

Počítacové metody v teoretické fyzice I

- informace o přednášce na utf.mff.cuni.cz/fuhoufek/

Úvod - reprezentace čísel v počítaci, zaokrouhlovací chyba, stabilita algoritmů a podmíněnost úlohy

- nejběžnější je reprezentace pomocí pohyblivé řádové čárky dle IEEE standardem (floating-point arithmetic)

npr. čísla $\pm 0,01$ není reprezentováno přesně, neboť

$$\pm 0,01 \text{ (DEC)} = \pm 0,000\ 000\ 101\ 000\ 111\ 101\dots \text{ (BIN)}$$

$$= \pm \underbrace{1.01000111101\dots}_{t \text{ cifer, určuje přesnost (precision)}} \times 2^e \text{ (BIN)}$$

$$= \pm m \times \beta^{e-t+1} \quad [pozn. neplatí s anglickou accuracy]$$

m ... mantisa (significand)

β ... základ (base) - obvykle 2, někdy 10, 16

e ... exponent

- mantisa je celočíselná (ovšem může být někdy brána i jako desetinné číslo)

$$m = d_1 d_2 \dots d_t$$

kde d_i nabývají hodnot $0 \leq d_i < \beta$, ovšem $d_1 \neq 0$

- číslo 0,01 je ve 32-bitovém IEEE standardu = single precision, neboť v REAL*4 (4 byty), reprezentováno takto:

$$0,01 \approx \underbrace{[0|01111100]}_{\substack{\text{znaménkový} \\ \text{bit } s: (-1)^s}} \underbrace{[0|0100011]}_{\substack{127 + e \\ (8 \text{ bitů})}} \underbrace{[11010111]}_{\substack{\text{mantisa } m \text{ bez první} \\ 1 \\ (23 \text{ bitů})}} \underbrace{[00001010]}_{\substack{\sim 7 \text{ platných míst} \\ 52 \text{ bitů} \\ \sim 16 \text{ platných míst} \\ v desítkové soustavě}}$$

v double precision: 11 bitů

- rozsah exponentů je $e_{\min} \leq e \leq e_{\max}$

příčinou pro single precision je $e_{\min} = -126$ a $e_{\max} = 127$
 $(2^{-126} \approx 1,18 \cdot 10^{-38}, 2^{127} \approx 3,4 \cdot 10^{38})$

a pro double precision $e_{\min} = -1022$ a $e_{\max} = 1023$
 $(2^{-1022} \approx 2,2 \cdot 10^{-308}, 2^{1023} \approx 1,8 \cdot 10^{308})$

- pozn. speciální exponenty $e = -127$ (samé nuly)

a $e = 128$ (samo jedničky) se používají na tzv. subnormální čísla a těž pro $\pm \infty$ a NaN

- pokud reálné číslo $x \in \mathbb{R}$ leží v rozsahu systému F
čísel s pohyblivou řádovou čárkou, pak

$$f_1(x) = x(1+\delta), \text{ kde } |\delta| < u = \frac{\epsilon_M}{2}$$

kde $f_1(x)$ znázorňuje číslo x vyjádřené nejbližším číslem s pohyblivou řádovou čárkou, tj. $f(x) \in F$
a u je tzv. unit roundoff $u = \frac{1}{2} \beta^{1-t}$
a ϵ_M je tzv. strojové epsilon.

Maximální relativní chyba vyjádření lib. čísla x v rozsahu F
je tedy dána u .

- je-li $f_1(x) = x + \Delta x = x(1 + \epsilon_x)$, hovoříme

o Δx jako o zaokrouhlovací chybě

a o ϵ_x jako o relativní zaokrouhlovací chybě

Pozn: protože je reprezentace čísel pomocí pohyblivé řádové čárky zdaleka nejběžnější, není jediná. Existují i další systémy jako např. pevná řádová čárka (fixed-point) nebo logaritmický číselný systém (logarithmic number system), a pod. (viz wikipedia nebo Higham: Accuracy and stability ...)

