

Vorlesungsskript

Falk Jonatan Strube

Vorlesung von Herrn Meinhold

16. November 2015





Inhaltsverzeichnis

I.	. Elementare Grundlagen I. Aussagen und Grundzüge der Logik					
1.						
2.	Mengen	1				
3.	Zahlen	1				
	3.1. Gruppen, Ringe, Körper	1				
	3.2. Zahlentheorie	2				
	3.3. Reelle Zahlen	Ę				
	3.3.1. Algebraische Struktur					
	3.3.2. Zahlendarstellung im Computer	7				
	3.3.3. Ordnungsstruktur					



Teil I. Elementare Grundlagen

- 1. Aussagen und Grundzüge der Logik
- 2. Mengen
- 3. Zahlen
- 3.1. Gruppen, Ringe, Körper
 - Gegeben sei eine Menge M und eine zweistellige Operation \circ (d.h. Abb. von $M \times M$ in M) Bezeichnung: (M, \circ) , analog $(M, \circ, *)$
 - Die Operation \circ heißt *kommutativ*, wenn $a \circ b = b \circ a$ für alle $a, b \in M$.
 - Die Operation \circ heißt *assoziativ*, wenn $(a \circ b) \circ c = a \circ (b \circ c)$ für alle $a, b, c \in M$.

Def. 1:

 (M, \circ) heißt *Gruppe*, wenn gilt:

- 1.) Die Operation ∘ ist assoziativ
- 2.) Es gibt genau ein *neutrales Element* $e \in M$ mit $a \circ e = e \circ a = a$ (für alle $a \in M$)
- 3.) Es gibt zu jedem $a \in M$ genau ein *inverses Element* a^{-1} mit $a \circ a^{-1} = a^{-1} \circ a = e$
- 4.) Eine Gruppe heißt *ABELsch*, wenn zusätzlich folgendes gilt:
 ∘ ist kommutativ

Def. 2:

 $(M, \oplus, *)$ heißt *Ring*, wenn gilt:

- 1.) (M, \oplus) ist eine ABELsche Gruppe.
- 2.) Die Operation * ist assoziativ.
- 3.) Es gelten für beliebige $a, b, c \in M$:

$$a*(b\oplus c)=(a*b)\oplus(a*c)$$
 $(a\oplus b)*c=(a*c)\oplus(b*c)$ (Distributivgesetze)

- 4.) Ein Ring heiß kommutativer Ring, wenn gilt:
 - * ist kommutativ

Def. 3:

 $(M, \oplus, *)$ heißt *Körper*, wenn gilt:

- 1.) $(M, \oplus, *)$ ist ein Ring (mit dem neutralen Element E_0 für die Operation \oplus)
- 2.) $(M \setminus \{E_0\}, *)$ ist eine ABELsche Gruppe (mit dem neutralen Element E_1 für die Operation *)



3.2. Zahlentheorie

- Eine natürliche Zahl p > 1, die nurch durch 1 und sich selbst teilbar ist heißt *Primzahl*.
- ullet Jede natürliche Zahl n>1 ist entweder eine Primzahl, oder sie lässt sich als Produkt von Primzahlen schreiben.

Diese sogenannte Primfaktorzerlegung ist bis auf die Reihenfolge der Faktoren eindeutig.

Def. 4:

Zwei natürliche zahlen aus \mathbb{N}^* heißen *teilerfremd*, wenn sie außer 1 keine gemeinsamen teiler besitzen

- Es sei $a \in \mathbb{Z}$ und $m \in \mathbb{N}^*$. Dann gibt es eine eindeutige Darstellung der Gestalt $a = q \cdot m + r$ mit $0 \le r < m$ und $q \in \mathbb{Z}$. Bezeichnung: $m \dots$ Modul $m \in \mathbb{Z}$. (kleinste nichtnegative) Rest modulo $m \in \mathbb{Z}$ mod(a, m))
- Zur Erinnerung: a und b seien ganze Zahlen, $m \in \mathbb{R}^*$, dann $a \equiv b \pmod{m}$ [a kongruent $b \bmod m$]

```
\Leftrightarrow a und b haben den gleicher Rest modulo\ m \Leftrightarrow a-b ist durch m teilbar (d.h. \exists k \in \mathbb{Z} \quad a-b=k\cdot m)
```

Satz 1:

Es sei $a \equiv b \pmod{m}$, $c \equiv d \pmod{m}$, dann gilt: $a + c \equiv b + d \pmod{m}$ und $a \cdot c \equiv b \cdot d \pmod{m}$ (d.h. in Summen und Produktenn darf jede Zahl durch einen beliebigen Vertreter der gleichen Restklasse ersetzt werden).

