

Numerisches Programmieren, Übungen

Musterlösung 1. Übungsblatt: Zahlendarstellung, Rundungsfehler

1) Umrechnung von Zahlen

a) Zahldarstellung in allgemeinem System: $\sum_{i=0}^N r_i \text{Basis}^i$ mit $r_i \in \{0, 1, \dots, \text{Basis} - 1\}$

	dezimal	binär	trinär	hexadezimal
Basis	10	2	3	16
darzust. Zahl	19	10011	201	13
darzust. Zahl	47	101111	1202	2f
darzust. Zahl	511	111111111	200221	1ff

b) Schriftliches Dividieren der Brüche im binären System analog zum dezimalen:

$$-\frac{1}{7} = -1_2 : 111_2 \text{ (binär) :}$$

$$\begin{array}{r} - \quad 1 \quad 0 \quad 0 \quad 0 \quad : \quad 1 \quad 1 \quad 1 = -0.00\overline{100} \quad \text{bzw.} \quad -0.\overline{001} \\ -) \quad \quad 1 \quad 1 \quad 1 \\ \hline \quad \quad \quad 1 \quad 0 \quad 0 \quad 0 \end{array}$$

$$\frac{1}{10} = 1_2 : 1010_2 \text{ (binär) :}$$

$$\begin{array}{r} \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad : \quad 1 \quad 0 \quad 1 \quad 0 = 0.000\overline{1100} \quad \text{bzw.} \quad 0.0\overline{0011} \\ -) \quad \quad 1 \quad 0 \quad 1 \quad 0 \\ \hline \quad \quad \quad 1 \quad 1 \quad 0 \quad 0 \\ -) \quad \quad \quad 1 \quad 0 \quad 1 \quad 0 \\ \hline \quad \quad \quad \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \end{array}$$

Alternative:

$$\text{i) } \begin{array}{ccccccc} & \cdot 2 & \cdot 2 & \cdot 2 & -1 \mid \cdot 2 & & \\ \frac{1}{7} & \frac{2}{7} & \frac{4}{7} & \frac{8}{7} & \frac{2}{7} & & \\ \hline 0, & 0 & 0 & 1 & & & \end{array}$$

$$\text{ii) } \begin{array}{cccccccc} & \cdot 2 & \cdot 2 & \cdot 2 & \cdot 2 & \frac{-1}{2} & \frac{-1}{2} & \\ \frac{1}{10} & \frac{2}{10} & \frac{4}{10} & \frac{8}{10} & \frac{16}{10} & \frac{12}{10} & \frac{4}{10} & \\ \hline 0, & 0 & 0 & 0 & 1 & 1 & & \end{array}$$

2) Binärdarstellung von ganzen Zahlen

Beachte: 1 Byte = 8 Bits!

Wert 1. Bit	-2^{n-1}
Wert 2. Bit	2^{n-2}
kleinste Zahl x_{min}	-2^{n-1}
größte Zahl x_{max}	$2^{n-1} - 1$

$$4 \text{ Bytes} = 32 \text{ Bits} \rightarrow [-2^{31}, 2^{31} - 1] = [2.147.483.648, 2.147.483.647]$$

3) Gleitkomma-Zahlen

a) Schriftliches Dividieren (wie in Aufg. 1): $-\frac{11}{10} = -1011_2 : 1010_2$ (binär) :

$$\begin{array}{r}
 1011 : 1010 = -1.000\overline{1100} \quad \text{bzw.} \quad -1.000\overline{11} \\
 -) 1010 \\
 \hline
 1000 \\
 -) 1010 \\
 \hline
 1100 \\
 -) 1010 \\
 \hline
 10000 \\
 \hline
 100000
 \end{array}$$

Damit erhält man insgesamt: $-1.000\overline{11} \cdot 2^0$

Alternative:

