

# Predavanje 1

## Uvod v numerične metode

### Prvi sklop izročkov

Fakulteta za računalništvo in informatiko  
Univerza v Ljubljani

5. oktober 2020

1/25

## Viri in obveznosti

### Viri:

- ▶ Bojan Orel, Osnove numerične matematike, Založba FE in FRI.
- ▶ Bor Plestenjak: Razširjen uvod v numerične metode, DMFA založništvo.

### Obveznosti:

- ▶ Predavanja: 3 ure na teden
- ▶ Vaje: laboratorijske vaje vsak teden po dve uri.
- ▶ 3 domače naloge
- ▶ Pisni izpit iz teorije

2/25

# Vsebina predmeta

1. Računanje in vloga napak pri numerični matematiki.
2. Sistemi linearnih enačb
  - ▶ Gausova eliminacija in LU razcep
  - ▶ Delno/kompletno pivotiranje - cena in občutljivost
  - ▶ Iterativne metode - Jacobijeva in Seidlova iteracija
3. Reševanje nelinearnih enačb
  - ▶ Bisekcija
  - ▶ Sekantna metoda in regula falsi
  - ▶ Tangentna metoda
  - ▶ Metoda fiksne točke
4. Aproksimacija in interpolacija
  - ▶ Lagrangeov in Newtonov interpolacijski polinom
  - ▶ Aproksimacija po metodi najmanjših kvadratov
5. Numerično odvajanje in integriranje
  - ▶ Trapezna metoda
  - ▶ Simpsonova metoda
  - ▶ Rombergova metoda
6. Diferencialne enačbe
  - ▶ Eulerjeva metoda
  - ▶ Runge-Kutta metode

3/25

## Numerično in simbolno računanje

### Numerično računanje:

- ▶ Takoj v formulo vstavljamo **števila**
- ▶ Pridemo do numeričnega rezultata - **numerične rešitve**

### Simbolno računanje:

- ▶ **simboli** predstavljajo števila
- ▶ izraz preoblikujemo s simbolnim računanjem do novega simbolnega izraza - **analitična rešitev**

### Primer

- ▶ *Numerično:*

$$\frac{(17.36)^2 - 1}{17.36 + 1} = 16.36; \quad 0.25, 0.33333 \dots (?), 3.14159 \dots (?)$$

- ▶ *Simbolno:*

$$\frac{x^2 - 1}{x + 1} = x - 1; \quad \frac{1}{4}, \frac{1}{3}, \pi, \tan 83$$

4/25

# Metoda, algoritem in implementacija

**Metoda**. . . matematična konstrukcija, s katero rešujemo problem

**Algoritem**. . . koraki metode

**Implementacija**. . . zapis algoritma v izbranem jeziku

'Dober' v NM. . . majhna sprememba  $\Rightarrow$  majhna napaka

## Vprašanja:

- ▶ Ali je problem občutljiv?
- ▶ Ali je metoda 'dobra'?
- ▶ Ali je algoritem robusten - deluje na širokem spektru problemov?
- ▶ Ali je implementacija hitra - časovna in prostorska zahtevnost?

5/25

## Občutljivost problema

Tipi problemov (glede na spremembo rezultata ob majhni spremembi začetnih podatkov):

- ▶ **občutljivi** oz. **slabo pogojeni** problemi.
- ▶ **neobčutljivi** oz. **dobro pogojeni** problemi.

Občutljivost je odvisna le od narave problema in ne od izbrane numerične metode.

### Primer (Presečišča premic)

#### 1. Sistem in njegova perturbacija

$$\begin{aligned}x + y = 2 &\rightarrow x + y = 1.9999 \\x - y = 0 &\rightarrow x - y = 0.0002\end{aligned}$$

ima rešitvi  $x = y = 1$  oz.  $x = 1.00005$  in  $y = 0.99985$ . Začetni sistem je neobčutljiv, saj je šlo za spremembo za isti velikostni razred.

