

Vorlesungsmitschriften

# Wintersemester 20/21

Universität Bonn \*  
Semester 1

Fabian

30.01.2022

Online aufrufbar auf

<https://github.com/git-fabus/lecutres/blob/main/notes.pdf>.

Veränderungsvorschläge und Verbesserungen bitte an [git-fabus@uni-bonn.de](mailto:git-fabus@uni-bonn.de).

\*Letzte Änderung: (None)

Commit: (None)

# Vorwort

Dies sind meine persönlichen Vorlesungsmitschriften. Das Dokument ist in 3 große Teilbereiche unterteilt, dazu zählen:

- Algorithmische Mathematik 1
- Analysis 1
- Lineare Algebra 1

Das Dokument wird stetig verbessert und es werden nach und nach immer mehr Aufgaben, Lösungen etc. erscheinen.

# Inhaltsverzeichnis

<b>I. Algorithmische Mathematik</b>	<b>1</b>
<b>Zahlendarstellung am Computer</b>	<b>3</b>
1. Zahlensystem . . . . .	3
2. Vorzeichen-Betrag-Darstellung . . . . .	5
3. Komplementdarstellung . . . . .	6
4. Fest-Komma-Darstellung . . . . .	10
5. Gleitkommadarstellung . . . . .	10
<b>Fehleranalyse</b>	<b>13</b>
6. Rechnerarithmetik . . . . .	13
7. Vorwärts- und Rückwärtsanalyse . . . . .	16
8. Kondition und Stabilität . . . . .	17
<b>Dreitermrekursion</b>	<b>20</b>
9. Theoretische Grundlagen . . . . .	20
10. Miller-Algorithmus . . . . .	24
<b>Sortieren</b>	<b>26</b>
11. Mergesort . . . . .	28
12. Quicksort . . . . .	30
13. Untere Schranke für das Sortierproblem . . . . .	32
<b>Graphen</b>	<b>35</b>
14. Grundlagen . . . . .	35
15. Zusammenhang . . . . .	37
16. Zyklische Graphen . . . . .	39

**Teil I**

# **Algorithmische Mathematik 1**

---

## Einführung

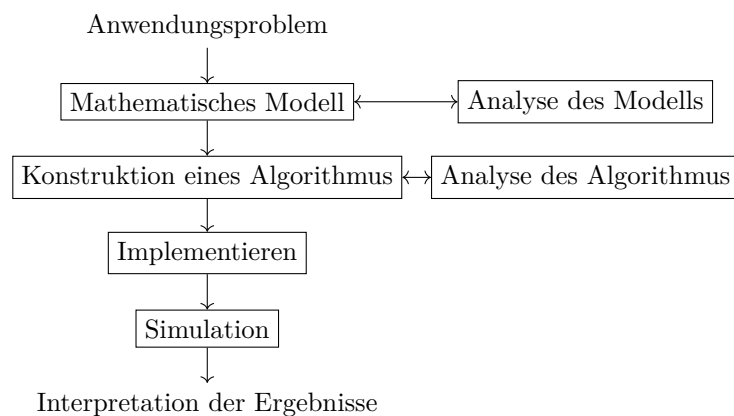
### Algorithmische Mathematik -Was ist das?

Gegenstand der Alma ist die Konstruktion und Analyse effizienter Algorithmen zur Lösung mathematischer Problemstellungen 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.

Diese Probleme können wie folgt aussehen:



**0.1 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

# Zahlendarstellung am Computer

## 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.

**1.1 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*

**1.2 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:

**1.3 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:*  $z < b$  hat die eindeutige Darstellung  $z_0 = z$  und  $z_i = 0$  sonst.

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

Betrachte

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

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

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

Bemerke  $z_{n-1} = 0$ , da

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

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:

#### 1.4 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. 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



### 3. Komplementdarstellung

**3.6 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

**3.7 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)

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

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

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

**3.10 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:

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

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

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

□

### 3.11 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

$$K_b(z) = b^n - z. \quad (3.1)$$

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.

### 3.12 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:

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

### 3.13 Beispiel. Der darstellbare Bereich.

- Sei  $b = 10, n = 2$ .  
Dann impliziert (3.1), 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  $b^n$ ).

**3.14 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  $x \geq 0, y < 0$ :

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

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

□

**3.15 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.

- $y \neq 0$ : (3.1) 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:

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

□

**3.16 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.

- 3.17 Beispiel.** In n-stelliger Binärarithmetik 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.

## 4. Fest-Komma-Darstellung

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

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

- 4.19 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.

- 4.20 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.

- 4.21 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.

- 4.22** Der relative Fehler dieses Runden ist:

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

Analog für  $z_2$ .

Solche Fehler machen jede Rechnung unbrauchbar.

## 5. Gleitkommadarstellung

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

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

mit einer Mantisse 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.

**5.24 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!

**5.25 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

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

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

**5.26 Bemerkung.** Das Hidden Bit ist hilfreich

- Für  $b = 2$  kann die erste Ziffer  $m_1$  der Matisse 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.

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

- 52 Bit für die Matisse 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.

**5.28 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  $z \in \mathbb{R}$ .

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

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

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

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

$$\begin{aligned} \varepsilon &= \frac{|x - \tilde{x}|}{|x|} \\ &= \frac{|x_1.x_2x_3 \dots x_t x_{t+1} \dots \cdot b^e - x_1x_2x_3 \dots x_t \cdot b^e|}{|x|} \\ &= \frac{|0.x_{t+1} \dots| \cdot |b^{e+1-t}|}{|x|} \end{aligned}$$

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  $\varepsilon_{mach}$  höchsten halb so groß ist wie der Abschneidefehler folgt die Behauptung.  $\square$

Die Maschinengenauigkeit  $\varepsilon_{mach}$  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.

# Fehleranalyse

## 6. Rechnerarithmetik

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

$$\begin{aligned}x &= 2.5684 \cdot 10^0 = 2.56840000 \\y &= 3.2791 \cdot 10^{-3} = 0.0032791\end{aligned}$$

Es gilt:

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

**6.2 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  $\circ \in \{+, -, \cdot, \div\}$  durch  $\boxdot \in \{\boxplus, \boxminus, \boxtimes, \boxdiv\}$  definiert durch:

$$x \boxdot y := rd(x \circ y) \tag{6.1}$$

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

**6.3 Bemerkung.** Für  $|x|, |y|, |x \circ y| \in [z_{min}, z_{max}]$  impliziert (6.1), dass

$$\frac{|x \boxdot y - x \circ y|}{|x \circ y|} = \frac{rd(x \circ y - x \circ y)}{|x \circ y|} \leq \varepsilon_{mach}$$

Das bedeutet, dass  $\boxdot$  im Computer bestmöglich umgesetzt ist.

**6.4 Beispiel.** Betrachte  $F = F(10, 5, -4, 5)$

- Setze:
  - $a = 0.98765$



- $b = 0.012424$
- $c = -0.0065432$

Dann gilt:

$$(a + b) + c = a + (b + c) = 0.9925208$$

Numerisch gilt:

$$(0.98765 \boxplus 0.012424) \boxminus 0.0065432 = rd(0.9935568) = 0.99356$$

und

$$0.98765 \boxplus (0.012424 \boxminus 0.0065432) = rd(0.9935308) = 0.99353$$

- Setze
  - $a = 4.2832$
  - $b = -4.2821$
  - $c = 5.7632$

Dann gilt

$$(a + b) \cdot c = a \cdot c + b \cdot c = 0.006339520000001$$

Numerisch gilt:

- 6.5 Bemerkung.** Mathematisch äquivalente Algorithmen auf Fließkommazahlen können je nach Implementierung zu wesentlich unterschiedlichen Ergebnissen führen, selbst wenn die Eingangszahlen exakt dargestellt werden.

## 6.1. Auslöschung

Unglücklicherweise pflanzen sich numerische Fehler, zum Beispiel durch Rundung im Verlauf eines Algorithmus fort.

- 6.6 Lemma** (Fehlerfortpflanzung). Es seien  $x, y \in \mathbb{R}$  mit Datenfehlern  $\Delta x, \Delta y \in \mathbb{R}$  behaftet, welche

$$\left| \frac{\Delta x}{x} \right|, \left| \frac{\Delta y}{y} \right| \ll 1$$

erfüllen. Für  $\circ \in \{+, -, \cdot, \div\}$  gilt für den fortgepflanzten Fehler

$$\Delta(x \circ y) := (x + \Delta x) \circ (y + \Delta y) - x \circ y,$$

dass

$$\begin{aligned} \frac{\Delta(x \pm y)}{x \pm y} &= \frac{x}{x \pm y} \frac{\Delta x}{x} \pm \frac{y}{x \pm y} \frac{\Delta y}{y} \\ \frac{\Delta(x \cdot y)}{xy} &\approx \frac{\Delta x}{x} + \frac{\Delta y}{y} \\ \frac{\Delta(\frac{x}{y})}{\frac{x}{y}} &\approx \frac{\Delta x}{x} - \frac{\Delta y}{y} \end{aligned}$$

Hierbei bedeutet " $\approx$ ", dass Terme mit  $(\Delta x)^2, (\Delta y)^2, \Delta x \Delta y$  vernachlässigt werden.

