

Übungen zur Vorlesung Numerik I

Sommersemester 2010

PD Dr. Thorsten Hüls
Dipl.-Math. Denny Otten

Übungsblatt 2
22.04.2010

Abgabe: Donnerstag, 29.04.2010, 10:00 Uhr in das Postfach des jeweiligen Tutors.

Mo.-Tutorium: Paul Voigt, paulvoigt@web.de, Postfach 195 in V3-128

Di.-Tutorium: Denny Otten, dotten@math.uni-bielefeld.de, Postfach 44 in V3-128

Mi.-Tutorium: Ingwar Petersen, ipeterse@math.uni-bielefeld.de, Postfach 227 in V3-128

Aufgabe 4: (Rundungsfehler der Gleitkommaarithmetik)

Sei $b \in \mathbb{N}$, $b \geq 2$ gerade, $G(\tau, b)$ die Menge der Gleitkommazahlen zur Basis b mit Mantissenlänge τ und sei $rd : \mathbb{R} \rightarrow G(\tau, b)$ die übliche Rundungsvorschrift.

Man zeige für alle $t \in \mathbb{R}$:

$$|t - rd(t)| \leq |t - g| \quad \forall g \in G(\tau, b). \quad (1)$$

Hinweis: Es empfiehlt sich, folgendermaßen vorzugehen:

- (i) Zeigen Sie: Gilt (1) für $t > 0$, so auch für $t < 0$.
- (ii) Wählen Sie für $t \in \mathbb{R}$ eine geeignete Darstellung der Form $t = b^n \sum_{j=1}^{\infty} \beta_j b^{-j}$.
- (iii) Zeigen Sie: Zu jedem $g \in G(\tau, b)$ gibt es ein $g' \in G(\tau, b) \cap [b^{n-1}, \infty)$, so dass $|g' - t| \leq |g - t|$.
- (iv) Zeigen Sie für $g, g' \in G(\tau, b) \cap [b^{n-1}, \infty)$ mit $g \neq g'$, dass $|g - g'| \geq b^{n-\tau}$.
- (v) Beweisen Sie nun (1) mit Hilfe von Skript, Satz 2.3 für alle $g \in G(\tau, b) \cap [b^{n-1}, \infty)$ mit $g \neq rd(t)$.

(6 Punkte)

Aufgabe 5: (Berechnung von Konditionszahlen)

(a) Berechnen Sie die relative komponentenweise Konditionszahl $\widehat{\kappa}(F, x)$ für die folgenden Auswertungsprobleme:

(i) $F : \mathbb{R} \rightarrow \mathbb{R}$ mit $F(x) = ax^n \sqrt{bx}$, $x > 0$, $a, b \in \mathbb{R}$ mit $b > 0$, $n \in \mathbb{N}_0$.

(ii) $F : \mathbb{R}^2 \rightarrow \mathbb{R}$ mit $F \left(\begin{array}{c} x_1 \\ x_2 \end{array} \right) = \frac{ax_1^n}{bx_2^m}$, $x_2 \neq 0$, $a, b \in \mathbb{R}$ mit $b \neq 0$, $n, m \in \mathbb{N}$.

(iii) $F : \mathbb{R} \rightarrow \mathbb{R}^2$ mit $F(x) = \left(\begin{array}{c} \operatorname{sech}(x) \\ \tanh(x) \end{array} \right)$

(b) Betrachten Sie die aus der Aufgabe 2 bekannte Iteration

$$x^{(n)} = f_{a,b}(x^{(n-1)}), \quad n = 1, 2, \dots$$

mit

$$f_{a,b} : \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad \text{mit} \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto f_{a,b} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} := \begin{pmatrix} ax_1 \\ -(b - a^2)x_1^2 + bx_2 \end{pmatrix}$$

und dem Anfangswert $x^{(0)} = \begin{pmatrix} c \\ c^2 \end{pmatrix}$.

- (1) Bestimmen Sie die relative komponentenweise Konditionszahl $\hat{\kappa}$ des Auswertungsproblems $(f_{a,b}, \begin{pmatrix} c \\ c^2 \end{pmatrix})$. Setzen Sie nun $a = \frac{1}{2}$, $b = 2$, $c = 2$ und berechnen Sie die
- (2) Kondition $\tilde{\kappa}(0)$ der Gesamtiteration, also des Auswertungsproblems

$$\left(\underbrace{f \circ \cdots \circ f}_{n-\text{mal}}, \begin{pmatrix} c \\ c^2 \end{pmatrix} \right)$$

Hinweis: Vgl. Skript, Formel (2.22) sowie die explizite Darstellung der Folge $(x^{(n)})_{n \in \mathbb{N}_0}$ aus Aufgabe 2.

Wie verhält sich für $n \rightarrow \infty$?

(6 Punkte)

Aufgabe 6: (Programmieraufgabe, Kondition von Auswertungsproblemen)

Man betrachte noch einmal das Beispiel zur Berechnung der Lösung $F(p, q) = -p + \sqrt{p^2 + q}$ der quadratischen Gleichung $t^2 + 2pt - q = 0$ für $q > 0, p \in \mathbb{R}$.

Zwei Algorithmen sind gegeben durch

Algorithmus 1: $F^1(p, q) = (p, \sqrt{p^2 + q})$, $F^2(y_1, y_2) = y_2 - y_1$.

Algorithmus 2: $F^1(p, q) = (p, q, \sqrt{p^2 + q})$, $F^2(y_1, y_2, y_3) = \frac{y_2}{y_1 + y_3}$.

In der Vorlesung wurden die Konditionszahlen bereits im Fall $p > 0$ untersucht. Bestimmen Sie nun für $p < 0 < q$

- (1) (a) die relative Kondition des Auswertungsproblems $\hat{\kappa}(F, (p, q))$.
- (1) (b) die relativen Konditionszahlen $\hat{\kappa}_1 = \hat{\kappa}(F^1, (p, q))$ und $\hat{\kappa}_2 = \hat{\kappa}(F^2, F^1(p, q))$ zum Algorithmus 1.
- (1) (c) die relativen Konditionszahlen $\hat{\kappa}_1$ und $\hat{\kappa}_2$ zum Algorithmus 2.
- (1) Welcher Algorithmus ist in diesem Fall gutartig?
- (1) (d) Berechnen Sie mit Hilfe der Rundungsfehler GUI die Werte der beiden Algorithmen zur Mantissenlänge 8 für $q = 1, p = -10, -100, -1000$. Interpretieren Sie die Ergebnisse.