Bsp. 1:

- a) $307 + 598 \equiv 1 + (-2) \equiv -1 \equiv 5 \pmod{6}$
- b) $307 \cdot 598 \equiv 1 \cdot (-2) \equiv -2 \equiv 4 \pmod{6}$
- c) $598^6 \equiv (-2)^6 \equiv 64 \equiv 4 \pmod{6}$
- Man wählt aus jeder Restklasse den kleinsten nichtnegativen Vertreter
 - \sim Menge von Resten $modulo\ m$: $\mathbb{Z}_m := \{0, 1, ..., m-1\}$
 - \sim "modulare Arithmetik": Operation \oplus und \odot für Zahlen aus \mathbb{Z}_m erklärbar, in dem für das Ergebnis jeweils der kleinste nichtnegative Rest $modulo\ m$ gewählt wird (vgl. Satz 1)

z.B.
$$\mathbb{Z}_7 = \{0, 1, ..., 6\}, \quad 5 \oplus 4 = 2$$
, da $5 + 4 \equiv 9 \equiv 2 \pmod{7}$ $5 \odot 6 = 2$, da $5 \cdot 6 \equiv 30 \equiv 2 \pmod{7}$

Falls keine Verwechselung zu befürchten ist, wird die übliche Schreibweise + und \cdot anstelle von \oplus und \odot verwendet.

Def. 5:

Wenn es zu $c \in \mathbb{Z}_m$ eine Zahl $d \in \mathbb{Z}_m$ gibt, mit $c \cdot d \equiv 1 \pmod{m}$ (bzw. $c \odot d \equiv 1$), so heißt d die *(multiplikative) modulare Inverse* zu c in \mathbb{Z}_m . Bezeichnung: $d = c^{-1}$

Bsp. 2:

$$c=3\in\mathbb{Z}_7$$
, wegen $3\cdot 5\equiv 1 \pmod{7}$ ist (in \mathbb{Z}_7) $3^{-1}=5$.

Satz 2: Zu $a \in \mathbb{Z}_m, a \neq 0$, gibt es genau dann eine modulare Inverse in \mathbb{Z}_m , wenn a und m teilerfremd sind (ggT(a,m)=1).



Satz 3: Es sei p eine Primzahl. Dann ist $(\mathbb{Z}_m, \oplus, \odot)$ ein Körper. Bemerkung: Falls m keine Primzahl ist, so ist $(\mathbb{Z}_m, \oplus, \odot)$ ein kommutativer Ring.

EUKLIDischer Algorithmus

- Verfahren zur Ermittlung des größten gemeinsamen Teilers t zweier positiver natürlicher Zahlen, t = ggT(a,b).
- In erweiterter Form bietet der Algorithmus eine Möglichkeit zur Bestimmung der modularen Inversen von a zum Modul m (mit a < m und a, m teilerfremd).

Satz 4: (EUKLIDischer Algorithmus)

Es seien $a, b \in \mathbb{N}^*, a > b$. Man bildet die endliche Folge

 $r_0 := b, \ r_1 = mod \ (a,b), \ r_2 = mod \ (r_0,r_1),..., \ r_n = mod \ (r_{n-2},r_{n-1}),$ Abbruch falls $r_n = 0$.

In diesem Fall gist $ggT(a,b) = r_{n-1}$ (letzter nicht verschwindender Rest).

Bezeichnung: j-te Division ... $r_{j-1} = q_j \text{ Rest } r_j$ (j = 1, ..., n) (dabei $r_1 := a$).