$\frac{11}{10} \cdot 1 = \frac{11}{10} \geq 1$	1
	.
$\frac{1}{10} \cdot 2 = \frac{2}{10} < 1$	0
$\frac{2}{10} \cdot 2 = \frac{4}{10} < 1$	0
$\frac{4}{10} \cdot 2 = \frac{8}{10} < 1$	0
$\frac{8}{10} \cdot 2 = \frac{16}{10} \geq 1$	1
$\frac{6}{10} \cdot 2 = \frac{12}{10} \geq 1$	1
$\frac{2}{10} \cdot 2 = \frac{4}{10} < 1$	0
$\frac{4}{10} \cdot 2 = \frac{8}{10} < 1$	0
$\frac{8}{10} \cdot 2 = \frac{16}{10} \geq 1$	1
$\frac{6}{10} \cdot 2 = \frac{12}{10} \geq 1$	1
$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$	\vdots

b) Rundungsregel (korrektes Runden) für eine Zahl $x = (-1)^\nu \cdot 2^e \cdot 1.x_1x_2 \dots x_{t-1} | x_t x_{t+1} x_{t+2} \dots$:

Nr.	Fallbed.	Rundungsvorschrift
Einfacher Fall: Abrunden		
1)	$x_t = 0$	$rd(x) = (-1)^\nu \cdot 2^e \cdot 1.x_1x_2 \dots x_{t-1}$
Einfacher Fall: Aufrunden		
2)	$x_t = 1 \wedge$ $ x_t x_{t+1} x_{t+2} \dots \neq 1000000000 \dots$	$rd(x) = (-1)^\nu \cdot 2^e \cdot (1.x_1x_2 \dots x_{t-1} + 2^{-(t-1)})$
Runden zur näheren geraden Mantissen-Zahl:		
3)	$x_{t-1} x_t x_{t+1} \dots = 0 1000000000 \dots$	$rd(x) = (-1)^\nu \cdot 2^e \cdot 1.x_1x_2 \dots x_{t-1}$
4)	$x_{t-1} x_t x_{t+1} \dots = 1 1000000000 \dots$	$rd(x) = (-1)^\nu \cdot 2^e \cdot (1.x_1x_2 \dots x_{t-1} + 2^{-(t-1)})$

Erläuterungen zu den verschiedenen Fällen:

Fall 1) und 2) stellen den normalen Fall dar, dass eine reelle Zahl nicht genau in der Mitte zwischen den beiden nächsten Maschinenzahlen liegt; es wird zum näheren Nachbarn gerundet. Fall 3) und 4) bewirken ein Runden mit Mantissenende $x_{t-1} = 0$ für reelle Zahlen, die genau zwischen zwei Maschinenzahlen liegen. In Fall 2) und 4) wird mit $+2^{-(t-1)}$ das letzte Bit um eins erhöht (und eventuell ein Übertrag durchgeführt).

Frage: Wie kann eine Floating Point Einheit auf der CPU unendlich viele Stellen zum Runden ausrechnen (z. B. bei periodischen Nachkommastellen)?

Das ist nicht nötig. Z. B. muss bei der Division lediglich bekannt sein, ob ein x_{t+1} ungleich 0 ist. Diese Information ist über den Rest erhältlich, der ab dem Berechnen von der Stelle x_{t+1} übrig bleibt.

Beispiel:

Es stehen zwei Mantissenbits zur Verfügung, also: $1, x_1 x_{t-1} | x_t x_{t+1} x_{t+2} \dots$

Damit sind die folgenden Zahlen exakt darstellbar:

$$\begin{aligned}
 1,00_2 &= 1_{10} \\
 1,01_2 &= 1,25_{10} \\
 1,10_2 &= 1,5_{10} \\
 1,11_2 &= 1,75_{10}
 \end{aligned}$$

Beispiele für alle vier Fälle aus der obigen Tabelle:

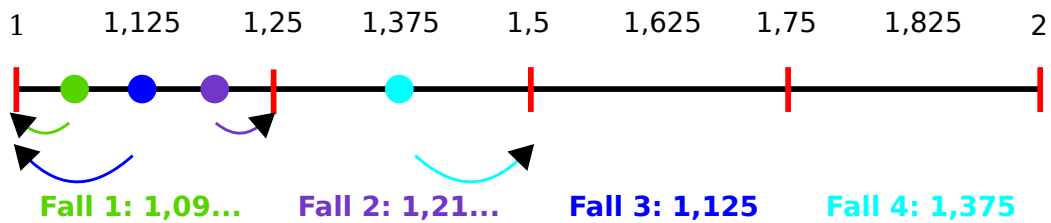


Abbildung 1: Veranschaulichung von Rundungen auf dem Zahlenstrahl

Nr.	Binärzahl	Dezimalzahl	ab-/aufrunden
1)	1,00 011 ₂	1,09375 ₁₀	abrunden
2)	1,00 111 ₂	1,21875 ₁₀	aufrunden
3)	1,00 100 ₂	1,125 ₁₀	abrunden
4)	1,01 100 ₂	1,375 ₁₀	aufrunden

Abbildung 1 veranschaulicht die vier Fälle auf dem Zahlenstrahl.