(6 Punkte)

AUFGABE 4: (Rundungsfehler der Gleitkommarithmetik)

Sei $b \in \mathbb{N}$ gerade mit $b \geq 2$, $G(\tau, b)$ die Menge der Gleitkommazahlen zur Basis b mit Mantissenlänge τ und $rd: \mathbb{R} \rightarrow G(\tau, b)$ die übliche Rundungsvorschrift. Dann gilt:

$$\forall t \in \mathbb{R} \quad \forall g \in G(\tau, b) : |t - rd(t)| \leq |t - g| \quad (4.1)$$

Beweis:

1. 1. Fall: ($t = 0$). Trivialeweise gilt:

$$|t - rd(t)| = 0 \leq |g| = |t - g| \quad \forall g \in G(\tau, b)$$

\uparrow
 $t = 0$
 $rd(0) = 0$

2. Fall: ($t < 0$). Angenommen die Aussage (4.1) sei für $t > 0$ bereits gezeigt, d.h. es gilt

$$\forall t \in \mathbb{R} \text{ mit } t > 0 \quad \forall g \in G(\tau, b) : |t - rd(t)| \leq |t - g|, \quad (4.2)$$

dann erhalten wir aus (4.2), " $g \in G(\tau, b) \Rightarrow -g \in G(\tau, b)$ " und $rd(t) = -rd(-t)$, denn

$$\begin{aligned} rd(t) &= \begin{cases} b^n(0, \beta_1 \beta_2 \dots \beta_\tau) & , \text{ falls } 0 \leq \beta_{\tau+1} \leq \frac{b}{2} - 1 \\ b^n[(0, \beta_1 \beta_2 \dots \beta_\tau) + b^{-\tau}] & , \text{ falls } \frac{b}{2} \leq \beta_{\tau+1} \leq b-1 \end{cases} \\ &= - \begin{cases} -b^n(0, \beta_1 \beta_2 \dots \beta_\tau) & , \text{ falls } 0 \leq \beta_{\tau+1} \leq \frac{b}{2} - 1 \\ -b^n[(0, \beta_1 \beta_2 \dots \beta_\tau) + b^{-\tau}] & , \text{ falls } \frac{b}{2} \leq \beta_{\tau+1} \leq b-1 \end{cases} \\ &= -rd(-t) \quad \forall t \in \mathbb{R} \end{aligned}$$

die Behauptung

$$|t - rd(t)| = |t + rd(-t)| = |(-t) - rd(-t)|$$

\uparrow
 $rd(t) = -rd(-t)$

$$\leq |(-t) - (-g)| = |g - t| = |t - g| \quad \forall t \in \mathbb{R} \text{ mit } t < 0 \text{ und } g \in G(\tau, b)$$

\uparrow
 $+ < 0 \Rightarrow - > 0$
 $\& (4.2) \& "g \in G(\tau, b) \Rightarrow -g \in G(\tau, b)"$

Damit ist die Aufgabe gelöst, insoweit wir (4.2) gezeigt haben. (d.h. wir behandeln nur den Fall $t > 0$).

2. Sei $t \in \mathbb{R}$ mit $t > 0$ beliebig, dann gilt zunächst

$$\exists_{n \in \mathbb{Z}} : t \in [b^{n-1}, b^n]. \quad (4.3)$$

Weiter können wir t darstellen als (siehe Skript (2.2))

$$t = b^n \sum_{j=1}^{\infty} \beta_j b^{-j}, \quad \beta_j \in \{0, \dots, b-1\}, \beta_1 \neq 0, j \in \mathbb{N} \quad (4.4)$$

Für die Eindeutigkeit der Gleitkommadarstellung (4.4) genügt es, wenn wir $\beta_j \neq b-1$ für unendlich viele $j \in \mathbb{N}$ fordern. Fortsetzung des Beweises folgt in Schritt 5.

3. HILFSAUSSAGE 1: Sei $t \in \mathbb{R}$ mit $t > 0$ und $t = b^n \sum_{j=1}^{\infty} \beta_j b^{-j}$ wie in (4.4). Dann gilt

$$\forall g \in G(\tau, b) \exists g' \in G(\tau, b) \cap [b^{n-1}, \infty[: |g' - t| \leq |g - t|.$$

Beweis:

1. Fall: ($g \geq b^{n-1}$). Wähle $g' := g$, dann gilt wegen $g \geq b^{n-1}$, dass $g' \in G(\tau, b) \cap [b^{n-1}, \infty[$ und es folgt

$$|g - t| = |g' - t|$$

Def. von g'

2. Fall: ($g < b^{n-1}$). Wähle $g' := b^{n-1}$, dann gilt $g' \in G(\tau, b) \cap [b^{n-1}, \infty[$. Daraus sowie aus der Darstellung von t erhalten wir

$$g < b^{n-1} = g' \leq t \Rightarrow g - t < g' - t \leq 0$$

$$\Rightarrow |g - t| > |g' - t|$$

4. HILFSAUSSAGE 2:

$$\forall g, g' \in G(\tau, b) \cap [b^{n-\tau}, \infty[\text{ mit } g \neq g' : |g - g'| \geq b^{n-\tau}$$

Beweis:

$$\begin{aligned} \text{Seien } g, g' \in G(\tau, b) \cap [b^{n-\tau}, \infty[\text{ mit } g \neq g', \text{ d.h.} \\ g = b^{m_1} \sum_{j=1}^{\tau} \alpha_j^1 b^{-j} \quad , \quad m_1 \geq n, \alpha_j^1 \in \{0, \dots, b-1\}, \alpha_1^1 \neq 0, j=1, \dots, \tau \\ g' = b^{m_2} \sum_{j=1}^{\tau} \alpha_j^2 b^{-j} \quad , \quad m_2 \geq n, \alpha_j^2 \in \{0, \dots, b-1\}, \alpha_n^2 \neq 0, j=1, \dots, \tau \end{aligned}$$

Offensichtlich gilt nun $g \cdot b^{\tau-n} \in \mathbb{Z}$ und $g' \cdot b^{\tau-n} \in \mathbb{Z}$, dann

$$g \cdot b^{\tau-n} = \underbrace{b^{m_1-n}}_{\in \mathbb{N}} \cdot \underbrace{\sum_{j=1}^{\tau} \alpha_j^1 b^{\tau-j}}_{\in \mathbb{N} \text{ (denn: } b \in \mathbb{N} \text{ mit } b \geq 2 \text{ & } m_1-n \in \mathbb{Z} \text{ mit } m_1-n \geq 0)} \quad \Rightarrow \quad g \cdot b^{\tau-n} \in \mathbb{Z} \quad (\text{sogar No})$$

$\in \mathbb{N}$
(denn: $b \in \mathbb{N}$
mit $b \geq 2$
 $\& m_1-n \in \mathbb{Z}$
mit $m_1-n \geq 0$)

$\in \mathbb{N}$
(denn: $b \in \mathbb{N}$
mit $b \geq 2$
 $\& \tau-j \in \mathbb{N}_0, j=1, \dots, \tau$)

$\in \mathbb{N}$
abgeschlossen
bzgl. + und .

$g' \cdot b^{\tau-n}$: analog

und daraus erhalten wir

$$|g - g'| \cdot b^{\tau-n} \geq 1.$$

↑
Zeile zuvor

Multiplication mit $b^{n-\tau}$ liefert die Behauptung

$$|g - g'| \geq b^{n-\tau}. \blacksquare$$

5. Für die Darstellung von t in (4.4) gilt nach der Definition der Rundungsvorschrift rd

$$rd(t) \in G(\tau, b) \cap [b^{n-\tau}, \infty[. \quad \text{deswegen und} \quad (4.5)$$

Sei nun $g \in G(\tau, b) \cap [b^{n-\tau}, \infty[$ mit $g \neq rd(t)$, dann sind wegen (4.5) die Voraussetzungen von Hilfsaussage 2 erfüllt und es gilt

$$|g - rd(t)| \geq b^{n-\tau} \quad (4.6)$$

Andererseits liefert Satz 2.3 aus der Vorlesung

$$|t - rd(t)| \leq \frac{1}{2} \cdot b^{n-\tau} \quad (4.7)$$

Beweisv. Satz 2.3

Aus (4.6) und (4.7) erhalten wir insgesamt

$$b^{n-\tau} \leq |g - rd(t)| \leq |g - t| + |t - rd(t)| \leq |g - t| + \frac{1}{2} \cdot b^{n-\tau} \quad (4.8)$$

