

Vorlesungsmitschrift

Algorithmische Mathematik 1

Wintersemester 2021/2022 *
Universität Bonn

Fabian

30.01.2022

Online aufrufbar auf <https://git-fabus.github.io>.
Veränderungsvorschläge und Verbesserungen bitte an git-fabus@uni-bonn.de.

*Letzte Änderung: (None)

Commit: (None)

Preface

Das folgende Skript sind meine persönlichen Notizen für die Vorlesung *Algorithmische Mathematik 1*, welche von Dr. Jürgen Dölz an der Universität Bonn gehalten wurden. Diese Mitschriften wurden initial als Klausurvorbereitung verfasst. Seit dem wurde das Skript ständig überarbeitet und soll nun als Vorlesungsskript für den jeweiligen Vortragenden der Alma 1 dienen.

Inhaltsverzeichnis

1	Einführung	1
2	Zahlendarstellung am Computer	2
3	Fehleranalyse	12

1 Einführung

Algorithmische Mathematik -Was ist das?

Gegenstand der Alma ist die Konstruktion und ANalyse effizienter Algorithmen zur Lösung mathematischer Problemstellugnen mit Hilfe des Computers.

Damit liegt sie im Bereich der Angewandten Mathematik. Konkrete Problemstellungen ergeben sich oft aus technischen Problemen, Naturwissenschaften, Medizin etc.

Man kann hierbei annehmen, dass verschiedene Problemstellungen aus der Anwendungen oft zu ähnlichen oder gleichen mathematischen Modellen zurückgeführt werden können.

1.1 Beispiel. Ein typisches Problem ist das lösen eines linearen Gleichungssystems.

Im Rahmen dieser Vorlesung konzentrieren wird uns auf die Teilbereiche

- a) Numerik
 - b) Diskrete Mathematik
 - c) Statistik
- der Angewandten Mathematik

2 Zahlendarstellung am Computer

2.1 Zahlensystem

Die Darstellung von Zahlen basiert auf sogenannten *Zahlensystemen*. Diese Zahlensysteme unterscheiden sich in der Wahl des zugrundeliegenden Alphabets. Unsere Zahl entspricht dann einem Wort, bestehend aus Elementen des Alphabets.

2.1 Definition (Alphabet). Es sei $\mathbb{N} = \{1, 2, 3, \dots\}$ und $b \in \mathbb{N}$.
Wir bezeichnen mit \sum_b das Alphabet des *b-adischen Zahlensystems*

2.2 Beispiel. Verschiedene Zahlensysteme

- a) Dezimalsystem: $\sum_{10} = \{0, 1, \dots, 9\}$
wie $(384)_{10}$
- b) Dual bzw. Binärsystem $\sum_2 = \{0, 1\}$
wie $(42)_{10} = (101010)_2$
- c) ...

Um die Zahlendarstellung sinnvoll zu nutzen ergibt sich folgendes:

2.3 Satz. Seien $b, n \in \mathbb{N}, b > 1$.

Dann ist jede ganze, nicht-negative Zahl z mit $0 \leq z \leq b^n - 1$ *eindeutig* als Wort der Länge n über \sum_b darstellbar durch:

$$z = \sum_{i=0}^{n-1} z_i b^i$$

mit $z_i \in \sum_b$ für alle $i = 0, 1, \dots, n-1$.

Wir schreiben vereinfachend:

$$z = (z_{n-1}, \dots, z_0)_b$$

Beweis. • Induktion

◇ *Induktionsanfang:* $\underline{z < b}$ hat die eindeutige Darstellung $z_0 = z$ und $z_i = 0$ sonst.

◇ *Induktionsschritt:* $\underline{z - 1 \rightarrow z \geq b}$

Betrachte

$$z = \left\lfloor \frac{z}{b} \right\rfloor \cdot b + (z \bmod b)$$

Da $\hat{z} < z$ besitzt die eindeutige Darstellung

$$\hat{z} = (\hat{z}_{n-1}, \dots, \hat{z}_0)_b$$

Bemerke $\hat{z}_{n-1} = 0$, da

$$(\hat{z}_{n-1} b^{n-1}) b \leq z \leq b^n - 1$$

Also ist z_b eine n-stellige Darstellung von z in b-adischen Zahlensystem

- Eindeutigkeit wird durch Widerspruch gezeigt.