- vliv zaokrouhlovacích chyb na výsledek

- při výpočtech na počítaci je nutné mít na paměti,

že zaokrouhlovací chyba je vždy přítomna

- čísla na vstupu jsou zaokrouhlena

- výsledek lib. aritmetické operace je zaokrouhlen

Pozn: dobrý systém by měl mít implementovány tyto operace tak, že výsledek bude

$$f_1(x \text{ op } y) = (x \text{ op } y)(1+\delta), |\delta| < u$$

pro $\text{op} = +, -, *, /$ a druhou odmocninu tedy

- numerická matematika se z velké části zabývá právě vlivem zaokrouhlovacích chyb na správnost výpočtu

1) pozor! - některé axiomy reálné aritmetiky
nejsou v systémech s pohyblivou řádovou čárkou
splňeny

napří: výsledek závisí na pořadí scítanců

a to i tehdy, jsou-li všechny scítance kladné

Pr. uvažujme jednoduchou sumu (viz notebook roundoff.errors.html)

$$S_n = 1 + \sum_{k=1}^n \frac{1}{k^2+k} = 1 + \sum_{k=1}^n \left(\frac{1}{k} - \frac{1}{k+1} \right) = 2 - \frac{1}{n+1}$$

- scítame-li odpoedu, objeví se zaokrouhlovací chyby,
zvláště pro velká n
- scítame-li od zadu (od nejménších příspěvků po největší),
zaokrouhlovací chyby se neprojeví

nebo: rovnice nemusí mít jednoznačné řešení

případně nemusí mít žádne řešení

(i když v přesné aritmetice existuje jediné)

Pr. rovnice $1+x=1$ má řešení $x=0$

avšem v double precision bude rovnice splňena
pro lib. $x \lesssim 10^{-17}$, které patří do systému čísel double prec.

Pr. rovnice $f(x)=0$ většinou nemá v daném systému
čísel s pohyblivou řádovou čárkou žádne řešení,
tže hledat pouze řešení, pro které je $|f(x)|$ nejménší,
případně menší než jisté ϵ , když nám stačí
určitá přesnost

2) odečítání blízkých (velkých) čísel - cancellation problem

- operace odečítání může vést do numerických výpočtů
značné chyby, jsou-li odečítaná čísla blízka

- máme-li $f_1(x_1) = x_1 + \Delta x_1$ pak $f_1(x_1 \pm x_2) = x_1 \pm x_2 + \Delta x_1 \pm \Delta x_2 + \epsilon(x_1 \pm x_2)$

a $f_1(x_2) = x_2 + \Delta x_2$ a relativní chyba je kde $\epsilon \leq \epsilon_M$

$$\delta = \left| \frac{f_1(x_1 \pm x_2) - (x_1 \pm x_2)}{x_1 \pm x_2} \right| \leq \frac{|\Delta x_1| + |\Delta x_2|}{|x_1 \pm x_2|} + |\epsilon|$$

- pokud je např. $x_1 \approx 1, x_2 \approx 1, \Delta x_1 \approx 10^{-8}$ a $\Delta x_2 \approx 10^{-8}$
a v přesné aritmetice $x_1 - x_2 \approx 10^{-16}$, pak $\delta \approx 2 \cdot 10^{-4}$
- často se s tímto problémem nedá moc dělat,
ale někdy ze znalosti problému lze změnit výpočet
tak, že se odečítání lze vyhnout

[Príklad: výraz $\sqrt{x+1} - \sqrt{x}$ lze vypočítat též pomocí

$\frac{1}{\sqrt{x+1} + \sqrt{x}}$ a pro velká x pak velká zokrouhlující chyba nevzniká]

[Príklad: řešení kvadratické rovnice $ax^2 + bx + c = 0$

$$\text{spojenečné pomocí } x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

může být zatíženo chybou, pokud $b^2 \gg 4ac$,

neboť pak $\sqrt{b^2 - 4ac} \approx |b|$ a jeden z kořenů

dostaneme odečítání dvou blízkých čísel,

lze obejít tak, že tento kořen spočteme jako $x_2 = \frac{c}{ax_1}$

3) součet řady „velkých“ čísel, jehož výsledek je „malý“

- jde o skryté odečítání blízkých čísel, neboť (smearing problem)
až mezi výsledky můžou být blízce
- opět lze někdy obejít jiným přístupem k danému problému

[Príklad: výpočet e^x pomocí Taylorovy řady $1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots$

pro $x > 0$ více méně v pořádku, jen pro větší x
máme pomalou konvergenci

ovšem pro $x < 0$ výpočet selže (výsledek je malé číslo
jako součet velkých s opačnými znaky)

trik: pro $x < 0$ lze psát $e^x = \frac{1}{e^{-x}}$

- často nastává u řad typu binomického rozvoje

Prí. výpočetní příklad ilustrující řešení přes rekurentní vztahy:

$$\text{integrál } I_n(a) = \int_0^1 \frac{x^n}{x+a} dx = (\text{trik } x^n = [(x+a)-a]^n)$$

$$= \sum_{k=0}^{n-1} (-a)^k \binom{n}{k} \left[(1+a)^{n-k} - a^{n-k} \right] \frac{1}{n-k} + (-a)^n \ln \frac{a+1}{a}$$