(4.6) (4.7) Ergänzen & Δ 's-Ungl.

$$-|g - t| \Rightarrow |t - rd(t)| \leq \frac{1}{2} \cdot b^{n-\tau} \leq \frac{1}{2} |g - rd(t)|$$

$$\stackrel{(4.7)}{=} \frac{1}{2} (|g - t| + |t - rd(t)|) = \frac{1}{2} |g - t| + \frac{1}{2} |t - rd(t)|$$

(2) Ergänzen & Δ 's-Ungl.

$$-\frac{1}{2} |t - rd(t)| \Rightarrow |t - rd(t)| \leq |g - t| \quad \forall g \in G(\tau, b) \cap [b^{n-\tau}, \infty[\quad (4.9)$$

Sei jetzt $g \in G(\tau, b)$ beliebig, dann gilt nach Hilfsaussage 1

$$\exists g' \in G(\tau, b) \cap [b^{n-\tau}, \infty[: |g' - t| \leq |g - t| \quad (4.10)$$

Daraus und aus (4.9) mit g' erhalten wir

$$|t - rd(t)| \leq |g' - t| \leq |g - t| \quad (4.9) \quad (4.10)$$

und es folgt (4.2) und somit (4.1). ■

AUFGABE 5 : zu(a):

Sei $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$. Wir berechnen in dieser Aufgabe die relative Komponentenweise Kondition des Auswertungsproblems (F, x) , die durch (siehe: Skript (2.17))

$$\hat{K}(F, x) := \lim_{\delta \rightarrow 0} \|L(\delta)\|_\infty = \max_{i=1, \dots, m} \sum_{j=1}^n \left| \frac{\partial F_i}{\partial x_j}(x) \right| \cdot \frac{|x_j|}{|F_i(x)|} \quad (5.1)$$

gegeben ist.

zu(i): Betrachte

$$F: \mathbb{R}_+^n \rightarrow \mathbb{R} \text{ mit } F(x) = a \cdot x^n \cdot \sqrt{b \cdot x} = a \cdot b^{\frac{1}{2}} \cdot x^{n+\frac{1}{2}}, \quad \mathbb{R}_+^n = \{x \in \mathbb{R}^n \mid x > 0\}$$

wobei $a, b \in \mathbb{R}$ mit $a, b > 0$ und $n \in \mathbb{N}_0$. Wegen

$$\frac{\partial F}{\partial x}(x) = a \cdot b^{\frac{1}{2}} \cdot (n + \frac{1}{2}) \cdot x^{n-\frac{1}{2}} \quad (5.2)$$

erhalten wir aus (5.1) die Kondition

$$\begin{aligned} \hat{K}(F, x) &= \left| \frac{\partial F}{\partial x}(x) \right| \cdot \frac{|x|}{|F(x)|} \stackrel{(5.2)}{=} |a \cdot b^{\frac{1}{2}} \cdot (n + \frac{1}{2}) \cdot x^{n-\frac{1}{2}}| \cdot \frac{|x|}{|a \cdot b^{\frac{1}{2}} \cdot x^{n+\frac{1}{2}}|} \\ &= |n + \frac{1}{2}| = n + \frac{1}{2}. \end{aligned}$$

zu(ii): Betrachte

$$F: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ mit } F(x_2) = \frac{a \cdot x_1^n}{b \cdot x_2^m} = \frac{a \cdot b^{-m} \cdot x_1^n \cdot x_2^{-m}}{b^m}, \quad x_2 \neq 0$$

wobei $a, b \in \mathbb{R}$ mit $b \neq 0$ und $n, m \in \mathbb{N}$. Wegen

$$\frac{\partial F}{\partial x_1}(x) = a \cdot b^{-1} \cdot n \cdot x_1^{n-1} \cdot x_2^{-m} \quad (5.3)$$

$$\frac{\partial F}{\partial x_2}(x) = a \cdot b^{-1} \cdot (-m) \cdot x_1^n \cdot x_2^{-(m+1)}$$

erhalten wir aus (5.1) die Kondition

$$\begin{aligned} \hat{K}(F, x) &= \sum_{j=1}^2 \left| \frac{\partial F}{\partial x_j}(x) \right| \cdot \frac{|x_j|}{|F(x)|} \stackrel{(5.1)}{=} \\ &= \left| a \cdot b^{-1} \cdot n \cdot x_1^{n-1} \cdot x_2^{-m} \right| \cdot \frac{|x_1|}{|a \cdot b^{-1} \cdot n \cdot x_1^{n-1} \cdot x_2^{-m}|} + \left| a \cdot b^{-1} \cdot (-m) \cdot x_1^n \cdot x_2^{-(m+1)} \right| \cdot \frac{|x_2|}{|a \cdot b^{-1} \cdot (-m) \cdot x_1^n \cdot x_2^{-(m+1)}|} \\ &= |n| + |-m| = |n| + |m| = n + m. \end{aligned}$$

zu(b): Betrachte die Funktion F aus Aufgabe 2

$$F: \mathbb{R}^2 \rightarrow \mathbb{R}^2 \text{ mit } F(x_2) = \begin{pmatrix} a \cdot x_1 \\ -(b-a^2)x_1^2 + b x_2 \end{pmatrix} =: \begin{pmatrix} F_1(x_2) \\ F_2(x_2) \end{pmatrix}$$

wobei $a, b \in \mathbb{R}$. Wegen

$$\frac{\partial F_1}{\partial x_1}(x) = a \quad \frac{\partial F_2}{\partial x_1}(x) = -2(b-a^2)x_1 \quad (5.5)$$

$$\frac{\partial F_1}{\partial x_2}(x) = 0 \quad \frac{\partial F_2}{\partial x_2}(x) = b$$

erhalten wir aus (5.1) die Kondition

$$\begin{aligned} \hat{K}(F, x) &= \max_{(5.1)} \left\{ \left| a \right| \cdot \frac{|x_1|}{|a \cdot x_1|} + \left| b \right| \cdot \frac{|x_2|}{|a \cdot x_1|} \right\} \\ &\quad + \left| -2(b-a^2)x_1 \right| \cdot \frac{|x_1|}{|-(b-a^2)x_1^2 + b x_2|} + \left| b \right| \cdot \frac{|x_2|}{|-(b-a^2)x_1^2 + b x_2|} \end{aligned}$$

$$= \max \left\{ 1, \frac{\left| -2(b-a^2)x_1^2 + b x_2 \right|}{\left| -(b-a^2)x_1^2 + b x_2 \right|} \right\} = \frac{\left| -2(b-a^2)x_1^2 \right| + \left| b x_2 \right|}{\left| -(b-a^2)x_1^2 + b x_2 \right|}.$$

$$\left| -2(b-a^2)x_1^2 + b x_2 \right| \leq \left| -2(b-a^2)x_1^2 \right| + \left| b x_2 \right| \quad \Delta's Ungl.$$

$$\Rightarrow 1 \leq \frac{\left| -2(b-a^2)x_1^2 \right| + \left| b x_2 \right|}{\left| -(b-a^2)x_1^2 + b x_2 \right|} \quad (5.6)$$

→ Fortsetzung folgt nach
(a)(iii)