Angenommen: Es gibt zwei verschiedene Darstellungen

$$z = (z_{n-1}^{(2)}, \dots, z_0^{(2)})_b = (z_{n-1}^{(1)}, \dots, z_0^{(1)})_b$$

Sei $m \in \mathbb{N}$ der größte Index mit $z_m^{(1)} \neq z_m^{(2)}$,

Ohne Beschränkung der Allgemeinheit kann gesagt werden $z_m^{(1)} > z_m^{(2)}$.

Dann müssten die Stellen $0, 1, \dots, m-1$ den niedrigeren Wert von $z_m^{(2)}$ kompensieren.

Die größte mit diesen Stellen darstellbare Zahl ist aber:

$$\sum_{i=0}^{m-1} (b-1)b^i = (b-1) \sum_{i=0}^{m-1} b^i = b^m - 1.$$

Da b^m der kleinstmögliche Wert der fehlende Stelle m ist, kann diese aber nicht kompensiert werden. Dies ist ein Widerspruch.

□

Durch den Beweis ergibt sich sofort ein Algorithmus zur Umwandlung einer Zahl in ein anderes Zahlensystem:

2.4 Beispiel. Umwandlung von $(1364)_{10}$ in das Oktalsystem.

- $1364 = 170 \cdot 8 + 4$
- $170 = 21 \cdot 8 + 2$
- ...
- $0 \cdot 8 + 2$

$\implies (2524)_8$

Daten: Dezimalzahl $z \in \mathbb{N}_0$, Basis $b \in \mathbb{N}$
Ergebnis: b-adische Darstellung $(z_{n-1}, \dots, z_0)_b$
 Initialisiere $i = 0$
solange $z > 0$ **tue**
 $z_i = z \bmod b$
 $z = \lfloor \frac{z}{b} \rfloor$
 $i = i + 1$
Ende

Algorithmus 1: Bestimmung der b-adischen Darstellung

Beobachtung Wir sehen, dass das Honor-Schema nur eine Schleife, Additionen und Multiplikationen benötigt. Diese Operationen können wir am Computer durchführen. Im Gegensatz dazu steht die Potenz b^i in modernen Programmiersprachen zwar zur Verfügung, wird aber im Hintergrund oft auf Multiplikationen zurückgeführt. Man überprüft leicht, dass das Honor-Schema weniger Multiplikationen benötigt und somit schneller ist.

2.2 Vorzeichen-Betrag-Darstellung

Um auch Zahlen mit Vorzeichen am Computer darstellen zu können, betrachten wir im Folgenden die Vorzeichen-Betrag-Darstellung für Binärzahlen. Das Binäralphabet besteht nur aus 0 und 1, welche wir auch als *Bits* bezeichnen. Bei einer Wortlänge von n Bits wird das erste Bit als Vorzeichen verwendet, die restlichen $n - 1$ -Bits für den Betrag der Zahl. Da die 0 die Darstellung $+0$ und -0 besitzt, können wir insgesamt $2^n - 1$ Zahlen darstellen.

2.5 Beispiel. Für $n = 3$

Bitmuster	Dezimaldarstellung
000	+0
001	+1
...	...
100	-0
...	...
111	-3

Aber: Diese Darstellung am Computer ist unpraktisch, da die vier Grundrechenarten auf Hardwareebene typischerweise mit Hilfe von Addition und Zusatz-Logik umgesetzt werden.

Lösung: Komplementdarstellung

2.3 Komplementdarstellung

2.6 Definition ((b-1)-Komplement). Sei $z = (z_{n-1} \dots z_1 z_0)_b$ eine n -stellige b -adische Zahl. Das $(b-1)$ -Komplement $K_{b-1}(z)$ ist definiert als:

$$K_{b-1} = (b - 1 - z_{n-1}, \dots, b - 1 - z_0)_b$$

Geben wir hierzu direkt ein paar Beispiele an

2.7 Beispiel. Komplemente

- $K_9((325)_{10}) = (674)_{10}$ (9er-Komplement im 10-er System)
- $K_1((10110)_2) = (01001)_2$ (1er-Komplement im 2er System)

