

1.5 Ersatzarithmetik

Durch Runden oder Abschneiden gelingt es, reelle Zahlen x mit $x_{\min} \leq |x| \leq x_{\max}$ in ein gegebenes Gleitpunkt Zahlensystem $F(b, t, l_{\min}, l_{\max})$ abzubilden. Deshalb werden die Grundoperationen $o \in \{+, -, *, /\}$ oft durch

$$x \tilde{o} y = rd(x \circ y) \quad \text{für } x, y \in F, x_{\min} \leq |x \circ y| \leq x_{\max} \quad (1.15)$$

oder

$$x \tilde{o} y = tc(x \circ y) \quad \text{für } x, y \in F, x_{\min} \leq |x \circ y| \leq x_{\max} \quad (1.16)$$

auf Rechnern realisiert (beide Division soll $y \neq 0$ sein).

Theorem 1.7

Benähtigt der durch (1.15) bzw. (1.16) definierten Ersatzoperationen $\tilde{+}, \tilde{-}, \tilde{*}, \tilde{/}$ in F abgeschlossen, d.h. im Ergebnis dieser Operation erhält man Elemente aus F .

Außerdem gilt die Beziehung bzw. Darstellung

$$x \tilde{o} y = (x \circ y) (1 + \epsilon) \quad \text{mit } |\epsilon| \leq \kappa \epsilon_{ps}, \quad (1.17)$$

wobei κ im Fall von (1.15) gleich 1 und im Fall von (1.16) gleich 2 ist (ϵ heißt Darstellungsefehler)

Beweis

Die Abgeschlossenheit von F bezüglich \tilde{o} folgt aus Th. 1.3

Die Darstellung (1.16) ergibt sich im Falle von (1.15) aus

$$\frac{|rd(x \circ y) - (x \circ y)|}{|x \circ y|} \leq \epsilon_{ps}$$

also aus (1.14)

1.6 Fehlerakkumulation

Wir betrachten Zahlen $x, y \in \mathbb{R}$. Durch eine evtl. Rundung erhalten wir mit

$$\begin{aligned} rd(x) &= x + \Delta x \in F & \text{mit } \frac{|\Delta x|}{|x|} &\leq \varepsilon \\ rd(y) &= y + \Delta y \in F & \text{mit } \frac{|\Delta y|}{|y|} &\leq \varepsilon \end{aligned}$$

Zahlen aus einem Gleitpunkt-Zahlensystem F
 $\tilde{\circ}$, sei nun Multiplikation oder Division.
Mit (1.15) und (1.17) erhält man

$$\begin{aligned} (x + \Delta x) \tilde{\circ} (y + \Delta y) &= (x(1 + \tau_x)) \tilde{\circ} (y(1 + \tau_y)), \quad |\tau_x|, |\tau_y| \leq \varepsilon \\ &= (x \circ y) \circ ((1 + \tau_x) \circ (1 + \tau_y)) \circ (1 + \kappa), \quad |\kappa| \leq \varepsilon \\ &= (x \circ y) (1 + \beta), \end{aligned}$$

wobei man bemerkt, dass

$$(1 + \tau_x) \circ (1 + \tau_y) (1 + \kappa) = 1 + \beta$$

mit einem β mit der Eigenschaft $|\beta| \leq \frac{3\varepsilon}{1-3\varepsilon}$
gilt (Beweis: Plato). Damit ergibt sich
für die Multiplikation / Division das

Theorem 1.8

Zu einem Gleitpunkt-Zahlensystem $F(b, t, e_{\min}, e_{\max})$
seien die Zahlen $x, y \in \mathbb{R}$ und $\Delta x, \Delta y \in \mathbb{R}$ gegeben mit

$x + \Delta x \in F, y + \Delta y \in F, \frac{|\Delta x|}{|x|} \leq \varepsilon, \frac{|\Delta y|}{|y|} \leq \varepsilon$ mit $\varepsilon \leq \frac{1}{4}$
0 steht für die Grundoperation \times bzw. $/$ und für
 $x \circ y$ soll $x_{\min} \leq |x \circ y| \leq x_{\max}$ gelten.

