

Numerische Mathematik für das Lehramt

Skriptum zur Vorlesung im SS 2009 †

PD Dr. Markus Neher
Universität Karlsruhe (TH)

Forschungsuniversität · gegründet 1825

Institut für Angewandte und Numerische Mathematik

10. Juli 2009

†©2007-2009 by Markus Neher. Dieses Skriptum ist urheberrechtlich geschützt. Weiterverbreitung und Einsatz in anderen Lehrveranstaltungen (auch von Teilen des Skriptums) ist ohne vorherige schriftliche Genehmigung des Autors untersagt. Insbesondere ist es nicht gestattet, das Skriptum oder Teile davon in elektronischer Form im Internet zugänglich zu machen.

Inhaltsverzeichnis

1	Einführung	7
1.1	Symbolisches und numerisches Rechnen	8
1.2	Gleitpunktzahlen	9
1.3	Gleitpunktarithmetik	10
1.4	Algorithmen	11
1.5	Maple	15
1.5.1	Schleifen in Maple	15
1.5.2	Auswahanweisungen in Maple	16
1.5.3	Prozeduren in Maple	18
1.6	Landau-Symbole	22
1.7	Ziele dieser Vorlesung	23
2	Iterationsverfahren	25
2.1	Fixpunktiteration	27
2.2	Lokale Fixpunktsätze	31
2.3	Relaxation	32
2.4	Das Newton-Verfahren	34
2.5	Verwandte Iterationsverfahren	36
2.5.1	Vereinfachtes Newton-Verfahren	36
2.5.2	Sekantenverfahren	36
2.5.3	Das Bisektionsverfahren	36
2.6	Vektor- und Matrixnormen	38
2.7	Iterationsverfahren im \mathbb{R}^n	40
3	Lineare Gleichungssysteme	47
3.1	Gauß-Elimination ohne Zeilentausch: <i>LU</i> -Zerlegung	47
3.1.1	Eliminationsmatrizen	48
3.1.2	<i>LU</i> -Zerlegung	49
3.1.3	Aufwand des Gauß-Algorithmus	52
3.2	Gauß-Elimination mit Zeilentausch: <i>PALU</i> -Zerlegung	53

3.2.1	Permutationsmatrizen	53
3.2.2	Pivotisierung	56
3.3	Fehlerabschätzungen für lineare Gleichungssysteme	57
3.4	Iterationsverfahren für lineare Gleichungssysteme	60
3.4.1	Fixpunktiteration für lineare Gleichungssysteme	60
3.4.2	Gesamt- und Einzelschrittverfahren	62
3.4.3	Die Methode des steilsten Abstiegs	66
3.4.4	Das cg-Verfahren	70
3.4.5	Vorkonditionierung beim cg-Verfahren	74
3.5	Die QR -Zerlegung einer Matrix	75
3.5.1	Orthogonale Matrizen	76
3.5.2	Gram-Schmidt-Orthogonalisierung	76
3.5.3	Reduzierte QR -Zerlegung	78
3.5.4	Householder-Matrizen	79
3.5.5	QR -Zerlegung durch Householder-Transformationen	81
3.6	Über- und unterbestimmte lineare Gleichungssysteme	84
3.6.1	Überbestimmte lineare Gleichungssysteme	84
3.6.2	Unterbestimmte lineare Gleichungssysteme	87
4	Das Eigenwertproblem für Matrizen (entfällt 2009)	89
5	Approximation und Interpolation	91
5.1	Approximation mit Taylor-Polynomen	93
5.1.1	Approximationsaufgaben zum Restglied	95
5.1.2	Anwendung: Berechnung transzendenter Funktionen	96
5.2	Polynom-Interpolation	97
5.2.1	Einführung	97
5.2.2	Polynom-Interpolation	98
5.2.3	Interpolationsfehler der Polynom-Interpolation	102
5.2.4	Polynom-Interpolation mit Tschebyscheff-Stützstellen	104
5.2.5	Hermite-Interpolation	107
5.3	Spline-Interpolation	107
5.3.1	Stückweise lineare Interpolation	108
5.3.2	Kubische Spline-Interpolation	108
5.3.3	Minimierungseigenschaft	112
5.3.4	Interpolationsfehler der kubischen Spline-Interpolation	113
5.3.5	Anwendung: Spline-Interpolation geschlossener Kurven	114
5.4	Trigonometrische Interpolation	115
5.4.1	Fourier-Reihen	115