2.8 Definition (b-Komplement). Das b -Komplement einer b -adischen Zahl $z \neq 0$ ist definiert als

$$K_b(z) = K_{b-1}(z) + (1)_b$$

2.9 Beispiel. • $K_{10}((325)_{10}) = (674)_{10} + (1)_{10} = (675)_{10}$

2.10 Lemma. Für jede n -stellige b -adische Zahl z gilt:

- i) $z + K_{b-1}(z) = (b-1, \dots, b-1)_b = b^n - 1$
- ii) $K_{b-1}(K_{b-1}(z)) = z$

Ist außerdem $z \neq 0$ so gilt:

- iii) $z + K_b(z) = b^n$
- iv) $K_b(K_b(z)) = z$

Beweis. Hilfssatz

(i) Durch nachrechnen:

$$\begin{aligned}
 z + K_{b-1}(z) &= (z_{n-1} \dots z_0)_b + (b-1 - z_{n-1}, \dots, b-1 - z_0)_b \\
 &= \sum_{i=0}^{n-1} z_i b^i + \sum_{i=0}^{n-1} (b-1 - z_i) b^i \\
 &= \sum_{i=0}^{n-1} (b-1) b^i = (b-1, \dots, b-1)_b \\
 &= (b-1) \sum_{i=0}^{n-1} b^i \\
 &= (b-1) \left(\frac{b^n - 1}{b-1} \right) \\
 &= b^n - 1
 \end{aligned}$$

- (ii) per Definition
- (iii) Nachrechnen:

$$z + K_b(z) = z + K_{b-1}(z) + 1 = b^n - 1 + 1 = b^n$$

- Definiere $\hat{z} = K_b(z) = K_{b-1}(z) + (1)_b > 0$ und rechne

$$z + K_b(z) = b^n = \hat{z} + K_b \hat{z} + K_b(K_b(z)) \implies \text{Behauptung.}$$

□

2.11 Bemerkung. Modifikation

- Die 3. Aussage gilt auch für $z = 0$, falls man dann bei der Addition von 1 die Anzahl der Stellen erweitert.
- die 4. Aussage gilt für $z = 0$, falls überall im Beweis modulo b^n gerechnet wird.

Außerdem impliziert die 3. Aussage des Lemmas, dass

$$K_b(z) = b^n - z. \quad (2.1)$$

Dies können wir geschickt zum Darstellen der b^n verschiedenen ganzen Zahlen z mit

$$-\left\lfloor \frac{b^n}{2} \right\rfloor \leq z \leq \left\lceil \frac{b^n}{2} \right\rceil - 1$$

nutzen. Diesen Bereich nennt man darstellbaren Bereich.

2.12 Definition (b-Komplement-Darstellung). Die b-Komplement-Darstellung $(z)_{K_b} = (z_{n-1} \dots z_0)_b$ einer Zahl $z \in \mathbb{Z}$ im darstellbaren Bereich ist definiert als:

$$(z)_{K_b} = \begin{cases} (z)_b & \text{falls } z \geq 0 \\ (K_b(|z|))_b & \text{falls } z < 0 \end{cases}$$

2.13 Beispiel. Der darstellbare Bereich.

- Sei $b = 10$, $n = 2$.
Dann impliziert (2.1), dass

$$K_{10}(50) = 10^2 - 50 = 50$$

$$K_{10}(49) = 100 - 49 = 51.$$

Der darstellbare Bereich ist nun

$$-50 \leq z \leq 49$$

und hat konkrete Darstellungen:

Darstellung	Zahl
0	+0
1	+1
...	...
49	+49
50	-50
...	...
99	-1

- Sei $b = 2$, $n = 3$ Der darstellbare Bereich ist $-4 \leq z \leq 3$.

Bitmuster	Dezimaldarstellung
000	0
001	1
...	...
100	-4
...	...
111	-1

Wir betrachten nun Addition und Subtraktion zweier Zahlen in b-Komplement-Darstellung. Hierzu bezeichne $(x)_{K_b} \oplus (y)_{K_b}$ die ziffernweise Addition der Darstellungen von x und y mit Übertrag ("schriftlich rechnen"), wobei ein eventueller Überlauf auf die $(n+1)$ -te Stelle vernachlässigt wird (wir rechnen also immer mit modulo b^n).

