

## Stopnja občutljivosti

Stopnjo občutljivosti merimo z razmerjem med velikostjo spremembe rezultata in velikostjo spremembe podatkov.

*Zgled:* Naj bo  $f : \mathbb{R} \rightarrow \mathbb{R}$  zvezna in odvedljiva funkcija. Zanima nas razlika med  $f(x)$  in  $f(x + \delta x)$ , kjer je  $\delta x$  majhna motnja.

Absolutna občutljivost: Iz ocene

$$|f(x + \delta x) - f(x)| \approx |f'(x)| \cdot |\delta x|,$$

sledi, da je  $|f'(x)|$  *absolutna občutljivost*  $f$  v točki  $x$ .

Relativna občutljivost: Iz ocene

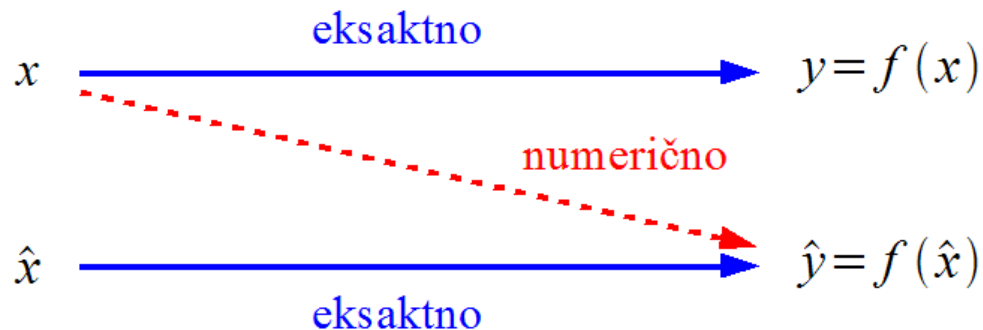
$$\frac{|f(x + \delta x) - f(x)|}{|f(x)|} \approx \frac{|f'(x)| \cdot |x|}{|f(x)|} \cdot \frac{|\delta x|}{|x|}$$

sledi, da je  $\frac{|f'(x)| \cdot |x|}{|f(x)|}$  *relativna občutljivost*  $f$  v točki  $x$ .

## 1.4 Stabilnost metode

Pri računskem procesu pravimo, da je *stabilen* oz. *nestabilen*, ločimo pa *direktno* in *obratno stabilnost*. S tem se ukvarja *analiza zaokrožitvenih napak*.

- **direktna analiza:** Iz  $x$  namesto  $y = f(x)$  izračunamo  $\hat{y}$ . Če je razlika med  $y$  in  $\hat{y}$  majhna (absolutno oz. relativno), je proces direktno stabilen (absolutno oz. relativno), sicer pa nestabilen.
- **obratna analiza:** Iz  $x$  namesto  $y = f(x)$  izračunamo  $\hat{y}$ . Sedaj se vprašamo, za koliko moramo spremeniti argument  $x$  v  $\hat{x}$ , da bo  $f(\hat{x}) = \hat{y}$ . Če je razlika med  $x$  in  $\hat{x}$  majhna (absolutno oz. relativno), je proces obratno stabilen (absolutno oz. relativno), sicer pa nestabilen.



## Občutljivost, stabilnost in natančnost

Algoritem je stabilen, če so rezultati, ki jih vrne, relativno neobčutljivi na motnje, ki se pojavijo zaradi zaokrožitvenih napak med samim računanjem.

Obratno stabilen algoritem tako vrne točno rešitev bližnjega problema.

Če je problem občutljiv, se točna rešitev bližnjega problema lahko zelo razlikuje od točne rešitve začetnega problema in izračunani rezultat je nenatančen.

Nenatančnost je tako lahko posledica:

- uporabe stabilnega algoritma na občutljivem problemu,
- uporabe nestabilnega algoritma na neobčutljivem problemu.

Natančnost je zagotovljena, kadar neobčutljiv problem rešimo s stabilno numerično metodo.

## 1.5 Tri vrste napak pri numeričnem računanju

Računamo vrednost funkcije  $f : X \rightarrow Y$  pri danem  $x$ . Numerična metoda vrne približek  $\hat{y}$  za  $y$ , razlika  $D = y - \hat{y}$  pa je *celotna napaka* približka.

Izvori napake so:

- nenatančnost začetnih podatkov,
- napaka numerične metode,
- zaokrožitvene napake med računanjem.

## Celotno napako lahko razdelimo na tri dele

**Neodstranljiva napaka:** Namesto z  $x$  računamo s približkom  $\bar{x}$  in namesto  $y = f(x)$  izračunamo  $\bar{y} = f(\bar{x})$ . Neodstranljiva napaka je  $D_n = y - \bar{y}$ .