*Beweis.* Übung. □

Für  $\cdot, \div$  addieren bzw. subtrahieren sich die Fehler.

**Achtung:** Ist  $|x \pm y|$  wesentlich kleiner als  $|x|$  oder  $|y|$ , kann der relative Fehler massiv verstärkt werden. Dieses Phänomen heißt Auslöschung.  
Bei der Konstruktion von Algorithmen sollte Auslöschung möglichst vermieden werden.

### 6.7 Beispiel. Betrachte

$$x^2 - 2px + q = 0 \tag{6.2}$$

mit Lösungen:

$$x_{1,2} = p \pm \sqrt{p^2 - q}$$

**Daten:**  $p, q \in \mathbb{R}$  so ,dass (6.2) lösbar ist

**Ergebnis:** Nullstellen  $x_1, x_2$  von (6.2).

$$d = \sqrt{p \cdot p - q}$$

$$x_1 = p + d$$

$$x_2 = p - d$$

**Algorithmus 2:** naive Nullstellenberechnung

Mit  $p = 100, q = 1$  in dreistelliger, dezimaler Gleitkommaarithmetik ergibt sich:

$$d = \sqrt{10000 - 1} = \sqrt{rd(9999)} = \sqrt{10000} = 100$$

$$x_1 = 100 + 100 = 200$$

$$x_2 = 100 - 100 = 0$$

Die exakten Werte sind  $x_1 \approx 199.99, x_2 \approx 0.00500$ , welche in dreistelliger, dezimaler Gleitkommaarithmetik als  $x_1 = 200, x_2 = 0$  dargestellt werden.

Gemäß 5.30 ist die Maschinengenauigkeit:

$$\varepsilon_{mach} = \frac{1}{2}b^{1-t} = \frac{1}{2}10^{1-3} = 0.005$$

Der relative Fehler  $\frac{|0-0.005|}{|0.005|} = 1$  ist also zu 100 Prozent falsch.

Wir können die Genauigkeit von  $x_2$  erhöhen, indem wir den Wurzelsatz von Vieta  $x_1 x_2 = q$  verwenden.

**Daten:**  $p, q \in \mathbb{R}$  so ,dass (6.2) lösbar ist.

**Ergebnis:** Nullstellen  $x_1, x_2$  (6.2).

$$d = \sqrt{p \cdot p - q}$$

**wenn  $q \geq 0$  dann**

$$x_1 = p + d$$

$$x_2 = p - d$$

**Ende**

**Ende**

$$x_2 = \frac{q}{x_1}$$

**Algorithmus 3:** verbesserte Nullstellenberechnung

Wir erhalten nun  $d = 100, x_1 = 200, x_2 = \frac{1}{200} = 0.005$

## 7. Vorwärts- und Rückwärtsanalyse

Abstrakt gesehen entspricht das Lösen eines Problems dem Auswerten einer Funktion  $f$ . Beim numerischen Auswerten von  $f$  können verschiedene Fehler passieren:

**7.8 Definition** (Fehlerarten). Sei  $D \subset \mathbb{R}$  eine Menge von Eingangsdaten und  $W \subset \mathbb{R}$  die Menge der möglichen Ergebnisse. Wir unterscheiden folgende Fehlerarten bei der Auswertung von  $f: D \rightarrow W$ .

- Datenfehler Typischerweise sind die Eingangsdaten  $x \in D$  nicht exakt, sondern mit einem Datenfehler  $\Delta x$  behaftet. Die gestörten Eingangsdaten  $\tilde{x} = x + \Delta x$  produzieren einen Fehler  $f(x + \Delta x) - f(x)$ .
- Verfahrensfehler Exakte Verfahren enden bei exakter Rechnung nach endlich vielen Operationen mit dem exakten Ergebnis. Näherungsverfahren enden in Abhängigkeit bestimmter Kriterien mit einer Näherung  $\tilde{y}$  für die Lösung  $y \in W$ .
- Rundungsfehler Verursacht durch Maschinenzahlen und Rundungen während der Arithmetik.

Ein Teilgebiet der Numerik versucht diese Fehler durch eine Fehleranalyse zu quantifizieren. Hierzu verfolgt man die Auswirkungen von allen Fehlern, die in den einzelnen Schritten vorkommen können.

**7.9 Definition** (Analysearten). Bei der *Vorwärtsanalyse* wird der Fehler von Schritt zu Schritt verfolgt und der akkumulierte Fehler für jedes Teil-Ergebnis abgeschätzt. Bei der *Rückwärtsanalyse* geschieht die Verfolgung des Fehlers hingegen so, dass jedes Zwischenergebnis als exakt berechnetes Ergebnis zu gestörten Daten interpretiert wird, d.h., der akkumulierte Fehler wird als Datenfehler interpretiert.

**7.10 Beispiel.** Betrachte  $f(x, y) = x + y$

- Vorwärtsanalyse:  $\boxed{f} = x \boxplus y = (x + y)(1 + \varepsilon)$
- Rückwärtsanalyse:  $x \boxplus y = x(1 + \varepsilon) + y(1 + \varepsilon) = f(x(1 + \varepsilon), y(1 + \varepsilon))$

mit  $|\varepsilon| \leq \varepsilon_{mach}$

In der Praxis ist die Vorwärtsanalyse kaum durchführbar. Für die meisten Algorithmen ist, wenn überhaupt, nur eine Rückwärtsanalyse bekannt.

**7.11 Beispiel.** Fortsetzung von Beispiel 6.7

Führe Rückwärtsanalyse durch zu Algorithmus 2 für die Betrags-kleinere Nullstelle  $x_2$  durch: Dann ist

$$f(p, q) = p - \sqrt[2]{p^2 - q}$$

mit  $p > 0$ . Für  $|\varepsilon_i| \leq \varepsilon_{mach}, i = 1, 2, 3, 4$  betrachten wir:

$$\begin{aligned} \left( p - \sqrt[2]{(p^2(1+\varepsilon_1) - q)(1+\varepsilon_2)(1+\varepsilon_3)} \right) (1+\varepsilon_4) &= p(1+\varepsilon_4) - \sqrt[2]{(p^2(1+\varepsilon_1) - q)(1+\varepsilon_2)(1+\varepsilon_3)^2(1+\varepsilon_4)^2} \\ &= p^2(1+\varepsilon_1)(1+\varepsilon_3)^2(1+\varepsilon_4)^2 - q(1+\varepsilon_2)(1+\varepsilon_3)^2(1+\varepsilon_4)^2 \\ &= \dots \\ &= p(1+\varepsilon_4) - \sqrt[2]{p^2(1+\varepsilon_4)^2 - q(1+\varepsilon_7)} \\ &= f(p(1+\varepsilon_4), q(1+\varepsilon_7)) \end{aligned}$$

Die Abschätzung für  $\varepsilon_7$  explodiert, falls  $0 < |q| \ll 1 < p$ . Dies war in Beispiel 6.7 der Fall.

## 8. Kondition und Stabilität

Gegeben sei eine stetige und differenzierbare Funktion

$$f: \mathbb{R} \rightarrow \mathbb{R}, x \mapsto y = f(x)$$

Für fehlerhafte Daten  $x + \Delta x$  mit kleinen Fehler  $\Delta x$  gilt:

$$\frac{f(x + \Delta x) - f(x)}{\Delta x} \approx f'(x).$$

Für den absoluten Datenfehler  $\Delta y$  gilt:

$$\Delta y = f(x + \Delta x) - f(x) \approx f'(x) \cdot \Delta x$$

und für den relativen Fehler:

$$\frac{\Delta y}{y} = \frac{f'(x) \cdot \Delta x}{f(x)} = \frac{f'(x)x}{f(x)} \cdot \frac{\Delta x}{x}$$

**8.12 Definition** (Konditionszahlen). Die Zahl

$$K_{abs} = |f'(x)|$$

heißt *absolute Konditionszahl* des Problems  $x \mapsto f(x)$ . Für  $f(x) \cdot x \neq 0$  heißt

$$K_{rel} = \left| \frac{f'(x) \cdot x}{f(x)} \right|$$

die entsprechende *relative Konditionszahl*. Ein Problem heißt schlecht konditioniert, falls eine dieser beiden Konditionszahlen deutlich größer als 1 sind. Ansonsten heißt das Problem gut konditioniert.