#### 2. Sistem in njegova perturbacija

$$\begin{aligned}x + 0.99y = 1.99 &\rightarrow x + 0.99y = 1.9899 \\0.99x + 0.98y = 1.97 &\rightarrow 0.99x + 0.98y = 1.9701\end{aligned}$$

ima rešitvi  $x = y = 1$  oz.  $x = 2.97$  in  $y = -0.99$ . Začetni sistem je občutljiv, saj je majhna sprememba začetnih podatkov povzročila veliko spremembo rezultata.

6/25

- ▶ Neskončne procese nadomestimo s končnimi (Taylorjeva vrsta).
- ▶ Neskončno razsežne prostore nadomestimo s končno razsežnimi (funkcijske enačbe rešujemo v prostoru polinomov).
- ▶ Diferencialne enačbe nadomestimo z algebraičnimi.
- ▶ Nelinearne probleme nadomestimo z linearnimi.
- ▶ Matrike nadomestimo z enostavnejšimi.

## Numerična matematika - zakaj?

- ▶ Znanost, ki temelji na matematičnih izračunih, je neposredno odvisna od NM.
- ▶ Nekatero katastrofo so se zgodile zaradi slabega numeričnega računanja: (<http://www-users.math.umn.edu/~arnold//disasters/>)
  - ▶ **Nesreča Misije Patriot, Dharan, Savdska Arabija, 25.2.1991**, 28 žrtev: **slaba analiza zaokrožitvenih napak**.
  - ▶ **Eksplozija rakete Ariana 5, Francoska Gvajana, 4.6.1996**: **posledica prekoračitve obsega števil**. [https://www.youtube.com/watch?v=PK\\_yguLapgA](https://www.youtube.com/watch?v=PK_yguLapgA)  
<https://www.youtube.com/watch?v=W3YJeoYgozw>  
<https://www.arianespace.com/vehicle/ariane-5/>
  - ▶ **Potop ploščadi Sleipner A, Stavanger, Norveška, 12.8.1991**, milijarda dolarjev škode: **nenatančna obdelava obremenitev pri reševanju PDE-jev**.  
<https://www.youtube.com/watch?v=eGdiPs4THW8>

# Predstavljliva števila

Števila shranjujemo v obliki

$$x = \pm(0.d_1d_2d_3 \dots d_m)_\beta \times \beta^e,$$

kjer je

- ▶  $\beta$  naravno število,
- ▶  $d_1d_2d_3 \dots d_m$  mantisa,  $e$  eksponent,
- ▶  $e$  celo število,

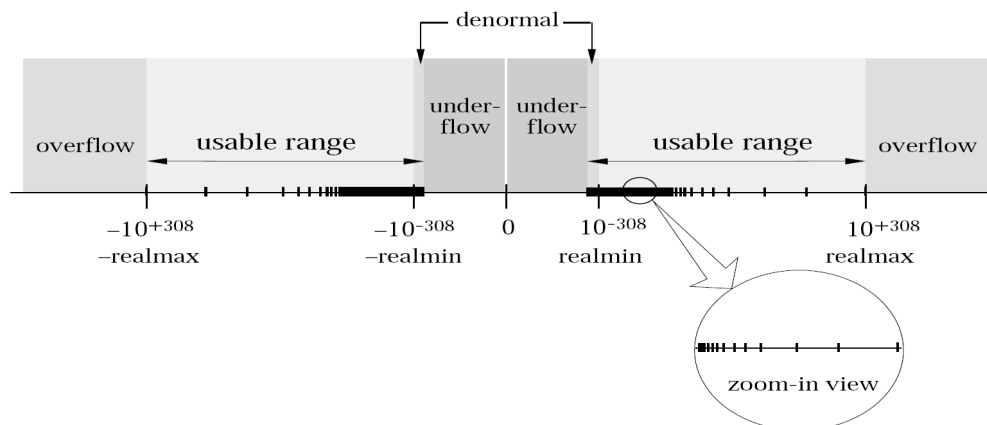
## Primer (baza 10)

- ▶ 1000.12345 zapišemo kot  $+(0.100012345)_{10} \times 10^4$ .
- ▶ 0.000812345 zapišemo kot  $+(0.812345)_{10} \times 10^{-3}$ .

