

LANA 2016 Demo 11-1

NA 1.6 Undersök felfortplantning av det relative felet i x vid evaluering av $\sin(x)$.

(a) Ge en gräns för det absoluta felet.

Lösning Använd Taylors formel: $\sin(x+\delta x) = \sin(x) + \cos(y)\delta x$ för något y mellan x och $x+\delta x$. Alltså

$$\frac{|\sin(x) - \sin(x+\delta x)|}{\text{absolut fel i utdata}} = |\cos(y)\delta x| = \underbrace{|\cos(y)|}_{\leq 1} \underbrace{|\delta x|}_{\text{absolut fel i indata}} \leq \underbrace{|\delta x|}_{\approx 1}.$$

(b) Ge en gräns för det relative felet.

Lösning Enligt ovan:

$$\frac{|\sin(x) - \sin(x+\delta x)|}{|\sin(x)|} = \frac{|\cos(y)|}{|\sin(x)|} |\delta x| = \frac{|\cos(y)||x|}{|\sin(x)|} \frac{|\delta x|}{|x|} \approx \frac{|\cos(x)||x|}{|\sin(x)|} \frac{|\delta x|}{|x|}$$

relativt fel i indata

(c) Ge en gräns för konditionstal.

Lösning Enligt definitionen av konditiontal K (relativt fel ut / relativt fel in) får vi $K = \frac{|\cos(y)||x|}{|\sin(x)|} \approx \frac{|\cos(x)||x|}{|\sin(x)|}$.

(d) För vilka x är evalueringen illakonditionerad, dvs. $K \gg 1$?

Lösning K är stort om $|x|/|\sin(x)|$ är stort, alltså $|\sin(x)| \approx 0$ om $|x| \gg 0$, alltså $x = n\pi$, $n = \pm 1, \pm 2, \pm 3, \dots$

NA 1.7 (a) Bestäm framåt- och bakåtfel för approximationen $\sin(x) \approx x$ med $x = 0.1, 0.5$ resp. 1.0 .

Lösning Framåtfel := |"exakt utdata" - "approximation"| = $|\sin(x) - x|$.

Bakåtfel := $|x - \hat{x}|$ där \hat{x} är sådan indata som exakt reproducerar approximationen, alltså här

$$\sin(\hat{x}) = x \Rightarrow \hat{x} = \arcsin(x).$$

Vi får

x	0.1	0.5	1.0
$ \sin(x) - x $	0.000167	0.0206	0.159
$ x - \arcsin(x) $	0.000167	0.0236	0.571

(b) Gamma som (a), men approximation $\sin(x) \approx x - \frac{x^3}{3!}$.

Lösning Här får vi istället $\hat{x} = \arcsin(x - \frac{x^3}{3!})$.

x	0.1	0.5	1.0
$ \sin(x) - x + \frac{x^3}{3!} $	0.000000833	0.000259	0.00814
$ x - \arcsin(x + \frac{x^3}{3!}) $	0.000000837	0.000255	0.0119

LANA 2016 Demo 11-2

NA1.16. Visa att det relativia felet vid flyttalsräkning av $x - y$ begränsas av $\mu + 2\mu(1+\mu) \frac{\max\{|x|, |y|\}}{|x-y|}$. Relatera felet till kancellation.

Lösning Avrundningshelen uppfyller $f(x * y) = (x * y)(1 + \delta_*)$, $f_l(z) = z(1 + \delta_z)$ mellan $|\delta_*|, |\delta_z| \leq \mu$, där $f_l(\cdot)$ representerar flyttalsapproximation, x, y, z är tal, $*$ är någon räkneoperation.

Vi får alltså

$$\begin{aligned} \text{relativt fel} &= \frac{|(x-y) - f_l(f_l(x) - f_l(y))|}{|x-y|} = \frac{|(x-y) - (f_l(x) - f_l(y))(1 + \delta_{-1})|}{|x-y|} \\ &= \frac{|(x-y) - (\alpha(1 + \delta_x) - \gamma(1 + \delta_y))(1 + \delta_{-1})|}{|x-y|} \\ &= \frac{|(x-y)\delta_{-1} + \alpha\delta_x - \gamma\delta_y(1 + \delta_{-1})|}{|x-y|} \leq \text{högelolikheten} \\ &\leq |\delta_{-1}| + \frac{|\alpha||\delta_x|(1 + \delta_{-1})}{|x-y|} + \frac{|\gamma||\delta_y|(1 + \delta_{-1})}{|x-y|} \\ &\leq \mu + \mu(1+\mu) \frac{|x|+|y|}{|x-y|} \leq \mu + 2\mu(1+\mu) \frac{\max\{|x|, |y|\}}{|x-y|}. \end{aligned}$$

Kancellation uppstår genom skadlig interaktion mellan avrundning av tal och approximation av räkneapproximationer. För att avgöra vilken del av felet som uppstår genom kancellation jämför i uttrycket ovan med det som erhålls då x och y är exakt lagrade flyttal, dvs. $f_l(x) = x$, $f_l(y) = y$. Då får vi

$$\begin{aligned} \text{relativt fel} &= \frac{|(x-y) - f_l(x-y)|}{|x-y|} = \frac{|(x-y) - (\alpha-y)(1 + \delta_{-1})|}{|x-y|} \\ &= \frac{|(x-y)\delta_{-1}|}{|x-y|} = |\delta_{-1}| \leq \mu. \end{aligned}$$

Vi sluter oss till att den andra termen i den första uppsättningen (som blir stor då $|x-y| \approx 0$ men $|x|, |y| \gg 0$) kommer från kancellation.