- použití této řady v double precision dají např. pro $a=10$ a $n \geq 10$ řešení se značnými chybami, nebo vyjdou přímo odvídne nesmysly ($I_n < 0$) (koeficienty jsou např. pro $n=10$ a $k=5$ řádu 10^{11} a pritom výsledek leží v $0 < I_n < 1$!)
- řešení tohoto problému spočívá ve využití rekurentního vztahu

$$I_n = \int_0^1 \frac{x^{n-1}(x+a-a)}{x+a} dx = \int_0^1 x^{n-1} dx - a \int_0^1 \frac{x^{n-1}}{x+a} dx = \frac{1}{n} - a I_{n-1}$$

ovšem dopředná rekurence, kdy startujeme s $I_0 = \ln \frac{a+1}{a}$, není stabilní!

- libovolná chyba (izanokrouhlkovací) na začátku bude v následujících krocích násobena faktorem $(-a)^n$ a výsledky pro $n \geq 15$ znehodnotí úplně

- co to zkoušit obráceně? $I_{n-1} = \frac{1}{a} \left(\frac{1}{n} - I_n \right)$
- ovšem jak začít? Ize libovolně, neboť počítací chyba bude potlačována faktorem $\frac{1}{a^n}$, po k iteracích, např. pro $n=10$ a $a=10$ zvolíme $I_{n+k} = 1$ pro vhodně velké $k \geq n + \log \frac{1}{\epsilon_M}$ a dostaneme I_n s přesností ϵ_M

• stabilita algoritmu

Řekneme, že algoritmus výpočtu veličiny $f=f(x)$, závisející na vstupních datech x , je stabilní, pokud jeho provedením v aritmetice s pohyblivou rádovou čárkou dá pro všechna x výsledek \tilde{f} , který se od přesného f liší s relativní chybou rádu $O(\epsilon_M)$, tj. při změně ϵ_M lze chybu odhadnout výrazem $c\epsilon_M$, kde konstanta c nezávisí na x .

- v praxi navíc požadujeme, aby pro konkrétní ϵ_M (např. pro $\epsilon_M \sim 10^{-16}$ v double prec.) nebylo $c \sim \epsilon_M^{-1}$ (tj. $\sim 10^{16}$), neboť pak by chyba byla srovnatelná s výsledkem!

Pr. při výpočtu $\sqrt{x+1} - \sqrt{x}$ lze pro velké kladné x odhadnout absolutní chybu jako $2\sqrt{x}\epsilon_M$ a tedy relativní

$$\text{chyba bude } \delta \leq \frac{2\sqrt{x}\epsilon_M}{\sqrt{x+1}-\sqrt{x}} = 2\epsilon_M \sqrt{x} (\sqrt{x+1} + \sqrt{x}) \sim 4\epsilon_M x$$

a nelze odhadnout výrazem $c\epsilon_M$ (pro velká x roste)

\Rightarrow výpočet je nestabilní

Oršen pro $\frac{1}{\sqrt{x+1} + \sqrt{x}}$ dostaneme pro odhad relativní

$$\text{chyby jmenovatele } \delta \leq \frac{2\sqrt{x}\epsilon_M}{\sqrt{x+1} + \sqrt{x}} \sim \epsilon_M \text{ a pro celý výraz tedy také } \epsilon_M \Rightarrow \text{stabilní výpočet}$$

• podmíněnost úlohy

- v některých úlohách, bez ohledu na použitý algoritmus, může výsledek silně záviset na vstupních datech
- klasickým příkladem takové úlohy je určení kořenů polynomu zadaného jeho koeficienty (vstupní data)

[Př. polynom $p(x) = x^{10} - 10x^9 + 45x^8 - \dots - 10x + 1 = (x-1)^{10}$

má desetinásobný kořen $x=1$

a tisk $q(x) = p(x) - 1^{10}$ (pouze a_0 se změnilo o 1^{-10})

má kořeny $x_k = 1 + \frac{1}{10} e^{\frac{2\pi i k}{10}}$ pro $k=0, \dots, 9$

a tedy změna výsledku je $\sim 10^{-1}$ (v přesné aritmetice!)

- obecně uvažujme problém, který řešíme, jako funkci $f: X \rightarrow Y$, kde X je normovaný vektorový prostor vstupních dat a Y je normovaný vektorový prostor řešení
- stav, že f je obecně (a též obvykle) nelineární.

Řekneme, že úloha (problém) je dobré podmíněná (well-conditioned),

když malá změna vstupních dat x vede

k malým změnám výstupních dat $f(x)$

(př. je např. výpočet vlastního integrálu, či určení vlastních čísel a vektorů hermitovské matice)

Naopak úloha je špatně podmíněná, (ill-conditioned)

když malá změna x vede k velké změně $f(x)$.

