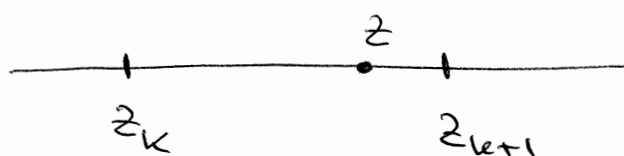


Betrachten wir Gleitpunktzahlen (normalisiert)
der Form

$$z_{\#} = m \cdot b^e \quad \text{mit } m > 0 \text{ und } m \geq 0,1$$

Dann gilt $\left| \frac{\lfloor z \rfloor - z}{z} \right| \leq \text{eps} = \begin{cases} b^{1-p} \text{ f. Abschn.} \\ \frac{1}{2} b^{1-p} \text{ f. Runden} \end{cases}$

für alle $z \in \mathbb{R} \cap \mathbb{F}_N$ (d.h. für alle reelle Zahlen z , die im Bereich der normal. Maschinenzahlen liegen.)



Seien z_k und z_{k+1} Maschinenzahlen mit

$$z_k = m_k b^e \quad \text{und} \quad z_{k+1} = m_{k+1} b^e, \quad \text{dann}$$

gilt $|\lfloor z \rfloor - z| \leq |z_{k+1} - z_k| = b^{-p} \cdot b^e$

$$\Rightarrow \left| \frac{\lfloor z \rfloor - z}{z} \right| \leq \frac{b^{-p} \cdot b^e}{m_k \cdot b^e} \leq \frac{b^{-p}}{0,1} = \frac{b^{-p}}{b^{-2}} =$$

$$= \underline{b^{1-p} \text{ f\"ur das Abschn.}}$$

$$\text{und } \left| \frac{\lfloor z \rceil - z}{z} \right| \leq \underline{\frac{1}{2} b^{1-p} \text{ f\"ur das Runden}}$$

2

Der wesentliche Unterschied zwischen den Fest- und Gleitkommazahlen besteht darin, daß für die Festkommazahlen der absolute Abstand zwischen ihnen konstant ist, während für die Gleitkommazahlen der relative Abstand konstant ist (im normal. Bereich.)

Schätzung des relativen Gesamtrechnungsfehlers mit Hilfe der "(1+ε)"-Technik.

$$\tilde{y} = \tilde{\ln}(1 \oplus x) \quad , \quad x \text{ Maschinenzahl}$$

$$1 \oplus x = (1+x)(1+\xi_1) \quad , \quad |\xi_1| \leq \epsilon_{ps}$$

$$\tilde{\ln}(1 \oplus x) = \ln(1 \oplus x)(1+\xi_2) \quad , \quad |\xi_2| \leq \epsilon_{ps}$$

Wir streben die Form an:

$$\tilde{\ln}(1 \oplus x) = \ln(1+x)(1+\delta) \quad ;$$

wobei δ dann der Gesamtrechnungsfehler ist.

Rechnung:

$$\tilde{\ln}(1 \oplus x) = \underbrace{\ln((1+x)(1+\xi_1))}_A (1+\xi_2).$$

Wir sehen uns zunächst ^A den Term A an!

$$A = \ln((1+x)(1+s_1)) = \ln((1+x) + s_1(1+x))$$

Wir haben also $\ln(a+b)$ mit b viel kleiner als $a \Rightarrow$ Linearisierung:

$$\begin{aligned} \ln(a+b) &= \ln(a) + (\ln(a+b))' \big|_{b=0} \cdot b + O(b^2) \\ &\approx \ln(a) + \frac{1}{a} \cdot b \end{aligned}$$

$$\text{d.h. } a = 1+x, b = s_1(1+x) \Rightarrow$$

$$\begin{aligned} \ln((1+x) + s_1(1+x)) &\approx \ln(1+x) + s_1(1+x) \frac{1}{(1+x)} \\ &= \ln(1+x) + s_1; \end{aligned}$$

Damit haben wir

$$\begin{aligned} \tilde{y} &= (\ln(1+x) + s_1)(1+s_2) = \\ &= \ln(1+x) + s_2 \ln(1+x) + s_1 + \underbrace{s_1 s_2}_{\approx 0} \\ &\doteq \ln(1+x) + s_1 + s_2 \ln(1+x) = \\ &= \ln(1+x) \left(1 + \underbrace{\frac{s_1}{\ln(1+x)}}_{\text{sehr groß für } x \approx 0} + s_2 \right) \end{aligned}$$

Der Gesamtfehler

$$\delta = \frac{\delta_1}{\ln(1+x)} + \delta_2 = \underbrace{\frac{1}{\ln(1+x)}}_{\text{großer Verstärkungsfaktor für den Additionsfehler}} \cdot \delta_1 + 1 \cdot \delta_2.$$

Dieser Fehler wird verglichen mit dem rel. Datenfehler

$$\tilde{y} = y(\tilde{x}) = y(x(1+\delta)), \text{ dabei ist } \delta \text{ der relative Datenfehler.}$$

$$\Rightarrow \text{Kondition } \tilde{y} = y(1 + \tilde{\delta}) \text{ mit}$$

$$\tilde{\delta} = K_{y \leftarrow x} \cdot \delta \approx 1 \cdot \delta,$$

vgl. Rechnungen auf meinen Seiten zur Vo. am 24.10.2005.

Damit sind die Verstärkungsfaktoren bei den RF viel größer als bei der Kondition und der Algorithmus instabil.

Auslöschung! Wegen $\ln(1+x) \approx x$ für kleine x bedeutet die Berechnung, daß wir die Operationen $(1 \oplus x) \ominus 1$ ausführen!