

# O predstavitvi podatkov v računalniku: decimalna števila



JURE SLAK

→ Z računalniki lahko dandanes izvajamo kompleksne matematično-fizikalne izračune in simulacije. V marcu 2020 je omrežje Folding@Home, ki se uporablja za izračune zlaganja proteinov, preseglo 1.5 exaFLOPS-a, to je več kot 1,500,000,000,000,000 računskih operacij na sekundo. Že v enoti sami, ki jo uporabljamo za merjenje hitrosti FLOPS, se skriva glavni podatkovni tip, ki stoji za vsemi simulacijami. FLOPS namreč pomeni *floating point operation per second* in pove, koliko računskih operacij z decimalnimi števili lahko naredimo v eni sekundi. V tem prispevku se bomo poglobili v to, kako decimalna števila sploh predstavimo v računalniku in kako z njimi računamo.

## Decimalna ali realna števila

Velikokrat se pogovorno reče, da v računalniku hranimo realna števila. Nekateri programski jeziki, npr. različne verzije SQL, tudi uporabljajo besedo `real` za oznako tipa. Vendar kljub temu v računalniku nekaterih realnih števil ne moremo predstaviti zelo enostavno. Iracionalna števila, kot npr.  $\sqrt{2}$  ali  $\pi$ , običajno le aproksimiramo. Ravno število  $\pi$  je v programskem jeziku C definirano kot

```
#define M_PI 3.14159265358979323846,
```

torej »le« na 20 decimalnk, precej manj kot trenutni slovenski rekord 3333 decimalnk, ki jih je znal na pamet povedati zmagovalec zadnjega  $\pi$ -dneva. Kot bomo videli, računalnik uporablja le racionalna števila omejene natančnosti – morebitna realna števila

so ustrezno zaokrožena. To nam omogoča enostavnost računanja in hitrost, natančnost pa ni največja. A brez skrbi, večinoma so decimalna števila, ki jih uporablja računalnik, povsem dovolj natančna.

## Fiksna in plavajoča pika

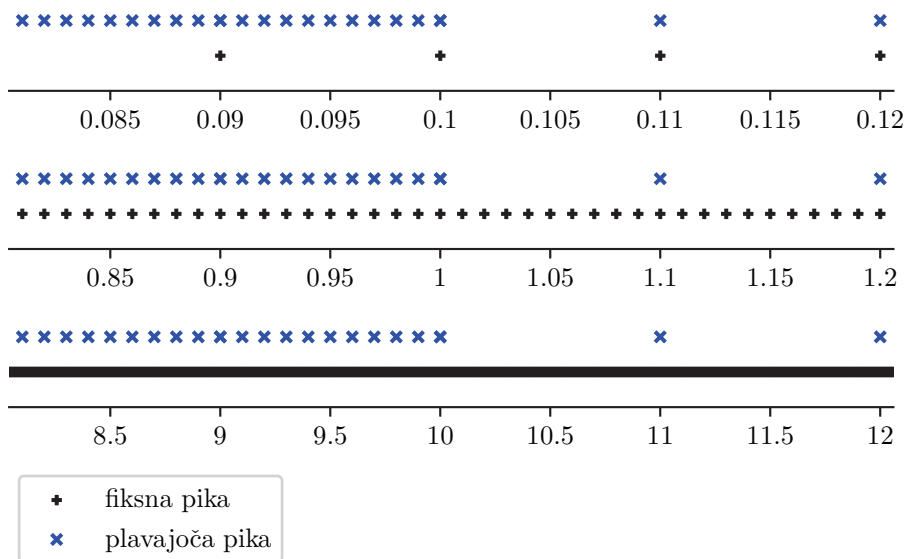
FLOPS se v prvih dveh črkah nanaša na operacije s plavajočo piko.<sup>1</sup> Predstavitev decimalnih števil s plavajočo piko je ena izmed dveh osnovnih načinov predstavitve števil. Druga, enostavnejša možnost je predstavitev s fiksno decimalno piko. Pri slednji imamo na voljo nekaj mest za del pred piko in nekaj mest za decimalno piko. Denimo, da imamo dve mesti pred in dve mesti po piki. Števila, ki jih lahko zapišemo, so torej od 00.00 do 99.99. Interval od [0, 100] smo enakomerno pokrili s 10000 števili na razdalji 0.01. Tem številom rečemo *predstavljiva*, vseh ostalim pa *nepredstavljiva*. Vsako realno število, s katerim želimo računati, moramo zaokrožiti; ponavadi izberemo najbližje predstavljivo število. Za število  $x$  bomo z  $\text{fl}(x)$  označili predstavljivo število, kjer zaokrožimo  $x$ .