9/25

## Prekoračitev in podkoračitev

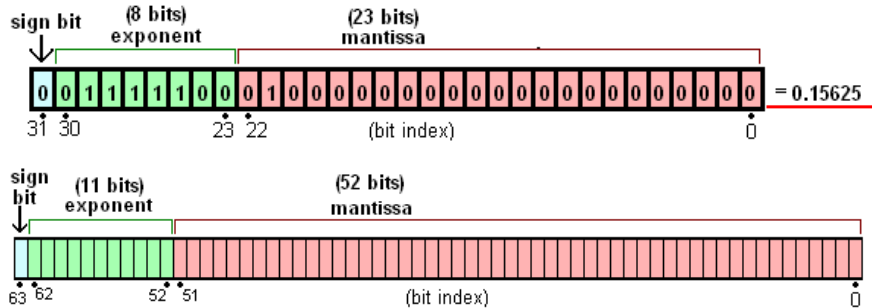
### Floating Point Number Line



- ▶ izračuni preblizu 0 lahko povzročijo **podkoračitev**
- ▶ preveliki izračuni lahko povzročijo **prekoračitev**
- ▶ prekoračitev je v splošnem hujši problem

10/25

- ▶ IEEE Enojna natančnost: števila so predstavljena z 32 biti.
- ▶ IEEE Dvojna natančnost: števila so predstavljena z 64 biti.



## Zaokrožitvene napake

- ▶ Vseh realnih števil ne moremo predstaviti v strojni aritmetiki  $\Rightarrow$  **zaokrožujemo** in delamo **zaokrožitvene napake**.
- ▶ IEEE standard... **zaokroži do najbližjega predstavljenega števila**
- ▶ Naj bosta  $x_- \leq x \leq x_+$  najbližji predstavljeni števili števila  $x$ :

$$\text{fl}(x) = \begin{cases} x_-, & \text{če je } x - x_- < x_+ - x, \\ x_+, & \text{če je } x - x_- \geq x_+ - x. \end{cases}$$

- ▶ Kako velika je napaka? Recimo, da je  $x$  bližje  $x_-$ :

$$x = (0.1b_2b_3 \dots b_{24}b_{25}b_{26})_2 \times 2^m$$

$$x_- = (0.1b_2b_3 \dots b_{24})_2 \times 2^m, \quad x_+ = \left( (0.1b_2b_3 \dots b_{24})_2 + 2^{-24} \right) \times 2^m$$

$$|x - x_-| \leq \frac{|x_+ - x_-|}{2} = 2^{m-25}, \quad \left| \frac{x - x_-}{x} \right| \leq \frac{2^{m-25}}{1/2 \times 2^m} \leq 2^{-24} =: \frac{1}{2} \epsilon.$$

- ▶  $\epsilon := 2^{1-24}$  ... **strojni epsilon**

- ▶  $u := \frac{1}{2} \epsilon$  ... **osnovna zaokrožitvena napaka**

# Strojni $\epsilon$

Vzemimo  $x = 1.0$  in dodajajmo  $\frac{1}{2}, \frac{1}{4}, \dots, 2^{-i}$ :

↓	←	23 bitov										→		
1	1	0	0	0	0	0	0	0	0	0	0	0	e	e
1	0	1	0	0	0	0	0	0	0	0	0	0	e	e
1	0	0	1	0	0	0	0	0	0	0	0	0	e	e
.....														
1	0	0	0	0	0	0	0	0	0	0	0	1	e	e
1	0	0	0	0	0	0	0	0	0	0	0	0	e	e

- ▶  $\text{fl}(1 + 2^{-23}) \neq 1$ , toda  $\text{fl}(1 + 2^{-24}) = 1$
- ▶ Enojna natančnost:  $2^{-23}$ ; Dvojna natančnost:  $2^{-52}$ .

```
1 >> eps
2 ans = 2.2204e-16
3
4 >> eps('single')
5 ans = 1.1921e-07
```

