

# Můžeme věřit své vlastní kalkulačce?

Edita Pelantová, katedra matematiky FJFI, ČVUT

Miloslav Znojil, Ústav jaderné fyziky, AV ČR

Abstrakt. V článku studujeme vliv zaokrouhlování, které počítač i kalkulačka během výpočtu nutně musejí provádět, na správnost výsledku. Na jedné konkrétní úloze demonstrujeme, jak použití výpočetní techniky bez uvažování o její přesnosti může způsobit chyby libovolné velikosti.

I v běžných matematických výpočtech se často setkáváme s iracionálními čísly. Výpočet obsahu kruhu vyžaduje použití Ludolfova čísla  $\pi$ , stanovení výšky úroků vyžaduje logaritmování, kde základem je Eulerovo číslo  $e$ . Je dávno známo, že obě zmíněné konstanty jsou iracionální čísla, to znamená, že jejich zápis v desítkové soustavě není konečný – jak je tomu např. u čísla  $\frac{1}{5} = 0,2 -$  ani od jistého místa periodický – jak je tomu u čísla  $\frac{1}{6} = 0,1666\dots$ . Zápisy čísel  $\pi$  a  $e$  mají tvar

$$\pi = 3,1415926535897932384626433832795028841971693993751\dots,$$

$$e = 2,7182818284590452353602874713526624977572470936999\dots$$

Výpočty, které provádíme pomocí kalkulačky či počítače, jsou omezeny pouze na racionální čísla. Počet platných číslic, se kterými kalkulačka a počítač pracují, je limitován a závisí na typu přístroje. Velikosti čísel, se kterými pracuje člověk, nejsou v principu omezeny, ale ruku na srdce, kolikamístná čísla byste byli ochotni násobit s tužkou na papíře vy?

Nemožnost pracovat s přesnou hodnotou Ludolfova čísla a Eulerova čísla nás většinou netrápí. Při školních výpočtech si vystačíme s přibližnou hodnotou čísla  $\pi \doteq 3,14$  nebo  $\pi = \frac{22}{7}$ . Chyba, která vznikne při výpočtu obsahu kruhu, je tak menší než jedno promile.

Je-li však úloha složitější a vyžaduje větší počet operací, může chyba přerůst všechny meze. Konstrukce mnohých technických zařízení vyžaduje komplikovaný výpočet. Nebudeme dopodrobna popisovat úlohu ze skutečné praxe - zabíhali bychom do podrobností nedůležitých pro vysvětlení vzniku chyb. Následující vymyšlený příběh dobře ilustruje podstatu problému. Pro úplnost dodejme, že matematické jádro příběhu bylo inspirováno přednáškou Francouze Jeana-Michela Mullera.

Otce právě narozeného syna, pana Nešetřila, upoutá reklama ve výloze banky se sloganem "Iracionálně ke štěstí". Banka nabízí rodičům, aby založili pro narozené dítě účet, na který vloží  $e$  korun, tedy iracionální částku. Banka slibuje, že po každém roce odečte z účtu jednu korunu jako poplatek za vedení účtu a vynásobí zbytek počtem let od založení účtu. V den 25. narozenin banka dítěti vyplatí jmění, které pro ně rodiče našetřili.

Pan Nešetřil se zamyslí, zda by neměl už teď pamatovat na štěstí svého syna. Rozhodne se o nabídce uvažovat a začne počítat: po prvním roce je na účtě  $p_1 = e - 1$  korun, po druhém roce  $p_2 = 2(p_1 - 1) = 2(e - 2)$  korun, po třetím roce  $p_3 = 3(p_2 - 1)$  atd. Protože má po ruce mobil, začne počítat na kalkulačce. Hodnotu Eulerova čísla si otec nepamatuje ani přibližně, ale ve výloze banky je jako dekorace uvedeno číslo  $e$  s více než sto místy. Kalkulačka dovolí však natukat pouze 9 platných míst. Proto rozvášňný otec správně zaokrouhlí a počítá částky  $p_n$ . Když mu mobil ukáže  $p_{25} = 0,239 \times 10^{17}$ , je celý bez sebe a spěchá oznámit svůj plán manželce. Ta, i když zrovna kojí, je ještě rozvášněnější (a taky ví, že jméno Nešetřil nedostala manželova rodina náhodou) a udělá kontrolní výpočet doma na kalkulačce počítače, která pracuje s přesností 16