5.4.2	Reelle diskrete Fourier-Transformation	116
5.5	Approximation nach der Methode der kleinsten Quadrate	118
5.6	Nichtlineare Ausgleichsrechnung	121
5.6.1	Newton-Verfahren	121
5.6.2	Gauß-Newton-Verfahren	123
6	Numerische Integration	125
6.1	Newton-Cotes-Formeln	125
6.2	Summierte Quadraturformeln	129
6.3	Extrapolation mit dem Romberg-Verfahren	130
6.4	Gauß-Quadratur	133
6.4.1	Orthogonale Polynome	135
6.4.2	Stützstellen und Gewichte bei der Gauß-Quadratur	137
7	Numerische Behandlung gewöhnlicher Differentialgleichungen (entfällt 2009)	141
8	Gleitpunktrechnung, Kondition, Stabilität	143
8.1	Gleitpunktzahlen	143
8.2	Gleitpunktarithmetik	147
8.2.1	Fehlerfortpflanzung bei den arithmetischen Grundoperationen	147
8.3	Die Kondition eines mathematischen Problems	151
8.3.1	Fehlerfortpflanzung	151
8.3.2	Stabilität	153
8.4	Beispiele	156

Kapitel 8

Gleitpunktrechnung, Kondition, Stabilität

8.1 Gleitpunktzahlen

Numerische Berechnungen werden in der Regel nicht mit Variablen, sondern mit Zahlen durchgeführt. Die Darstellung von Zahlen ist dabei von entscheidender Bedeutung. Bereits die Konstruktion von Algorithmen kann vom gewählten Zahlenformat abhängen. Beispielsweise ist der heute gebräuchliche Algorithmus zur Division von Zahlen im römischen Zahlensystem praktisch nicht durchführbar. Noch stärker werden die Rechengenauigkeit und die Rechengeschwindigkeit durch das Zahlenformat beeinflusst.

Auch die Anschaulichkeit von Zahlen hängt wesentlich von der gewählten Darstellung ab. So wird $\frac{355}{113}$ durch 3.1416 nur approximiert, aber unter der Dezimalnäherung können sich die meisten Menschen die Zahl besser vorstellen. Auch im Fall großer Zahlen wie der Avogadro-Konstante

$$N_A = 602214175141592653589793 \approx 6.022 \cdot 10^{23}$$

ist die Angabe des exakten Werts wenig hilfreich. Die wesentlichen Informationen über den Wert der Zahl lassen sich aus der Dezimalnäherung einfacher und schneller ablesen.

In Anwendungen sind gemessene oder berechnete Größen üblicherweise nur ungefähr bekannt. Dies gilt sogar für mathematische Konstanten wie π oder mathematische Symbole wie $\sqrt{2}$, die lediglich Namen von Zahlen darstellen, aber nichts über den Wert der repräsentierten Zahl aussagen. Daher ist es üblich, bei praktischen Rechnungen sogenannte Gleitpunktzahlen der Bauart

$$\pm m \cdot b^e \tag{8.1}$$

mit der Mantisse m , der Basis b und dem Exponenten e zu verwenden. Dabei ist die Basis b des b -adischen Zahlensystems eine natürliche Zahl größer als 1, der Exponent e eine ganze Zahl und die Mantisse m eine rationale Zahl mit endlicher b -adischer Darstellung. Beispiele solcher Zahlen sind

$$3.1416 \cdot 10^0, \quad -44.123 \cdot 5^3, \quad 0.0001101 \cdot 2^{-1011}.$$

Die angegebene Gleitpunktdarstellung (8.1) einer Zahl ist nicht eindeutig. Beispielsweise gilt

$$3.1416 \cdot 10^0 = 3141600.00 \cdot 10^{-6} = 0.00000000000000000000000031416 \cdot 10^{22}.$$

Beim Rechnen ist diese Mehrdeutigkeit unpraktisch. Deshalb fordern wir in (8.1) zusätzlich, dass – außer für die Zahl Null – die Mantisse im Intervall $[1, b)$ liegt. Durch eine Anpassung des

Exponenten (bei der im Gegenzug der Dezimalpunkt gleitet) existiert für jede Gleitpunktzahl $z \neq 0$ eine eindeutig bestimmte normalisierte Darstellung

$$z = \pm m_1.m_2m_3 \dots m_l \cdot b^e$$

mit $m_1, m_2, \dots, m_l \in \{0, 1, \dots, b-1\}$, $m_1 \neq 0$, $e \in \mathbb{Z}$. Die Menge aller normalisierten Gleitpunktzahlen zur Basis b mit fester Mantissenlänge $l \geq 1$ und Exponentenbereich $e_{\min} \leq e \leq e_{\max}$, vereinigt mit der Zahl Null, bildet das Gleitpunktsystem $S_{\text{norm}}(b, l, e_{\min}, e_{\max}) \subset \mathbb{Q}$.