$D_n$  je posledica napak začetnih podatkov.

Zgled: Računanje  $\sin(\pi/10)$  z osnovnimi operacijami v  $P(10, 4, -5, 5)$

Namesto z  $x = \pi/10$  računamo z  $\bar{x} = 0.3142 \cdot 10^0$

$$D_n = y - \bar{y} = \sin(\pi/10) - \sin(0.3142) = -3.9 \cdot 10^{-5}$$

**Napaka metode:** Namesto  $f$  računamo vrednost funkcije  $g$ , ki jo lahko izračunamo s končnim številom operacij. Namesto  $\bar{y} = f(\bar{x})$  tako izračunamo  $\tilde{y} = g(\bar{x})$ . Napaka metode je  $D_m = \bar{y} - \tilde{y}$ .

Pri sami numerični metodi pogosto neskončen proces nadomestimo s končnim (seštejemo le končno členov neskončne vrste, po končnem številu korakov prekinemo iterativno metodo).

Zgled: Namesto  $\sin(\bar{x})$  izračunamo  $g(\bar{x})$  za  $g(x) = x - x^3/6$ .

$$D_m = \bar{y} - \tilde{y} = 2.5 \cdot 10^{-5}$$

## Celotna napaka

**Zaokrožitvena napaka:** Pri računanju  $\tilde{y} = g(\bar{x})$  se pri vsaki računski operaciji pojavi zaokrožitvena napaka, tako da namesto  $\tilde{y}$  izračunamo  $\hat{y}$ . Sama vrednost  $\hat{y}$  je odvisna od vrstnega reda operacij in načina izračuna  $g(\bar{x})$ . Zaokrožitvena napaka je  $D_z = \tilde{y} - \hat{y}$ .

Zgled:  $D_z$  je odvisna je od vrstnega reda in načina računanja  $g(\bar{x})$ . Primer:

$$\begin{aligned}a_1 &= fl(\bar{x} \cdot \bar{x}) = fl(0.09872164) = 0.9872 \cdot 10^{-1} \\a_2 &= fl(a_1 \cdot \bar{x}) = fl(0.03101154) = 0.3101 \cdot 10^{-1} \\a_3 &= fl(a_2/6) = fl(0.0051683\dots) = 0.5168 \cdot 10^{-2} \\ \hat{y} &= fl(\bar{x} - a_3) = fl(0.309032) = 0.3090 \cdot 10^0\end{aligned}$$

$$D_z = \tilde{y} - g(\bar{x}) = 3.0 \cdot 10^{-5}$$

**Celotna napaka:** Končna napaka je  $D = D_n + D_m + D_z$ . Velja

$$|D| \leq |D_n| + |D_m| + |D_z|.$$

Zgled: Celotna napaka je  $D = D_n + D_m + D_z = 1.6 \cdot 10^{-5}$

## 1.6.1 Analiza zaokrožitvenih napak za produkt $n$ števil

Računamo produkt  $p = x_0 x_1 \cdots x_n$  predstavljenih števil  $x_0, x_1, \dots, x_n$ .

Eksaktni algoritem je:

$$p_0 = x_0$$

$$i = 1, \dots, n$$

$$p_i = p_{i-1} x_i$$

$$p = p_n$$

Dejanski algoritem pa:

$$\hat{p}_0 = x_0$$

$$i = 1, \dots, n$$

$$\hat{p}_i = \hat{p}_{i-1} x_i (1 + \delta_i), \quad |\delta_i| \leq u$$

$$\hat{p} = \hat{p}_n$$

Dobimo

$$\hat{p} = p(1 + \gamma) = p(1 + \delta_1) \cdots (1 + \delta_n).$$

Velja

$$(1 - u)^n \leq (1 + \gamma) \leq (1 + u)^n.$$

Ocenimo

$$(1 + u)^n = 1 + \binom{n}{1} u + \binom{n}{2} u^2 + \cdots = 1 + nu + \mathcal{O}(u^2),$$

$$(1 - u)^n \geq 1 - nu \quad (\text{indukcija})$$

Če je  $nu \ll 1$  ocenimo  $|\gamma| < nu$ . To pomeni, da je relativna napaka odvisna od števila množenj in da se z vsakim množenjem poveča za  $u$ . Računanje produkta  $n$  števil je direktno in obratno stabilno.