míst. Paní Nešetřilová provádí stejný výpočet a dostane  $p_{25} = -0.365 \times 10^{10}$ . Vyleká se a okamžitě telefonuje zpátky manželovi, že bankéři jsou vydřiduši a že jejich syn by po 25. narozeninách byl tak nanejvýš velkým dlužníkem, rozhodně ne boháčem. Pobouřená paní Nešetřilova manželovi vynadá a navrhuje banku žalovat pro klamavou reklamu. Pan Nešetřil se neodvažuje manželce odporovat, ale dřív, než podá žalobu, vezme tužku a papír a začne počítat v ruce. Vidí, že kalkulačkám věřit nelze. Když zjistí, že faktická částka, kterou by syn k narozeninám dostal, by byla kladná – něco kolem jedné koruny – koupí manželce iracionálně za posledních  $e$  korun kytku a spěchá domů, aby stihl vykoupit syna Bohuslava.

Pro čtenáře, který si bude provádět kontrolní výpočet, upřesněme, že uvedené částky  $p_{25}$  jsme získali na kalkulačce mobilu značky Motorola a na kalkulačce zabudované v příslušenství k Windows.

Rekonstruuje úvahy, které zkušený matematik pan Nešetřil s tužkou v ruce provedl.

Číslo  $e$  je definováno pomocí posloupnosti  $a_n = \left(1 + \frac{1}{n}\right)^n$ . Tato posloupnost je rostoucí, tj.  $a_n < a_{n+1}$ , a právě reálné číslo, ke kterému se s rostoucím  $n$  hodnoty  $a_n$  přibližují, bylo na počest Eulera označeno  $e$ . V matematickém jazyce se  $e$  nazývá limitou posloupnosti  $a_n$  a zapisujeme

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n .$$

Už samotný Euler věděl, že  $e$  lze vyjádřit i jako limitu jiné rostoucí posloupnosti, totiž posloupnosti

$$b_n = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots + \frac{1}{(n-1)!} + \frac{1}{n!} .$$

Tato posloupnost se pro naše účely bude hodit více. Pro zajímavost srovnáme prvních pár členů obou posloupností.

$n$	$a_n$	$b_n$
1	2	2
2	2,25	2,5
3	2,370370...	2,666666...
4	2,441406...	2,708333...
5	2,48832	2,716666...
⋮		
9	2,581174...	2,718281...
⋮		

Když srovnáme hodnoty  $a_n$  a  $b_n$  v tabulce s uvedenou hodnotou  $e$ , vidíme, že posloupnost  $b_n$  se ke své limitě přibližuje rychleji<sup>1</sup>.

Obecně lze ukázat nerovnosti

$$a_n < b_n < e .$$

Čtenáři, který by se chtěl ve větších podrobnostech věnovat číslu  $e$ , doporučujeme učebnici [2]. Pro ilustraci toho, jak lze od posloupnosti  $a_n$  dospět k posloupnosti  $b_n$ , odvodíme první ze dvou nerovností. Použijeme známou binomickou větu

$$(A + B)^n = A^n + \binom{n}{1} A^{n-1} B^1 + \binom{n}{2} A^{n-2} B^2 + \binom{n}{3} A^{n-3} B^3 + \dots + \binom{n}{n-1} A^1 B^{n-1} + B^n ,$$

kde koeficienty  $\binom{n}{k}$  jsou kombinační čísla definovaná předpisem  $\binom{n}{k} = \frac{n!}{(n-k)!k!}$ , nebo chcete-li, jsou to prvky z  $n$ -tého řádku Pascalova trojúhelníku. Binomickou větu použijeme na výpočet hodnoty  $a_n$ , kde za  $A$  dosadíme 1 a za  $B$  dosadíme  $\frac{1}{n}$ . Dostaneme

$$a_n = 1 + \binom{n}{1} \frac{1}{n} + \binom{n}{2} \frac{1}{n^2} + \binom{n}{3} \frac{1}{n^3} + \dots + \binom{n}{n-1} \frac{1}{n^{n-1}} + \frac{1}{n^n} .$$

<sup>1</sup>Např.  $b_9$  se shoduje s číslem  $e$  na prvních 6 místech za desetinnou čárkou, zatímco  $a_9$  na žádném.