2.14 Satz (Addition in b-Komplement-Darstellung). Seien x und y zwei n-stellige, b-adische Zahlen und x, y und $x + y$ im darstellbaren Bereich. Dann gilt:

$$(x + y)_{K_b} = (x)_{K_b} \oplus (y)_{K_b}$$

Beweis. Wir betrachten dazu mehrere Fälle.

- Fall $x, y \geq 0$:

$$\begin{aligned} (x)_{K_b} \oplus (y)_{K_b} &\stackrel{\text{Def.}}{=} ((x)_b + (y)_b) \mod b^n \\ &\stackrel{\text{Def.}}{=} (x + y) \mod b^n \\ &\stackrel{\text{Def.}}{=} (x + y)_{K_b} \end{aligned}$$

Der letzte Schritt ist möglich, da $(x + y)_{K_b}$ im darstellbaren Bereich sind.

- Fall $x, y < 0$

$$\begin{aligned} (x)_{K_b} \oplus (y)_{K_b} &\stackrel{\text{Def.}}{=} ((K_b(|x|))_b + (K_b(|y|))_b) \mod b^n \\ &= ((K_b(|x|)) + (K_b(|y|))) \mod b^n \\ &= (b^n - |x| + b^n - |y|) \mod b^n \\ &= (x + y) \mod b^n \\ &= (x + y)_{K_b} \end{aligned}$$

- Fall $x \geq 0, y < 0$:

$$\begin{aligned} (x)_{K_b} \oplus (y)_{K_b} &\stackrel{\text{Def.}}{=} ((x)_b + (K_b(|y|))_b) \mod b^n \\ &= (x + K_b(|y|)) \mod b^n \\ &= (x + b^n - |y|) \mod b^n \\ &= (x + y) \mod b^n \\ &= (x + y)_{K_b} \end{aligned}$$

- Fall $x < 0, y \geq 0$: Analog

□

2.15 Satz (Subtraktion in b-Komplement-Darstellung). Seien x und y n -stellige b -adische Zahlen und x, y und $x-y$ im darstellbaren Bereich. Dann gilt:

$$(x - y)_{K_b} = (x)_{K_b} \oplus (K_b(y))_{K_b}$$

Beweis. • Fall $y = 0$ nichts zu zeigen.

- $y \neq 0$: (2.1) impliziert:

$$-y = K_b(y) - b^n$$

mit modulo b^n -rechnen folgt, dass

$$(-y)_{K_b} = (K_b(y))_{K_b}$$

ist und somit:

$$\begin{aligned} (x - y)_{K_b} &= (x + (-y))_{K_b} \\ &= (x)_{K_b} \oplus (-y)_{K_b} \\ &= (x)_{K_b} \oplus (K_b(y))_{K_b} \end{aligned}$$

□

2.16 Beispiel. Für $b = 10$ und $n = 2$ ist der darstellbare Bereich $-50 \leq z \leq 49$

- $(20 + 7)_{K_{10}} = (20)_{K_{10}} \oplus (7)_{K_{10}} = (27)_{K_{10}} = 27$
- $28 - 5 = 28 + (-5) = (28)_{K_{10}} \oplus (95)_{K_{10}} = (23)_{K_{10}} = 23$
- $-18 - 20 = (-18) + (-20) = \dots = -28$

Die darstellbaren Zahlen kann man sich beim K_b -Komplement als Zahlenrad vorstellen.

Achtung Ein eventueller Überlauf bzw. Unterlauf wird im Allgemeinen nicht aufgefangen.

2.17 Beispiel. In n -stelliger Binäarithmetik ist die größte darstellbare Zahl $x_{max} = (011 \dots 1)_{K_2}$ gleich $2^{n-1} - 1$. Hingegen ist $x_{max} + 1 = (100 \dots 0)_{K_2}$ und wir als -2^{n-1} interpretiert.