**8.13 Beispiel.** Konditionierung der Addition und Multiplikation

- Für die Addition  $f(x) = x + a$  gilt

$$K_{rel} = \left| \frac{f'(x)x}{f(x)} \right| = \left| \frac{x}{x+1} \right|$$

$\implies K_{rel}$  ist groß, falls  $|x+a| \ll |x|$ .

- Für die Multiplikation  $f(x) = x \cdot a$  gilt

$$K_{rel} = \frac{|f'(x)x|}{|f(x)|} = \frac{|ax|}{|ax|} = 1$$

$\implies$  Die absolute Kondition ist schlecht, falls  $1 \ll q$ . Die relative Kondition ist immer gut.

**8.14 Definition.** Erfüllt die Implementierung eines Algorithmus  $\boxed{f}$  zur Lösung eines Problems  $x \mapsto f(x)$  die Abschätzung

$$\left| \frac{\boxed{f} - f(x)}{f(x)} \right| \leq C_V K_{rel} \varepsilon_{mach}$$

mit einem mäßig großen  $C_V > 0$ , so wird der Algorithmus *vorwärtsstabil* genannt. Ergibt die Rückwärtsanalyse  $\boxed{f}(x) = f(x + \Delta x)$  mit

$$\left| \frac{\Delta x}{x} \right| \leq C_R \varepsilon_{mach}$$

mit  $C_R > 0$  nicht zu groß, so heißt der Algorithmus  $\boxed{f}$  *rückwärtsstabil*

**8.15 Satz.** Jeder rückwärtsstabile Algorithmus ist auch vorwärtsstabil

*Beweis.* Übung. □

Oft ist Rückwärtsstabilität einfacher nachzuweisen.

**Faustregel** zu konditionierten Probleme

- Gut konditioniertes Problem + stabiler Algorithmus  
 $\implies$  Gute numerische Resultate.
- Schlecht konditioniertes Problem oder instabiler Algorithmus  
 $\implies$  Fragwürdige Ergebnisse.

**8.16 Beispiel.** Fortsetzung von Beispiel 7.11

Die Rückwärtsanalyse hat gezeigt: Falls  $0 < |q| \ll 1 < p$  wird der numerische Fehler untragbar. Unsere Abbildung ist:

$$f(q) = p - \sqrt[p]{p^2 - 1}$$

Als Konditionszahlen ergeben sich:

$$K_{abs} = |f'(q)| = \left| \frac{1}{2\sqrt[2]{p^2 - q}} \right| < 1$$

$$\begin{aligned} K_{rel} &= \left| \frac{f'(q)q}{f(q)} \right| \\ &= \left| \frac{q}{2\sqrt[2]{p^2 - q}(p - \sqrt[2]{p^2 - q})(p + \sqrt[2]{p^2 - q})} \right| \\ &= \frac{1}{2} \left| \frac{p + \sqrt[2]{p^2 - q}}{\sqrt[2]{p^2 - q}} \right| \\ &\approx \frac{1}{2} \left| \frac{p + p}{p} \right| \approx 1 \end{aligned}$$

$\Rightarrow$  Nullstellenberechnung ist ein gut konditioniertes Problem, aber Algorithmus 2 muss instabil sein

# Dreitermrekursion

## 9. Theoretische Grundlagen

**9.1 Definition.** Für gegebene  $p_0$  und  $p_1$  heißt eine Rekursion der Form

$$p_k = a_k p_{k-1} + b_k p_{k-2} + c_k, \quad k = 2, 3, \dots \quad (9.1)$$

mit  $b_k \neq 0$  eine *Dreitermrekursion*. Die zugehörige *Rückwärtsrekursion* ist

$$p_{k-2} = -\frac{a_k}{b_k} p_{k-1} + \frac{1}{b_k} p_k - \frac{c_k}{b_k}, \quad k = n, n-1, \dots \quad (9.2)$$

Ist  $b_k = 1$ , das heißt (9.1) und (9.2) gehen durch vertauschen von  $p_{k-2}$  und  $p_k$  auseinander hervor, so heißt die Rekursion symmetrisch. Gilt  $c_k = 0$  für alle  $k$ , so heißt die Rekursion homogen.

**9.2 Beispiel.** Die Fibonacci-Zahlen sind rekursiv definiert durch:

$$p_k = p_{k-1} + p_{k-2}$$

mit  $p_0 = 0, p_1 = 1$ . Es gilt  $a_k = b_k = 1$ .

Sei

$$p_k(x) = 2 \cos(x) p_{k-1}(x) - 1 \cdot p_{k-2}(x)$$

mit  $p_0 = 1$  und  $p_1 = \cos(x)$ . Man kann zeigen, dass

$$p_k(x) = 2 \cos(kx)$$

ist.

Die Chebychev-Polynome  $T_k$  erfüllen

$$T_k(x) = 2xT_{k-1} - T_{k-2}(x)$$

mit  $T_0(x) = 1$  und  $T_1(x) = x$

**Daten:** Koeffizienten von  $\{a_k\}_{k=2}^n, \{b_k\}_{k=2}^n, \{c_k\}_{k=2}^n$ , Startwerte  $p_0, p_1$   
**Ergebnis:** Werte von  $\{p_k\}_{k=2}^n$   
**für**  $k \leftarrow 2$  **bis**  $n$  **do**  
   $p_k = a_k p_{k-1} + b_k p_{k-2} + c_k$   
**Ende**

**Algorithmus 4:** Dreitermrekursion

**9.3 Beispiel.** Betrachte:  $p_0 = 1, p_1 = \sqrt[3]{2} - 1, p_k = -2p_{k-1} + p_{k-2}, k = 2, 3, \dots$

Algorithmus 5 liefert:

Die Oszillationen für höhere  $k$  sollten uns misstrauische machen und uns motivieren das Problem genauer anzuschauen.

Wir können homogene Dreitermrekursion schreiben als:

$$\begin{bmatrix} p_k \\ p_{k-1} \end{bmatrix} = \begin{bmatrix} a_k p_{k-1} + b_k p_{k-2} \\ p_{k-1} \end{bmatrix} =: \begin{bmatrix} a_k & b_k \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} p_{k-1} \\ p_{k-2} \end{bmatrix} = A_k \begin{bmatrix} p_{k-1} \\ p_{k-2} \end{bmatrix}$$

Rekursiv folgt:

$$\begin{bmatrix} p_k \\ p_{k-1} \end{bmatrix} = A_k \begin{bmatrix} p_{k-1} \\ p_{k-2} \end{bmatrix} = A_k A_{k-1} \begin{bmatrix} p_{k-2} \\ p_{k-3} \end{bmatrix} = \dots = A_k \dots A_2 \begin{bmatrix} p_1 \\ p_0 \end{bmatrix}$$

Offensichtlich gilt für alle  $\alpha, \beta \in \mathbb{R}$  und Startwerte  $p_0, p_1, q_0, q_1$  und  $B_k := A_k \dots A_2$ , dass

$$B_k \begin{bmatrix} \alpha p_1 + \beta q_1 \\ \alpha p_0 + \beta q_0 \end{bmatrix} = \alpha B_k \begin{bmatrix} p_1 \\ p_0 \end{bmatrix} + \beta B_k \begin{bmatrix} q_1 \\ q_0 \end{bmatrix} = \alpha \begin{bmatrix} p_k \\ p_{k-1} \end{bmatrix} + \beta \begin{bmatrix} q_k \\ q_{k-1} \end{bmatrix}$$

Das heißt, die Lösungsfolge  $\{p_k\}$  hängt *linear* von den Startwerten  $\begin{bmatrix} p_1 \\ p_0 \end{bmatrix} \in \mathbb{R}^2$  ab. Im Falle unseres Beispiels gilt:  $a_k = a = -2, b_k = b = 1$ . Das bedeutet, es gilt:

$$B_k = A^{k-1}, A = \begin{bmatrix} a & b \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} -2 & 1 \\ 1 & 0 \end{bmatrix}$$

Aus den Eigenwerten  $\lambda_1, \lambda_2$  der Matrix  $A$  kann man die Lösungen der homogenen Dreitermrekursion direkt angeben.