Obecný tvar sčítance v předchozím součtu můžeme upravit

$$\binom{n}{k} \frac{1}{n^k} = \frac{n!}{(n-k)!k!} \frac{1}{n^k} = \frac{1}{k!} \frac{n(n-1)(n-2)\dots(n-k+1)}{n^k}.$$

Poněvadž každý člen součtinu  $\frac{n(n-1)(n-2)\dots(n-k+1)}{n^k} = \frac{n}{n} \frac{n-1}{n} \frac{n-2}{n} \dots \frac{n-k+2}{n} \frac{n-k+1}{n}$  je nanejvýš jedna, lze odhadnout  $\binom{n}{k} \frac{1}{n^k} \leq \frac{1}{k!}$  a celkově

$$a_n < 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots + \frac{1}{(n-1)!} + \frac{1}{n!} = b_n.$$

Vzhledem ke zmíněné nerovnosti  $a_n < b_n < e$  má posloupnost  $b_n$  limitu  $e$ . Číslo  $e$  si tedy můžeme představit jako výsledek nekonečného součtu

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \frac{1}{5!} + \dots$$

Začněme nyní vyjádřovat částky  $p_n$  nespořené v bance,

$$\begin{aligned} p_1 &= 1 \cdot (e - 1) &= \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \frac{1}{5!} + \dots \\ p_2 &= 2 \cdot (p_1 - 1) &= 2 \left( \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \frac{1}{5!} + \dots \right) \\ p_3 &= 3 \cdot (p_2 - 1) &= 3 \cdot 2 \left( \frac{1}{3!} + \frac{1}{4!} + \frac{1}{5!} + \frac{1}{6!} + \dots \right) \\ p_4 &= 4 \cdot (p_3 - 1) &= 4 \cdot 3 \cdot 2 \left( \frac{1}{4!} + \frac{1}{5!} + \frac{1}{6!} + \frac{1}{7!} + \dots \right) \\ p_5 &= 5 \cdot (p_4 - 1) &= 5! \left( \frac{1}{5!} + \frac{1}{6!} + \frac{1}{7!} + \frac{1}{8!} + \dots \right) \\ &\vdots & \\ p_{24} &= 24 \cdot (p_{23} - 1) &= 24! \left( \frac{1}{24!} + \frac{1}{25!} + \frac{1}{26!} + \frac{1}{27!} + \dots \right) \\ p_{25} &= 25 \cdot (p_{24} - 1) &= 25! \left( \frac{1}{25!} + \frac{1}{26!} + \frac{1}{27!} + \frac{1}{28!} + \dots \right) \end{aligned}$$

Je jasné, že částka nespořená po 25. roce musí být kladná a větší než jedna.

Odhadněme, jak je opravdu veliká. Po zkrácení faktoriálů ve vyjádření  $p_{25}$  dostaneme

$$p_{25} = 1 + \frac{1}{26} + \frac{1}{26 \cdot 27} + \frac{1}{26 \cdot 27 \cdot 28} + \frac{1}{26 \cdot 27 \cdot 28 \cdot 29} + \dots < 1 + \frac{1}{26} + \frac{1}{26 \cdot 26} + \frac{1}{26 \cdot 26 \cdot 26} + \frac{1}{26 \cdot 26 \cdot 26 \cdot 26} + \dots$$