(hledání kořenů polynomu, řešení soustavy lin. rovnic

pro nevhodné matice či počáteční úloha pro dif. rovnice)

Podmíněnost se kvantifikuje buď pomocí

absolutního čísla podmíněnosti

$$\hat{\rho}(x) = \lim_{\varepsilon \rightarrow 0} \sup_{\|\delta x\| \leq \varepsilon} \frac{\|f(x+\delta x) - f(x)\|}{\|\delta x\|}, \text{ kde tedy } \delta x \text{ je malá změna} \\ \text{a } \delta f = f(x+\delta x) - f(x)$$

nebo pouze relativního čísla podmíněnosti

$$\alpha(x) = \limsup_{\varepsilon \rightarrow 0} \frac{\|\delta f\|}{\|\delta x\|} \leq \varepsilon \sqrt{\frac{\|\delta f\|}{\|\delta x\|}} / \sqrt{\frac{\|\delta x\|}{\|x\|}}$$

- pro diferencovatelné f je $\alpha(x)$ dánou pouze Jacobijnu

$$\text{zobr. } f(x) : J_{ij}(x) = \frac{\partial f_i}{\partial x_j} \Big|_x, \quad \delta f \approx J(x) \delta x$$

a tedy $\hat{\alpha}(x) = \|J(x)\|$, kde maticová norma $J(x)$
je indukovaná normami v X a Y

$$\text{a } \alpha(x) = \frac{\|J(x)\|}{\frac{\|f(x)\|}{\|x\|}}$$

- relativní číslo podmíněnosti je obvykle dležitější
při výpočtech v pohyblivé rádiové sítce, kde
pracujeme s relativní chybou ε

- problém je (typicky) dobré podmíněny pro $\alpha \leq 10^2$
a špatně podmíněny pro $\alpha \geq 10^{10}$

ovšem přesná hodnota α závisí na zvolené normě

Pr. výpočet $f(x) = \sqrt{x}$ je dobré podmíněna úloha,

$$\text{neboť } \alpha = \frac{\frac{1}{2\sqrt{x}}}{\frac{\sqrt{x}}{x}} = \frac{1}{2}$$

ovšem pro $f(x) = x_1 - x_2$ dostaneme $J = (1, -1)$ a pro $\|\cdot\|_\infty$

$$\text{normu pak } \alpha = \frac{2}{\max\{|x_1|, |x_2|\}} = \frac{2 \max\{|x_1|, |x_2|\}}{|x_1 - x_2|} \quad \text{což může být libovolně velké, pro } x_1 \sim x_2$$

Pozn.: $\|J\|_\infty = 2 = \text{maximum součtu rádkových prvků}$
(viz poznámky k vektorovým a maticovým normám)

Pr. pro hledání kořenů polynomu jde $\alpha \rightarrow \infty$

$$\text{neboť pro } \delta x = 10^{-10} \text{ je } \frac{\|\delta f\|}{\|\delta x\|} \approx 10^9$$

$$\text{pro } \delta x = 10^{-20} \text{ je } \frac{\|\delta f\|}{\|\delta x\|} \approx 10^{18}$$

atd.

tj. $\frac{\|\delta f\|}{\|\delta x\|}$ roste

nade všechny meze
a proto $\|f(x)\| \sim 10^0$ a $\|x\| \sim 10^0$

Normy na vekt. prostorech a normy matic

Def: Norma $\|\cdot\|$ je funkce zadaná $\mathbb{C}^m \rightarrow \mathbb{R}$ splňující pro $x, y \in \mathbb{C}^m$

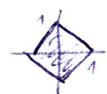
$$\alpha \in \mathbb{C} : 1) \|x\| \geq 0 \quad 2) \|x\| = 0 \text{ právě, když } x=0$$

$$2) \|x+y\| \leq \|x\| + \|y\|$$

$$3) \|\alpha x\| = |\alpha| \|x\|$$

p-normy na \mathbb{C}^n a jejich jednotkové koule $\{x \in \mathbb{C}^n, \|x\| \leq 1\}$.