Satz 5: (erweiterter EUKLIDischer Algorithmus)

Zusätzlich zur Folge (r_n) aus Satz 4 bilde man die Folgen

$$\begin{array}{ll} x_0=0,\; x_1=1,\; x_2=x_0-q_2x_1,...,\; x_j=x_{j-2}-q_jx_{j-1} & (j\leq n-1) \text{ und } \\ y_0=1,\; y_1=-q_1,\; y_2=y_0-q_2y_1,...,\; y_j=y_{j-2}-q_jy_{j-1} & (j\leq n-1) \\ \text{Dann gilt für alle } j=0,...,\; n-1: \boxed{r_j=x_j\cdot a+y_j\cdot b} \\ \text{Insbesondere gilt } \boxed{ggT(a,b)=x_{n-1}\cdot a+y_{n-1}\cdot b} \end{array}$$

Diskussion:

- 1.) Der Sinn der erweiterten EUKLIDischen Algorithmus besteht darin, in jedem Schrit den *Divisionsrest* r als linearkombination von a und b mit ganzzahligen Koeffizienten x und y darzustellen: $r = x \cdot a + y \cdot b$
 - Der Mechanismus wird am besten im Rechenschema des nachfolgenden Bsp. 4 deutlich.
- 2.) Sind c und m teilerfremd, $1 \le c < m$, d.h. ggT(m,c) = 1, so erhält man mit dem erweiterten EUKLIDischen Algorithmus (a = m, b = c) eine Darstellung in der Form $1 = x \cdot m + y \cdot c$. $y \cdot c \equiv 1 \pmod{m}$ und damit $c^{-1} \equiv y \pmod{m}$ (für die modulare Inverse muss eventuell noch der in \mathbb{Z}_m liegende, zu y kongruente, Wert gebildet werden!).

Bsp. 3:

Man ermittle den größten gemeinsamen Teiler t sowie das kleinste gemeinsame Vielfache v der Zahlen 132 und 84.

• Es genügt der "einfache" Algorithmus:

$$132:84=1$$
 Rest 48 $84:48=1$ Rest 36 $48:36=1$ Rest $12 \curvearrowright t=ggT(132,84)=\underline{\underline{12}}$ $36:12=3$ Rest $\boxed{0} \curvearrowright$ Ende.

•
$$v = \frac{a \cdot b}{t} = \frac{132 \cdot 84}{12} = \underline{924} = kgV(132, 84)$$



Bsp. 4:

Man ermittle die modulare Inverse von $\overbrace{11}^{b}$ zum Modul $\overbrace{25}^{a}$

 $\curvearrowright (-9) \cdot 11 \equiv 1 \pmod{25}$

 $\sim 11^{-1} \equiv -9 \equiv 16 \pmod{25}$, die Inverse von 11 in \mathbb{Z}_{25} ist 16.

Zu den Schritten:

- (1) $b = 0 \cdot a + 1 \cdot b$
- (2) mittleres Feld als Linearkombination
- (3) ab hier Rechnung links spaltenweise durchführen, dabei Faktoren a und b beibehalten.

EULERsche φ -Funktion, Satz von EULER

Def. 6:

Es sei $n \in \mathbb{N}^*$. Dann *EULERsche* φ -Funktion:

 $\varphi(n) := \text{Anzahl der zu } n \text{ teilerfremden Elemente aus } \{1, 2, ..., n\}.$ Eigenschaften der φ -Funktion:

- Es sei p eine Primzahl, dann ist $\varphi(p)=p-1$, $\varphi(p^k)=p^{k-1}(p-1)$ $(k\in\mathbb{N}^*)$
- Falls ggT(m,n)=1, so gilt $\varphi(m\cdot n)=\varphi(m)\cdot\varphi(n)$.
- Speziell: $n=p\cdot q$ (p,q Primzahlen), dann $\boxed{\varphi(n)=(p-1)\cdot (q-1)}$ (1).

Satz 6: (Satz von EULER)

Es sei ggT(a, n) = 1, dann gilt:

$$\boxed{a^{\varphi(n)} \equiv 1 \pmod{n}}$$
 (2).

RSA-Verschlüsselung

- Die Formeln (1) und (2) [siehe oberhalb] bilden die Grundlage für die sogenannte RSA-Verschlüsselung (RIVES, SHAMIR, ADLEMAN - 1978)
- Schlüsselerzeugung:
 - 1.) Man wählt (in der Praxis sehr große) Primzahlen d und q.
 - **2.)** $n := p \cdot q, m := \varphi(n) \stackrel{(1)}{=} (p-1)(q-1)$
 - 3.) e wird so gewählt, dass ggT(e, m) = 1
 - 4.) $d := e^{-1} \pmod{m}$ (modulare Inverse)
 - 5.) (n, e) ... öffentlicher Schlüssel (n, d) ... geheimer Schlüssel (geheim ist nur d) p, q und m werden nicht mehr benötigt, bleiben aber geheim!