Zu (iii): Betrachte die Funktion

$$F: \mathbb{R} \rightarrow \mathbb{R}^2 \text{ mit } F(x) = \begin{pmatrix} \operatorname{sech}(x) \\ \tanh(x) \end{pmatrix} =: \begin{pmatrix} F_1(x) \\ F_2(x) \end{pmatrix}$$

wegen

$$\frac{\partial F_1}{\partial x}(x) = -\operatorname{sech}(x) \cdot \tanh(x)$$

$$\frac{\partial F_2}{\partial x}(x) = 1 - \tanh^2(x)$$

(5.6)⁴

erhalten wir aus (5.1) die Kondition

$$\hat{\kappa}(F_1, x) \stackrel{(5.1)}{=} \max \left\{ \left| \frac{\partial F_1}{\partial x}(x) \right| \cdot \frac{|x|}{|F_1(x)|}, \left| \frac{\partial F_2}{\partial x}(x) \right| \cdot \frac{|x|}{|F_2(x)|} \right\}$$

$$\stackrel{(5.6)}{=} \max \left\{ \frac{|-\operatorname{sech}(x) \cdot \tanh(x)| \cdot |x|}{|\operatorname{sech}(x)|}, \frac{|1 - \tanh^2(x)| \cdot |x|}{|\tanh(x)|} \right\}$$

$$= \max \left\{ |x \cdot \tanh(x)|, \left| \frac{x(1 - \tanh^2(x))}{\tanh(x)} \right| \right\}$$

$$= \begin{cases} |x \cdot \tanh(x)| & , x \in]-\infty, -\operatorname{arctanh}\left(\frac{\sqrt{2}}{2}\right] \cup [\operatorname{arctanh}\left(\frac{\sqrt{2}}{2}\right), +\infty[\\ \left| \frac{x(1 - \tanh^2(x))}{\tanh(x)} \right| & , x \in [\operatorname{arctanh}\left(\frac{\sqrt{2}}{2}\right), \operatorname{arctanh}\left(\frac{\sqrt{2}}{2}\right)] \end{cases}$$

Fortsetzung zu (b):

Aus der Kondition des Auswertungsproblems (F_1, x) , die in (5.6) berechnet wurde, erhalten wir die Kondition des Auswertungsproblems $(F_1, (c_2))$, indem $(x_2) = (c_2)$ gesetzt wird:

$$\hat{\kappa}(F_1, (c_2)) \stackrel{(5.6)}{=} \frac{|-(b-a^2)c^2| + |bc^2|}{|-(b-a^2)c^2 + bc^2|} = \frac{|-2(b-a^2)| + |b|}{|a^2|} \quad (5.7)$$

Setze nun $a = \frac{1}{2}$, $b = 2$ und $c = 2$. Zur Berechnung der Kondition $\tilde{\kappa}(0)$ der Gesamtiteration wissen wir aus Aufgabe 2, Blatt 1, dass sich $x^{(n)}$ explizit darstellen lässt durch

$$x^{(n)} = f(x^{(n-1)}) = (f_0 \dots \circ f)(x^{(0)}) = \begin{pmatrix} a^n \cdot c \\ a^{2n} \cdot c^2 \end{pmatrix} = \begin{pmatrix} (\frac{1}{2})^n \cdot 2 \\ (\frac{1}{2})^{2n} \cdot 2^2 \end{pmatrix} \quad \begin{matrix} x^{(0)} = (c_2) \\ \text{Aufgabe 2, Blatt 1} \end{matrix} \quad \begin{matrix} a = \frac{1}{2} \\ c = 2 \end{matrix}$$

$$= \begin{pmatrix} (\frac{1}{2})^{n-1} \\ (\frac{1}{2})^{2(n-1)} \end{pmatrix}, n \in \mathbb{N} \quad (5.8)$$

Nach der Konditionsformel (2.22) im Skript, gilt für die Kondition der Gesamtiteration

$$\tilde{\kappa}(0) = 1 + \sum_{j=1}^n \prod_{v=j}^n \hat{\kappa}_v$$

Skript (2.22)

mit

$$\hat{\kappa}_v := \hat{\kappa}(F_1, x^{(n)}) \stackrel{(5.8)}{=} \hat{\kappa}(F_1, \begin{pmatrix} (\frac{1}{2})^{n-1} \\ (\frac{1}{2})^{2(n-1)} \end{pmatrix})$$

$$\stackrel{\cancel{22}}{=} \frac{\left| 2 - \frac{7}{4} \cdot \left(\frac{1}{2}\right)^{2(n-1)} \right| + \left| \left(\frac{1}{2}\right)^{2(n-1)} \cdot 2 \right|}{\left| \left(\frac{1}{2}\right)^2 \cdot \left(\frac{1}{2}\right)^{2(n-1)} \right|} = \frac{2 \cdot \frac{7}{4} + 2}{\frac{1}{4}} = \frac{\frac{22}{4} + 2}{\frac{1}{4}} = \frac{22}{1} \quad n \in \mathbb{N}$$

$\cancel{22}$

und demzufolge

$$\tilde{\kappa}(0) = 1 + \sum_{j=1}^n \prod_{v=j}^n \hat{\kappa}_v \xrightarrow{n \rightarrow \infty} \infty$$

$$= 1 + \sum_{j=1}^n 22^j = \frac{22^{n+1} - 1}{22 - 1}$$

AUFGABE 6:

ZU (a): Betrachte

$$F: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ mit } F\left(\begin{array}{c} p \\ q \end{array}\right) = -p + \sqrt{p^2 + q^2} \quad (6.1)$$

wobei $p, q \in \mathbb{R}$ mit $p < 0 < q$. Wegen

$$\frac{\partial F}{\partial p}(p, q) \stackrel{(6.1)}{=} -1 + \frac{1}{2}(p^2 + q^2)^{-\frac{1}{2}} \cdot 2p = -1 + \frac{p}{\sqrt{p^2 + q^2}} = \frac{p - \sqrt{p^2 + q^2}}{\sqrt{p^2 + q^2}} \quad (6.2)$$

$$\frac{\partial F}{\partial q}(p, q) \stackrel{(6.1)}{=} \frac{1}{2}(p^2 + q^2)^{-\frac{1}{2}} = \frac{1}{2\sqrt{p^2 + q^2}}$$

erhalten wir aus (5.1) (mit $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} p \\ q \end{pmatrix}$) die relative komponentenweise Kondition

$$\hat{\chi}(F_1 x) \stackrel{(6.1)}{=} \left| \frac{\partial F}{\partial p}(p, q) \right| \cdot \frac{|p|}{|F(p, q)|} + \left| \frac{\partial F}{\partial q}(p, q) \right| \cdot \frac{|q|}{|F(p, q)|}$$

$$\stackrel{(6.2)}{=} \left| \frac{1}{\sqrt{p^2 + q^2}} \right| \cdot \left[\frac{|p - \sqrt{p^2 + q^2}| \cdot |p|}{|p - \sqrt{p^2 + q^2}|} + \left| \frac{1}{2} \right| \cdot \frac{|q|}{|p - \sqrt{p^2 + q^2}|} \right]$$

$$\stackrel{p < 0}{=} \left| \frac{1}{\sqrt{p^2 + q^2}} \right| \cdot \left[\frac{(-p + \sqrt{p^2 + q^2}) \cdot (-p) \cdot 2 + q}{2 \cdot (-p + \sqrt{p^2 + q^2})} \right]$$