**9.4 Satz.** Seien  $\lambda_1, \lambda_2$  die Nullstellen des *charakteristischen Polynoms*

$$q(\lambda) = \lambda^2 - a\lambda - b \tag{9.3}$$

Dann ist die Lösung der homogenen Dreitermrekursion:

$$p_k = ap_{k-1} + bp_{k-2}, k = 2, 3, \dots$$

gegeben durch:

$$p_k = \alpha \lambda_1^k + \beta \lambda_2^k, k = 2, 3, \dots$$

mit  $\alpha, \beta \in \mathbb{R}$  Lösungen des linearen Gleichungssystems.

$$\alpha + \beta = p_0$$

$$\alpha \lambda_1 + \beta \lambda_2 = p_1$$

*Beweis.* • Induktion über  $k$

$$\diamond \text{ Induktionsanfang: } \underline{k=0}: p_0 = \alpha + \beta = \alpha \lambda_1^0 + \beta \lambda_2^0$$



◇ Induktionsschritt:  $k \rightarrow k+1$ :

$$\begin{aligned}
 p_{k+1} &= ap_k + bp_{k-1} \\
 &= a(\alpha\lambda_1^k + \beta\lambda_2^k) + b(\alpha\lambda_1^{k-1} + \beta\lambda_2^{k-1}) \\
 &= \alpha(a\lambda_1^k + b\lambda_1^{k-1}) + \beta(a\lambda_2^k + b\lambda_2^{k-1}) \\
 &= \alpha\lambda_1^{k-1}(a\lambda_1 + b) + \beta\lambda_2^{k-1}(a\lambda_2 + b) \\
 &= \alpha\lambda_1^{k+1} + \beta\lambda_2^{k+1}
 \end{aligned}$$

Da nach (9.3)  $\lambda_1^2$  und  $\lambda_2^2$  Nullstellen sind.

□

Wende 9.4 auf das Beispiel 9.3 an: Mit  $a = -2, b = 1$  hat das charakteristische Polynom (9.3) die Nullstellen:

$$\lambda_1 = \sqrt[2]{2} - 1, \lambda_2 = -\sqrt[2]{2} - 1$$

Aus

$$\alpha + \beta = p_0 = 1$$

$$\alpha\lambda_1 + \beta\lambda_2 = p_1 = \sqrt[2]{2} - 1$$

folgt  $\alpha = 1, \beta = 0$ . Der Satz 9.4 impliziert:

$$p_k = \lambda_1^k > 0, k = 0, 1, 2, \dots$$

in Beispiel 9.3

**Daten:** Koeffizienten a,b, Startwerte  $p_0, p_1$

**Ergebnis:** Werte  $\{p_k\}_{k=2}^n$

Initialisiere:

- $\lambda_1 = \frac{a}{2} + \sqrt{\frac{a^2}{4} + b}$
- $\lambda_2 = \frac{a}{2} - \sqrt{\frac{a^2}{4} + b}$
- $\beta = \frac{p_1 - \lambda_1 p_0}{\lambda_2 - \lambda_1}$
- $\alpha = p_0 - \beta$

**für**  $k \leftarrow 2$  **bis**  $n$  **tue**

  |  $p_k = \alpha\lambda_1^k + \beta\lambda_2^k$

**Ende**

**Algorithmus 5:** verbesserte Dreitermrekursion

Dieser Algorithmus liefert für Beispiel 9.3 folgende Grafik:

Dies ist in diesem Fall das richtige Ergebnis, da

$$p_k = \lambda_1^k = (\sqrt[2]{2} - 1)^k \ll 1$$

Was geht schief in dem vorherigen Algorithmus 5?

Betrachten wir unser Problem  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(p_0, p_1) = p_k$  mit gestörten Eingangsdaten:  $\hat{p}_0 = 1(1 + \varepsilon_0)$ ,  $\hat{p}_1 = \lambda_1(1 + \varepsilon_1)$ ,  $|\varepsilon_0|, |\varepsilon_1| \leq \varepsilon_{mach}$ .

Dies ergibt:

$$\begin{aligned}\tilde{\beta} &= \frac{\lambda_1(1 + \varepsilon_1) - \lambda_1(1 + \varepsilon_0)}{\lambda_2 - \lambda_1} = (\varepsilon_1 - \varepsilon_0) \frac{\lambda_1}{\lambda_2 - \lambda_1} \\ \tilde{\alpha} &= 1 + \varepsilon_0 \frac{\lambda_2 - \lambda_1}{\lambda_2 - \lambda_1} - (\varepsilon_1 - \varepsilon_0) \frac{\lambda_1}{\lambda_2 - \lambda_1} = 1 + \varepsilon_0 \frac{\lambda_2}{\lambda_2 - \lambda_1} - \varepsilon_1 \frac{\lambda_1}{\lambda_2 - \lambda_1}\end{aligned}$$

und somit

$$\begin{aligned}\tilde{p}_k &= \tilde{\alpha} \lambda_1^k + \tilde{\beta} \lambda_2^k \\ &= (1 + \varepsilon_0 \frac{\lambda_2}{\lambda_2 - \lambda_1} - \varepsilon_1 \frac{\lambda_1}{\lambda_2 - \lambda_1}) \lambda_1^k + (\varepsilon_1 - \varepsilon_0) \frac{\lambda_1}{\lambda_2 - \lambda_1} \lambda_2^k\end{aligned}$$

Der relative Fehler von  $\tilde{p}_k$  zu  $p_k = \lambda_1$  ist somit :

$$\begin{aligned}|\frac{\tilde{p}_k - p_k}{p_k}| &= |\varepsilon_0 \frac{\lambda_2}{\lambda_2 - \lambda_1} - \varepsilon_1 \frac{\lambda_1}{\lambda_2 - \lambda_1} + (\varepsilon_1 - \varepsilon_0) \frac{\lambda_1}{\lambda_2 - \lambda_1} (\frac{\lambda_2}{\lambda_1})^k| \\ &= |\varepsilon_0 + (\varepsilon_1 - \varepsilon_0) \frac{\lambda_1}{\lambda_2 - \lambda_1} \left( \left( \frac{\lambda_2}{\lambda_1} \right)^k - 1 \right)|\end{aligned}$$

**Beobachte:** Falls  $\frac{\lambda_2}{\lambda_1} > 1$  explodiert der relative Fehler für wachsende k. Dies war in Beispiel 9.3 der Fall.

Falls wir den Begriff der Konditionszahl aus der Definition 8.12 geeignet für Funktion  $\mathbb{R}^2 \rightarrow \mathbb{R}$  erweitern, sehen wir, dass das Problem schlecht konditioniert ist.

**9.5 Definition** (Minimallösung). Die Lösung  $\{p_k\}$  der Dreiterm-Rekursion (9.1) zu den Startwerten  $p_0$  und  $p_1$  heißt *Minimallösung*, falls für jede Lösung  $\{p_k\}$  zu den von  $p_0, p_1$  linear unabhängigen Startwerten  $q_0, q_1$  gilt, dass:

$$\lim_{k \rightarrow \infty} \frac{p_k}{q_k} = 0$$

Die Lösung  $\{p_k\}$  wird dominante Lösung genannt.

**9.6 Beispiel.** In Beispiel 9.3 ist  $\{p_k\}$  Minimallösung genau dann, wenn  $\beta = 0$ .

Die Minimallösung ist nur bis auf einen skalaren Faktor eindeutig. Deshalb normieren wir sie mit der zusätzlichen Bedingung:  $p_0^2 + p_1^2 = 1$ .

**Frage:** Wie können wir eine normierte Minimallösung berechnen, obwohl das Problem schlecht konditioniert ist?

## 10. Miller-Algorithmus

Betrachte die Rückwärtsrekursion zu

$$q_k = a_k q_{k-1} + b_k q_{k-2}, k = 2, 3, \dots$$

d.h.

$$q_{k-2} = -\frac{a_k}{b_k} q_{k-1} + \frac{1}{b_k} q_k, k = n, n-1, \dots, 2$$

mit Startwerten  $q_n = 0, q_{n-1} = 1$ . Für  $a_k = a, b_k = b$  besitzt das charakteristische Polynom

$$q(\mu) = \mu^2 + \frac{a}{b}\mu - \frac{1}{b}$$

Die Nullstellen:

$$\mu_{1,2} = \frac{-a \pm \sqrt{a^2 + 4b}}{2b} = \frac{1}{\lambda_{1,2}}$$

**Idee:** Wende Satz 9.4 "rückwärts" an.

Lösen wir

$$\alpha + \beta = 0$$

und

$$\alpha \mu_1 + \beta \mu_2 = q_{n-1} = 1$$

ergibt sich:

$$\alpha = \frac{\lambda_1 \lambda_2}{\lambda_2 - \lambda_1} \neq 0$$

$$\beta = -\frac{\lambda_1 \lambda_2}{\lambda_2 - \lambda_1} \neq 0$$

und

$$q_k = \alpha \mu_1^{n-k} + \beta \mu_2^{n-k} = \frac{\alpha}{\lambda_1^{n-k}} + \frac{\beta}{\lambda_2^{n-k}}$$