- Verschlüsselung: Klartext a teilerfremd zu n verschlüsseln mit e, d.h. $b :\equiv a^e \pmod{n}$ bilden $(b \dots$ Geheimtext)
- Entschlüsselung: Der Empfänger und Besitzer des geheimen Schlüssels bildet $b^d (mod\ m)$ und erhält $b^d \equiv a (mod\ n)$ denn $b^d \equiv (a^e)^d \equiv a^{ed} \equiv a^{1+k\cdot m} \equiv a \cdot \begin{pmatrix} a^{\varphi(n)} \end{pmatrix}^k \equiv a (mod\ n)$.
- Praktische Durchführung vgl Übungsaufgabe 2.4

3.3. Reelle Zahlen

 \mathbb{R} ... Menge der reellen Zahlen.

Auf \mathbb{R} existiert eine algebraische Struktur und eine Ordnungsstruktur.

3.3.1. Algebraische Struktur

 $(\mathbb{R}, +, \cdot)$ mit den arithmetischen Operationen + (Addition) und \cdot (Multiplikation) ist ein Körper.

Def. 7:

a.) $0! := 1, \ n! = n \cdot (n-1)!$ mit $n \in \mathbb{N}^*$ Fakultät (rekursive Funktion)

b.) Sei
$$\alpha \in \mathbb{R}, k \in \mathbb{N}^*$$
, dann sei $\left[\begin{pmatrix} \alpha \\ 0 \end{pmatrix} := 1, \begin{pmatrix} \alpha \\ k \end{pmatrix} := \frac{\alpha}{k} \begin{pmatrix} \alpha - 1 \\ k - 1 \end{pmatrix} \right]$ Binominalkoeffizient α über k . d.h. $\left[\begin{pmatrix} \alpha \\ k \end{pmatrix} = \frac{\alpha(\alpha-1)(\alpha-2)...(\alpha-k-1)}{k!} \right]$

Diskussion:

1.) Für
$$k,n\in\mathbb{N},\ 0\leq k\leq n$$
 gilt
$$\binom{n}{k}=\binom{n}{n-k}=\frac{n!}{k!(n-k)!}$$

2.) Für
$$k \in \mathbb{N}, \alpha \in \mathbb{R}$$
 gilt
$$\binom{\alpha}{k} + \binom{\alpha}{k+1} = \binom{\alpha+1}{k+1}$$



Stellenwertsysteme:

• Es sei k > 1 eine natürliche Zahl (die sogenannte Basis)

$$x = \underbrace{(x_p x_{p-1} ... x_1 x_0, \underbrace{x_{-1} x_{-2} ... x_{-q}}_{\text{Nachkomma}})_b }_{\text{Vorkomma}}$$

$$:= \underbrace{x_p \cdot b^p + x_{p-1} \cdot b^{p-1} + ... + x_1 \cdot b^1 + x_0 \cdot b^0}_{\text{Vorkomma}} + \underbrace{x_{-1} \cdot b^{-1} + x_{-2} \cdot b^{-2} + ... + x_{-q} \cdot b^{-q}}_{\text{Nachkomma}}$$

heißt Darstellung zur Basis b (*).

Bsp. 5:

- $b = 2 \dots$ Dual- oder Binärsystem (Ziffern $\{0, 1\}$)
- $b=3\dots$ Trialsystem
- $b = 10 \dots$ Dezimalsystem
- $b = 16 \dots$ Hexadezimalsystem (Ziffern $\{0, 1, 2, ..., 9, \underbrace{A}_{10}, \underbrace{B}_{11}, \underbrace{C}_{12}, \underbrace{D}_{13}, \underbrace{E}_{14}, \underbrace{F}_{15}\}$

z.B.
$$(47)_{10} = (101111)_2 = (1202)_3 = (57)_8 = (2F)_{16}$$

Übergang von einem Ziffernsystem zu einem anderen

z.B.
$$p = 3$$
, $q = 2$

$$x = x_3b^3 + x_2b^2 + x_1b^1 + x_0 + x_{-1}b^{-1} + x_{-2}b^{-2}$$

$$= (x_3b^2 + x_2x^1 + x_1)b + x_0 + (x_{-1} + x_{-2}b^{-1})b^{-1}$$

$$= ((x_3b^1 + x_2)b + x_1)b + x_0 + (x_{-1} + x_{-2}b^{-1})b^{-1}$$

Grundlage: fortgesetztes Klammern:

$$x = ((....((x_pb + x_{p-1})b + x_{p-2})b + ... + x_2)b + x_1)b + x_0 + ((...(x_{-q}b^{-1} + x_{-(q-1)})b^{-1} + ... + x_{-2})b^{-1} + x_{-1})b^{-1}$$