Na pravé straně se objevil nekonečný součet geometrické posloupnosti s kvocientem  $q = \frac{1}{26}$ , viz [1]. Protože součet prvních  $N$  členů této posloupnosti je roven  $\frac{1-q^N}{1-q}$  a pro  $N$  blížící se do nekonečna se naše  $q^N$  blíží k 0, dostaneme pro součet nekonečně mnoha členů součet  $\frac{1}{1-q}$ , a tedy celkový odhad

$$p_{25} < \frac{1}{1 - \frac{1}{26}} = \frac{26}{25} = 1,04.$$

Hodnotu  $p_{25}$  můžeme odhadnout zdola jednoduše součtem prvních dvou členů geometrické posloupnosti. Získáme tak dolní odhad

$$p_{25} > 1 + \frac{1}{26} = \frac{27}{26} = 1,038.$$

Vidíme, že skutečná výše uspořené peněz  $p_{25}$  se podstatně liší od výsledků, které dostaneme pomocí kalkulačky mobilu nebo kalkulačky počítače. Paradoxní je, že i když např. displeje různých mobilů umožňují pracovat se stejným počtem platných míst, můžete dostat různé výsledky (zkuste si výpočet  $p_{25}$  i s mobily svých spolužáků). Jak je to možné? Můžeme vůbec důvěřovat výsledkům, které dostaneme pomocí výpočetní techniky?

Prvním krokem k pochopení zdánlivého paradoxu je zjištění, že zabudované programy pro násobení a sčítání jsou různé. Mohou se lišit v tom, jak a kdy zaokrouhlují – zda před převodem do dvojkové soustavy, ve které většina algoritmů pracuje, nebo až po tomto převodu. Vyskytují se i kalkulačky, které pracují pouze v desítkové soustavě. Při výpočtu mohou dále používat více platných cifer, než kolik se jich nakonec ukáže na displeji, atp. Zjistit, jakou zaokrouhlovací taktiku vlastně váš počítač používá, je téměř nemožné.

I našeho čtenáře jistě napadne, že v předchozím příkladě by k vyřešení celé záhady mohlo pomoci použití algoritmu, který by počítal prostě na více platných cifer. Při práci v pohyblivé desetinné čárce dnes samozřejmě existuje možnost měnit přesnost výpočtů velmi snadno. Pro naše potřeby jsme proto použili programovací jazyk Maple, který byl vyvinut na univerzitě v kanadském Waterloo (proto název javorový (list) - v angličtině maple (leaf)). V následující tabulce

$D$	$p_{25}$
17	$-0,54848 \times 10^9$
18	$0,71967 \times 10^8$
19	$-0,55884 \times 10^7$
20	615990
21	-4457,9
22	-4457,9
23	195,37
24	40,262
25	-6,2713
26	1,4842
27	1,0189
28	1,0344
29	1,0406
30	1,0399
31	1,0399

symbol  $D$  znamená počet platných míst, které jsme programu Maple předepsali pro provádění výpočtu. Z hodnot  $p_{25}$  napočítaných Maplem uvádíme v tabulce pouze prvních pět platných míst<sup>2</sup>.

Při pohledu na naši tabulku se musíme zeptat: Podle jakého obecného pravidla máme tedy *předem* zvolit počet platných míst pro výpočet, abychom se vyvarovali velkých chyb? Nebo naopak: Co nám tedy při dané přesnosti výpočtu vlastně na displeji počítač ukazuje, když ne správné hodnoty?

Označme tyto hodnoty z displeje  $\tilde{p}_1, \tilde{p}_2, \tilde{p}_3, \dots, \tilde{p}_{25}$  a zkusme odvodit, s jakou přesností by musel počítač pracovat, aby se zobrazená částka  $\tilde{p}_{25}$  lišila od skutečné  $p_{25}$  o méně než korunu. Odchylku výsledku operace prováděné počítačem od výsledku operace, který bychom měli dostat při absolutně přesných výpočtech, budeme označovat písmenkem  $\varepsilon$ , a to bude mít v indexu příslušné pořadové číslo operace. Odchylka  $\varepsilon_n$  může být kladná nebo záporná podle toho, zda se zaokrouhlí nahoru nebo dolů. Tedy