Für  $|\lambda_1| < |\lambda_2|$  folgt, dass

$$\begin{aligned} p_k^{(n)} &:= \frac{q_k}{q_0} \\ &= \frac{\frac{\alpha}{\lambda_1^{n-k}} + \frac{\beta}{\lambda_2^{n-k}}}{\frac{\alpha}{\lambda_1^n} + \frac{\beta}{\lambda_2^n}} \\ &= \frac{\lambda_1^k + \frac{\beta}{\alpha} \lambda_2^k \left(\frac{\lambda_1}{\lambda_2}\right)^n}{1 + \frac{\beta}{\alpha} \left(\frac{\lambda_1}{\lambda_2}\right)^n} \\ &\rightarrow \lambda_1^k = p_k \end{aligned}$$

**Beobachtung:** Für großes  $n$  approximiert  $p_k^{(n)}$  die Minimallösung.

**Daten:**  $n$  genügend groß  $\{a_k\}_{k=2}^n, \{b_k\}_{k=2}^n$   
**Ergebnis:** Approximation von  $\{p_k^{(n)}\}_{k=0}^n$  um eine normierte Minimallösung  
 Setze  $\hat{p}_n = 0, \hat{p}_{n-1} = 1$   
**für**  $k \leftarrow n$  **bis** 2 **tue**  
      $\hat{p}_{k-2} = -\frac{a_k}{b_k} \hat{p}_{k-1} + \frac{1}{b_k} \hat{p}_k$   
**Ende**  
**für**  $k \leftarrow 0$  **bis**  $n$  **tue**  
      $p_k^{(n)} = \frac{\hat{p}_n}{\sqrt{\hat{p}_0^2 + \hat{p}_n^2}}$   
**Ende**

**Algorithmus 6:** Miller-Algorithmus

**10.7 Satz** (Miller-Algorithmus). Sei  $\{p_k\}_{k=0}^\infty$  Minimallösung der normierten, homogenen Dreitermrekursion. Dann gilt für die Lösung des Miller-Algorithmus:

$$\lim_{n \rightarrow \infty} p_k^{(n)} = p_k, k = 0, 1, 2, \dots$$

*Beweis.* Wir bemerken, dass sich jede Lösung der Dreitermrekursion als Linearkombination einer Minimallösung und einer dominanten Lösung schreiben lässt (Übung). Mit den richtigen  $\alpha$  und  $\beta$  folgt (Übung):

$$\begin{aligned} \hat{p}_k &= \alpha p_k + \beta q_k \\ &= \frac{p_k q_n - q_k p_n}{p_{n-1} q_n - q_{n-1} p_n} = \frac{q_n}{p_{n-1} q_n - q_{n-1} p_n} \left( p_k - \frac{p_n}{q_n} q_k \right) \end{aligned}$$

Aus (Normierungs-Konstante)

$$\begin{aligned} \hat{p}_0^2 + \hat{p}_1^2 &= \frac{q_n^2}{(p_{n-1} q_n - q_{n-1} p_n)^2} \left( p_0^2 + p_1^2 - 2 \frac{p_n}{q_n} (p_0 q_0 + p_1 q_1) + \frac{p_n^2}{q_n^2} (q_0^2 + q_1^2) \right) \\ &= \frac{q_n^2}{(p_{n-1} q_n - q_{n-1} p_n)^2} \left( 1 - 2 \frac{p_n}{q_n} + \frac{p_n^2}{q_n^2} \right) \end{aligned}$$

folgt:

$$\begin{aligned} p_k^{(n)} &= \frac{\hat{p}_k}{\sqrt{\hat{p}_0^2 + \hat{p}_1^2}} \\ &= \frac{p_k - \frac{p_n}{q_n} q_k}{\sqrt{1 - 2 \frac{p_n}{q_n} + \frac{p_n^2}{q_n^2}}} \rightarrow p_k \end{aligned}$$

□

Wir können also auch für schlecht konditionierte Probleme stabile Algorithmen finden.

# Sortieren

## 10.1. Das Sortierproblem

**Gegeben**  $n \in \mathbb{N}$  verschiedene Zahlen  $z_1, \dots, z_n \in \mathbb{R}$ .

**Gesucht** Permutation  $\pi_1, \dots, \pi_n$ , so dass  $z_{\pi_1} < \dots < z_{\pi_n}$

**10.1 Definition** (Permutation). Eine Permutation  $\pi$  von  $\{1, 2, \dots, n\}$  ist eine bijektive Abbildung von  $\{1, 2, \dots, n\}$  auf sich selbst. Wir schreiben  $\pi(k) = \pi_k$  für  $k = 1, \dots, n$

**10.2 Bemerkung.** Da wir die Zahlen  $z_1, \dots, z_n$  als verschieden annehmen ist das Sortierproblem eindeutig lösbar.

Probiere so lange alle möglichen Permutationen durch, bis die gewünschte Sortierung vorliegt

**Algorithmus 7:** Brute-Force

**10.3 Satz.** Es gibt  $n! = n(n-1) \dots 2 \cdot 1$  Permutationen der Menge  $\{1, 2, \dots, n\}$

*Beweis.* Wir haben

- $n$  Möglichkeiten die erste Zahl auszuwählen
- $n - 1$  Möglichkeiten die zweite Zahl auszuwählen
- ...
- 1 Möglichkeit die letzte Zahl auszuwählen

woraus die Behauptung folgt. □

Im schlimmsten Fall (worst case) muss der Algorithmus 7 also

$$(n-1) \cdot n!$$

Vergleiche durchführen. Da dies extrem aufwändig sein kann, betrachten wir im Folgenden einen Algorithmus, der die transitive Struktur

$$x < y \text{ und } y < z \implies x < z$$

der Ordnungsrelation ausnutzt.

**Daten:** Menge  $S_n = \{z_1, \dots, z_n\}$   
**Ergebnis:** Sortierte Menge  $S^\pi = \{z_{\pi_1}, \dots, z_{\pi_n}\}$   
**für**  $k \leftarrow n$  **bis** 1 **tue**  
     $z_{\pi_k} = \max(S_k)$   
     $S_{k-1} = S_k \setminus \{z_{\pi_k}\}$   
**Ende**

**Algorithmus 8:** Bubblesort

**10.4 Beispiel.** Die zu sortierende Menge ist:  $\{4, 1, 2\}$ .

- $S_3 = \{4, 1, 2\}$ ,  $k = 3$ ,  $z_{\pi_3} = 4$
- $S_2 = \{1, 2\}$ ,  $k = 2$ ,  $z_{\pi_2} = 2$
- $S_1 = \{1\}$ ,  $k = 1$ ,  $z_{\pi_1} = 1$

Die sortierte Liste ist:  $\{1, 2, 4\}$ .

Um den Aufwand von Algorithmen zu untersuchen, beschränken wir uns auf das asymptotische Verhalten für große  $n$ .

**10.5 Definition** (Landau-Notation). Wir schreiben

- $f(x) = \mathcal{O}(g(x))$ , falls Zahlen  $C > 0, x_0 > 0$  existieren, so dass:  $|f(x)| \leq Cg(x), \forall x > x_0$ .
- $f(x) = \Omega(g(x))$ , falls Zahlen  $C > 0, x_0 > 0$  existieren, so dass  $|f(x)| \geq Cg(x) \forall x > x_0$ .
- $f(x) = o(g(x))$ , falls für jedes  $c > 0$  ein  $x_0 > 0$  existiert, so dass  $|f(x)| \leq cg(x), \forall x > x_0$ .
- $f(x) = \Theta(g(x))$ , falls  $f(x) = \mathcal{O}(g(x)), g(x) = \mathcal{O}(f(x))$

**10.6 Bemerkung.** Genau dann wenn Beziehung der Landau-Notation

- $f(x) = \mathcal{O}(g(x)) \iff \limsup_{x \rightarrow \infty} \left| \frac{f(x)}{g(x)} \right| \leq C < \infty$
- $f(x) = \Omega(g(x)) \iff \liminf_{x \rightarrow \infty} \left| \frac{f(x)}{g(x)} \right| > 0$
- $f(x) = o(g(x)) \iff \lim_{x \rightarrow \infty} \left| \frac{f(x)}{g(x)} \right| = 0$
- $f(x) = \Theta(g(x)) \iff f(x) = \mathcal{O}(g(x)), f(x) = \Omega(g(x))$

**10.7 Beispiel.** Es gilt  $\sin(x) = \mathcal{O}(1)$ , da  $|\sin(x)| \leq 1$  für alle  $x \in \mathbb{R}$ . Bei Polynomen gilt die Laufzeit ist die höchste Potenz, sofern  $x \geq 1$ .