Beispiel 8.1 Das normalisierte dezimale Gleitpunktsystem $S_{\text{norm}}(10, 4, 0, 9)$ mit vier dezimalen Mantissenstellen und einer Dezimalstelle für den Exponenten umfasst die folgenden Zahlen:

$$\begin{aligned} &0, \\ &+1.000 \cdot 10^0, -1.000 \cdot 10^0, +1.001 \cdot 10^0, -1.001 \cdot 10^0, \dots, +9.999 \cdot 10^0, -9.999 \cdot 10^0, \\ &+1.000 \cdot 10^1, -1.000 \cdot 10^1, +1.001 \cdot 10^1, -1.001 \cdot 10^1, \dots, +9.999 \cdot 10^1, -9.999 \cdot 10^1, \\ &\vdots \\ &+1.000 \cdot 10^9, -1.000 \cdot 10^9, +1.001 \cdot 10^9, -1.001 \cdot 10^9, \dots, +9.999 \cdot 10^9, -9.999 \cdot 10^9. \end{aligned}$$

Die kleinste darstellbare positive normalisierte Gleitpunktzahl heißt *mininorm*, die größte heißt *maxreal*. In $S_{\text{norm}}(10, 4, 0, 9)$ gilt:

$$\text{mininorm} = 1.000 \cdot 10^0, \quad \text{maxreal} = 9.999 \cdot 10^9$$

Man beachte, dass die Gleitpunktzahlen in $S_{\text{norm}}(10, 4, 0, 9)$ nicht gleichabständig verteilt sind. Der Abstand der Nachbarzahlen $9.998 \cdot 10^9$ und $9.999 \cdot 10^9$ beträgt 10^6 , der Abstand der Nachbarzahlen $1.000 \cdot 10^0$ und $1.001 \cdot 10^0$ nur 10^{-3} . Normalisierte Gleitpunktzahlen besitzen aber die wertvolle Eigenschaft, dass der relative Fehler zweier benachbarter Gleitpunktzahlen zueinander ungefähr konstant ist.

Wir diskutieren nun, wie Zahlen möglichst effizient auf einem Computer gespeichert werden können. Halten wir die Basis b und die Mantissenlänge l fest, müssen die folgenden Informationen zu einer Gleitpunktzahl gespeichert werden:

1. Das Vorzeichen der Zahl,
2. die Mantisse,
3. das Vorzeichen des Exponenten,
4. die Ziffern des Exponenten.

Mit einem kleinen Trick kann man die separate Speicherung des Vorzeichens des Exponenten umgehen. Man speichert auf dem Rechner in der Praxis nur nichtnegative Exponenten (mit l_e Stellen zur Basis b), interpretiert die gespeicherte Zahl aber so, dass vom Exponenten ein fester Basiswert subtrahiert wird.

Beispiel 8.2 Gegeben sei das dezimale normalisierte Gleitpunktsystem $S_{\text{norm}}(10, 4, -5, 4)$ mit vier dezimalen Mantissenstellen und einer Dezimalstelle für den Exponenten nach dem folgenden Schema:

$$\boxed{\pm \mid e \mid m_1 \mid m_2 \mid m_3 \mid m_4} \quad e, m_1, m_2, m_3, m_4 \in \{0, 1, \dots, 9\}.$$

Vom gespeicherten Wert e wird bei der Umrechnung in eine Zahl der feste Basiswert 5 abgezogen. Normalisierte Gleitpunktzahlen besitzen also in $S_{\text{norm}}(10, 4, -5, 4)$ die Darstellung

$$\pm m_1.m_2m_3m_4 \cdot 10^{e-5} \quad (m_1 \neq 0, 0 \leq e \leq 9).$$

Die kleinste darstellbare positive normalisierte Gleitpunktzahl ist

$$\text{mininorm} = 1.000 \cdot 10^{-5}$$

(gespeichert wird $e = 0$), die größte ist

$$\text{maxreal} = 9.999 \cdot 10^4$$

(gespeichert wird $e = 9$).