$$\begin{aligned} \tilde{p}_1 &= e - 1 + \varepsilon_1 = p_1 + \varepsilon_1 \\ \tilde{p}_2 &= 2(\tilde{p}_1 - 1) + \varepsilon_2, \quad \text{a odtud} \quad \tilde{p}_2 = 2(p_1 + \varepsilon_1 - 1) + \varepsilon_2 = p_2 + 2\varepsilon_1 + \varepsilon_2 \\ \tilde{p}_3 &= 3(\tilde{p}_2 - 1) + \varepsilon_3, \quad \text{a odtud} \quad \tilde{p}_3 = 3(p_2 + 2\varepsilon_1 + \varepsilon_2 - 1) + \varepsilon_3 = p_3 + 3 \cdot 2\varepsilon_1 + 3\varepsilon_2 + \varepsilon_3 \\ \tilde{p}_4 &= 4(\tilde{p}_3 - 1) + \varepsilon_4, \quad \text{a odtud} \quad \tilde{p}_4 = p_4 + 4 \cdot 3 \cdot 2\varepsilon_1 + 4 \cdot 3\varepsilon_2 + 4\varepsilon_3 + \varepsilon_4 \end{aligned}$$

atd. Nakonec odvodíme

$$\tilde{p}_{25} = 25(\tilde{p}_{24} - 1) + \varepsilon_{25} = p_{25} + 25!\varepsilon_1 + \frac{25!}{2}\varepsilon_2 + \frac{25!}{2 \cdot 3}\varepsilon_3 + \frac{25!}{2 \cdot 3 \cdot 4}\varepsilon_4 + \dots + \varepsilon_{25}.$$

Rozdíl mezi skutečnou hodnotou  $p_{25}$  a vypočítanou hodnotou  $\tilde{p}_{25}$ , tj. celková

$$\text{chyba} = 25! \left( \varepsilon_1 + \frac{1}{2!}\varepsilon_2 + \frac{1}{3!}\varepsilon_3 + \frac{1}{4!}\varepsilon_4 + \dots + \frac{1}{25!}\varepsilon_{25} \right).$$

Kdybychom měli zaručeno, že velikosti  $\varepsilon$  nepřesáhnou hodnotu  $10^{-N}$  (to lze docílit tím, že budeme zaokrouhlovat na  $N$  platných desetinných míst), bude chyba v absolutní hodnotě menší než

$$|\text{chyba}| \leq 25! \left( 1 + \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \dots + \frac{1}{25!} \right) 10^{-N} < 25! 10^{-N} e.$$

<sup>2</sup>Na první pohled překvapí, že pro volbu přesnosti výpočtu  $D = 21$  a  $D = 22$  dostaneme stejnou (a přesto nesprávnou) hodnotu  $p_{25}$ . Čtenář si snad vysvětlí tento jev, když se podívá na 21. číslici za desetinnou čárkou v zápisu čísla  $e$

Bude-li tedy  $N$  alespoň tak velké, že  $25! 10^{-N} e < 1$ , budeme mít zaručeno, že chyba výsledku je menší než jedna. Nejmenší takové  $N$  je  $N = 26$ . Tedy pro naši úlohu je třeba počítat s 26 místy za desetinnou čárkou.

A co dodat nakonec? Naše úloha byla velice jednoduchá, v podstatě se v ní 25 krát násobilo. Skutečné úlohy z praxe konstruktérů letadel nebo jaderných elektráren, tvůrců modelů pro předpovědi počasí nebo pro plánování letů do kosmu jsou daleko složitější. Číslo  $10^{10}$  prováděných operací není nijak neobvyklý počet. Proto důsledná analýza numerických chyb je nezbytnou součástí každého špičkového softwarového produktu. Komerčně nabízené programy nám velice usnadňují práci, musíme si však být vědomi i jejich – ne vždy na první pohled zřejmých – omezení.

Traduje se, že věhlasný matematik Householder – zakladatel oboru numerická matematika – nikdy necestoval letadlem. Prý si byl až moc dobře vědom, kde všude mohlo při výpočtech spojených s konstrukcí letadel dojít k chybě.

### **Použitá literatura:**

1. O. Odvárko, Matematika pro gymnázia, Posloupnosti a řady, Prometheus, 1995.
2. S. Pošta, P. Pošta, Analýza v příkladech, skriptum ČVUT, 2009.