## Aritmetika v predstavljenih številih

- ▶  $\text{fl}(x) = x(1 + \delta)$ , kjer je  $|\delta| \leq u := \frac{1}{2}\beta^{1-m}$ .
- ▶ Za **predstavljeni** števili  $x, y$  in  $\odot \in \{+, -, *, /\}$  velja

$$\text{fl}(x \odot y) = (x \odot y)(1 + \delta).$$

- ▶  $(a + b) + c \neq a + (b + c)$ :

$$\begin{aligned} (a + b) + c &= \text{fl}(\text{fl}(a + b) + c) = \text{fl}((a + b)(1 + \delta_1) + c) \\ &= [(a + b)(1 + \delta_1) + c](1 + \delta_2) \\ &= [(a + b + c) + (a + b)\delta_1](1 + \delta_2) \\ &= (a + b + c) \left[ 1 + \frac{a + b}{a + b + c} \delta_1(1 + \delta_2) + \delta_2 \right] \end{aligned}$$

Podobno

$$a + (b + c) = (a + b + c) \left[ 1 + \frac{b + c}{a + b + c} \delta_3(1 + \delta_4) + \delta_4 \right].$$

Če pozabimo na člena  $\delta_1\delta_2$  in  $\delta_3\delta_4$ , dobimo

$$(a + b) + c = (a + b + c)(1 + \epsilon_3) \quad \text{kjer je} \quad \epsilon_3 \approx \frac{a + b}{a + b + c} \delta_1 + \delta_2,$$

$$a + (b + c) = (a + b + c)(1 + \epsilon_4) \quad \text{kjer je} \quad \epsilon_4 \approx \frac{b + c}{a + b + c} \delta_3 + \delta_4.$$

**Sklep:** Ko seštevamo  $n$  števil, je najbolje začeti z najmanjšim in prištevati večje.

## Napake pri numeričnem računanju

- ▶ **Neodstranjiva napaka**  $D_n \dots$  nenatančni začetni podatki.
- ▶ **Napaka metode**  $D_m \dots$  npr. neskončni proces aproksimiramo s končnim.
- ▶ **Zaokrožitvena napaka**  $D_z \dots$  računanje s približki in zaokroževanje.

Celotna napaka je

$$D = D_n + D_m + D_z.$$

Primer ( $\sin \frac{\pi}{10}$  računamo v  $\beta = 10$  in  $m = 4$ )

- ▶  $D_n: fl(\frac{\pi}{10}) = 0.3142 \cdot 10^0.$

$$|D_n| \approx \sin'(\frac{\pi}{10})|x - fl(x)| \leq \frac{1}{2} \cdot 10^{-4}.$$

- ▶  $D_m: \sin x \approx x - x^3/6.$

$$|D_m| \leq x^5/120 = 2.6 \cdot 10^{-5}.$$

- ▶  $D_z: fl(x - fl(fl(fl(x \cdot x) \cdot x)/6)).$

$$|D_z| \leq 3.0 \cdot 10^{-5}.$$

15/25

## Stabilnost metode

Stabilnost metode preverimo z **analizo zaokrožitvenih napak**.

Vrste napak ( $x$  naj bo točna vrednost,  $\bar{x}$  pa približek zanjo):

- ▶ Prva delitev:

- ▶ **Absolutna napaka**:  $\bar{x} - x.$

- ▶ **Relativna napaka**:  $\frac{\bar{x} - x}{x}.$

- ▶ Druga delitev:

- ▶ **Direktna napaka**: Absolutna in relativna napaka.

- ▶ **Obratna napaka**: Koliko je potrebno spremeniti začetne podatke, da dobimo izračunan rezultat.

Glede na napake, nas lahko zanima:

- ▶ **direktna stabilnost** metode.
- ▶ **obratna stabilnost** metode.

Velja

$$|\text{direktna napaka}| \approx \text{občutljivost} \times |\text{obratna napaka}|.$$

Izračunana vrednost je blizu pravi, če rešujemo neobčutljiv problem z obratno stabilno metodo.

16/25



## Katastrofalno odštevanje

- ▶ Glavni problem zaokrožitvene aritmetike... **odštevanje/seštevanje dveh približno enakih/nasprotnih števil**

$$a = x.xxxx\ xxxx\ xxx1 \overbrace{ssss \dots}^{\text{izguba}}$$

$$b = x.xxxx\ xxxx\ xxx0 \overbrace{tttt \dots}^{\text{izguba}}$$

$$\begin{array}{r} \text{Potem} \\ \begin{array}{r} \text{končna natančnost} \\ x.xxx\ xxxx\ xxx1 \\ -\ x.xxx\ xxxx\ xxx0 \\ \hline =\ 0.000\ 0000\ 0001 \end{array} \end{array} \underbrace{\hspace{10em}}_{\text{izguba natančnosti}} \underbrace{????\ ????}$$

- ▶ S ponavljanjem se napake seštevajo.