Dann gilt die Fehlerdarstellung

$$(x + \Delta x) \tilde{\circ} (y + \Delta y) = x \circ y + \eta \quad \text{mit } \frac{|\eta|}{|x \circ y|} \leq \frac{3\varepsilon}{1-4\varepsilon} \quad (1.18)$$

Die Darstellung (1.78) zeigt, dass die Multiplikation bzw. Division verhältnismäßig genau mit einem kleinen rel. Fehler ist.

Im Folgenden soll die Fehlerverstärkung bei der Hintereinanderausführung von Additionen in einem gegebenen Gleitpunkt-Zahlensystem F betrachtet werden. $\tilde{\Sigma}$ und Σ^{\sim} sollen für die Addition in F und die Summation von links nach rechts in F stehen.

Theorem 1.9

Zu $F(b, t, l_{\min}, l_{\max})$ seien $x_1, \dots, x_n \in \mathbb{R}$ und $\Delta x_1, \dots, \Delta x_n \in \mathbb{R}$ Zahlen mit

$$x_k + \Delta x_k \in F, \quad \frac{|\Delta x_k|}{|x_k|} \leq \varepsilon \quad \text{für } k=1, \dots, n$$

und es bezeichne

$$\tilde{S}_k := \sum_{j=1}^k (x_j + \Delta x_j), \quad S_k := \sum_{j=1}^k x_j, \quad k=1, \dots, n$$

die entsprechende Partialsumme (Summation von links nach rechts). Dann gilt

$$|\tilde{S}_k - S_k| \leq \underbrace{\left(\sum_{j=1}^k (1+\varepsilon)^{k-j} (2|x_j| + |S_j|) \right)}_{=: M_k} \varepsilon \quad \text{für } k=1, \dots, n \quad (1.19)$$

falls die Partialsummen (Notation $M_0 = 0$) innerhalb gewisser Schwärze liegen:

$$x_{\min} + (M_{k-1} + |x_k|) \varepsilon \leq |S_k| \leq x_{\max} - (M_{k-1} + |x_k|) \varepsilon \quad k=1, \dots, n$$

Beweis:

Vollr. Induktion über k

(1.19) gilt für $k=1$, denn $|\tilde{S}_1 - S_1| = |\Delta x_1| \leq 3|x_1|\varepsilon$

Wir nehmen jetzt an, dass (1.19) für ein $k \geq 1$ richtig ist

Mit der Notation

$$\Delta S_j := \tilde{S}_j - S_j \quad \text{für } j \geq 1, \quad \Delta S_0 = 0$$

Berechnet man mit einer gewissen Zahl $r_k \in \mathbb{R}, |r_k| \leq \epsilon$

$$\begin{aligned} \Delta S_k &= \tilde{S}_k - S_k = \tilde{S}_{k-1} \tilde{f}(X_k + \Delta X_k) - S_k \\ &= (S_{k-1} + \Delta S_{k-1}) \tilde{f}(X_k + \Delta X_k) - S_k \\ &= (S_k + \Delta S_{k-1} + \Delta X_k)(1 + r_k) - S_k \\ &= (1 + r_k) \Delta S_{k-1} + r_k S_k + (1 + r_k) \Delta X_k \end{aligned}$$

und damit

$$|\Delta S_k| \leq (1 + \epsilon) |\Delta S_{k-1}| + \epsilon (|S_k| + 2|X_k|) \quad (1.20)$$

Aus (1.20) und der Induktionsannahme folgt die Behauptung (1.19) des Theorems

Bemerkung 1.10

Der Faktor $(1 + \epsilon)^{n-j}$ in der Abschätzung (1.19) ist umso größer, je kleiner j ist. Daher ist es vorteilhaft beim Aufsummieren mit der betragsmäßig kleinen Zahlen zu beginnen. Dies gewährleistet zudem, dass die Partialsumme S_k betragsmäßig nicht unnötig anwachsen

Theorem 1.9 liefert mit (1.19) nur eine Abschätzung für den absoluten Fehler. Der relative Fehler

$\frac{|\tilde{S}_n - S_n|}{|S_n|}$ kann jedoch groß ausfallen, falls $|S_n|$ klein gegenüber $\sum_{j=1}^{n-1} (|X_j| + |S_j|) + |X_n|$ ist!