Mit einem zweiten Trick kann man die Menge der darstellbaren Gleitpunktzahlen vergrößern, ohne zusätzlichen Speicherplatz bereitstellen zu müssen. Besonders vorteilhaft ist es, dass dabei die Lücke zwischen der Null und der kleinsten positiven Gleitpunktzahl verkleinert wird. Dazu verzichtet man bei Zahlen, die mit dem kleinsten Exponentenwert gespeichert sind, auf die Normalisierungsbedingung $m_1 \geq 1$ und interpretiert die gespeicherten Ziffern der Mantisse als Nachkommastellen einer nicht normalisierten Gleitpunktzahl. Damit im Gleitpunktsystem keine Lücke entsteht, muss man dann auch die Exponentenanpassung geringfügig ändern.

Beispiel 8.3 Wir betrachten das dezimale Gleitpunktsystem $\mathcal{T} := S(10, 4, -5, 4)$ mit vier dezimalen Mantissenstellen und einer Dezimalstelle für den Exponenten. Für $1 \leq e \leq 9$ interpretieren wir die gespeicherte Zahl z wie in $S_{\text{norm}}(10, 4, -5, 4)$:

$$z = \pm m_1.m_2m_3m_4 \cdot 10^{e-5} \quad (m_1 \neq 0, 1 \leq e \leq 9).$$

Die kleinste darstellbare positive normalisierte Gleitpunktzahl ist nun

$$\text{mininorm} = 1.000 \cdot 10^{-4}$$

(gespeichert wird $e = 1$), die größte ist wie in $S_{\text{norm}}(10, 4, -5, 4)$ die Zahl

$$\text{maxreal} = 9.999 \cdot 10^4$$

(gespeichert wird $e = 9$).

Für $e = 0$ stellt der gespeicherte Wert die nicht normalisierte Gleitpunktzahl

$$z = \pm 0.m_1m_2m_3m_4 \cdot 10^{e-4} \quad (e = 0)$$

dar. Die Zahl Null wird durch $e = m_1 = m_2 = m_3 = m_4 = 0$ dargestellt. Die kleinste darstellbare positive (nicht normalisierte) Gleitpunktzahl ist

$$\text{minreal} := 0.0001 \cdot 10^{-4} = 1.0 \cdot 10^{-8},$$

die größte nicht normalisierte Gleitpunktzahl $0.9999 \cdot 10^{-4}$ schließt an mininorm an. Im Vergleich zu $S_{\text{norm}}(10, 4, -5, 4)$ ist die kleinste positive Gleitpunktzahl in \mathcal{T} um drei Zehnerpotenzen kleiner.

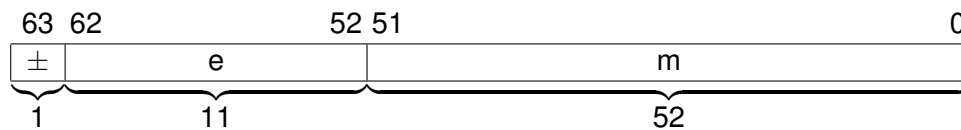
Man beachte insbesondere, dass bei der Neuinterpretation der Zahlen mit Exponentenwert $e = 0$ in \mathcal{T} im Vergleich zu $S_{\text{norm}}(10, 4, -5, 4)$ keine Zahl verloren geht. Die gespeicherten Mantissenwerte m_1, m_2, m_3, m_4 repräsentieren in $S_{\text{norm}}(10, 4, -5, 4)$ für $e = 0$ die Zahl

$$m_1.m_2m_3m_4 \cdot 10^{-5},$$

in \mathcal{T} die identische Zahl

$$0.m_1m_2m_3m_4 \cdot 10^{-4}.$$

Auf modernen Computern werden üblicherweise Gleitpunktzahlen zur Basis $b = 2$ (Dualzahlen, Binärzahlen) entsprechend dem IEEE 754 Double-Standard verwendet:



Eine solche Gleitpunktzahl besteht aus 64 Bits (1 Bit = 0 oder 1):

- dem Vorzeichenbit (Bit 63, “+” oder “-“)
- 11 Exponentenbits (Bits 52-62)
- 52 Mantissenbits (Bits 0-51)

Zur Umrechnung der gespeicherten Bitfolge in eine Zahl werden die folgenden Konventionen verwendet:

- Der Exponent $e = 2^{11} - 1 = 2047$ wird nur für Ausnahmefälle ($m = 0$: signed infinity; $m \neq 0$: NaN = not a number) benutzt.
- Für $1 \leq e \leq 2046$ werden vom Exponenten 1023 subtrahiert; für $e = 0$ werden 1022 subtrahiert.
- Für $1 \leq e \leq 2046$ werden die Mantissenbits als Nachkommastellen der Zahl

$$1.m \cdot 2^{(e-1023)}$$

interpretiert (normalisierte Gleitpunktzahl, bei der die führende 1 nicht gespeichert wird).