**10.8 Definition.** Der Aufwand eines Algorithmus ist die kleinste obere Schranke für das betrachtete Aufwandsmaß

Für uns ist der Speicherbedarf irrelevant und benutzen als Aufwandsmaß für den Rechen die Anzahl der benötigten Vergleiche.

In der  $k$ -ten Iteration des Bubblesort-Algorithmus 8 müssen  $k$  Vergleiche ausgeführt werden. Das heißt, der Gesamtaufwand des Algorithmus ist:

$$\sum_{k=1}^{n-1} k = \frac{n(n-1)}{2} = \mathcal{O}(n^2)$$

Dies ist eine drastische Verbesserung im Vergleich zu  $\mathcal{O}(n(n!))$  von Algorithmus 7

## 11. Mergesort

Zu erst wollen wir folgendes beobachten:

**11.9 Lemma.** Gegeben seien zwei sortierte Mengen

- $S_x = \{x_1 < \dots < x_m\}$
- $S_y = \{y_1 < \dots < y_n\}$

Dann lässt sich die Menge  $S = S_x \cup S_y$  mit linearem Aufwand sortieren. Genauer werden  $m + n - 1$  Vergleiche benötigt.

*Beweis.* Konstruktiv, durch den entsprechenden Algorithmus. □

**Daten:** Sortierte Mengen  $S_x$  und  $S_y$   
**Ergebnis:** Sortierte Menge  $S = S_x \cup S_y$   
 Initialisiere  $i = j = k = 1$   
**solange**  $i \leq m$  *und*  $j \leq n$  **tue**  
     **wenn**  $x_i < y_j$  **dann**  
          $z_k = x_i$   
          $i = i + 1$   
     **Ende**  
     **sonst**  
          $z_k = y_j$   
          $j = j + 1$   
     **Ende**  
      $k = k + 1$   
**Ende**  
**für**  $l \leftarrow 0$  **bis**  $m - i$  **tue**  
      $z_{k+l} = x_{k+l}$   
**Ende**  
**für**  $l \leftarrow 0$  **bis**  $n - j$  **tue**  
      $z_{k+l} = y_{k+l}$   
**Ende**

**Algorithmus 9:** Merge

Basierend auf dieser Beobachtung können wir eine divide-and-conquer Strategie angeben um Mengen der Länge  $n = 2^m, m \in \mathbb{N}$  zu sortieren.

**Daten:** Menge  $S = \{z_1, \dots, z_n\}$   
**Ergebnis:** Sortierte Menge  $S^\pi = \{z_{\pi_1} < \dots < z_{\pi_n}\}$   
**wenn**  $n = 1$  **dann**  
    |  $S^\pi = S$   
**Ende**  
**sonst**  
    |  
        • Sortiere  
           $L = \{z_1, \dots, z_{\frac{n}{2}}\}$   
           $R = \{z_{\frac{n}{2}+1}, \dots, z_n\}$  mittels Mergesort zu  $L^\pi$  und  $R^\pi$   
        • Sortiere  $L^\pi \cup R^\pi$  mittels Merge-Algorithmus 9 zu  $S^\pi$   
**Ende**

**Algorithmus 10:** Mergesort

**11.10 Beispiel.** Mergesort der Menge  $\{20, 7, 84, 31, 71, 42, 18, 10\}$

**11.11 Bemerkung.** Da Mergesort sich selbst aufruft, sprechen wir von einem *rekursiven Algorithmus*. Im Allgemeinen ist es schwierig zu beurteilen, ob solche Algorithmen terminieren. Im Fall von Mergesort ist der Fall jedoch klar, da die Rekursion im Fall  $n = 1$  abgebrochen wird.

**11.12 Satz.** Der Aufwand von Mergesort ist  $\mathcal{O}(n \log n)$

*Beweis.* Bezeichne  $A(n)$  den Aufwand für das Sortieren einer  $n = 2^m$  elementigen Menge mittels Mergesort. Dann gilt:

$$A(1) = 0$$

$$A(n) = n - 1 + 2A\left(\frac{n}{2}\right)$$

Auflösen der Rekursion ergibt:

$$\begin{aligned} A(n) &= n - 1 + 2A\left(\frac{n}{2}\right) \\ &= 2n - 1 - 2 + 4A\left(\frac{n}{4}\right) \\ &= \dots \\ &= mn - \sum_{i=0}^{m-1} 2^i \\ &= mn - \frac{1 - 2^m}{1 - 2} \\ &= (m - 1)n + 1 \end{aligned}$$

$m = \log_2(n) = \frac{\log(n)}{\log(2)}$  impliziert die Behauptung

□



$\log n \ll n$ , also ist die Verbesserung von  $\mathcal{O}(n^2)$  nach  $\mathcal{O}(n \log(n))$  signifikant. Man nennt das Wachstum auch beinahe linear.

- 11.13 Bemerkung.** Die Implementierung von Mergesort als in-place-Algorithmus ist je nach Datenstrukturen trickreich. Deshalb wird oft nicht in-place implementiert, weshalb für jeden "divide"-Schritt zusätzlicher Speicherplatz implementiert werden muss.

## 12. Quicksort

**Idee:** Divide-and-Conquer Strategie basierend auf dem Inhalt der zu sortierenden Liste.

**Daten:** Menge  $S = \{z_1, \dots, z_n\}$   
**Ergebnis:** Sortierte Menge  $S^\pi = \{z_{\pi_1}, \dots, z_{\pi_n}\}$   
 Wähle ein Pivot-Element  $x \in S$   
 Bestimme eine Permutation  $\pi$ , so dass  $x = z_{m_\pi}$   
**wenn**  $L = \{z_{\pi_1}, \dots, z_{\pi_{m-1}}\} \neq \emptyset$  **dann**  
     | Sortiere  $L$  zu  $L^\pi$  mittels Quicksort  
**Ende**  
**wenn**  $R = \{z_{\pi_{m+1}}, \dots, z_{\pi_n}\} \neq \emptyset$  **dann**  
     | Sortiere  $R$  zu  $R^\pi$  mittels Quicksort  
**Ende**  
 Vereinige  $S^\pi = R^\pi \cup \{x\} \cup L^\pi$

**Algorithmus 11:** Quicksort

- 12.14 Beispiel.** Beispiel anhand der gleichen Menge von Mergesort.

- 12.15 Lemma.** Im schlimmsten Fall ist der Aufwand von Quicksort  $\mathcal{O}(n^2)$

*Beweis.*  $A(n)$  aus 11.12 wird umso größer, je unterschiedlicher die Größe der beiden Teilprobleme  $A(m-1)$  und  $A(n-m)$  ist.  $A(n)$  ist also maximal für  $m=1$  oder  $m=n$ , also wenn das Pivot-Element das größte oder kleinste Element ist. Dann gilt:

$$A(n) = n - 1 + A(n - 1)$$

Der Rest erfolgt Analog zu der Aufwandserklärung von Bubblesort 8:

$$A(n) = \frac{n(n-1)}{2} = \mathcal{O}(n^2)$$

Damit ist der Aufwand gezeigt. □

- 12.16 Satz.** Alle Permutationen der Zahlen  $\{1, 2, \dots, n\}$  seien gleich wahrscheinlich. Dann benötigt Quicksort im Durchschnitt  $\mathcal{O}(n \log(n))$  Vergleiche zum sortieren von Zahlen.