## 1.6.2 Skalarni produkt - obratna analiza zaokrožitvenih napak

Imamo dva vektorja predstavljenih števil  $x = [x_1 \ \cdots \ x_n]^T$  in  $y = [y_1 \ \cdots \ y_n]^T$ , računamo pa

$$s = y^T x = \sum_{i=1}^n x_i y_i.$$

Eksaktni algoritem je:

$$s_0 = 0$$

$$i = 0, \dots, n$$

$$p_i = x_i y_i$$

$$s_i = s_{i-1} + p_i$$

$$s = s_n$$

Dejanski algoritem pa:

$$\hat{s}_0 = 0$$

$$i = 0, \dots, n$$

$$\hat{p}_i = x_i y_i (1 + \alpha_i), \quad |\alpha_i| \leq u$$

$$\hat{s}_i = (\hat{s}_{i-1} + \hat{p}_i)(1 + \beta_i), \quad |\beta_i| \leq u$$

$$\hat{s} = \hat{s}_n$$

Obratna analiza nam da  $\hat{s} = \sum_{i=1}^n x_i y_i (1 + \gamma_i)$ , kjer je

$$1 + \gamma_1 = (1 + \alpha_1)(1 + \beta_2) \cdots (1 + \beta_n)$$

in

$$1 + \gamma_i = (1 + \alpha_i)(1 + \beta_i) \cdots (1 + \beta_n), \quad i = 2, \dots, n.$$

Tako dobimo ocene  $|\gamma_1| \leq nu$  in  $|\gamma_i| \leq (n - i + 2)u$  za  $i = 2, \dots, n$ . To pomeni, da je  $\hat{s}$  točni skalarni produkt relativno malo zmotenih vektorjev  $x$  in  $y$ . Računanje skalarnega produkta je tako obratno stabilno.



## Skalarni produkt - direktna analiza zaokrožitvenih napak

Pri direktni analizi najprej izračunamo absolutno napako:

$$\hat{s} - s = \sum_{i=1}^n x_i y_i \gamma_i,$$

torej

$$|\hat{s} - s| \leq \sum_{i=1}^n |x_i| \cdot |y_i| \cdot |\gamma_i| \leq nu \sum_{i=1}^n |x_i| \cdot |y_i| = nu |y|^T |x|.$$

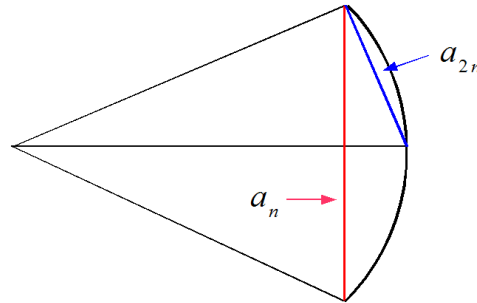
Dobimo

$$\left| \frac{\hat{s} - s}{s} \right| \leq \frac{|y|^T |x|}{|y^T x|} nu.$$

Če so vsi  $x_i y_i$  enakega predznaka, dobimo  $\left| \frac{\hat{s} - s}{s} \right| \leq nu$  in računanje je direktno stabilno, sicer pa imamo v primeru, ko sta vektorja skoraj pravokotna, lahko veliko relativno napako.

## 1.7.1 Poučni primeri - računanje števila $\pi$

$\pi$  je limita obsega  $S_n$  pravilnega mnogokotnika, včrtanega v krog s polmerom  $r = \frac{1}{2}$ . Naj bo  $a_n$  stranica pravilnega  $n$ -kotnika. Poiščimo zvezo med  $a_n$  in  $a_{2n}$ :



Velja

$$a_{2n} = \sqrt{\left(\frac{a_n}{2}\right)^2 + \left(\frac{1}{2} - \sqrt{\frac{1}{4} - \left(\frac{a_n}{2}\right)^2}\right)^2} = \sqrt{\frac{1 - \sqrt{1 - a_n^2}}{2}},$$

od tod pa iz  $S_n = na_n$  sledi

$$S_{2n} = 2na_{2n} = 2n\sqrt{\frac{1 - \sqrt{1 - \left(\frac{S_n}{n}\right)^2}}{2}}.$$

## Računanje števila $\pi$ , 2.del

Začnemo pri  $S_6 = 3$  in uporabljamo formulo  $S_{2n} = 2n \sqrt{\frac{1 - \sqrt{1 - \left(\frac{S_n}{n}\right)^2}}{2}}$ .

$n$	$S_n$	$n$	$S_n$
6	3.0000000	768	3.1430728
12	3.1058285	1536	3.1486604
24	3.1326292	3072	3.1374750
48	3.1393456	6144	3.1819806
96	3.1410186	12288	3.0000000
192	3.1414995	24576	4.2426405
384	3.1416743	49152	0.0000000

Formula odpove, saj pride do odštevanja skoraj enako velikih števil, napaka pa se množi z  $2n$ . V zgornji tabeli (enojna natančnost) so napačne decimalke rdeče.