~~$$\stackrel{q > 0}{=} \frac{1}{\sqrt{p^2 + q^2}} \cdot \left[-p + \frac{q}{2(-p + \sqrt{p^2 + q^2})} \right]$$~~

~~$$\cdot \left(\frac{p + \sqrt{p^2 + q^2}}{p + \sqrt{p^2 + q^2}} \right) \stackrel{= 1}{=} \frac{1}{\sqrt{p^2 + q^2}} \cdot \left[-p + \frac{q \cdot (p + \sqrt{p^2 + q^2})}{2 \cdot q} \right]$$~~

~~$$\stackrel{= 1}{=} \frac{1}{2} \cdot \underbrace{\frac{-p}{\sqrt{p^2 + q^2}}}_{\hookrightarrow 1} + \frac{1}{2} < \frac{1}{2} + \frac{1}{2} = 1$$~~

Begründung: $p < 0 \Rightarrow p^2 > 0$ und $-p > 0$

$$q > 0 \Rightarrow p^2 + q^2 > p^2 > 0$$

q > 0 zweite Zeile zuvor

f. i. streng monoton
wachsend $\Rightarrow \sqrt{p^2 + q^2} > \sqrt{p^2} = -p > 0$

$$\Rightarrow 1 = \frac{\sqrt{p^2 + q^2}}{\sqrt{p^2 + q^2}} > \frac{-p}{\sqrt{p^2 + q^2}} > 0$$

Bemerkung: Mit der Kondition für $p > 0$, die bereits in der Vorlesung gezeigt wurde erhalten wir insgesamt:

$$\hat{\chi}(F_1 x) \begin{cases} \leq 1, & p \leq 0 \\ \leq 2, & p > 0 \end{cases}$$

ZU (b): Betrachte

$$F^1: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ mit } F^1\left(\begin{array}{c} p \\ q \end{array}\right) = \left(\begin{array}{c} p \\ \sqrt{p^2 + q^2} \end{array}\right) =: \begin{pmatrix} F_1^1(p, q) \\ F_2^1(p, q) \end{pmatrix}$$

$$F^2: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ mit } F^2\left(\begin{array}{c} y_1 \\ y_2 \end{array}\right) = y_2 - y_1$$

wobei $p, q \in \mathbb{R}$ mit $p < 0 < q$. Wegen

$$\frac{\partial F_1^1}{\partial p}(p, q) = 1, \quad \frac{\partial F_1^1}{\partial q}(p, q) = 0, \quad \frac{\partial F_2^1}{\partial p}(p, q) = \frac{p}{\sqrt{p^2 + q^2}}, \quad \frac{\partial F_2^1}{\partial q}(p, q) = \frac{1}{2\sqrt{p^2 + q^2}} \quad (6.3)$$

erhalten wir aus (5.1)

$$\begin{aligned} \hat{K}_1 := \hat{K}(\mathbb{F}^1, (\frac{p}{q})) &= \max_{(5.1)} \left\{ \left| \frac{\partial \mathbb{F}^1}{\partial p}(p, q) \right| \cdot \frac{|p|}{|\mathbb{F}^1(p, q)|} + \left| \frac{\partial \mathbb{F}^1}{\partial q}(p, q) \right| \cdot \frac{|q|}{|\mathbb{F}^1(p, q)|} \right\} \\ &\quad \left| \frac{\partial \mathbb{F}^1_2}{\partial p}(p, q) \right| \cdot \frac{|p|}{|\mathbb{F}^1_2(p, q)|} + \left| \frac{\partial \mathbb{F}^1_2}{\partial q}(p, q) \right| \cdot \frac{|q|}{|\mathbb{F}^1_2(p, q)|} \} \\ &\stackrel{(6.3)}{=} \max \left\{ \underbrace{|1| \cdot \frac{|p|}{|p|}}_{=1} + \underbrace{|0| \cdot \frac{|q|}{|p|}}_{=0}, \left| \frac{p}{\sqrt{p^2+q}} \right| \cdot \frac{|p|}{|\mathbb{F}^1(p, q)|} + \left| \frac{1}{2\sqrt{p^2+q}} \right| \cdot \frac{|q|}{|\mathbb{F}^1(p, q)|} \right\} \\ &= \max \left\{ 1, \frac{p^2}{p^2+q} + \frac{q}{2(p^2+q)} \right\} \\ &= \max \left\{ 1, \underbrace{\frac{2p^2+q}{2p^2+2q}}_{< 1 \text{ (denn: } q < 2q \uparrow \atop q > 0)} \right\} = 1 \end{aligned}$$

Weiter erhalten wir wegen

$$\frac{\partial \mathbb{F}^2}{\partial y_1}(y_1, y_2) = -1, \quad \frac{\partial \mathbb{F}^2}{\partial y_2}(y_1, y_2) = 1 \quad (6.4)$$

die Konditionszahl

$$\begin{aligned} \hat{K}(\mathbb{F}^2, (\frac{p}{q})) &= \left| \frac{\partial \mathbb{F}^2}{\partial y_1}(y_1, y_2) \right| \cdot \frac{|y_1|}{|\mathbb{F}^2(y_1, y_2)|} + \left| \frac{\partial \mathbb{F}^2}{\partial y_2}(y_1, y_2) \right| \cdot \frac{|y_2|}{|\mathbb{F}^2(y_1, y_2)|} \\ &\stackrel{(6.4)}{=} |-1| \cdot \frac{|y_1|}{|y_2 - y_1|} + |1| \cdot \frac{|y_2|}{|y_2 - y_1|} \\ &= \frac{|y_1| + |y_2|}{|y_2 - y_1|} \end{aligned} \quad (6.5)$$

und somit

$$\begin{aligned} \hat{K}_2 := \hat{K}(\mathbb{F}^2, \mathbb{F}^1(\frac{p}{q})) &= \frac{|p| + \sqrt{p^2+q}|}{|1-p+\sqrt{p^2+q}|} = \frac{-p + \sqrt{p^2+q}}{-p + \sqrt{p^2+q}} = 1 \\ (6.5) \text{ mit } \left(\begin{array}{c} y_1 \\ y_2 \end{array} \right) = \mathbb{F}^1(\frac{p}{q}) &= \left(\begin{array}{c} p \\ -\sqrt{p^2+q} \end{array} \right) \quad \begin{array}{l} p < 0 \\ \Rightarrow -p > 0 \\ \Rightarrow |p| = -p \end{array} \end{aligned}$$

Bemerkung:

$$\begin{aligned} \hat{K}_1 = \hat{K}(\mathbb{F}^1, (\frac{p}{q})) &= \begin{cases} 1, & p \leq 0 \\ 1, & p \geq 0 \end{cases} \quad \begin{array}{l} \text{(siehe Aufgabe)} \\ \text{(siehe Vorlesung)} \end{array} \\ \hat{K}_2 = \hat{K}(\mathbb{F}^2, \mathbb{F}^1(\frac{p}{q})) &= \begin{cases} 1, & p \leq 0 \\ \geq \frac{4p^2}{q}, & p \geq 0 \end{cases} \quad \begin{array}{l} \text{(siehe Aufgabe)} \\ \text{(siehe Vorlesung)} \end{array} \end{aligned}$$

