineární eliptické úlohy získali v každém bodě $P \in N_\Omega$ síťovou rovnici

$$(124) \quad L_h U_P = f_P + q_P,$$

kde q_P je určeno okrajovými podmínkami jinými než dirichletovskými. Příkladem síťové rovnice s nenulovým členem q_P je rovnice (123). Množina N_Ω obsahuje všechny uzly ležící v oblasti Ω , ale může obsahovat i některé body na hranici Ω , v nichž je předepsána jiná okrajová podmínka než dirichletovská. Dále zavádíme množinu $N_{\partial\Omega}$ sestávající z bodů na hranici Ω , v nichž je předepsána Dirichletova okrajová podmínka. Body množiny $N_{\partial\Omega}$ mohou náležet mezi uzly sítě pokrývající oblast Ω nebo mohou být průsečíky síťových přímek s hranicí $\partial\Omega$ jako např. body A a B v obr. 4. Přibližné řešení $\{U_P\}$ v rovnici (124) představuje množinu hodnot v bodech množiny $N_\Omega \cup N_{\partial\Omega}$.

Činíme následující předpoklady o operátoru L_h a množinách N_Ω a $N_{\partial\Omega}$:

(P1) Pro každé $P \in N_\Omega$ definujeme množinu $N_P \subset (N_\Omega \cup N_{\partial\Omega}) \setminus \{P\}$ a předpokládáme, že

$$(125) \quad L_h U_P = c_P U_P - \sum_{Q \in N_P} c_{PQ} U_Q,$$

kde

$$(126) \quad c_{PQ} > 0 \quad \forall Q \in N_P, \quad c_P \geq \sum_{Q \in N_P} c_{PQ}.$$

(P2) Množina N_Ω je souvislá v tom smyslu, že

$$\forall P, Q \in N_\Omega \quad \exists P_1, \dots, P_m \in N_\Omega : \\ P_1 \in N_P, P_2 \in N_{P_1}, \dots, P_m \in N_{P_{m-1}}, Q \in N_{P_m}.$$

(P3) $\exists P \in N_\Omega, Q \in N_{\partial\Omega} : Q \in N_P$.

Lze snadno ověřit, že rovnice (108), (118), (121) a (123) splňují (125) a (126) (místo druhé nerovnosti v (126) platí ve všech čtyřech případech rovnost). Pro diskretizaci (108) Poissonovy rovnice je splněn též předpoklad (P2) o souvislosti množiny N_Ω . Předpoklad (P3) implikuje, že na části hranice oblasti Ω je předepsána Dirichletova okrajová podmínka. Pro diskretizace (108) i (121) je předpoklad (P3) zřejmě splněn.

Nyní zformulujeme a dokážeme diskrétní princip maxima pro výše uvedený operátor L_h .

Věta 10 *Nechť jsou splněny předpoklady (P1)–(P3). Pak pro libovolnou síťovou funkci $\{U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ platí*

$$(127) \quad L_h U_P \leq 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \max_{P \in N_\Omega} U_P \leq \max_{Q \in N_{\partial\Omega}} U_Q^+,$$

$$(128) \quad L_h U_P \geq 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \min_{P \in N_\Omega} U_P \geq \min_{Q \in N_{\partial\Omega}} U_Q^-,$$

kde $U_Q^+ = \max\{U_Q, 0\}$, $U_Q^- = \min\{U_Q, 0\}$. *Speciálně*

$$(129) \quad L_h U_P = 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \max_{P \in N_\Omega} |U_P| \leq \max_{Q \in N_{\partial\Omega}} |U_Q|.$$