- $e = 0, m = 0$ stellt die Zahl 0 dar.
- $e = 0, m \neq 0$ stellt die Zahl

$$0.m \cdot 2^{-1022}$$

dar (nicht normalisierte Gleitpunktzahl im Unterlaufbereich).

Die auf dem Rechner darstellbaren Zahlen heißen Maschinenzahlen.

Somit ergibt sich der folgende Bereich darstellbarer Zahlen:

- Größte darstellbare positive normalisierte Gleitpunktzahl (maxreal):

$$2^{1024} - 2^{1023-52} \approx 1.798 \cdot 10^{308}.$$

- Kleinste darstellbare positive normalisierte Gleitpunktzahl (mininorm):

$$2^{-1022} \approx 2.225 \cdot 10^{-308}.$$

- Kleinste darstellbare positive nicht normalisierte Gleitpunktzahl (minreal):

$$2^{-1022-52} \approx 4.941 \cdot 10^{-324}.$$

Das IEEE 754 Double-Format ist also gegeben durch

$$\text{IEEE 754 Double} = S(2, 53, -1022, 1024).$$

Für die praktischen Beispiele benutzen wir im Folgenden das wesentlich übersichtlichere normalisierte Gleitpunktsystem $\mathcal{S} := S_{\text{norm}}(10, 4, -9, 9)$.

8.2 Gleitpunktarithmetik

Ein endliches Gleitpunktsystem enthält nur endlich viele Zahlen. Nicht exakt darstellbare Zahlen werden zur Speicherung auf dem Computer gerundet. Dies gilt auch für die Ergebnisse arithmetischer Verknüpfungen von Maschinenzahlen, sofern sie nicht exakt darstellbar sind.

Bezeichnet \mathcal{R} die Menge der Maschinenzahlen, dann ist jede Rundung \diamond eine Abbildung von \mathbb{R} nach \mathcal{R} . Um eine Rundungsfehleranalyse zu ermöglichen, setzen wir für alle Rechnungen mit Gleitpunktzahlen die Existenz einer Konstante eps (der Maschinengenauigkeit) mit der folgenden Eigenschaft voraus: Für jedes x mit $|x| \in [\text{mininorm}, \text{maxreal}]$ gibt es ein ε mit

$$\diamond x = x(1 + \varepsilon), \quad |\varepsilon| < \text{eps}. \quad (8.2)$$

Für den relativen Rundungsfehler, der bei der Speicherung von x auf dem Rechner entsteht, gilt unter diesen Annahmen

$$\frac{|\diamond x - x|}{|x|} \leq \text{eps}.$$

Wir verwenden im Folgenden ausschließlich die Rundung $\square : \mathbb{R} \rightarrow \mathcal{R}$, die $x \in \mathbb{R}$ auf die nächstgelegene Maschinenzahl $\square x \in \mathcal{R}$ abbildet. Gibt es zwei Maschinenzahlen, die den gleichen Abstand von x besitzen, wird auf die betragsgrößere Zahl gerundet. Für \square ist (8.2) mit

$$\text{eps} := \frac{1}{2} \min\{x \geq \text{minreal} : 1 + x \in \mathcal{R}\}.$$

erfüllt. Im IEEE 754-Standard gilt beispielsweise $\text{eps} = 2^{-53}$, in S besitzt eps den Wert $\frac{1}{2} \cdot 10^{-3}$.

Der IEEE 754-Standard legt nicht nur ein Zahlenformat fest. Er verlangt auch, dass das Ergebnis jeder Elementaroperation auf dem Rechner mit Hilfe von \square zu runden ist (man beachte, dass die Summe, das Produkt oder der Quotient zweier Maschinenzahlen nicht wieder eine Maschinenzahl sein müssen). Sind x und y Maschinenzahlen und bezeichnet \boxtimes die Realisierung der arithmetischen Grundoperation \circ auf dem Rechner, muss für $\circ \in \{+, -, \cdot, /\}$ gelten:

$$\underbrace{x \boxtimes y}_{\in \mathcal{R}} := \square \left(\underbrace{x \circ y}_{\in \mathbb{R}} \right).$$