17/25

## Katastrofalno odštevanje

Iščemo rešitve kvadratne enačbe

$$x^2 + 2a^2x - q = 0$$

Rešitev z manjšo absolutno vrednostjo je

$$x_2 = \frac{-2a^2 + \sqrt{4a^4 + 4q}}{2} = -a^2 + \sqrt{a^4 + q}.$$

```
1 s := (a^2)^2
2 t := s + q
3 u := sqrt(t)
4 x2 := -a^2 + u
```

Če je  $a^2$  veliko večji od  $q$ , potem ima lahko korak 4 zgoraj veliko napako. Bolje:

$$x_2 = \left(-a^2 + \sqrt{a^4 + q}\right) \cdot \frac{a^2 + \sqrt{a^4 + q}}{a^2 + \sqrt{a^4 + q}} = \frac{q}{a^2 + \sqrt{a^4 + q}}.$$

```
1 s := (a^2)^2
2 t := s + q
3 u := sqrt(t)
4 v := a^2 + u
5 y := q/v
```

18/25

## Rešitve $x^2 + 2Ax - q = 0$

```
1 >> A = 10000;  
2 >> q = 1;  
3 >> y = -A+sqrt(A^2 + q)  
4 y = 5.0000000055588316e-05  
5  
6 >> y2 = q/(A+sqrt(A^2+q))  
7 y2 = 4.999999987500000e-05  
8  
9 >> x = y;  
10 >> x^2 + 2 * A * x - q  
11 ans = 1.361766321927860e-08  
12  
13 >> x = y2;  
14 >> x^2 + 2 * A * x - q  
15 ans = -1.110223024625157e-16
```

19/25

## Seštevanje in odštevanje nista direktno stabilni operaciji

$x, y \in \mathbb{R}$ . Računamo približek  $\bar{p}$  za  $p = x \pm y$ .

$$\begin{aligned}\bar{p} &= \text{fl}(\text{fl}(x) \pm \text{fl}(y)) = \text{fl}(x(1 + \delta_1) \pm y(1 + \delta_2)) \\ &= (x(1 + \delta_1) \pm y(1 + \delta_2))(1 + \delta_3) \\ &= (x(1 + \delta_1)(1 + \delta_3) \pm y(1 + \delta_2)(1 + \delta_3)) \\ &= (x \pm y + x(\delta_1 + \delta_3 + \delta_1\delta_3) \pm y(\delta_2 + \delta_3 + \delta_2\delta_3))\end{aligned}$$

kjer je  $|\delta_i| \leq u$ . Relativna napaka je

$$\frac{|\bar{p} - p|}{|p|} \leq \frac{|x(\delta_1 + \delta_3 + \delta_1\delta_3) \pm y(\delta_2 + \delta_3 + \delta_2\delta_3)|}{|x \pm y|}.$$

Torej:

$$x + y \approx 0 \quad \Rightarrow \quad \frac{|\bar{p} - p|}{|p|} \text{ veliko.}$$

20/25

## Množenje je direktno stabilna operacija

$x, y \in \mathbb{R}$ . Računamo približek  $\bar{p}$  za  $p = x \cdot y$ .

$$\begin{aligned}\bar{p} &= \text{fl}(\text{fl}(x) \cdot \text{fl}(y)) = \text{fl}(x(1 + \delta_1) \cdot y(1 + \delta_2)) \\ &= x(1 + \delta_1) \cdot y(1 + \delta_2)(1 + \delta_3) \\ &= xy(1 + \delta_1 + \delta_2 + \delta_3 + \text{produkti več } \delta),\end{aligned}$$

kjer je  $|\delta_i| \leq u$ . Relativna napaka je

$$\frac{|\bar{p} - p|}{|p|} \leq \frac{|xy||\delta_1 + \delta_2 + \delta_3 + \mathcal{O}(u^2)|}{|xy|} = |\delta_1 + \delta_2 + \delta_3 + \mathcal{O}(u^2)|.$$

Torej:

Relativna napaka  $\frac{|\bar{p} - p|}{|p|}$  ni odvisna od velikosti produkta  $xy$ .