Seštevanje in odštevanje predstavljivih števil je enostavno in natančno, lahko pa pride do prekoračitve ali podkoračitve; to pomeni, da rezultat leži izven razpona možnih vrednosti. Pri množenju in deljenju pridemo do novih težav. Če pomnožimo 0.5 in 0.41, pri točnem računanju dobimo 0.205, kar pa ni predstavljivo; rezultat je potrebno zaokrožiti na najbližje predstavljivo število. V našem primeru imamo dve izbiri: zaokrožimo lahko na 0.20 ali na 0.21. Čeprav je v resničnem življenju pogosto, da polovice zaokrožamo navzgor, v računalništvu pona-

<sup>1</sup>Pogosto se jim v slovenščini reče tudi števila s plavajočo vejico, saj je vejica v slovenščini ločilo, ki se uporablja za označevanje decimalnih mest. Vendar je v računalništvu precej bolj pogosta pika, zato jo bomo uporabljali tudi v tem prispevku.







SLIKA 1.

Primerjava pogostosti predstavljivih števil na treh različnih delih realne osi. Števila s fiksno piko imajo enako absolutno natančnost ne glede na lokacijo, števila s premično piko pa postajajo čedalje bolj redka, toda ohranjajo enako relativno natančnost. Na grafih se spreminjajo enote, zato števila s fiksno piko izgledajo, kot da so čedalje bolj gosta, števila s premično piko pa se zdijo enako gosta.

kar je popolnoma točno. Če bi dodali malo več decimalk in izračunali  $0.053 \times 0.082$ , bi izračun potekal tako:

$$\begin{aligned}
 \blacksquare \quad & 0.53 \cdot 10^{-1} \times 0.82 \cdot 10^{-1} = \text{fl}(0.4346 \cdot 10^{-1}) \\
 & = 0.43 \cdot 10^{-2}.
 \end{aligned}$$

V tem primeru rezultat ni točen, je pa precej boljši kot pri fiksni piki. Njegova relativna napaka je približno 1%.

Poučno je tudi, kaj se zgodi, če izračunamo npr.  $100 + 0.1$ . Dobimo

$$\blacksquare \quad 0.1 \cdot 10^3 \times 0.1 = \text{fl}(0.1001 \cdot 10^3) = 0.1 \cdot 10^3.$$

Rezultat je zopet točno 100, saj vrednost 0.1 ni bila dovolj velika, da bi jo upoštevali, in se je zgubila pri zaokroževanju. Še vedno pa je to znotraj zaokrožitvene napake: 0.1 predstavlja le 0.1% od 100, kar je močno znotraj dovoljene 5% napake.

### Dejanska števila v računalniku

Pri delu z decimalnimi števili v računalniku ne uporabljamo le dveh decimalnih mest, kot smo jih mi do sedaj, a principi računanja kljub temu ostanejo enaki. Delo z decimalnimi števili predpisuje standard IEEE 754. Glavni tip, ki ga najpogosteje uporabljamo, se imenuje `double`. Velik je 64 bitov in se imenuje dvojna natančnost – tako ime ima, ker je

dvakrat večji od 32-bitnega tipa `single`, ki predstavlja enojno natančnost. Decimalna števila so shranjena v dvojiškem sistemu, ne v desetiškem. Primer števila bi bilo npr.

$$\blacksquare \quad 0.101101 \times 2^3,$$

kar pretvorimo v

$$\begin{aligned}
 \blacksquare \quad & (1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4} + 0 \cdot 2^{-5} + 1 \cdot 2^{-6}) \\
 & \times 2^3 = 5.625.
 \end{aligned}$$

Izmed 64 bitov, ki so na voljo, jih je 11 rezerviranih za eksponent, 52 za mantiso (decimalke), in 1 za predznak števila (+ ali –). 11-bitni prosto za eksponent nam omogoča eksponente od –1022 do 1023, najmanjši in največji eksponent –1023 in 1024 pa sta rezervirana za posebna števila, o katerih bomo več povedali kasneje. V dvojiškem sistemu velja omeniti tudi posebno obliko normaliziranega zapisa. Če število zapišemo v običajnem normaliziranem zapisu, bo oblike  $0.x$ , kjer je  $x$  neničelna številka. Toda v dvojiškem je ta številka lahko le 1, zato je nepotrebno, da jo shranjujemo, in se raje dogovorimo, da kot normalizirano obliko za dvojiška decimalna števila vzamemo  $1.x$ , kjer je  $x$  poljubna, 0 ali 1. Tako najbolje izkoristimo 52 dvojiških decimalk, ki jih imamo na voljo, kar je enako približno 16 desetiškim decimalkam natančnosti. V tem sistemu je osnovna zaokrožitvena napaka enaka  $1.11 \cdot 10^{-16}$ .



→ Decimalno število s 64 biti je v računalniku predstavljeno kot:

$$\blacksquare s \underbrace{eee\dots eee}_{11 \text{ bitov za eksponent}} \underbrace{mmm\dots mmm}_{52 \text{ bitov za mantiso}},$$

kjer  $s$  predstavlja en bit za predznak.