Arithmetische Grundoperationen können auf dem Rechner dann folgendermaßen beschrieben werden:

$$x \boxtimes y = (x \circ y)(1 + \varepsilon) \quad \text{mit} \quad |\varepsilon| \leq \text{eps}. \quad (8.3)$$

Bei der Stabilitätsanalyse von Algorithmen nehmen wir im Folgenden an, dass eine ähnliche Eigenschaft auch für die üblichen auf dem Rechner verfügbaren Standardfunktionen (Wurzel, Exponentialfunktion, trigonometrische Funktionen, ...) gilt. Ist f eine Standardfunktion, \tilde{f} die auf dem Rechner implementierte Prozedur zur Berechnung von f und x eine Maschinenzahl, für die $f(x)$ im Intervall $[\text{mininorm}, \text{maxreal}]$ liegt, soll die Genauigkeitsbedingung

$$\tilde{f}(x) = f(x)(1 + \varepsilon) \quad \text{mit} \quad |\varepsilon| \leq \text{eps}$$

erfüllt sein.

8.2.1 Fehlerfortpflanzung bei den arithmetischen Grundoperationen

Wir haben bereits im Abschnitt 1.3 bemerkt, dass in Gleitpunktarithmetik sowohl das Assoziativgesetz bezüglich Addition und Multiplikation als auch das Distributivgesetz verletzt werden

können. In diesem Abschnitt untersuchen wir, wie sich vorhandene Fehler (in den Maschinenzahlen \tilde{x} und \tilde{y}) in den arithmetischen Grundoperationen (zwischen \tilde{x} und \tilde{y}) fortpflanzen. In der Diskussion der Genauigkeit von Gleitpunktberechnungen ist diese Frage von großer praktischer Bedeutung. Die Forderung (8.3) regelt nur, wie groß der bei einer einzelnen Gleitpunktoperation neu entstehende Fehler höchstens sein darf, macht aber keine Aussage darüber, ob und wie in den beteiligten Maschinenzahlen bereits enthaltene (z.B. durch Rundung bei der Speicherung oder in früheren Rechenoperationen entstandene) Fehler verstärkt werden.

Gegeben seien zwei Zahlen $x, y \in \mathbb{R}$ sowie ihre Gleitpunktnäherungen

$$\tilde{x} := \square x = x(1 + \varepsilon_x), \quad \tilde{y} := \square y = y(1 + \varepsilon_y), \quad |\varepsilon_x|, |\varepsilon_y| \leq \text{eps}.$$

Dann gilt für die arithmetischen Verknüpfungen von \tilde{x} und \tilde{y} :

1. Multiplikation:

$$\tilde{x}\tilde{y} = x(1 + \varepsilon_x)y(1 + \varepsilon_y) = xy(1 + \varepsilon_x + \varepsilon_y + \varepsilon_x\varepsilon_y) \approx xy(1 + \varepsilon_x + \varepsilon_y) =: xy(1 + \varepsilon_{xy})$$

mit

$$\varepsilon_{xy} \approx \varepsilon_x + \varepsilon_y, \quad |\varepsilon_{xy}| \lesssim 2 \text{eps}.$$

Die übliche (und intuitive, aber mathematisch eigentlich unsinnige) Notation \lesssim bedeutet hier, dass die Relation \leq höchstens in der Größenordnung von eps^2 verletzt wird. Die relativen Fehler von x und y werden bei der Multiplikation also ungefähr addiert. Weiter gilt für das gerundete Ergebnis der Multiplikation

$$\square(\tilde{x}\tilde{y}) = xy(1 + \varepsilon_{xy})(1 + \varepsilon) \approx xy(1 + \varepsilon_x + \varepsilon_y + \varepsilon) =: xy(1 + \varepsilon_{\square(xy)})$$

mit $|\varepsilon_{\square(xy)}| \lesssim 3 \text{eps}$.

2. Division:

$$\frac{\tilde{x}}{\tilde{y}} = \frac{x(1 + \varepsilon_x)}{y(1 + \varepsilon_y)} = \frac{x}{y}(1 + \varepsilon_x)(1 - \varepsilon_y + \varepsilon_y^2 - \dots) \approx \frac{x}{y}(1 + \varepsilon_x - \varepsilon_y) =: xy(1 + \varepsilon_{\frac{x}{y}})$$

mit $|\varepsilon_{\frac{x}{y}}| \lesssim 2 \text{eps}$.