21/25

## Deljenje je direktno stabilna operacija

$x, y \in \mathbb{R}$ . Računamo približek  $\bar{p}$  za  $p = x/y$ .

$$\begin{aligned}\bar{p} &= \text{fl}(\text{fl}(x)/\text{fl}(y)) = \text{fl}(x(1 + \delta_1)/y(1 + \delta_2)) \\ &= x(1 + \delta_1)/(y(1 + \delta_2))(1 + \delta_3) \\ &= x/y(1 + \delta_1)(1 - \delta_2 + \delta_2^2 - \dots)(1 + \delta_3), \\ &= x/y(1 + \delta_1 - \delta_2 + \delta_3 + \text{produkti več } \delta),\end{aligned}$$

kjer je  $|\delta_i| \leq u$ . Relativna napaka je

$$\frac{|\bar{p} - p|}{|p|} \leq \frac{|x/y||\delta_1 - \delta_2 + \delta_3 + \text{produkti več } \delta|}{|x/y|} = |\delta_1 - \delta_2 + \delta_3 + \mathcal{O}(u^2)|$$

Torej:

Relativna napaka  $\frac{|\bar{p} - p|}{|p|}$  ni odvisna od velikosti kvocienta  $x/y$ .

22/25

## Računanje skalarnega produkta - obratno, a ne direktno stabilno

Računamo skalarni produkt s vektorjev  $x = (x_1, \dots, x_n)$  in  $y = (y_1, \dots, y_n)$ .

```
1  $\bar{s} = x_1 \cdot y_1$ 
2 for  $i = 2$  to  $n$  do
3    $\bar{s} = \bar{s} + x_i \cdot y_i$ 
4 end
5 return  $\bar{s}$ 
```

Velja:

$$\bar{s} = \sum_{i=1}^n x_i \cdot y_i \cdot \underbrace{(1 + 3\delta_i) \cdot \prod_{j=i}^n (1 + \epsilon_j)}_{1 + \gamma_i}.$$

Ocenimo:

$$(1 - 3u)^{n-i+2} \leq 1 + \gamma_i \leq (1 + 3u)^{n-i+2}.$$

Z indukcijo:

$$1 - (n - i + 2)3u \leq (1 - 3u)^{n-i+2}.$$

Z razvojem po binomski formuli:

$$(1 + 3u)^{n-i+2} \leq 1 + (n - i + 2)3u + \mathcal{O}(u^2).$$

Sledi:

$$|\gamma_i| \leq (n - i + 2)3u \leq 3nu.$$

23/25

## Računanje skalarnega produkta - obratno, a ne direktno stabilno

Relativna napaka je

$$\frac{|s - \bar{s}|}{|s|} \leq \frac{3nu \sum_{i=1}^n |x_i| |y_i|}{|x^T y|}.$$

Če sta vektorja blizu pravokotnosti, je relativna direktna napaka lahko zelo velika. Je pa obratno stabilno, saj lahko vzamemo

$$\bar{x}_i = x_i, \quad \bar{y}_i = y_i(1 + \gamma_i) \quad \text{za } i = 1, \dots, n.$$

24/25

# Računanje vrednosti polinoma - obratno, a ne direktno stabilno

Računamo vrednost polinoma

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

v neki točki  $x = x_0$ .

```
1  $\bar{p} = a_n$ 
2 for  $i = 1 : n$ 
3    $\bar{p} = \bar{p} \cdot x + a_{n-i}$ 
4 end
5 return  $\bar{p}$ 
```

Velja

$$\bar{p} = \sum_{i=0}^n a_{n-i} x^{n-i} \underbrace{\prod_{j=i+1}^n (1 + 3\delta_j) \prod_{k=i}^n (1 + \epsilon_j)}_{1 + \gamma_i}.$$

Torej je

$$|\gamma_i| \leq 2(n - i + 1)3u \leq 6nu$$

in računanje vrednosti je obratno, ne pa direktno stabilno.