2.4 Fest-Komma-Darstellung

2.18 Definition. Bei der Festkommadarstellung einer n -stelligen Zahl werden k Vorkomma und $n-k$ Nachkommastellen definiert:

$$z = + - (z_{k-1} \dots z_0 . z_{-1} \dots z_{k-n})_b = + - \sum_{i=k-n}^{k-1} z_i b^i$$

2.19 Beispiel. Im Zehner System wie gehabt.

Im Binärsystem ergibt sich folgendes: $(101.01)_2 = 2^2 + 2^0 + 2^{-2} = 5.25$

Achtung Im Gegensatz zur Darstellung ganzer Zahlen können bereits bei der Konvertierung von Dezimalzahlen in das b-adische Zahlensystem Rundungsfehler auftreten.

2.20 Beispiel. $(0.8)_{10} = (0.110\overline{1100})_2$

Das größte Problem der Festkommadarstellung ist allerdings, dass der darstellbare Bereich stark eingeschränkt ist und schlecht aufgelöst ist, da der Abstand zwischen zweier Zahlen immer gleich ist.

2.21 Beispiel. Die kleinste darstellbare Zahl in Fixkommadarstellung ist

$$z_1 = (0 \dots 0.0 \dots 01)_b$$

Die zweitkleinste Zahl ist $z_2 = 2z_1$. Wir wollen $x = \frac{z_1 + z_2}{2}$ in Fixkommadarstellung darstellen, müssen wir entweder zu z_1 abrunden oder zu z_2 aufrunden.

2.22 Der relative Fehler dieses Runden ist:

$$\frac{|x - z_1|}{|x|} = \frac{1}{3}$$

Analog für z_2 .

Solche Fehler machen jede Rechnung unbrauchbar.

2.5 Gleitkommadarstellung

2.23 Definition (Gleitkommadarstellung). Die Gleitkommastellung einer Zahl $z \in \mathbb{R}$ ist gegeben durch:

$$z = + - m \cdot b^e$$

mit einer Matisse m , dem Exponenten e und der Basis b .

Der Einfachheit halber nehmen wir an, dass die Basis für alle Zahlen gleich ist, obwohl sie prinzipiell verschieden sein könnte.

2.24 Beispiel. Die Gleitkommadarstellungen von:

- $(-384.753)_{10} = -3.84753 \cdot 10^2$
- $(0.00042)_{10} = 4.2 \cdot 10^{-4}$
- $(1010.101)_2 = (1.010101)_2 \cdot 2^3$

Achtung: Die Gleitkommadarstellung ist nicht eindeutig!

2.25 Definition. Die Matisse m heißt normalisiert, falls $m = m_1.m_2m_3 \dots m_t$ mit $1 \leq m_1 \leq b$.

Um die eindeutige, normalisierte Gleitkommadarstellung im Computer zu speichern, müssen wir festlegen, wie viele Stellen für Matisse und Exponent zur Verfügung gestellt werden. Dies resultiert in der Menge

$$F = F(b, t, e_{min}, e_{max}) = \{z = + - m_1.m_2m_3 \dots m_t \cdot b^e \mid e_{min} \leq e \leq e_{max}\}$$

Da 0 nicht in dieser Form dargestellt werden kann, reserviert man dafür eine spezielle Ziffernfolge in Matisse und Exponent.

2.26 Bemerkung. Das Hidden Bit ist hilfreich

- Für $b = 2$ kann die erste Ziffer m_1 der Mantissee weggelassen werden (Hidden Bit)
- Um den Exponenten besser vergleichen zu können, verwendet man oft die Exzess- oder Bias-Darstellung. Durch Addition der Exzesser $|e_{\min}| + 1$ wird der Exponent auf den Bereich $1, 2, \dots, |e_{\min}| + e_{\max} + 1$ transformiert.

2.27 Beispiel (IEEE 754 Standard). Betrachte binäre Gleitkommazahlen mit 64 Bit auf dem Computer.

- 52 Bit für die Mantissee in Hidden Bit Darstellung
- 11 Bit für den Exponenten mit $e_{\min} = -1022$ und $e_{\max} = 1023$ gespeichert in Exzessdarstellung.
- 1 Bit als Vorzeichen.

Weitere Fälle können Online nachgelesen werden.

Genauigkeit der Gleitkommadarstellung

Da die Menge aller Zahlen in $F = F(b, t, e_{\min}, e_{\max})$ endlich ist, müssen wir Zahlen in $\mathbb{R} \setminus F$ geeignet annähern.