LANA 2016 Demo 11 - 3.

NA 1.18 Beskriv hur $\|x\|_2 := \sqrt{\sum_{i=1}^n x_i^2}$ kan beräknas utan risk för "overflow" eller "underflow".

Lösning Risken med "overflow" (för stora tal sätts till ∞ i flyttalssystem) är att ett x_i kvadreras och avrundas till ∞ , även om $\sqrt{x_i^2}$ kunde lagras i flyttalssystemet.

Risken med "underflow" (för små tal avrundas till 0 i flyttalssystem) är att ett x_i kvadreras och avrundas till 0, även om $\sqrt{x_i^2}$ inte är försäkrat och kunde lagras i flyttalssystemet.

Båda problemen kan avrundas genom att normalera beräkningen: Sätt $m := \max_i |x_i|$ och beräkna $\|x\|_2 = m \sqrt{\sum_{i=1}^n (\frac{x_i}{m})^2}$. Eftersom $|x_i/m| \leq 1$ kan undviks "overflow", och termen $(\frac{x_i}{m})^2$ försummas bara om de faktiskt är försäkbara i förhållande till det största elementet x_i med $|x_i|=m$.

NA 1.22 Låt $x, y \geq 0$ vara intilliggande flyttal i IEEE-DP (dubbel precision).

(a) Vilket är det minsta möjliga avståndet mellan x och y ?

Lösning IEEE-DP är flyttalssystemet $(\beta, t, l, u) = (2, 53, -1022, 1023)$, alltså tal på formen $(a_0 \cdot 2^0 + a_1 \cdot 2^{-1} + \dots + a_{52} \cdot 2^{-52}) \cdot 2^e$ med $-1022 \leq e \leq 1023$ och $a_i \in \{0, 1\}$, $i = 0, \dots, 52$, med $a_0 = 1$ om $e > -1022$. Om $e = -1022$ tillåts $a_0 = 0$ så att mindre tal kan representeras (gradvis underspill, "gradual underflow").

Talen i ett flyttalssystem är som näst längst från 0, så minsta avståndet är det mellan det minsta representerbara talet och 0 (eller näst minsta utan gradvis underspill), här $(0 \cdot 2^0 + 0 \cdot 2^{-1} + \dots + 0 \cdot 2^{-52} + 1 \cdot 2^{-52}) \cdot 2^{-1022} = -(0 \cdot 2^0 + \dots + 0 \cdot 2^{-52}) \cdot 2^{-1022} = 2^{-1022-52} = 2^{-1074}$.

(b) Vilket är det största möjliga avståndet mellan x och y ?

Lösning Talen blir glesare ju större de är, så största avståndet är det mellan det största och det näst största talet. Alltså $(1 \cdot 2^0 + \dots + 1 \cdot 2^{-52}) \cdot 2^{1023} - (1 \cdot 2^0 + \dots + 1 \cdot 2^{-52} + 0 \cdot 2^{-52}) = 2^{1023-52} = 2^{971}$.

LANA 2016 Demo 11-4

NA 1.25 Betrakta flyttalsystemet $(\beta, t, u, \lambda) = (10, 3, -98, 98)$.

(a) Bestäm UFL i detta system.

Lösning UFL (= "under flow limit", underrillsgränsen) = "minsta tal som kan representeras" = $(1 \cdot 10^0 + 0 \cdot 10^{-1} + 0 \cdot 10^{-2}) \cdot 10^{-98} = 10^{-98}$.

Notera att vi inte antar gradvis underrill, dvs. första niffran är alltid nollbild.

(b) Beräkna $x - y$ i detta system då $x = 6.87 \cdot 10^{-97}$ och $y = 6.81 \cdot 10^{-97}$.

Lösning $x - y = (6.87 - 6.81) \cdot 10^{-97} = 0.06 \cdot 10^{-97} < \text{UFL}$
 $\Rightarrow \text{fl}(x - y) = 0$.

(c) Vad ger beräkningen i (b) om gradvis underrill tillåts?

Lösning $x - y = 0.06 \cdot 10^{-97} = 0.60 \cdot 10^{-98}$ kan nu representeras i flyttalsystemet, så $\text{fl}(x - y) = 0.60 \cdot 10^{-98}$.