Kadar imamo nestabilen postopek, nam ne pomaga niti računanje z večjo natančnostjo. Prava rešitev je preurediti postopek tako, da se med računanjem ne izgublja natančnost.

## Računanje števila $\pi$ , 3.del

Za stabilno računanje je potrebno formulo preurediti. Stabilna oblika je

$$S_{2n} = 2n \sqrt{\frac{\left(1 - \sqrt{1 - \left(\frac{S_n}{n}\right)^2}\right) \left(1 + \sqrt{1 - \left(\frac{S_n}{n}\right)^2}\right)}{2 \left(1 + \sqrt{1 - \left(\frac{S_n}{n}\right)^2}\right)}} = S_n \sqrt{\frac{2}{1 + \sqrt{1 - \left(\frac{S_n}{n}\right)^2}}}.$$

Sedaj dobimo pravilne rezultate:

$n$	$S_n$	$n$	$S_n$
6	3.0000000	768	3.1415839
12	3.1058285	1536	3.1375904
24	3.1326287	3072	3.1415920
48	3.1393502	6144	3.1415925
96	3.1410320	12288	3.1415925
192	3.1414526	24576	3.1415925
384	3.1415577	49152	3.1415925

## 1.7.2 Seštevanje Taylorjeve vrste za $e^{-x}$

Vemo, da je

$$e^{-x} = \sum_{n=0}^{\infty} (-1)^n \frac{x^n}{n!}$$

in da vrsta konvergira za vsak  $x \in \mathbb{C}$ . Če pa to vrsto seštevamo numerično po vrsti, potem za  $x > 0$  ne dobimo najboljših rezultatov.

$x$	$e^{-x}$	vrsta	relativna napaka
1	$3.678795 \cdot 10^{-1}$	$3.678794 \cdot 10^{-1}$	$1.4 \cdot 10^{-7}$
2	$1.353353 \cdot 10^{-1}$	$1.353353 \cdot 10^{-1}$	$2.3 \cdot 10^{-7}$
3	$4.978707 \cdot 10^{-2}$	$4.978702 \cdot 10^{-2}$	$9.3 \cdot 10^{-7}$
4	$1.831564 \cdot 10^{-2}$	$1.831531 \cdot 10^{-2}$	$1.8 \cdot 10^{-5}$
5	$6.737947 \cdot 10^{-3}$	$6.737477 \cdot 10^{-3}$	$7.0 \cdot 10^{-5}$
6	$2.478752 \cdot 10^{-3}$	$2.477039 \cdot 10^{-3}$	$6.9 \cdot 10^{-4}$
7	$9.118820 \cdot 10^{-4}$	$9.139248 \cdot 10^{-4}$	$2.2 \cdot 10^{-3}$
8	$3.354626 \cdot 10^{-4}$	$3.485951 \cdot 10^{-4}$	$3.9 \cdot 10^{-2}$
9	$1.234098 \cdot 10^{-4}$	$1.799276 \cdot 10^{-4}$	$4.6 \cdot 10^{-1}$
10	$4.539992 \cdot 10^{-5}$	$-7.266693 \cdot 10^{-5}$	$2.6 \cdot 10^{-0}$

## Seštevanje Taylorjeve vrste za $e^{-x}$ , 2. del

Razlog je, da zaporedje členov vrste alternira, poleg tega pa po absolutni vrednosti nekaj časa naraščajo, preden začnejo padati proti 0. Ko so člani največji, se zameglijo majhne decimalke, ki ostanejo netočne do konca računanja.

$n$	$a_n$	$s_n$	$n$	$a_n$	$s_n$
0	1.000000	1.000000	20	41.103188	13.396751
1	-10.000000	-9.000000	21	-19.572947	-6.176195
2	50.000000	41.000000	22	8.896794	2.720599
3	-166.666672	-125.666672	23	-3.868171	-1.147572
4	416.666687	291.000000	24	1.611738	0.464166
5	-833.333374	-542.333374	25	-0.644695	-0.180529
6	1388.888916	846.555542	26	0.247960	0.067430
7	-1984.127075	-1137.571533	27	-0.091837	-0.024407
8	2480.158936	1342.587402	28	0.032799	0.008392
9	-2755.732178	-1413.144775	29	-0.011310	-0.002918
10	2755.732178	1342.587402	30	0.000380	0.000852
11	-2505.211182	-1162.623779	31	-0.001216	-0.000364
12	2087.676025	925.052246	32	0.000380	0.000016
13	-1605.904663	-680.852417	33	-0.000115	-0.000099
14	1147.074707	466.222290	34	0.000034	-0.000065
15	-764.716492	-298.494202	35	-0.000010	-0.000075
16	477.947815	179.453613	36	0.000003	-0.000072
17	-281.145782	-101.692169	37	-0.000001	-0.000073
18	156.192108	54.499939	38	0.000000	-0.000073
19	-82.206375	-27.706436	39	-0.000000	-0.000073