Bsp. 6: Übergang Dezimalsystem → anderes System

- ganzer Teil: fortgesetzte Division durch b und Restabspaltung liefert b-Ziffern in der Reihenfolge x_0, x_1, x_2, \dots
- gebrochener Teil: fortgesetzte Multiplikation mit b und Abspaltung des ganzzahligen Anteils liefert b-Ziffern in der Reihenfolge $x_{-1}, x_{-2}, ...$

z.B. Dezimalsystem \rightarrow Hexadezimalsystem (b = 16) x = 435, 9

ganzer Teil:

(**)

$$435: 16 = 27$$
 Rest $3 \rightarrow x_0$
 $27: 16 = 1$ Rest $11 \rightarrow x_1$
 $1: 16 = \boxed{0}$ Rest $1 \rightarrow x_2$

gebrochener Teil:

$$0,9\cdot 16=0.4$$
 $+14$ $\to x_{-1}$ $0,4\cdot 16=\boxed{0.4}$ $+6$ $\to x_{-2}$ (Periode, da gleicher "Nachkommarest") $\to x=(1B3,E\overline{6})_{16}$



Bsp. 7: Übergang anderer Systeme → Dezimalsystem

Entweder direkte Auswertung von (*) (v.a. beim Dualsystem \rightarrow Addition von 2er-Potenzen) *oder* Klammern in (**) von innen nach außen berechnen (zweckmäßig HORNER-Schema).

- ullet ganzer Teil: beginnend mit x_p ABB1
- $\bullet \ \ {\rm gebrochener\ Teil:\ beginnend\ mit\ } x_{-q} \\ \ \ {\rm ABB\ 2} \\$

$$x = (1E2, B8)_{16}$$

ganzer Teil:

gebrochener Teil:

Bsp. 8: Hexadezimalsystem ↔ Dualsystem

- 4 Dualziffern entsprechen einer Hexadezimalziffer ($2^4=16$) \sim 4er Gruppen von Dualziffern ab Komma bilden.
 - a.) $(A8C, B2)_{16} = (1010\ 1000\ 1100,\ 1011\ 1011\ 001(0))_2$
 - **b.)** $((0)11011110, 101(0))_2 = (6E, A)_{16}$

3.3.2. Zahlendarstellung im Computer

- 1.) Ganze Binärzahlen in Zweierkomplementdarstellung (n Bit, meist n = 8, 16, 32, 64)
 - Bsp.: n = 8 $(100)_{10} = (64)_{16}$ 0 1 1 0 0 1 0 2^{7} 2^{6} 2^5 2^{4} 2^{3} 2^{1} 2^{0}

MSB: most siginficant bit (LSB: least significant bit)

- Um auch negative Zahlen darstellen zu können, wird das MSB als Vorzeichen reserviert. Negative Zahlen -a ($1 \le a \le 2^{n-1}$) werden im sogenannten Zweierkomplement $\overline{a} := 2^n a$ dargestellt ($\curvearrowright \overline{a} \ge 2^{n-1} \curvearrowright MSB = 1$)
- Nightnegtavie Zahlen $0 \le a \le 2^{n-1} 1$ werden unverändert dargestellt (MSB = 0)
- Damit Darstellung ganzer Zahlen von -2^{n-1} bis $2^{n-1}-1$ (Anzahl 2^n) möglich, n=8:-128bis127.
- $\bullet \ \ Umwandlung \ negativer \ Zahlen \rightarrow Zweierkomplement$



Bsp. 9: n=8, umzuwandeln sei -100 (dezimal) zwei Möglichkeiten:

Bemerkung: Das Zweierkomplement der positiven Zahl 100 ist die positive Zahl $156 = \overline{100}$, diese wird wegen MSB = 1 als negative Zahl -100 interpretiert.