*Beweis.* Sei  $\Pi$  die Menge aller Permutationen von  $\{1, 2, \dots, n\}$  und sei  $A(\pi), \pi \in \Pi$  die Anzahl Vergleiche um eine Permutation  $\pi$  mittels Quicksort zu sortieren. Der durchschnittliche Aufwand ist:

$$\bar{A}(n) = \frac{1}{n!} \sum_{\pi \in \Pi} A(\pi)$$

O.B.d.A: Sei das erste Element das Pivotelement. Definiere:

$$\Pi_k = \{\pi \in \Pi | \pi_1 = k\}, k = 1, \dots, n$$

mit

$$|\Pi_k| = (n-1)!, k = 1, \dots, n$$

Für  $k$  fix und  $\pi \in \Pi_k$  teilt Quicksort im ersten Aufruf in zwei Mengen

$$\begin{aligned} \pi_{<} &= \{\text{Permutationen von } 1, 2, 3, \dots, k-1\} \\ \pi_{>} &= \{\text{Permutationen von } k+1, \dots, n\} \end{aligned}$$

Analog zu 11.12 folgt

$$A(\pi) = n-1 + A(\pi_{<}) + A(\pi_{>})$$

und

$$\sum_{\pi \in \Pi_k} A(\pi) = (n-1)!(n-1) + \sum_{\pi \in \Pi_k} A(\pi_{<}) + \sum_{\pi \in \Pi_k} A(\pi_{>}) \quad (12.1)$$

Wenn  $\pi$  alle Permutationen aus  $\Pi_k$  durchläuft, entstehen für  $\Pi_{<}$  alle Permutationen von  $\{1, 2, \dots, k-1\}$ . Dabei tritt jede Permutation genau  $\frac{(n-1)!}{(k-1)!} = \frac{|\Pi_k|}{|\Pi_{<}|}$  mal auf.

$$\Rightarrow \sum_{\pi \in \Pi_k} A(\pi_{<}) = \frac{(n-1)!}{(k-1)!} \sum_{\pi \in \Pi_k} A(\pi_{<}) = (n-1)! \bar{A}(k-1) \quad (12.2)$$

Analog:

$$\sum_{\pi \in \Pi_k} A(\pi_{>}) = (n-1)! \bar{A}(n-k) \quad (12.3)$$

Zusammensetzen:

$$\begin{aligned} \bar{A}(n) &= \frac{1}{n!} \sum_{\pi \in \Pi} A(\pi) \\ &= \frac{1}{n!} \sum_{k=1}^n \sum_{\pi \in \Pi_k} A(\pi) \end{aligned}$$

Unter Verwendung der Gleichungen 12.1, 12.2, 12.3 folgt weiter:

$$\begin{aligned}
 &= \frac{1}{n!} \sum_{k=1}^n (n-1)! (n-1 + \bar{A}(k-1) + \bar{A}(n-k)) \\
 &= n-1 + \frac{1}{n} \sum_{k=1}^n \bar{A}(k-1) + \bar{A}(n-k) \\
 &= n-1 + \frac{2}{n} \sum_{k=0}^{n-1} \bar{A}(k)
 \end{aligned}$$

mit Startwerten  $\bar{A}(0) = \bar{A}(1) = 0$ .

Per Induktion folgt:

$$\begin{aligned}
 \bar{A}(n) &= 2(n+1) \sum_{i=1}^n \frac{1}{i} - 4n \\
 &\leq 2(n+1) \left(1 + \int_1^n \frac{1}{x} dx\right) - 4n \\
 &= 2(n+1)(1 + \log(n)) - 4n \\
 &= 2(n+1) \log(n) - 2(n-1) \\
 &= \mathcal{O}(n \log(n))
 \end{aligned}$$

□

## 13. Untere Schranke für das Sortierproblem

**Frage:** Gibt es einen Sortieralgorithmus, der schneller ist als  $\mathcal{O}(n \log(n))$ ?

**13.17 Satz.** Jedes deterministische Sortierverfahren, das auf paarweisen Vergleichen basiert und keine Vorkenntnisse über die zu sortierende Menge hat, benötigt im schlimmsten Fall mindestens  $\log_2(n!)$  Vergleiche zum Sortieren von  $n$  verschiedenen Zahlen.

Bevor wir den Beweis machen gibt es noch einige Hinweise, die wir für den Beweis benötigen.

**13.18 Bemerkung.** Es gilt:

- $\log_2(n!) = \Theta(n \log(n))$
- $\log_2(n) \leq \log_2(n^n) = n \log_2(n) = \frac{1}{\log(2)} n \log n$
- $\log(n!) \geq \log\left(\frac{n}{2} \dots \frac{n}{2}\right) \geq \log\left(\left(\frac{n}{2}\right)^{\frac{n}{2}}\right) = \frac{n}{2} \log\left(\frac{n}{2}\right) = \frac{n}{2} \log(n) - n \log(2)$

Zum Beweis benötigen wir einen Entscheidungsbaum. Dieser kann jedem Sortieralgorithmus mit paarweisen Vergleichen zugeordnet werden, und illustriert das Verhalten des Algorithmus:

- Innere Knoten des Entscheidungsbaums sind Vergleiche im Algorithmus.
- Ein Weg von der Wurzel zu einem Blatt entspricht der zeitlichen Abfolge der Vergleiche. Ein Weitergeben nach rechts illustriert einen richtigen Vergleich, nach links einen falschen.
- Die  $n!$  Blätter des Baumes stellen die Permutationen der zu sortierenden Menge dar, die jeweils zur Abfolge der Vergleiche gehört.

**13.19 Beispiel.** Baum für  $\{z_1, z_2, z_3\}$ :

*Beweis.* 13.17

Idee: Finde eine untere Schranke für die Tiefe des Baumes. Dies ist die Mindestanzahl der im schlimmsten Fall nötigen Vergleiche. Im besten Fall ist ein Baum ausbalanciert, d.h. alle Blätter haben die gleiche Tiefe  $m$  und wir haben  $2^m$  Blätter. Da wir  $n!$  Permutationen haben (also Blätter), muss  $2^m = n!$  gelten. Es folgt  $m \log_2(n!)$   $\square$

Satz 13.17 sagt, dass Mergesort in der asymptotischen Laufzeit optimal ist.

**Frage:** Können wir etwas Ähnliches für Quicksort zeigen?

**13.20 Satz.** Alle Permutationen der Zahlen  $z_1, \dots, z_n$  seien gleich wahrscheinlich. Dann benötigen Sortierverfahren wie in Satz 13.17 im Mittel  $\log_2(n!)$  Vergleiche.

*Beweis.* Die mittlere Höhe  $\overline{H}(z)$  des Entscheidungsbaums  $\tau$  ist gegeben als der Mittelwert der Tiefen  $t_\tau$  aller seiner Blätter  $v$ . Sei  $B(\tau)$  die Menge der Blätter und  $\beta(\tau) = |B(\tau)|$ . Dann gilt:

$$\overline{H}(\bar{c}) = \frac{1}{\beta(\bar{c})} \sum_{y \in B(\bar{c})} t_{\bar{c}}(v)$$

Zu zeigen ist, dass:

$$\overline{H}(\bar{c}) \geq \log_2(\beta(\bar{c}))$$

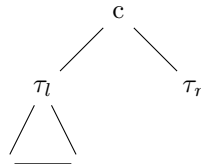
Vollständige Induktion über die Höhe des Entscheidungsbaums  $\bar{c}$ .

◇ Induktionsverankerung:  $H(\bar{c}) = 0$

$\bar{c}$  besteht nur aus seiner Wurzel.  $\implies \overline{H}(\bar{c}) = \frac{1}{1} \sum_{y \in B(\bar{c})} t_{\bar{c}}(v) = 0 = \log_2(1)$

◇ Induktionsschritt:  $H(\bar{c}) \leq h - 1 \implies H(\bar{c}) = h$

Seien  $\tau_l, \tau_r$  die Wurzeln des linken bzw. rechten Teilbaums zum Wurzelknoten von  $\bar{c}$ :



**Beobachte:**

- $H(\tau_l), H(\tau_r) < h$
- $B(\tau) = B(\tau_l) \cup B(\tau_r)$
- $B(\tau_l) \cap B(\tau_r) = \emptyset$
- $\beta(\bar{c}) = \beta(\tau_l) + \beta(\tau_r)$
- $\beta(\tau_l), \beta(\tau_r) \geq 1$
- $t_\tau(v) = t_{\tau_l}(v) + 1$
- $t_\tau(v) = t_{\tau_r}(v) + 1$

Rechne:

$$\begin{aligned} \overline{H}(\tau) &= \frac{1}{\beta(\tau)} \sum_{v \in B(\tau)} t_\tau(v) \\ &= \dots \\ &\geq 1 + \frac{1}{b} \min_{x \in [0, b]} (x \log_2(x) + (b-x) \log_2(b-x)) \end{aligned}$$

Da

$$\begin{aligned} f'(x) &= \log_2(x) - \log_2(b-x) \\ &= \log_2\left(\frac{x}{b-x}\right) = 0 \end{aligned}$$

genau dann, wenn  $x = \frac{b}{2}$ , ist dies die einzige Möglichkeit, wo das Minimum in  $(0, b)$  angenommen werden kann. Die anderen Möglichkeiten sind die Intervallrandpunkte, was ausgeschlossen ist. Also gilt:

$$\begin{aligned} \overline{H}(\tau) &\geq 1 + \log_2\left(\frac{b}{2}\right) \\ &= 1 + \log_2(\beta(\tau)) - 1 \\ &= \log_2(\beta(\tau)) \end{aligned}$$

□

# Graphen

## 14. Grundlagen

**14.1 Definition (Graph).** Ein Graph ist ein Paar  $G = (V, E)$ , bestehend aus endlichen Mengen  $V$  von Knoten (vertices) und  $E$  von Kanten (edges). Die Kanten beschreiben die Verbindungen zwischen den Knoten und sind gerichtet oder ungerichtet.

