

L^AT_EX- Grundlagen

von

Dominik Wille

22 Oktober 2013

Freie Universität Berlin
Zentraleinrichtung für Datenverarbeitung
Betriebsysteme und Programmieren

Dozent:
Dr. Herbert Voß

1 ...

1.1 ...

1.2 ...

1.3 Rundungsfehler

Rechenoperationen mit reellen Zahlen im Computer → Rundungsfehler.

1.3.1 Gleitkommaarithmetik

Im Vergleich zum Festpunktformat: geringerer Speicherplatzbedarf.

n -stellige Gleitkommazahl, Basis B :

$$x = \pm (0, z_1, z_2, \dots, z_n)_B \cdot B^E = \pm \sum_{i=1}^n z_i \cdot B^{-i} = 0 \quad (1)$$

(Normalisierte Gleitkommadarstellung) Exponent: $E \in \mathbb{Z} : m \leq E \leq M$

Bsp: $+1234,567 = +(0,1234567)_{10} \cdot 10^4$

($B = 10, n = 7$) Die Werte n, B, m, M maschinenabhängig (Hardware und Compiler) Übliche Basen:

- $B = 2$ (Dualzahlen, im Computer)
- $B = 8$ (Oktalzahlen)
- $B = 10$ (Dezimal)
- $B = 16$ (Hexadezimal)

Bsp: binäre Darstellung:

$$(5,0625)_{10} = 0,50625 \cdot 10^1 \quad (2)$$

$$= 1 \cdot 2^2 + 1 \cdot 2^0 + 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} \quad (3)$$

$$= (101,0001)_2 = (0,1010001)_2 \cdot 2^3 \quad (4)$$

manche Zahlen lassen sich nur schwer als Dualzahlen darstellen:

- $(3)_{10} = (11)_2$ geht
- $(0,3)_{10} = 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + \dots = (0,010011001\dots)_2$ geht nicht

Genauigkeit der Darstellung

23 Stellen $111111111111111111111111 = 2^{23} - 1 = 8.388.608$

\Rightarrow 6 Ziffern können unterschieden werden.

52 Stellen $2^{52} = 4.503.599.627.370.496$

\Rightarrow 15 Stellen können unterschieden werden. Die größte darstellbare Zahl entspricht der größten Maschienenzahl.

$$x_{max} = (0, [B-1][B-1] \dots [B-1])_B \cdot B^M = (1 - B^{-n}) \cdot B^M \quad (5)$$

kleinste darstellbare Zahl

$$x_{min} = (0, 1000000)_B \cdot B^m = (1 - B^{-n}) \cdot B^{m-1} \quad (6)$$

\Rightarrow Die Menge der Maschienenzahlen ist endlich

Bsp:

$$x_{max} + x_{max} = \infty$$

$$x_{min} \cdot B^{-1} = 0$$

1.3.2 Rundungsfehler

Beim Runden einer Zahl x wird eine Näherung $rd(x)$ unter den Maschienenzahlen geliefert, so dass der absolute Fehler $|x - rd(x)|$ minimal ist, der unvermeidbare Fehler ist der Rundungsfehler. Eine n -stellige Dezimalzahl im Gleitkommaformat

$$x = \pm(0, z_1, \dots, z_n)_{10} = rd(x) \quad (7)$$

hat einen maximalen absoluten Fehler :

$$|x - rd(x)| \leq 0,000\ldots005 \cdot 10^E \quad (8)$$

$$= 0,5 \cdot E^{E-n} \quad (9)$$

, für allgemeine Basis B :

$$|x - rd(x)| \leq \frac{B \cdot 1}{2 \cdot B} B^{E-n} = \frac{1}{2} B^{E-n} \quad (10)$$

Rundungsfehler werden durch die Rechnung getragen!

n -stellige Gleitkommaarithmetik:

jede einzelne Rechenoperation $(+, -, \times, \div)$ wird auf $n + 1$ Stellen genau berechnet und dann auf n Stellen gerundet. Jedes Zwischenergebnis, nicht Endergebnis!

Bsp:

rechne $2590 + 4 + 4$ in 3-stelliger dez. G.P.A.

links 1. $2590 + 4 \rightarrow 2590$

2. $2590 + 4 \rightarrow 2590$

rechts 1. $4 + 4 \rightarrow 10$

2. $2590 + 10 \rightarrow 2600$

\Rightarrow Rechenwege unterscheiden sich!

Regel: beim Addieren Summanden in der Reihenfolge aufsteigender Beträge addieren.

Maß für der Rechenzeit eines Computers: “flops” floating point operations per second
(typisch Multiplikation oder Division)

(top500.org) 1 Tiake-2 3 Mio Cores, 54.000 T Flops, 17 MW

relative Fehler wichtiger als absoluter Fehler:

Näherung \tilde{x} zu exaktem wert x , rel. fehler $E = \left| \frac{x-\tilde{x}}{\tilde{x}} \right| \approx \frac{x-\tilde{x}}{\tilde{x}}$ für duale rechnungen am Computer $B=2 \rightarrow E_{max} = 2^{-n}$

E_{max} wird auch maschinenzahlgenauigkeit genannt, und gibt die kleinste potentielle Zahl an, für die gilt $|E_{max}|$; E_{max} kann aus Rechenergebnissen errechnet werden (ÜB1)