2.) (am schnellsten): Rechts (beim LSB) beginnend alle Ziffern bis einschließlich der ersten 1 übernehmen (unverändert lassen), für alle höherwertigen Ziffern Z das Einerkomplement 1-z bilden:

Rückumwandlung (Zahl mit $\overline{M}SB = 1 \rightarrow$ neg. Zahl) analog:

$$\overline{156} = 256 - 156 = 100 \rightarrow -100$$

Die Subtraktion wird damit auf die Addition des Zweierkomplements zurückgeführt.

Bsp. 10:
$$a = 64 - 100 = 64 + (-100)$$

 $64 = 2^6 = 0100\,0000$
 $-100 = 1001\,1100 +$
 $Summe = 1101\,1100$ Ergebins negativ

 $36 = 0010\,0100$ dargestellst ist aber $\overline{z} = 2^n - z$

$$\sim$$
 Ergebnis: $a = -36 = -z$

Ein Überlauf (Ergebnis $\geq 2^{n-1}$ oder $< -2^{n-1}$) ensteht in folgenden Fällen (\rightarrow ERROR!)

	a	b	a+b	
MSB	0	0	1	(MSB = 0 erwartet!)
MSB	1	1	0	(MSB = 1 erwartet!)

Bemerkung: Für die Handrechnung (z.B. 2-5=:a) kleinere zahl von größerer Subtrahieren a=-(5-2), dabei genügt es für n die Binärstellenzahl des Minuenden $(5)_{10}=(101)_2$ also n=3 zu verwenden. Es wird dabei ausschließlich mit nicht-negativen Zahlen gerechnet $(0,1,...,2^n-1)$:

$$(5-2)_{10} = ((5+2^{n}-2)-2^{n})_{10} = (5+\overline{2}-2^{n})_{10}$$

$$(2)_{10} = (010)_{2} \curvearrowright \overline{2} = (110)_{2}$$

$$(5)_{10} = (101)_{2}$$

vordere Stelle 2ⁿ ignorieren

2.) Gleitkommasysteme

$$x = v \cdot m \cdot b^e$$
 dabei

$$\bullet \ v = (-1)^V \ \dots \ \textit{Vorzeichen} \begin{cases} V = 0 & \text{positive Zahl} \\ V = 1 & \text{negative Zahl} \end{cases}$$

• *m* ... *Mantisse*, Stellenzahl *p*, die Mantisse heißt *normalisiert* falls sie folgende Gestalt besitzt:

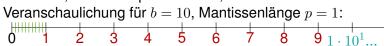


 $m_1, m_2, ..., m_p$ oder $0, m_1, m_2, ..., m_p$ mit $m \neq 0$. Dabei sind $m_1, m_2, ..., m_p$ die Ziffern zur

• $e \dots$ Exponent, ganzzahlig $e_{min} \le e \le e_{max}$.

In jedem Gleitkommasystem sind nur endlich viele Zahlen darstellbar, die Menge der reellen zahlen ist aber überabzählbar (unendlich).

Abstände, wächst Exponent um k, so wachsen die Abstände auf b^k -fache!)



Exponent 0: $1 \cdot 10^{0}, 2 \cdot 10^{0}, ..., 9 \cdot 10^{0}$

Exponent -1: $1 \cdot 10^{-1}$, $2 \cdot 10^{-1}$, ..., $9 \cdot 10^{-1}$

Exponent 1: $1 \cdot 10^1, 2 \cdot 10^1, ..., 9 \cdot 10^1$

Rundung: Zahlen, die nicht in diesem "Raster" enthalten sind, werden auf dei nächstgelegene darstellbare Gleitkommazahl gerundet. Falls die Zahl genau in der Mitte zwischen zwei darstellbaren Zahlen liegt, wird auf die gerade Zahl gerundet (Bsp. $3,75 \rightarrow 3,8$ oder $4,65 \rightarrow 4,6$ bei Rundung auf eine Stelle nach dem Komma).

Numerische Probleme beim Rechnen mit Gleitkommazahlen

 Kommutativ-, Assoziativ- und Distributivgesetze gelten im allgemeinen nicht mehr. Ursachen sind bspw. Ziffernauslöschung bei der Subtraktion fast gleicher Zahlen, Addition oder Subtraktion von Zahlen unterschiedlicher Größenordnung.