- Für *gerichtete* Graphen oder *Digraphen* gilt:

$$E \subset \{(v, w) \in V \times V | v \neq w\}$$

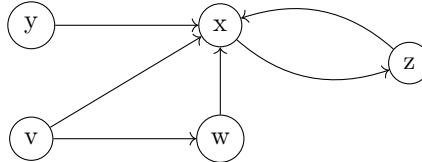
Ist  $e = (v, w) \in E$ , so heißt  $v$  der Anfangsknoten und  $w$  der Endknoten.

- Für *ungerichtete* Graphen ist  $E$  eine Menge von ungeordneten Paaren  $\{v, w\} \subset E, v \neq w$ , das heißt:

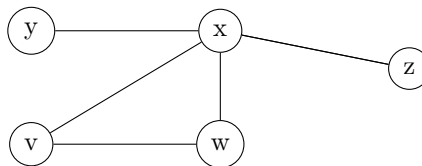
$$E \subset \{X \subset V | |X| = 2\}$$

**14.2 Beispiel.** Sei  $V = \{v, w, x, y, z\}$

- $G = (V, E), E = \{(v, w), (v, x), (w, x), (x, z), (z, x), (y, x)\}$  ist ein Digraph, der wie folgt illustriert werden kann:



- Der gleiche Graph sieht ungerichtet wie folgt aus:



Wir bemerken insbesondere, dass  $\{v, w\} = \{w, v\}$  im ungerichteten Graphen ist.

**14.3 Bemerkung.** In Digraphen werden manchmal auch Kanten der Form  $\{v, v\}, v \in V$  zugelassen.

**14.4 Definition.** Sei  $G = (V, E)$  ein Graph. Ein *Weg* oder *Pfad* von  $v = v_0$  nach  $w = v_r$  in  $G$  ist eine Kantenfolge

$$\pi = v_0, v_1, \dots, v_r$$

mit  $r \geq 1$  und

- $(v_i, v_{i+1}) \in E, i = 0, \dots, r - 1$  im Falle von Digraphen.
- $\{v_i, v_{i+1}\} \in E, i = 0, \dots, r - 1$  im Falle von ungerichteten Graphen.

Der Weg heißt *einfach*, falls

- $v_0, \dots, v_r$  paarweise verschieden sind oder
- $v_0, \dots, v_r$  paarweise verschieden mit  $v_0 = v_r$ .

Die *Länge* von  $\pi$  ist gegeben als  $|\pi| = r$ , ist also die Anzahl Knoten, die in  $\pi$  durchlaufen werden. Ein Knoten  $w$  heißt von Knoten  $v$  *erreichbar*, falls ein Weg von  $v$  nach  $w$  existiert.

**14.5 Beispiel.** Bezogen auf Beispiel 14.2 ist:

- $\pi_1 = v, w$  ein einfacher Weg der Länge 1
- $\pi_2 = v, w, x$  ein einfacher Weg der Länge 2
- $\pi_3 = v, w, x, z, x, z$  ein Weg der Länge 5

Im ungerichteten Fall ist

- $\pi_4 = v, w, x, v$  ein einfacher Weg der Länge 3.

In beiden Fällen ist  $z$  von  $v$  erreichbar, im ungerichteten Fall gilt dies auch umgekehrt. Im gerichteten Fall ist  $v$  nicht von  $z$  erreichbar.

**14.6 Definition.** Sei  $G = (V, E)$  ein Digraph und  $v \in V$ . Wir definieren:

- die Menge der (direkten) *Nachfahren* von  $v$ :

$$post(v) = \{w \in V | (v, w) \in E\}$$

- die Menge der (direkten) *Vorfahren* von  $v$ :

$$pre(v) = \{w \in V | v \in post(w)\}$$

- die Menge der *erreichbaren Knoten* von  $v$ :

$$post^*(v) = \{w \in V | \text{es gib einen Weg von } v \text{ nach } w\}$$

- die Menge aller Knoten die  $v$  erreichen können:

$$pre^*(v) = \{w \in V | v \in post^*(w)\}$$

- die Nachbarn als die Menge  $post(v) \cup pre(v)$  von  $v$
- und den Knotengrad  $deg(v)$ , welcher der Anzahl seiner Kanten entspricht.

**14.7 Bemerkung.** Mit den offensichtlichen Modifikationen kann die Definition 14.6 auch auf ungerichtete Graphen angewendet werden. Wir erhalten in diesem Fall für  $v \in V$ , dass

$$\begin{aligned} post(v) &= pre(v) \\ post^*(v) &= pre^*(v) \end{aligned}$$

**14.8 Beispiel.** Für den Digraphen aus Beispiel 14.2 gilt:

- $post(v) = \{w, x\}$
- $post^*(v) = \{w, x, z\}$
- $pre(v) = \emptyset$
- $pre^*(v) = \emptyset$
- $deg(v) = 2$

Für den ungerichteten Graphen gilt geändert:

- $pre(v) = \{w, x\}$
- $pre^*(v) = \{w, x, y, z\}$
- $post^* = \{w, x, y, z\}$

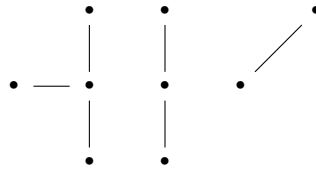
## 15. Zusammenhang

**15.9 Definition.** Sei  $G = (V, E)$  ein ungerichteter Graph  $\emptyset \neq C \subset V$

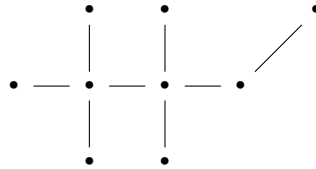
- Die Menge  $C$  heißt *zusammenhängend*, falls je zwei Knoten  $v, w \in C, v \neq w$ , voneinander erreichbar sind, das heißt,  $v \in post^*(w), w \in post^*(v)$ .  
Ist  $G$  ein Digraph, so heißt  $C$  zusammenhängend, falls  $C$  im zugrundeliegenden ungerichteten Graphen zusammenhängend ist.
- Eine zusammenhängende Knotenmenge heißt *Zusammenhangskomponente*, falls sie *maximal* ist. Das bedeutet, es gibt keine weitere zusammenhängende Knotenmenge  $D \subset V$  mit  $C \subsetneq D$ .
- Ein Graph heißt *zusammenhängend*, falls  $V$  zusammenhängend ist.

**15.10 Beispiel.** Dies ist ein unzusammenhängender Graph mit drei Zusammenhangskomponenten:





Analog hierzu ist der folgende Graph zusammenhängend:



Wir können beobachten, dass die Zusammenhangskomponenten eines ungerichteten Graphen die Äquivalenzklassen der Knotenmenge  $V$  unter der Äquivalenzrelation

$$v = w \iff \{v\} \cup \text{post}^*(v) = \{w\} \cup \text{post}^*(w)$$

ist.

Insbesondere zerfällt  $G$  in paarweise disjunkte Zusammenhangskomponenten  $C_1, \dots, C_r$  mit

$$V = \bigcup_{i=1}^r C_i$$

$$E = \bigcup_{i=1}^r E_i$$

wobei  $E_i = E \cap \{X \subset C_i : |X| = 2\}$ .

**15.11 Satz.** Sei  $G = (V, E)$  ein ungerichteter Graph mit  $n = |V| \geq 1$  Knoten und  $m = |E|$  Kanten. Ist  $G$  zusammenhängend, so folgt für die Anzahl Knoten und Kanten, dass

$$m \geq n - 1$$

*Beweis.*

□

Per vollständiger Induktion:

◇ IV  $n=1$

$$m = 0 = n - 1$$

◇ IS  $n \geq 3$

Wähle  $v \in V$  mit:

$$\deg(v) = \min_{v \in V} \deg(v) =: k$$

Da  $G$  zusammenhängend ist, muss  $k > 0$  gelten.

$k \geq 2$

$$\begin{aligned}
 2m &= 2|E| \\
 &= \sum_{w \in V} \deg w \\
 &\geq 2|V| \\
 &= 2n \qquad \qquad \qquad \implies m \geq n \geq n-1
 \end{aligned}$$

$k = 1$

Sei  $G'(V', E')$  derjenige Graph, der durch das Streichen von  $v$  und der zugehörigen Kante entsteht.  $G'$  ist zusammenhängend, weil  $G$  zusammenhängend ist. Die Induktionsannahme impliziert:

$$m - 1 = |E'| \geq |V'| - 1 = (n - 1) - 1 = n - 2 \implies m \geq n - 1$$

## 16. Zyklische Graphen