Analog erhält man

$$\square\left(\frac{\tilde{x}}{\tilde{y}}\right) = xy(1 + \varepsilon_{\square(\frac{x}{y})}), \quad |\varepsilon_{\square(\frac{x}{y})}| \lesssim 3 \text{eps}.$$

3. Addition:

$$\begin{aligned} \tilde{x} + \tilde{y} &= x(1 + \varepsilon_x) + y(1 + \varepsilon_y) = x + y + x\varepsilon_x + y\varepsilon_y \\ &= (x + y)\left(1 + \frac{x}{x + y}\varepsilon_x + \frac{y}{x + y}\varepsilon_y\right) =: (x + y)(1 + \varepsilon_{x+y}). \end{aligned}$$

Haben x und y gleiche Vorzeichen, dann gilt

$$\left|\frac{x}{x + y}\right| \leq 1, \quad \left|\frac{y}{x + y}\right| \leq 1, \quad |\varepsilon_{x+y}| \leq 2 \text{eps}$$

(die relativen Fehler von x und y werden höchstens addiert). Eine problematische Situation tritt jedoch ein, wenn x und y verschiedene Vorzeichen besitzen. Im Fall $x \approx -y$ kann ε_{x+y} beliebig groß werden. Die Subtraktion ungefähr gleich großer Zahlen bezeichnet man als Auslöschung, weil sich die führenden Stellen in den Mantissen der Gleitpunktzahlen \tilde{x} und \tilde{y} gegenseitig auslöschen. Bei der Exponentenanpassung zur Normalisierung der Gleitpunktdifferenz wird die Mantisse von hinten mit Nullen ausgefüllt. Im gleichen Maß geht (relative) Genauigkeit verloren. Die Subtraktion annähernd gleich großer gerundeter Zahlen sollte daher möglichst vermieden werden.

Die beschriebene Fehlerfortpflanzung kann sich bereits in einfachsten Berechnungen negativ bemerkbar machen. Wir illustrieren dies im folgenden Beispiel.

Beispiel 8.4 Einfluss von Rundungsfehlern auf das Ergebnis einer Gleitpunktrechnung.

1. Nullstellen quadratischer Gleichungen: $x^2 - 2px + q = 0 \iff x_{1/2} = p \pm \sqrt{p^2 - q}$.

Im Fall $p^2 \gg q$ gilt $\sqrt{p^2 - q} \approx |p|$, so dass bei der Berechnung der betragskleineren Nullstelle Auslöschung auftritt.

Zahlenbeispiel: Berechnung der Nullstellen von $x^2 - 40x + 14 = 0$ in S . Korrekte Ziffern werden schwarz, fehlerhafte Ziffern blau dargestellt. In S gilt:

$$\square(p^2 - q) = \square(4.000 \cdot 10^2 \ominus 1.400 \cdot 10^1) = 3.860 \cdot 10^2,$$

$$\square\sqrt{3.860 \cdot 10^2} = \square 19.646 \dots = 1.965 \cdot 10^1,$$

$$x_1 = 2.000 \cdot 10^1 \ominus 1.965 \cdot 10^1 = 3.965 \cdot 10^1,$$

$$x_2 = 2.000 \cdot 10^1 \oplus 1.965 \cdot 10^1 = \square 0.35 = 3.500 \cdot 10^{-1}.$$

Die zwei letzten Mantissenstellen in x_2 wurden bei der Normalisierung mit Nullen aufgefüllt. Die auf vier Dezimalstellen gerundeten Nullstellen des Polynoms sind

$$x_1 = \square 39.646 \dots = 3.965 \cdot 10^1, \quad x_2 = \square 0.35311 \dots = 3.531 \cdot 10^{-1}.$$

Obwohl sich die Größenordnungen von p^2 und q nicht einmal besonders stark unterscheiden, tritt bei der Berechnung von x_2 bereits ein relativer Fehler von ca. 0.9% auf. Dieser Fehler wächst umso stärker an, je größer p^2 im Vergleich zu q ist.

2. Numerische Differentiation: $f'(x) \approx \frac{f(x+h) - f(x)}{h}$.

Bei der numerischen Differentiation treten zwei Fehler auf: der analytische Fehler in der Größenordnung von h (Satz von Taylor) und die Auslöschung durch Differenzbildung benachbarter Funktionswerte. Um den analytischen Fehler zu minimieren, sollte ein kleiner Wert für h gewählt werden. Für kleine Werte von h vergrößert sich jedoch der Effekt der Auslöschung.