Bsp. 11:

- 1.) Man berechne (a+b)+c und a+(b+c) in einem System mit 3-Stelliger Mantisse: $a = 3,73 \cdot 10^6, b = -3,71 \cdot 10^6 \text{ und } c = 6,42 \cdot 10^3$
 - $-a+b=0.02\cdot 10^6=2.00\cdot 10^4$ (Normalisierung) $c=0,642\cdot 10^4=0,64\cdot 10^4$ (Exponentenangleichung und Rundung) $(a+b)+c=2,64\cdot 10^4=26400$
 - $-c = 0.00642 \cdot 10^6 = 0.01 \cdot 10^6$ (Exponentenangleichung und Rundung) $b+c=-3.70\cdot 10^6$ $a + (b + c) = 0.03 \cdot 10^6 = 3.00 \cdot 10^4 = 30000$
 - **–** exakter Wert: a + b + c = 26420
- 2.) Aufgabe der numerischen Mathematik ist es, die unvermeidlichen Genauigkeitsverluste beim Rechnen mit Maschinenzahlen durch optimale Organisation (Reihenfolge) der Rechnung und Fehleranalyse in Grenzen zu halten.
- 3.) Gleitkommaformat IEEE 754 (single precision, 32 Bit)

$$x = v \cdot m \cdot b^e = (-1)^V \cdot 1, m_2 m_3 ... m_2 4 \cdot 2^{E-B}$$
 (b = 2, Binärsystenm)

- Vorzeichen $V=0 \curvearrowright \text{positiv}, V=1 \curvearrowright \text{negativ}$ (1 Bit)
- Mantisse m_1 im Binärsystem stets gleich 1. \sim nur Abspeicherung von $M = m_2 m_3 ... m_2 4$ (23 Bit)
- Exponent: Abgespeichert wird E := e + Bmit dem sogenannten Biaswert B=127 (Bias = Verzerrung) als nichtnegative 8-stellige Binärzahl, $e_{min}=-126\;(E=1),\;e_{max}=127\;(E=254=(1111\;1110)_2),\;{\rm die\;Grenzf\"{a}lle}$



 $E=(0000\ 0000)_2$ und $E=(1111\ 1111)_2$ sind für Sonderfälle (0, ∞ , nichtdefinierte Werte) vorgesehen.

Abspeicherung in der Reihenfolge VEM (32 Bit)

Bsp. 12: Umwandlung einer Dezimalzahl in das IEEE 754-Format (32-Bit), x=435,9 (vgl. Bsp. 6)

- 1.) Konvertierung in Dualzahl (unter Verwendung von Bsp. 6/Hexadez.) $x=1B3, E\overline{6})_{16}=(1\ 1011\ 0011, 1100\ 0110\ 0110...)_2$
- 2.) Normalisierte Gleitkommadarstellung, Mantisse mit 23 Stellen nach dem Komma, Kommaverschiebung um 8 Stellen.

$$x = 1, \underbrace{1011\ 0011\ 1110\ 0110\ 0110\ 011}_{M} (0\ 0110...))_{2} \cdot 2^{8}$$
 (Abrundung!)

- 3.) Exponent $e = 8 \curvearrowright E = e + B = 8 + 127 = 135 = \underbrace{(1000\ 0111)_2}_{E}$
- 4.) Vorzeichenbit V=0 (da x positiv)

Bsp. 13: IEEE 754→ Dezimalzahl

1 1000 0011 0111 1100 0000 0000 0000 0000

1.)
$$E = (1000\ 00111)_2 = 131 \land e = E - B = 131 - 127 = 4$$

2.)
$$V=1$$
 \curvearrowright negativ, normalisierte Mantisse $1,M$ $\curvearrowright x=-(1,011111)_2\cdot 2^4=-(10111,11)_2$ $\curvearrowright x=-23,75$

Bemerkung:

- 1.) Neben dem single-Format gibt es in IEEE 754 das double-Format (54 Bit, V=1Bit, E=11Bit, M=52 Bit, B=1023) sowie das erweiterbare Format
- 2.) Zahlbereiche single: $1,401 \cdot 10^{-45} \dots 3,403 \cdot 10^{38}$, double: $4,941 \cdot 10^{-324} \dots 1,798 \cdot 10^{308}$