Zu (c): Betrachte

$$\mathbb{F}^1: \mathbb{R}^2 \rightarrow \mathbb{R}^3 \quad \text{mit} \quad \mathbb{F}^1\left(\begin{array}{c} p \\ q \end{array}\right) = \left(\begin{array}{c} p \\ q \\ \frac{p}{\sqrt{p^2+q}} \end{array}\right)$$

$$\mathbb{F}^2: \mathbb{R}^3 \rightarrow \mathbb{R} \quad \text{mit} \quad \mathbb{F}^2\left(\begin{array}{c} y_1 \\ y_2 \\ y_3 \end{array}\right) = \frac{y_2}{y_1+y_3}$$

wobei $p, q \in \mathbb{R}$ mit $p < 0 < q$. Wegen

$$\begin{aligned} \frac{\partial \mathbb{F}^1_1}{\partial p}(p, q) &= 1, \quad \frac{\partial \mathbb{F}^1_2}{\partial p}(p, q) = 0, \quad \frac{\partial \mathbb{F}^1_3}{\partial p}(p, q) = \frac{p}{\sqrt{p^2+q}} \\ \frac{\partial \mathbb{F}^1_1}{\partial q}(p, q) &= 0, \quad \frac{\partial \mathbb{F}^1_2}{\partial q}(p, q) = 1, \quad \frac{\partial \mathbb{F}^1_3}{\partial q}(p, q) = \frac{1}{2\sqrt{p^2+q}} \end{aligned} \quad (6.6)$$

erhalten wir aus (5.1)

$$\begin{aligned}
 \hat{F}_1 &:= \hat{F}(F^1, (p, q)) = \max_{(p, q)} \left\{ \left| \frac{\partial F_1^1}{\partial p} (p, q) \right| \cdot \frac{|p|}{|F_1(p, q)|} + \left| \frac{\partial F_1^1}{\partial q} (p, q) \right| \cdot \frac{|q|}{|F_1(p, q)|}, \right. \\
 &\quad \left| \frac{\partial F_1^2}{\partial p} (p, q) \right| \cdot \frac{|p|}{|F_1^2(p, q)|} + \left| \frac{\partial F_1^2}{\partial q} (p, q) \right| \cdot \frac{|q|}{|F_1^2(p, q)|}, \\
 &\quad \left. \left| \frac{\partial F_1^3}{\partial p} (p, q) \right| \cdot \frac{|p|}{|F_1^3(p, q)|} + \left| \frac{\partial F_1^3}{\partial q} (p, q) \right| \cdot \frac{|q|}{|F_1^3(p, q)|} \right\} \\
 &= \max_{(p, q)} \left\{ \underbrace{|1| \cdot \frac{|p|}{|p|}}_{=1} + \underbrace{|0| \cdot \frac{|q|}{|p|}}_{=0}, \underbrace{|0| \cdot \frac{|p|}{|q|}}_{=0} + \underbrace{|1| \cdot \frac{|q|}{|q|}}_{=1}, \right. \\
 &\quad \left. \left| \frac{p}{\sqrt{p^2+q}} \right| \cdot \frac{|p|}{\sqrt{p^2+q}} + \left| \frac{1}{2\sqrt{p^2+q}} \right| \cdot \frac{|q|}{\sqrt{p^2+q}} \right\} \\
 &= \max \left\{ 1, 1, \underbrace{\frac{2p^2+q}{2p^2+2q}}_{< 1 \text{ (da } q \uparrow \text{, } q > 0\text{)}} \right\} = 1
 \end{aligned}$$

Weiter erhalten wir wegen

$$\frac{\partial F^2}{\partial y_1}(y_1, y_2, y_3) = \frac{-y_2}{(y_1+y_3)^2}, \quad \frac{\partial F^2}{\partial y_2}(y_1, y_2, y_3) = \frac{1}{y_1+y_3}, \quad \frac{\partial F^2}{\partial y_3}(y_1, y_2, y_3) = \frac{-y_2}{(y_1+y_3)^2} \quad (6.7)$$

die Konditionstahl

$$\begin{aligned}
 & \hat{\chi}(F^2, \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}) = \left| \frac{\partial F^2}{\partial y_1} (y_1, y_2, y_3) \right| \cdot \frac{|y_1|}{|F^2(y_1, y_2, y_3)|} + \left| \frac{\partial F^2}{\partial y_2} (y_1, y_2, y_3) \right| \cdot \frac{|y_2|}{|F^2(y_1, y_2, y_3)|} \\
 & \quad + \left| \frac{\partial F^2}{\partial y_3} (y_1, y_2, y_3) \right| \cdot \frac{|y_3|}{|F^2(y_1, y_2, y_3)|} \\
 & \stackrel{(6.7)}{=} \left| \frac{-y_2}{(y_1+y_3)^2} \right| \cdot \frac{|y_1| \cdot |y_1+y_3|}{|y_2|} + \left| \frac{1}{y_1+y_3} \right| \cdot \frac{|y_2| \cdot |y_1+y_3|}{|y_2|} \\
 & \quad + \left| \frac{-y_2}{(y_1+y_3)^2} \right| \cdot \frac{|y_3| \cdot |y_1+y_3|}{|y_2|} \\
 & = \left| \frac{y_1}{y_1+y_3} \right| + 1 + \left| \frac{y_3}{y_1+y_3} \right| = 1 + \frac{|y_1| + |y_3|}{|y_1+y_3|} \quad (6.7)
 \end{aligned}$$

and sonit

$$\text{d} \text{ sonst } \hat{x}_2 := \hat{x}(T^2, T^1(q)) = 1 + \frac{|p| + \sqrt{|p^2 + q|}}{|p + \sqrt{p^2 + q}|} = 1 + \frac{-p + \sqrt{p^2 + q}}{|p + \sqrt{p^2 + q}|}$$

$$(6.8) \text{ mit } \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = T^{-1}\left(\begin{pmatrix} p \\ q \end{pmatrix}\right) \quad p < 0$$

$$= 1 + \frac{(-p + \sqrt{p^2+q})^2}{|q|} = 1 + \frac{(-p + \sqrt{p^2+q})^2}{q}$$

$$\cdot \begin{pmatrix} -p + \sqrt{p^2+q} \\ -p + \sqrt{p^2+q} \end{pmatrix} \geq \frac{(-p)^2 + 2 \cdot (-p) \cdot \sqrt{p^2+q} + p^2+q}{q} \geq \frac{4p^2}{q}$$

$$(-p)^2 + 2 \cdot (-p) \cdot \underbrace{\sqrt{p^2+q}}_{\geq (-p) > 0} + p^2 + q \underbrace{> 0}_{> 0} > 4(-p)^2 = 4p^2$$

Bemerkung:

$$\hat{f}_1 = \hat{f}(\mathbb{F}, \binom{p}{q}) = \begin{cases} 1 & , p \leq 0 \quad (\text{siehe Aufgabe}) \\ 1 & , p \geq 0 \quad (\text{siehe Vorlesung}) \end{cases}$$

$$\hat{X}_2 = \hat{X}(F^2, F^1(\frac{p}{q})) = \begin{cases} \geq \frac{4p^2}{q}, & p \leq 0 \quad (\text{siehe Aufgabe}) \\ \leq 2, & p \geq 0 \quad (\text{siehe Vorlesung}) \end{cases}$$

Welcher Algorithmus ist in diesem Fall gutartig?

Der Algorithmus 1 ist

- gutartig $\Leftrightarrow p \leq 0$
- nicht-gutartig $\Leftrightarrow p > 0$.