Důkaz. Nechť platí levá strana implikace (127) a označme $M_\Omega := \max_{P \in N_\Omega} U_P$ a $M_{\partial\Omega} := \max_{Q \in N_{\partial\Omega}} U_Q$. Je-li $M_\Omega \leq 0$, pak (127) triviálně platí. Nechť tedy $M_\Omega > 0$ a předpokládejme, že $M_\Omega > M_{\partial\Omega}$. Nechť $P^* \in N_\Omega$ je takový, že $M_\Omega = U_{P^*}$. Pak

$$(130) \quad M_\Omega = U_{P^*} \leq \frac{1}{c_{P^*}} \sum_{Q \in N_{P^*}} c_{P^*Q} U_Q \leq M_\Omega,$$

kde první nerovnost plyne z toho, že $L_h U_{P^*} \leq 0$, a druhá z (126). Vztah (130) může platit pouze tehdy, pokud v něm platí rovnost, a z (126) pak plyne, že $U_Q = M_\Omega$ pro všechna $Q \in N_{P^*}$. Ze souvislosti množiny N_Ω formulované v předpokladu (P2) pak plyne, že $U_P = M_\Omega$ pro všechna $P \in N_\Omega$. Konečně z předpokladu (P3) plyne, že existuje $Q \in N_{\partial\Omega}$ takové, že $U_Q = M_\Omega$, což je spor s předpokladem, že $M_\Omega > M_{\partial\Omega}$. Je tedy $M_\Omega \leq M_{\partial\Omega}$, tj. platí (127).

Platí-li levá strana implikace (128), pak $\{-U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ splňuje levou stranu implikace (127), a tudíž

$$-\min_{P \in N_\Omega} U_P = \max_{P \in N_\Omega} (-U_P) \leq \max_{Q \in N_{\partial\Omega}} (-U_Q)^+ = \max_{Q \in N_{\partial\Omega}} (-U_Q^-) = -\min_{Q \in N_{\partial\Omega}} U_Q^-,$$

tj. platí (128).

Platí-li levá strana implikace (129), pak $\{U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ i $\{-U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ splňují levou stranu implikace (127), a tudíž pro libovolné $P \in N_\Omega$ platí

$$U_P \leq \max_{Q \in N_{\partial\Omega}} U_Q^+, \quad -U_P \leq \max_{Q \in N_{\partial\Omega}} (-U_Q)^+ \quad \Rightarrow \quad |U_P| \leq \max_{Q \in N_{\partial\Omega}} |U_Q|,$$

což dává (129). □

Důsledek 1 *Nechť jsou splněny předpoklady (P1)–(P3) a navíc*

$$(131) \quad c_P = \sum_{Q \in N_P} c_{PQ} \quad \forall P \in N_\Omega.$$

Pak pro libovolnou síťovou funkci $\{U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ platí

$$(132) \quad L_h U_P \leq 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \max_{P \in N_\Omega} U_P \leq \max_{Q \in N_{\partial\Omega}} U_Q,$$

$$(133) \quad L_h U_P \geq 0 \quad \forall P \in N_\Omega \quad \Rightarrow \quad \min_{P \in N_\Omega} U_P \geq \min_{Q \in N_{\partial\Omega}} U_Q.$$

Důkaz. Nechť platí levá strana implikace (132) a označme $M := \min_{Q \in N_\Omega \cup N_{\partial\Omega}} U_Q$. Buď $\tilde{U}_P := U_P - M$ pro všechna $P \in N_\Omega \cup N_{\partial\Omega}$. Pak $\tilde{U}_P \geq 0$ pro všechna $P \in N_\Omega \cup N_{\partial\Omega}$ a díky (131) je $L_h \tilde{U}_P = L_h U_P \leq 0$ pro všechna $P \in N_\Omega$. Použitím (127) získáme

$$\left(\max_{P \in N_\Omega} U_P \right) - M = \max_{P \in N_\Omega} \tilde{U}_P \leq \max_{Q \in N_{\partial\Omega}} \tilde{U}_Q^+ = \max_{Q \in N_{\partial\Omega}} \tilde{U}_Q = \left(\max_{Q \in N_{\partial\Omega}} U_Q \right) - M,$$

z čehož plyne (132). Implikace (133) se získá přechodem k $\{-U_P\}_{P \in N_\Omega \cup N_{\partial\Omega}}$ a aplikací implikace (132), analogicky jako bylo získáno (128) z (127) v důkazu věty 10. □

Stejně jako ve spojitém případě plyne z diskrétního principu maxima formulovaného ve větě 10 jednoznačnost řešení úlohy

$$(134) \quad L_h U_P = f_P + q_P \quad \forall P \in N_\Omega,$$

$$(135) \quad U_P = g_P \quad \forall P \in N_{\partial\Omega},$$

kde $g_P := g(P)$ udává Dirichletovu okrajovou podmínku v bodě P . Jelikož úloha (134)–(135) je lineární a konečně rozměrná, plyne z jednoznačnosti též existence řešení, jak ukážeme v následující větě.