### 1.7.3 Računanje $I_{10}$

Integrale  $I_n = \int_0^1 x^n e^{x-1} dx$ ,  $n = 0, 1, \dots$ , lahko računamo rekurzivno preko formule

$$I_n = x^n e^{x-1} \Big|_0^1 - n \int_0^1 x^{n-1} e^{x-1} dx = 1 - nI_{n-1},$$

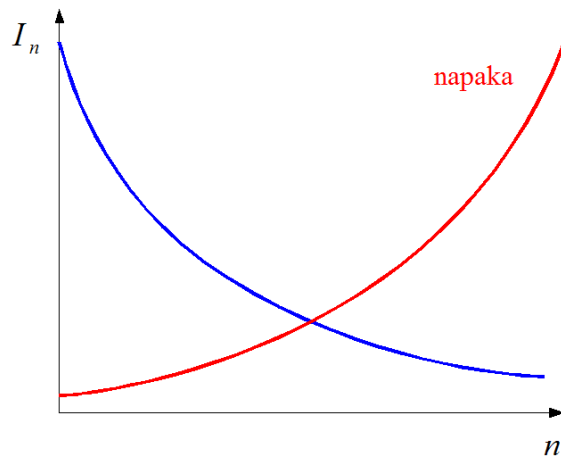
saj poznamo začetno vrednost  $I_0 = 1 - e^{-1}$ . V enojni natančnosti dobimo

$n$	$I_n$	$n$	$I_n$
0	0.6321205	7	0.1124296
1	0.3678795	8	0.1005630
2	0.2642411	9	0.0949326
3	0.2072767	10	0.0506744
4	0.1708932	11	0.4425812
5	0.1455340	12	-4.3109741
6	0.1267958	13	57.0426636

Razlog je v formuli  $I_n = 1 - nI_{n-1}$ . Napaka pri členu  $I_{n-1}$  se pomnoži z  $n$  in torej po absolutni vrednosti hitro narašča, točne vrednosti  $I_n$  pa padajo.

## Računanje $I_{10}$ , 2. del

Če računamo v obratni smeri:  $I_{n-1} = \frac{1-I_n}{n}$ , se napaka v vsakem koraku deli z  $n$ . Če začnemo pri nekem dovolj velikem členu, lahko z začetnim  $I_n = 0$  izračunamo vse začetne člene dovolj natančno. Če začnemo z  $I_{26} = 0$  tako dobimo (v enojni natančnosti) vse člene od  $I_{12}$  do  $I_0$  na vse decimalke točno.



$$I_n = 1 - nI_{n-1}$$

$n$	$I_n$	$n$	$I_n$
0	0.6321205	8	0.1009320
1	0.3678795	9	0.0916123
2	0.2642411	10	0.0838771
3	0.2072766	11	0.0773522
4	0.1708934	12	0.0717733
5	0.1455329	13	0.0669477
6	0.1268024	⋮	⋮
7	0.1123835	26	0.0000000



## 1.7.4 Seštevanje številske vrste

Znano je, da velja

$$\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6} = 1.644934066848 \dots$$

Kako bi to sešteli, če tega ne bi vedeli?

- Prištevamo člene, dokler se vsota ne spreminja več. V enojni natančnosti tako dobimo  $1.64472532 \dots$ , vrednost pa dosežemo pri  $k = 4096$ . V tem koraku delni vsoti  $\approx 1.6$  prištevamo  $2^{-24}$ .
- Seštevamo v obratnem vrstnem redu od malih cifer proti velikim. Izkaže se, da bi za približek  $1.64493406$  z 8 točnimi decimalkami morali sešteti  $10^9$  členov!

Seveda nobeden izmed zgornjih dveh načinov ni primeren, se pa da z ustrezno numerično metodo dobiti dovolj natančen približek z uporabo relativno malo členov.

## Zgled za nelinearno enačbo

Imamo  $150m$  dolgo tračnico, ki je na obeh koncih trdno vpeta. Zaradi velike vročine se tračnica raztegne za  $1cm$  in se dvigne v obliki krožnega loka. Na 8 decimalk točno izračunaj, kolikšna je maksimalna oddaljenost od tal.