3.3.3. Ordnungsstruktur

- Durch ≤ (auch ≤) ist auf ℝ eine vollständige Ordnungsrelation erklärt.
- Verträglichkeit mit der algebraischen Struktur (für alle $x, y, z \in \mathbb{R}$):

$$(1) x \le y \quad \Rightarrow \quad x + z \leqq y + z$$

(2)
$$(x \le y) \land (z \ge 0)$$
 \Rightarrow $x \cdot z \le y \cdot z$ $(x \le y) \land (z \le 0)$ \Rightarrow $x \cdot z \ge y \cdot z$

Für die strikte Ordnung < gilt:

$$(x < y) \land (z < 0) \quad \Rightarrow \quad x \cdot > y \cdot z$$

Def. 8:

Sei x eine reele Zahl. Dann heißt $|x|:=\begin{cases} x & \text{für } x\geq 0 \\ -x & x<0 \end{cases}$ der (absolute) Betrag von x.



ABB 5

Vorzeichenfunktion
$$sgn(x) := \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

Diskussion: Es gilt:

- 1.) |a-b| = "Abstand der Zahlen a und b auf der Zahlengeraden" ABB 6 (speziell: |a| = "Abstand von a zum Ursprung")
- $2.) \ a = |a| \cdot sgn(a)$

3.)
$$|a| = |-a|, |ab| = |a| \cdot |b|, \left| \frac{a}{b} \right| = \frac{|a|}{|b|}$$
 (falls $b \neq 0$)

4.) $|a \pm b| \le |a| + |b|$ (*Dreiecksungleichung*) für alle $a, b \in \mathbb{R}$

Lösung von Ungleichung

Bsp. 14: (Ungleichung mit Beträgen)

Gesucht sei die Lösungmenge L der reellen Zahlen, die die Ungleichung $|x-1| < 3 + \frac{1}{2}x$ (*) erfüllen.

- ullet kritische Stelle(n): Nullstellen des Terms innerhalb der Betragsstriche d.h. $x=1 \curvearrowright$ Fallunterscheidung ABB 7
- 1. Fall: x-1 < 0 d.h. x < 1 in (*): $-(x-1) < 3 + \frac{1}{2}x \Leftrightarrow$ $-\frac{3}{2}x < 2 \Leftrightarrow$ $\underbrace{x > -\frac{4}{3}}_{\frown} L_1 = \left\{x | (x < 1) \land x > \frac{4}{3}\right\} = \left(-\frac{4}{3}, 1\right)$
- 2. Fall $x-1 \ge 0$ d.h. $x \ge 1$ in (*): $x-1 < 3 + \frac{1}{2}x \Leftrightarrow \frac{1}{2}x < 4 \Leftrightarrow \frac{\underline{x} < 8}{\frown L_2} = \{x | (x \ge 1) \land (x < 8)\} = [1, 8)$
- $\bullet \Rightarrow L = L_1 \cup L_2 = \underbrace{\left(-\frac{4}{3}, 8\right)}_{}$



Bsp. 15: (Ungleichung mit gebrochen rationalen Termen) $\frac{x}{x+1} < 1$ (*)

- kritische Stelle(n): Nenner-Nullstellen, hier: x=-1. ABB 7 (äquiv. mit -1)
- 1. Fall: x < -1 in (*): $\Leftrightarrow x > x + 1 \Leftrightarrow 0 > 1$ (Widerspruch) $L_1 = \emptyset$.
- 2. Fall: x > -1 in (*): $\Leftrightarrow x < x+1 \Leftrightarrow 0 < 1$ (wahre Aussage) $L_2 = \{x|x > -1\} = (-1, \infty)$
- $\bullet \Rightarrow L = L_1 \cup L_2 = \underline{(-1, \infty)}$

Bsp. 16: (quadratische Ungleichungen)

$$x^2 + 3x < 10 \Leftrightarrow$$

$$\left(x + \frac{3}{3}\right)^2 - \frac{9}{4} < 10 \Leftrightarrow \text{(vereinfacht durch quadratische Ergännzung)}$$

$$\left(x + \frac{3}{2}\right)^2 < \frac{49}{4} \Leftrightarrow$$

$$\left|x + \frac{3}{2}\right| < \frac{7}{2} \Leftrightarrow \text{(Äquivalenz siehe Übung)}$$

$$\frac{-7}{2} < x + \frac{3}{2} < \frac{7}{2} \Leftrightarrow$$

$$-5 < x < 2$$

$$L = (-5, 2)$$