Für unsere Zahl $-\frac{11}{10} = -1.00011 \cdot 2^0$ aus Teilaufgabe i) bedeutet das:

Bit-Verwendungszweck	gesetztes Bit für $-1.00011 \cdot 2^0$
Vorzeichen (1 Bit) ('-' wird zu '1')	1
Exponent (8 Bits) ($0 + 127 = 127$)	0 1 1 1 1 1 1 1
Mantisse (23 Bits)	0 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 1

Das letzte Bit der Mantisse ergibt sich nach der Rundungsregel (s.o.).

c) Zwei Beispiele:

- Die Zahl $x = 1 + 2^{-30}$ ist nicht exakt darstellbar (zuwenig Stellen). Die Zahl wird als $rd(x) = 1$ abgespeichert (vgl. Rundungsregeln).

$$\text{Absoluter Fehler: } f_{abs} := |x - rd(x)| = 2^{-30}$$

$$\text{Relativer Fehler: } f_{rel} := |f_{abs}/x| < 2^{-30}$$

- $x = \frac{11}{10} = 1.00011 \cdot 2^0$.

$$\text{Absoluter Fehler: } f_{abs} = 1.1001 \cdot 2^{-26}$$

$$\text{Relativer Fehler: } f_{rel} = 1.01110100 \dots \cdot 2^{-26} \approx f_{abs}$$

d) Definition: Maschinengenauigkeit = Die größte positive Zahl ε_{Ma} , so dass $1 \oplus \varepsilon_{Ma} = 1$.

Die kleinste darstellbare Zahl größer als 1 ist $1 + 2^{-23}$ (23 Bits echt für Mantissen-Nachkommastellen frei, da Normierung '1.' nicht extra gespeichert werden muss!). Folglich liegt ε in einer Größenordnung von 2^{-24} .

4) Fehlerabschätzung von Zeitschrittweiten

Floating Point Darstellung von 1/60:

Wir starten mit einfacher Division:

$$\begin{array}{r} 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0 : 1\ 1\ 1\ 1\ 0\ 0 = 0.000001\overline{0001} \text{ bzw. } 0.000001 \\ -) \quad 1\ 1\ 1\ 1\ 0\ 0 \\ \hline \quad \quad \quad 1\ 0\ 0\ 0\ 0\ 0\ 0 \\ -) \quad \quad \quad 1\ 1\ 1\ 1\ 0\ 0 \\ \hline \quad \quad \quad \quad \quad 1\ 0\ 0 \end{array}$$

Für die Mantisse stehen 23 Bit zur Verfügung. Die erste '1' an der 6-ten Nachkommastelle wird allerdings nicht explizit abgespeichert. D. h. erst ab der 30-ten Nachkommastelle können die Ziffern nicht mehr abgespeichert werden:

	0,000001	0001	0001	0001	0001	0001	0001	000[1
Nr. der Nachkommastelle:	7	11	15	19	23	27		

Mit der Normalisierung und dem nicht-Abspeichern der führenden '1' kann die Zahl deshalb bis auf die 29. Stelle genau abgespeichert werden.

Wie wird die Zahl dann als floating-point Zahl abgespeichert?

Sign: 0 für positive Zahlen

Exponent: $-6+127 = 0b1111001$

Mantisse: $0b000100010001000100010001001$ mit nach oben gerundeter letzten Stelle!