Věta 11 *Nechť jsou splněny předpoklady (P1)–(P3). Pak pro libovolnou pravou stranu má úloha (134)–(135) právě jedno řešení.*

Důkaz. Úloze (134)–(135) odpovídá soustava lineárních algebraických rovnic se čtvercovou maticí pro neznámé hodnoty $\{U_P\}_{P \in N_\Omega}$. Úloha je tedy jednoznačně řešitelná právě tehdy, když homogenní soustava má pouze triviální řešení. Je-li však pravá strana v (134)–(135) nulová, dostáváme z (129) ihned, že $U_P = 0$ pro všechna $P \in N_\Omega$. □

6.6 Odhad chyby aproximace

V této části odvodíme obecné odhady chyby aproximace pro úlohu (134)–(135), přičemž zobecníme odhad chyby aproximace provedený na začátku této kapitoly pro diskretizaci Poissonovy rovnice.

Nechť u je řešení spojité úlohy, kterou aproximujeme pomocí diskretizace (134)–(135). Jako obvykle zavádíme chybu diskretizace

$$\varepsilon_P := L_h u_P - f_P - q_P \quad \forall P \in N_\Omega,$$

kde $u_P = u(P)$. Předpokládáme, že přibližné řešení splňuje stejné Dirichletovy okrajové podmínky jako řešení přesné, tj. $U_P = u(P)$ pro všechna $P \in N_{\partial\Omega}$. Pak chyba aproximace $e_P := U_P - u_P$ opět splňuje

$$L_h e_P = -\varepsilon_P \quad \forall P \in N_\Omega, \quad e_P = 0 \quad \forall P \in N_{\partial\Omega}.$$

Věta 12 *Nechť jsou splněny předpoklady (P1)–(P3) a necht' Φ je nezáporná síťová funkce definovaná na $N_\Omega \cup N_{\partial\Omega}$ splňující*

$$L_h \Phi_P \leq -1 \quad \forall P \in N_\Omega.$$

Pak pro chybu aproximace odpovídající diskretizaci (134)–(135) platí odhad

$$\max_{P \in N_\Omega} |e_P| \leq \left(\max_{Q \in N_{\partial\Omega}} \Phi_Q \right) \left(\max_{P \in N_\Omega} |\varepsilon_P| \right).$$

Důkaz. Buď $D := \max_{P \in N_\Omega} |\varepsilon_P|$. Pak

$$L_h(D \Phi_P \pm e_P) = D L_h \Phi_P \mp \varepsilon_P \leq -D \mp \varepsilon_P \leq 0 \quad \forall P \in N_\Omega,$$

a tudíž podle (127) platí pro libovolné $P \in N_\Omega$

$$\pm e_P \leq D \Phi_P \pm e_P \leq \max_{Q \in N_{\partial\Omega}} (D \Phi_Q \pm e_Q)^+ = D \max_{Q \in N_{\partial\Omega}} \Phi_Q,$$

z čehož plyne dokazovaný odhad. □

Zatímco odhad chyby diskretizace je poměrně snadný, nalezení srovnávací funkce Φ nemusí být vždy jednoduché. Tato funkce není samozřejmě určena jednoznačně a cílem je nalézt takové Φ , aby $\max_{Q \in N_{\partial\Omega}} \Phi_Q$ bylo co nejmenší.

Pro úlohu (107) a diskretizaci (108)–(109) plyne z věty 12 odhad $|U_P - u_P| \leq C h^2$ pro libovolné $P \in N_\Omega$. Bude-li však mít Ω zakřivenou hranici, získáme odhad $|\varepsilon_P| \leq C h^2$ pouze v regulárních uzlech, zatímco v uzlech u hranice budeme mít obecně jen $|\varepsilon_P| \leq C h$ (viz vztah (120)). Aplikací věty 12 pak dostaneme pouze $|U_P - u_P| \leq C h$ pro $P \in N_\Omega$. Tento odhad však není optimální a lze ho zlepšit použitím následující věty.