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



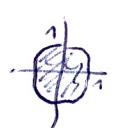
$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$



$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2}$$



$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$



vážené p-normy $\|x\|_W = \|Wx\|_p = \left(\sum_{i=1}^n |w_i x_i|^p \right)^{1/p}$, kde W je diagonální matice

- lze zobecnit na lib. nesingularní matice W

Maticová norma indukovaná vekt. normou

zobrazuje matici $A \in \mathbb{C}^{n \times n}$ jeho zobrazení z vekt. prostoru s normou $\|\cdot\|_{(n)}$ do prostoru s normou $\|\cdot\|_{(m)}$

pak induk. matic. norma $\|A\|_{(n,m)}$ je největší C

pro které platí $\|Ax\|_{(m)} \leq C \|x\|_{(n)}$

neboli

$$\|A\|_{(n,m)} = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|_{(m)}}{\|x\|_{(n)}} = \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|_{(n)}=1}} \|Ax\|_{(m)}$$

Pr. $A = \begin{pmatrix} 1 & 2 \\ 0 & 2 \end{pmatrix}$ - shodné A zobrazi $e_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ na $\begin{pmatrix} a_{11} \\ a_{21} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ a $e_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ na $\begin{pmatrix} a_{12} \\ a_{22} \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$

1-norm



\xrightarrow{A}



$$\|A\|_1 = 4$$

2-norm



\xrightarrow{A}

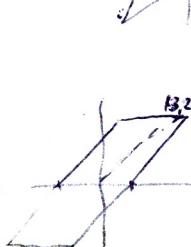


$$\|A\|_2 = \sqrt{\frac{1}{2} + \frac{7}{2\sqrt{65}}} \approx 2,9208$$

oo-norm



\xrightarrow{A}



$$\|A\|_\infty = 3 \quad (\text{vidíme zde, že } \|A\|_2 = \max_{1 \leq i \leq 2} \sigma_i \text{ kde } \sigma_i \text{ jsou singularní čísla})$$

Př. pro lib. diagonální matici $D = \begin{pmatrix} d_1 & & \\ & d_2 & \\ & & \ddots & d_m \end{pmatrix}$ se jednotková koule

ztransf. na hyperelipsu s polosazemi ~~p~~ je délce $|d_i|$

$$\text{a tedy } \|D\|_2 = \max_{1 \leq i \leq m} |d_i| = \|D\|_P$$

absolutní hodnota

~~platí~~, tedy $\|A\|_1 = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$ je součet sloupcových prvků

$$= \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \text{ kde } a_{ij} \text{ je } j\text{-tý sloupec maticy}$$

$$\text{neboť } \|Ax\|_1 = \left\| \sum_{j=1}^n x_j a_{ij} \right\|_1 \leq \sum_{j=1}^n |x_j| \|a_{ij}\|_1 \leq \max_{1 \leq j \leq n} \|a_{ij}\|_1, \text{ což platí i pro } \|x\|_1 = 1$$

a vždykdy $x = e_j$, pro které je $\|a_{ij}\|_1$

maximální, dostatečné rovnost a tedy

$$\|A\|_1 = \sup_{\|x\|_1=1} \|Ax\|_1 = \max_{1 \leq j \leq n} \|a_{ij}\|_1$$

$$Ax = \begin{pmatrix} a_1 & | & a_2 & | & \dots & | & a_n & | & x_1 \\ & & & & & & & & x_2 \\ & & & & & & & & \vdots \\ & & & & & & & & x_n \end{pmatrix}$$

• podobně $\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$ je řádek

$$= \max_{1 \leq i \leq n} \|a_i^*\|_1, \text{ kde } a_i^* \text{ je } i\text{-tý řádek}$$

$$\text{neboť } \|Ax\|_\infty = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \text{ maximálnost násobku}$$

pro jeden vektoru
 $x = (\pm 1, \pm 1, \dots, \pm 1)$

Pozn. měli by studenti znít z LA

Hölderova nerovnost

$$|xy| \leq \|x\|_p \|y\|_q \text{ kde } \frac{1}{p} + \frac{1}{q} = 1$$

a spec. případ Cauchy-Schwarzkova nerovnost

$$|xy| \leq \|x\|_2 \|y\|_2$$

- Pro odhad ^{normy} součtu dvou matic platí

$$\|AB\|_{(e_i)} \leq \|A\|_{(e_i)} \|B\|_{(-i)}$$

ale obecně zde není rovnost, např. $\|A^n\| \leq \|A\|^n$, ale pro $n \geq 2$ obecně neplatí $\|A^n\| = \|A\|^n$

Obecné matice normy - splňují stejně podmínky jako normy vektoru

jejich normy (jako hlyba - měli man rozdíly)
vektorských prostorů

- nejpoužívanější je Frobeniova norma (též Hilbertova-Schmidtova norma)

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}}, \text{ tj. jde o 2-normu v max. prostoru vektorů.}$$

i zde platí odhad pro součin

$$\|AB\|_F^2 \leq \|A\|_F \|B\|_F$$

$$\text{(bez-důkaz, že } \|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2}$$