Der Algorithmus 2 ist

- gutartig $\Leftrightarrow p \geq 0$
- nicht-gutartig $\Leftrightarrow p < 0$.

zu (d): Starte zunächst die Rundungsfehler GUI

zu Algorithmus 1:

- ①: Wähle „Mantissenlänge 8“ unter allgemeine Parameter
- ②: Wähle unter Standardbeispiele die „quad. Gleichung I“
- ③: Unter Funktionen gebe im Feld „Anfangswerte“ -1^{10} (bzw. -1^{100} , -1^{1000}) ein.
- ④: Wähle den Button Berechnen.

zu Algorithmus 2:

- ①: siehe oben
- ②: Wähle unter Standardbeispiele die „quad. Gleichung II“
- ③: siehe oben
- ④: siehe oben

Screenshots:

- ①: Gebe im Terminal „Gimp“ ein.
- ②: In Gimp wähle in der Menüleiste „File → Acquire → Screenshot“ und Klicke auf „Snap“
- ③: Wähle nun das Fenster aus, vor dem ein Screenshot gemacht werden soll.
- ④: Wähle in Gimp nun „drucken“ und speichere das Bild.

Interpretation der Ergebnisse:

Wie wir im theoretischen Teil dieser Aufgabe festgestellt haben, ist für $p < 0$ der Algorithmus 1 gutartig und der Algorithmus 2 nicht-gutartig. Dieses Ergebnis wird von den numerischen Untersuchungen bestätigt: Dies erkennen wir, wenn wir uns den zugehörigen Fehler ansehen: ($q=1$, Mantissenlänge 8)

	Algorithmus 1	Algorithmus 2
$p = -10$	10^{-7}	10^{-4}
$p = -100$	10^{-7}	10^{-3}
$p = -1000$	10^{-10}	10^{-4}

Größenordnung des Fehlers

wie wir feststellen ist der Fehler beim Algorithmus 1 erheblich kleiner.

Algorithmus 1:

$p = -10$



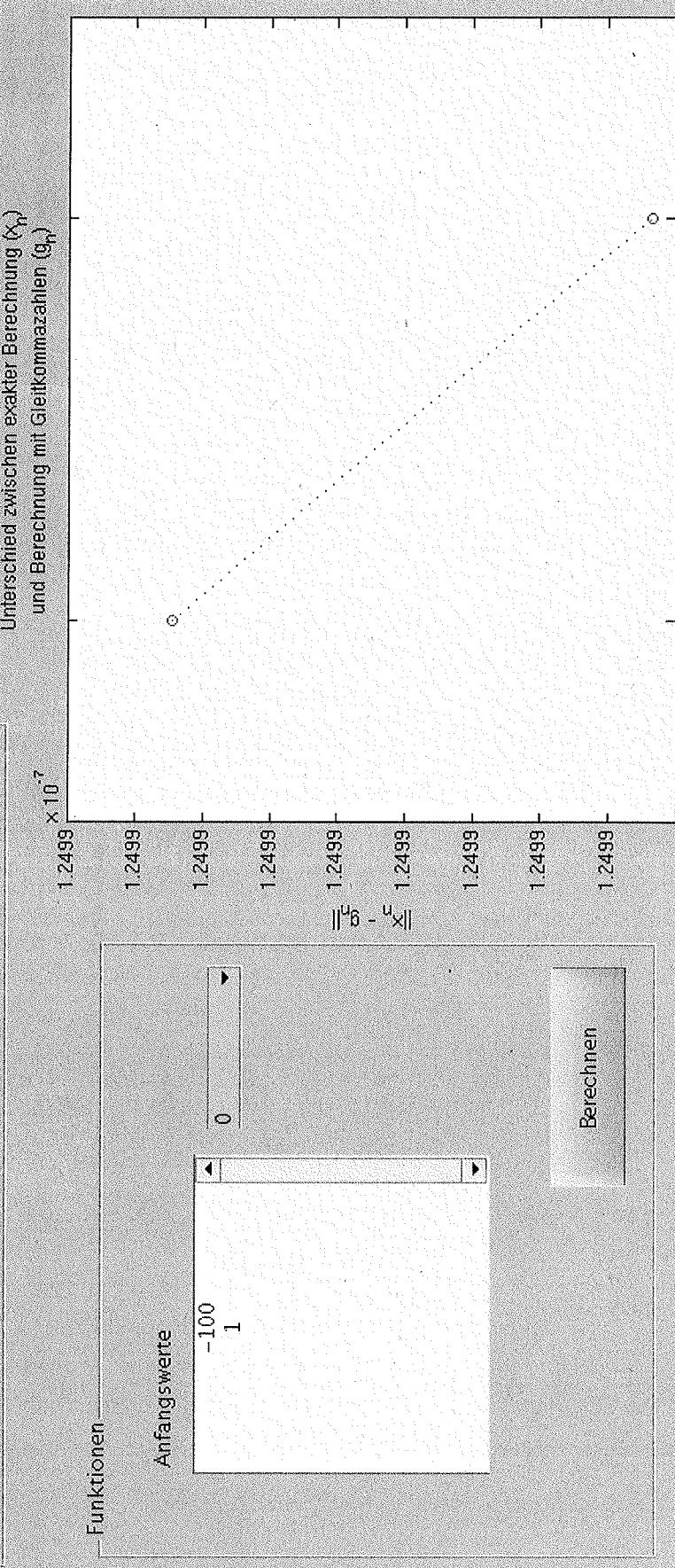
Algorithmus 1:

$P = 100$

Rundungsfehler

allgemeine Parameter	
maximale Anzahl an Schritten:	2
Ausgabe alle	1 Schritte.
iterierte Funktion	Benutzerdef. Fkt.

Standardbeispiele	
harmonische Reihe	
quad. Gleichung I	
quad. Gleichung II	



Wertetabelle

Algorithmus 1:

$p = -1000$

Rundungsfehler

allgemeine Parameter	
maximale Anzahl an Schritten:	2
Mantissaenlaenge:	8
Ausgabe alle Schritte	1
iterierte Funktion	
Benutzerdef. Fkt.	

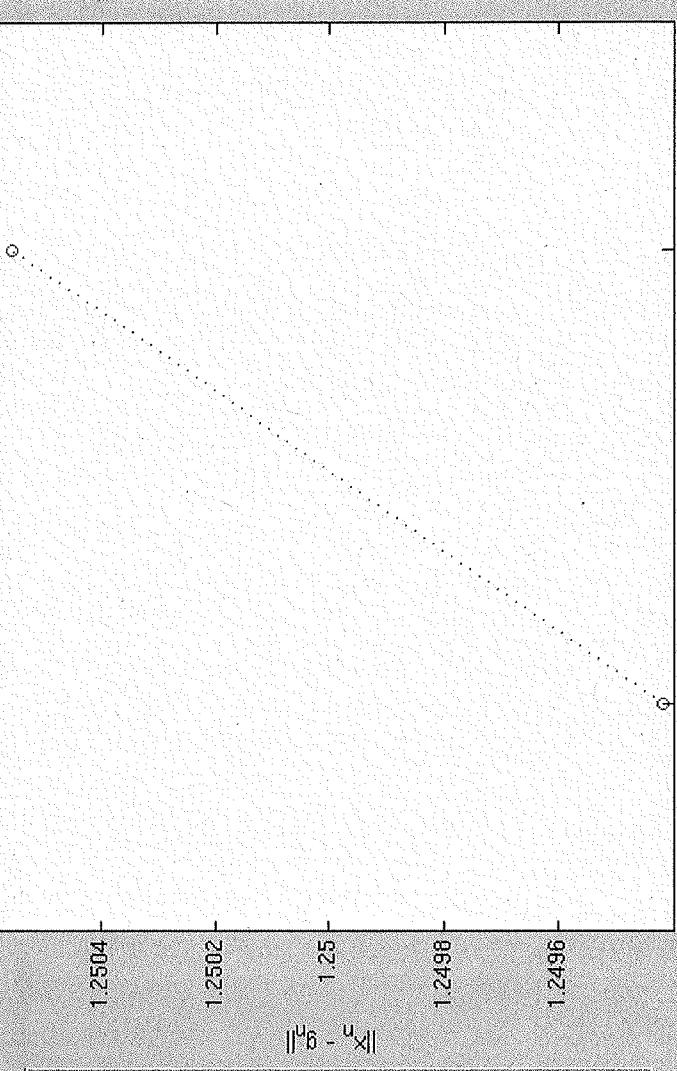
Standardbeispiele	
harmonische Reihe	
quad. Gleichung I	
quad. Gleichung II	