0	01111001	000100010001000100010001001
Sign	Exponent	Mantissa

Floating point number: $(1) \cdot 2^{-6-23} \cdot 8947849 = 0.0166666675359010696411132\dots$

Kleines Beispielprogramm für Zweifler:

```
#include <iostream>

int main()
{
    float l = (1.0/60.0);
    std::cout << l << std::endl;
    std::cout << (void*)(unsigned long*)&l << std::endl;
}
/*
Output:
0.0166667
0x3c888889 ( $\hat{=}$  111100100010001000100010001001)
*/
```

Wie kann man das Problem umgehen?

- Positionsänderung relativ zur Startposition berechnen.
- Als Zeitschrittweite eine 2er-Potenz wählen (z.B. 1/64).
- Rechengenauigkeit erhöhen.
- Fehler ignorieren, wenn die Simulationsdauer klein genug ist.

5) Ermittlung von π nach Archimedes

a) Quadrat:

$$s = \sqrt{2}$$

$$U = 4\sqrt{2} \approx 5.6568$$

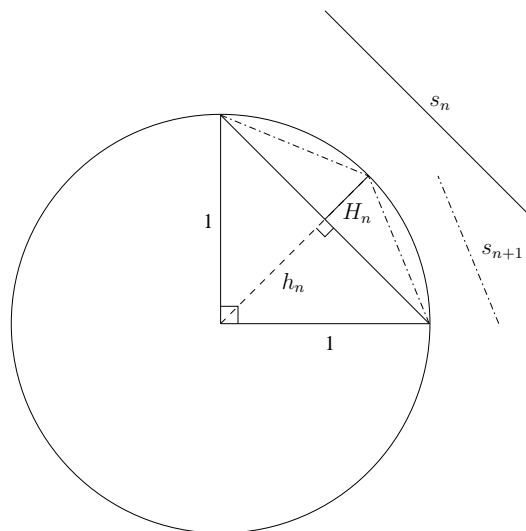


Abbildung 2: Definition der rekursiven Größen.

b) Es gilt mit den Definitionen aus Abbildung 2 (+Satz von Pythagoras!):

$$h_n^2 + \left(\frac{s_n}{2}\right)^2 = 1$$

$$H_n = 1 - h_n$$

$$H_n^2 + \left(\frac{s_n}{2}\right)^2 = s_{n+1}^2$$

und damit

$$\begin{aligned}
 s_{n+1} &= \sqrt{H_n^2 + \left(\frac{s_n}{2}\right)^2} = \sqrt{(1-h_n)^2 + \left(\frac{s_n}{2}\right)^2} \\
 &= \sqrt{\left(1 - \sqrt{1 - \left(\frac{s_n}{2}\right)^2}\right)^2 + \left(\frac{s_n}{2}\right)^2} \\
 &= \sqrt{\left(1 - 2\sqrt{1 - \left(\frac{s_n}{2}\right)^2} + 1 - \left(\frac{s_n}{2}\right)^2\right) + \left(\frac{s_n}{2}\right)^2} \\
 &= \sqrt{2 - 2\sqrt{1 - \left(\frac{s_n}{2}\right)^2}} = \sqrt{2 - \sqrt{4 - s_n^2}}.
 \end{aligned}$$

Der gesamte Umfang des 2^n -Ecks ergibt sich somit zu $U_n = 2^n \cdot s_n$ und damit also $\pi_n = U_n/2 = 2^{n-1} \cdot s_n$.

- c) Effekt wie in Tabelle auf Angabe: Vermeintlich höhere Genauigkeit durch mehr Rekursionen verbessert das Ergebnis nicht sondern zerstört den gesamten Wert!
 Problem: Auslöschung!

Algebraische Umformung zur Vermeidung der Auslöschung:

$$\begin{aligned}
 s_{n+1} &= \sqrt{2 - \sqrt{4 - s_n^2}} = \sqrt{2 - \sqrt{4 - s_n^2}} \cdot \frac{\sqrt{2 + \sqrt{4 - s_n^2}}}{\sqrt{2 + \sqrt{4 - s_n^2}}} \\
 &= \frac{|s_n|}{\sqrt{2 + \sqrt{4 - s_n^2}}}
 \end{aligned}$$

Vergleich der absoluten Fehler der beiden Formeln in Bezug auf echte Lösung π in semi-logarithmischer Skala (erstellt mit matlab-Programm archimedes.m aus www):