Poglejmo si konkreten primer. Bitni zapis števila 4269.6842 je enak

$$\blacksquare 0 \ 10000001011 \ 0000101011011010111100100 \ 11110111011001011111101100,$$

kjer so vrinjeni presledki za lažje branje. Prvi bit je enak 0, kar nam pove, da je število pozitivno. Sledi eksponent; če 10000001011 pretvorimo v desetiško, dobimo 1035. Vendar so eksponenti zamaknjeni, saj nimajo razpona od 0 do 2048, temveč od  $-1023$  do 1024. Z upoštevanjem zamika je iskani eksponent enak  $1035 - 1023 = 12$ . Zapisano število je tako enako

$$\blacksquare 1.000010101101101011110010011110111011 \ 001011111101100 \times 2^{12},$$

kar je enako

$$\blacksquare 1000010101101.101011110010011110111011 \ 001011111101100$$

oz. približno 4269.6841999999996915. Kot vidimo, rezultat ni točno 4269.6842, temveč je za približno  $3.09 \cdot 10^{-13}$ , kar (relativno gledano) ustreza osnovni zaokrožitveni napaki.

### Posledice zaokroževanja

Zaokroževanje na najbližjo vrednost pomeni, da običajna pravila računanja ne držijo več. Če v računalniku izračunamo  $(a + b) + c$ , to ni več nujno enako kot  $a + (b + c)$ . Poglejmo primer: vzemimo  $a = 0.88$ ,  $b = 0.56$  in  $c = 0.13 \cdot 10^1$ . Če izračunamo  $a + b + c$  od leve proti desni, dobimo

$$\begin{aligned} \blacksquare a + b + c &= 0.88 + 0.56 + 0.13 \cdot 10^1 \\ &= \text{fl}(1.44) + 0.13 \cdot 10^1 \\ &= 0.14 \cdot 10^1 + 0.13 \cdot 10^1 \\ &= \text{fl}(0.27 \cdot 10^1) = 0.27 \cdot 10^1. \end{aligned}$$

Če pa izračunamo  $a + b + c$  od desne proti levi, dobimo

$$\begin{aligned} \blacksquare a + b + c &= 0.88 + 0.56 + 0.13 \cdot 10^1 \\ &= 0.88 + \text{fl}(1.86) = 0.88 + 0.19 \cdot 10^1 \\ &= \text{fl}(2.78) = 0.28 \cdot 10^1. \end{aligned}$$

Razlika je majhna, vendar števili nista enaki. To je tudi eden izmed razlogov, da decimalna števila redko neposredno primerjamo, ali imajo popolnoma enako vrednost. Že majhne razlike v načinu izračuna namreč lahko prinesejo napake pri zadnjih decimalkah. Poslužimo se raje primerjanja s *toleranco*: namesto da bi pogledali, ali je  $a = b$ , pogledamo, ali je absolutna vrednost razlike med  $a$  in  $b$  manjša od tolerance  $t$ :  $|a - b| \leq t$ , za npr.  $t = 0.00001$ . Z izbiro vrednosti  $t$  lahko tudi določimo, kako velike napake so še sprejemljive.

Še ena zanimivost se pojavi pri računanju aritmetične sredine dveh števil. V matematiki smo navajeni, da aritmetična sredina  $\frac{a+b}{2}$  dveh števil  $a$  in  $b$  leži natančno na sredini med številoma. Pri decimalnih številih temu ni tako: vzemimo npr.  $a = 0.21$  in  $b = 0.24$ . Njuna izračunana aritmetična sredina je

$$\begin{aligned} \blacksquare \frac{a + b}{2} &= \frac{1}{2}(0.21 + 0.24) = \frac{1}{2} \text{fl}(0.43) = \\ &= \frac{1}{2}(0.43) = \text{fl}(0.215) = 0.22, \end{aligned}$$

kar ni točen rezultat 0.215, toda je (eden izmed) najboljših možnih približkov.

Oglejmo si še en primer izračuna aritmetične sredine, tokrat za  $a = 0.66$  in  $b = 0.67$ :

$$\begin{aligned} \blacksquare \frac{a + b}{2} &= \frac{1}{2}(0.66 + 0.67) = \frac{1}{2} \text{fl}(1.33) = \\ &= \frac{1}{2}(0.13 \cdot 10^1) = \text{fl}(0.65) = 0.65. \end{aligned}$$

Tokrat smo dobili, da je sredina števil 0.66 in 0.67 enaka 0.65, kar seveda leži izven intervala  $[0.66, 0.67]$ ! Če pogledamo izračun, vidimo, da je bil glavni krivec za izgubo natančnosti to, da smo zašli v prevelika števila. Pri predstavitvi 1.33 smo lahko obdržali le dve mesti in smo bili prisiljeni zadnjo zavreči, da smo shranili 1.3.

Izračun bi lahko popravili tako, da bi ga namesto  $\frac{a+b}{2}$  napisali kot  $a + \frac{b-a}{2}$ . V tem primeru ne bi prišlo