Status

bereit

Unterschied zwischen exakter Berechnung (X_n) und Berechnung mit Gleitkommazahlen (y_n)

1.2506×10^{-10}



Wertertabelle

Anfangswerte	-1000	0
	1	
Funktionen	$\frac{1}{x}$	x^2
Berechnen		

Algorithmus 2:

$$P = -10$$

Rundungsfehler

allgemeine Parameter

maximale Anzahl an Schritten: 2

Mantissaenlaenge: 8

Ausgabe alle 1 Schritte.

iterierte Funktion

Benutzerdef. Fkt.

Standardbeispiele

harmonische Reihe

quad. Gleichung I

quad. Gleichung II

Status

bereit

Unterschied zwischen exakter Berechnung (x_n) und Berechnung mit Gleitkommazahlen (q_n)

$\times 10^{-4}$

Anfangswerte

-10	0	►
1		

Berechnen



Wertetabelle

Algorithmus 2:

$$d = -100$$

Rundungsfehler

allgemeine Parameter

maximale Anzahl an Schritten:

Mantissenlänge:

Ausgabe alle Schritte.

Iterierte Funktion

Benutzerdef. Fkt.

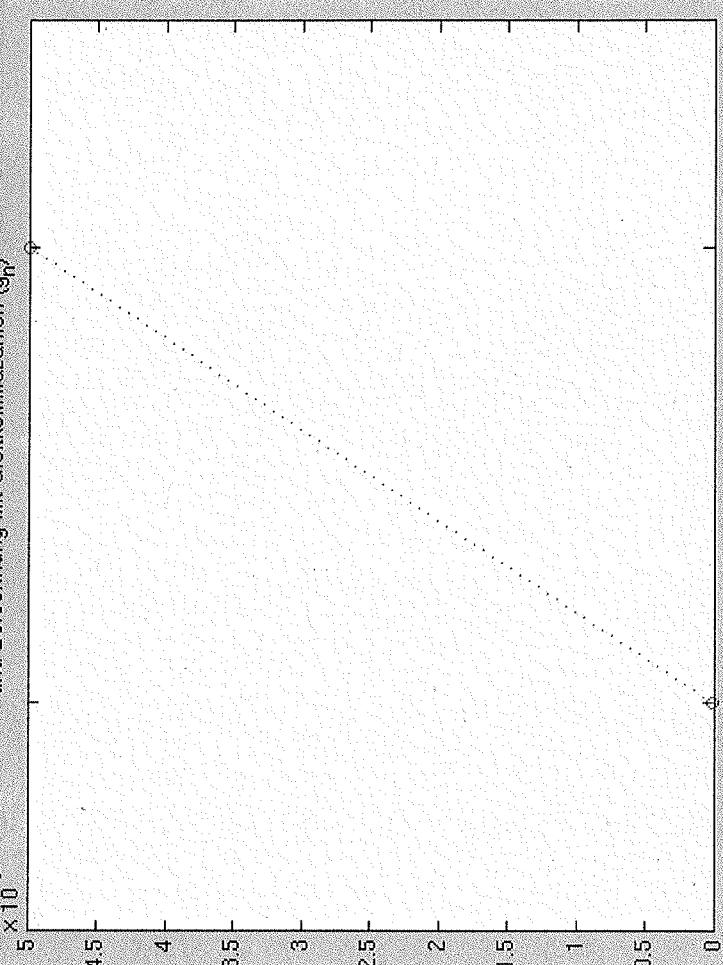
Standardbeispiele

Status

harmonische Reihe

quad. Gleichung I

quad. Gleichung II



Funktionen

Anfangswerte

-100	0
1	

Berechnen

Wertetabelle

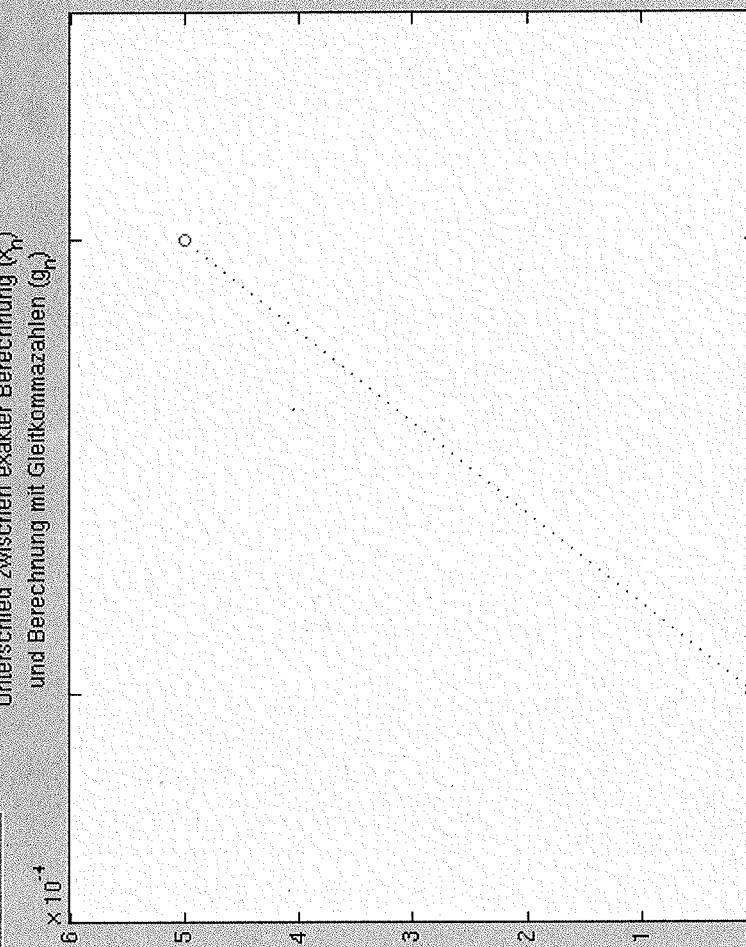
2

Algorithmus 2:

$p = -1000$

Rundungsfehler

allgemeine Parameter	
maximale Anzahl an Schritten:	2
Ausgabe alle	1
Mantissaenlaenge:	8
iterierte Funktion	
Benutzerdef. Fkt.	
Standardbeispiele	
harmonische Reihe	
quad. Gleichung I	
quad. Gleichung II	



Funktionen	$\frac{1}{x}$
Anfangswerte	-1000
Berechnen	0