Věta 13 *Nechť $N_\Omega = N_1 \cup N_2$, kde $N_1 \cap N_2 = \emptyset$. Necht' jsou splněny předpoklady (P1)–(P3) a necht' Φ je nezáporná síťová funkce definovaná na $N_\Omega \cup N_{\partial\Omega}$ splňující*

$$L_h \Phi_P \leq -C_1 < 0 \quad \forall P \in N_1, \quad L_h \Phi_P \leq -C_2 < 0 \quad \forall P \in N_2.$$

Nechť chyba diskretizace odpovídající úloze (134)–(135) splňuje

$$|\varepsilon_P| \leq D_1 \quad \forall P \in N_1, \quad |\varepsilon_P| \leq D_2 \quad \forall P \in N_2.$$

Pak pro chybu aproximace platí odhad

$$\max_{P \in N_\Omega} |e_P| \leq \left(\max_{Q \in N_{\partial\Omega}} \Phi_Q \right) \max \left\{ \frac{D_1}{C_1}, \frac{D_2}{C_2} \right\}.$$

Důkaz. Označíme-li

$$D := \max \left\{ \frac{D_1}{C_1}, \frac{D_2}{C_2} \right\},$$

můžeme důkaz provést analogicky jako pro větu 12. \square

Příklad 2 Uvažujme úlohu (107) s $\Omega := \{(x, y) \in \mathbb{R}^2; x^2 + y^2 < 1\}$. Oblast Ω pokryjeme rovnoměrnou čtvercovou sítí s krokem h a v regulárních uzlech použijeme schéma (108), zatímco v neregulárních uzlech použijeme pro druhé derivace aproximace typu (120). Aplikací věty 13 odvoďte odhad chyby aproximace, který je druhého řádu v h .

Řešení: Nechť N_1 je množina uzlů z N_Ω , jejichž všechny sousední uzly leží v Ω a $N_2 := N_\Omega \setminus N_1$. Podle (113) a (120) platí

$$|\varepsilon_P| \leq K_1 h^2 \quad \forall P \in N_1, \quad |\varepsilon_P| \leq K_2 h \quad \forall P \in N_2.$$

Nechť $\Psi(x, y) = x^2 + y^2$ a položme

$$\Phi_P = M_1 \Psi(P) \quad \forall P \in N_\Omega, \quad \Phi_P = M_1 \Psi(P) + M_2 \quad \forall P \in N_{\partial\Omega},$$

kde M_1, M_2 jsou kladné konstanty, které budou určeny později. Pak

$$L_h \Phi_P = M_1 L_h \Psi_P = M_1 (-\Delta \Psi)_P = -4 M_1 \quad \forall P \in N_1.$$

Je-li $P \in N_2$, uvažujme např. situaci z obr. 4 a aproximaci $u_{xx}(P)$ pomocí výrazu v (120). Jelikož chyba této aproximace závisí na třetích derivacích u , je tato aproximace pro kvadratické funkce přesná, a proto získáme

$$\frac{\Phi_A + \theta \Phi_W - (1 + \theta) \Phi_P}{\frac{1}{2} \theta (1 + \theta) h^2} = M_1 \Psi_{xx}(P) + \frac{M_2}{\frac{1}{2} \theta (1 + \theta) h^2} \geq 2 M_1 + \frac{M_2}{h^2}.$$

Obdobně postupujeme v ostatních situacích, které mohou nastat, což vede ke vztahu

$$L_h \Phi_P \leq -4 M_1 - \frac{M_2}{h^2} \leq -\frac{M_2}{h^2} \quad \forall P \in N_2.$$

Aplikací věty 13 dostáváme

$$\max_{P \in N_\Omega} |e_P| \leq (M_1 + M_2) \max \left\{ \frac{K_1 h^2}{4 M_1}, \frac{K_2 h^3}{M_2} \right\}.$$

Položíme-li $M_1 = \frac{1}{4} K_1 h^2$ a $M_2 = K_2 h^3$, získáme

$$\max_{P \in N_\Omega} |e_P| \leq \frac{1}{4} K_1 h^2 + K_2 h^3,$$

tj. ukázali jsme, že chyba aproximace je skutečně druhého řádu přesnosti.