2.28 Definition. Die Rundung ist eine Abbildung $rd: \mathbb{R} \rightarrow F$ mit

- $rd(a) = a, a \in F$
- $rd(z), z \in \mathbb{R}$ ist gegeben so, dass $|z - rd(z)| = \min_{a \in F} (z - a)$ für alle $z \in \mathbb{R}$.

2.29 Definition (Maschinengenauigkeit). Den maximalen relativen Rundungsfehler $\varepsilon_{\text{mach}}$ für $z_{\min} \leq |z| \leq z_{\max}$ nennt man Maschinengenauigkeit. Die Stellen der Mantissee heißen signifikante Stellen. Wenn die Mantissee t -stellig ist, spricht man von t -stelliger Arithmetik.

2.30 Satz (Maschinengenauigkeit von F). Die Maschinengenauigkeit $\varepsilon_{\text{mach}}$ für $F = F(b, t, e_{\min}, e_{\max})$ ist:

$$\varepsilon_{\text{mach}} = \frac{1}{2} b^{1-t}.$$

Beweis. Betrachte relativen Rundungsfehler ε , der vom Abscheiden nicht signifikanten Stellen herrührt. Sei $z_{\min} < x < z_{\max}$ und \tilde{x} die durch Abschneiden entstehende Zahl. Dann gilt:

$$\begin{aligned} \varepsilon &= \frac{|x - \tilde{x}|}{|x|} \\ &= \frac{|x_1 x_2 x_3 \dots x_t x_{t+1} \dots \cdot b^e - x_1 x_2 x_3 \dots x_t \cdot b^e|}{|x|} \\ &= \frac{|0.x_{t+1} \dots| \cdot |b^{e+1-t}|}{|x|} \end{aligned}$$

Da $|0.x_{t+1} \dots| < 1$ und $b^e \leq |x|$ gilt:

$$\varepsilon < \frac{b^{e+1-t}}{b^e} = b^{1-t}$$

Da der Rundungsfehler ε_{mach} höchsten halb so groß ist wie der Abschneidefehler folgt die Behauptung. \square

Die Maschinengenauigkeit ε_{mach} ist die wichtigste Größe zur Beurteilung der Genauigkeit von Gleitkommarechnungen am Computer. Sie gibt Aufschluss zur Anzahl signifikanter Stellen zur Basis b . Die zugehörige Anzahl s Stellen im Dezimalsystem erhält man durch das Auflösen von:

$$\varepsilon_{mach} = \frac{1}{2}b^{1-t} = \frac{1}{2}10^{1-s}$$

nach s . Es folgt daher:

$$s = \lfloor 1 + (t - 1) \log_{10}(b) \rfloor$$

Für den IEEE 754 Standard mit $b = 2, t = 53$ erhält man $s = 16$ signifikante Stellen.

3 Fehleranalyse

3.1 Rechnerarithmetik

3.1 Beispiel. Betrachte $F = F(10, 5, -4, 5)$ und Maschinenzahlen

$$x = 2.5684 \cdot 10^0 = 2.56840000$$

$$y = 3.2791 \cdot 10^{-3} = 0.0032791$$

Es gilt:

$$\left. \begin{array}{l} x + y = 2.5716792 \\ x - y = 2.5651209 \\ x \cdot y = 0.00842204044 \\ \frac{x}{y} = 783.2637004 \end{array} \right\} \notin F$$

3.2 Bemerkung. Die Menge $F(b, t, e_{min}, e_{max})$ ist nicht abgeschlossen bezüglich der Grundrechenarten und können somit im Allgemeinen nicht im Computer implementiert werden.

Lösung Wir runden das Ergebnis und implementieren so eine Pseudoarithmetik. Das bedeutet wir ersetzen $o \in \{+, -, \cdot, /\}$ durch \boxed{o} , definiert durch $x \boxed{o} y := rd(x \circ y)$.

$$s = t \tag{3.2}$$

Auf Hardwareebene wird üblicherweise mit einer längeren Mantisse gearbeitet und dann normalisiert und gerundet. Dies entspricht dem IEEE 754 Standard.

3.3 Bemerkung. Für $|x|, |y|, |x \circ y| \in [z_{min}, z_{max}]$ impliziert (3.2), dass