Zahlenbeispiel: $f(x) = \sqrt{x}$, $x = 1$: $f'(1) = \frac{1}{2} = 5.000 \cdot 10^{-1}$. In S erhält man:

$$h = 1 : \quad \square(\sqrt{2} \oplus 1) \ominus 1 = (\square 1.4142 \dots \oplus 1) = 1.414 - 1 = 4.140 \cdot 10^{-1}$$

(der analytische Fehler ist sehr groß),

$$h = 0.1 : \quad \square(\sqrt{1.1} \oplus 1) \ominus 0.1 = (\square 1.0488 \dots \oplus 1) \ominus 0.1 = (1.049 \oplus 1) \ominus 0.1$$

$$= 4.900 \cdot 10^{-1} \quad (\text{der analytische Fehler ist groß}),$$

$$h = 0.01 : \quad \square(\sqrt{1.01} \oplus 1) \ominus 0.01 = (\square 1.00498 \dots \oplus 1) \ominus 0.01$$

$$= (1.005 \oplus 1) \ominus 0.01 = 5.000 \cdot 10^{-1}$$

(der analytische Fehler ist klein, die Auslöschung ist noch moderat),

$$h = 0.001 : \quad \square(\sqrt{1.001} \oplus 1) \ominus 0.001 = (\square 1.00049 \dots \oplus 1) \ominus 0.001$$

$$= (1.000 \oplus 1) \ominus 0.001 = 0$$

(der analytische Fehler ist sehr klein, aber die Auslöschung wirkt katastrophal).

Man kann zeigen, dass bei der numerischen Differentiation die Wahl $h = 2\sqrt{\text{eps}/f''(x)}$ optimal ist. Für größere Werte von h überwiegt der analytische Fehler, für kleinere Werte von h dominiert die Auslöschung.

Auslöschung kann man mit Mitteln der Rechnerarithmetik bekämpfen, indem man auf dem Computer Register bereitstellt, in denen Zwischenergebnisse von Berechnungen mit höherer Genauigkeit (längeren Mantissen) gespeichert werden. Derartige Register sind allerdings teuer und stehen dem Anwender nicht auf jedem Rechner zur Verfügung. Manchmal ist es aber möglich, Auslöschung dadurch zu vermeiden, dass man den auszuwertenden Ausdruck geeignet umformt.

Beispiel 8.5

1. Im Fall der quadratischen Gleichung gilt nach dem Satz von Vieta:

$$x_1 x_2 = q.$$

Man kann also x_1 wie in Beispiel 8.4 berechnen und anschließend $x_2 := \frac{q}{x_1}$ setzen. Im Vergleich zur obigen Rechnung erfordert dies eine zusätzliche Division. Dafür wird die Auslöschung vermieden. Im obigen Beispiel:

$$x_2 = 1.400 \cdot 10^1 \boxminus 3.965 \cdot 10^1 = \boxminus 3.5308 \dots \cdot 10^{-1} = 3.531 \cdot 10^{-1}.$$

2. Geht man von exakten Maschinenzahlen x und y aus, ist die Differenz $x - y$ im Fall $x \approx y$ ebenfalls eine Maschinenzahl. Die Differenz $x - y$ wird dann fehlerfrei berechnet, obwohl Auslöschung auftritt. Kritisch bleibt dagegen die Berechnung von Ausdrücken der Bauart $f(x) - f(y)$, wenn die Funktionswerte gerundet berechnet werden. Einige typische Umformungen zur Bekämpfung von Auslöschung sind:

$$(a) \sqrt{x} - \sqrt{y} = \frac{x - y}{\sqrt{x} + \sqrt{y}}.$$

$$\text{Z.B. } x = 4.48, y = 4.47:$$

$$\boxminus \sqrt{x} \boxminus \boxminus \sqrt{y} = \boxminus 2.1166 \dots \boxminus \boxminus 2.1142 \dots = 2.117 \boxminus 2.114 = 3.000 \cdot 10^{-3},$$

$$(x \boxminus y) \boxminus (\boxminus \sqrt{x} \boxplus \boxminus \sqrt{y}) = 0.01 \boxminus 4.231 = \boxminus 0.0023635 \dots = 2.364 \cdot 10^{-3}.$$

$$(b) \ln x - \ln y = \ln \frac{x}{y} = \ln \left(1 + \frac{x - y}{y} \right) \approx \frac{x - y}{y} \quad (\text{Satz von Taylor, für } x \approx y).$$

$$\text{Z.B. } x = 4.48, y = 4.47:$$

$$(\boxminus \ln x) \boxminus (\boxminus \ln y) = \boxminus 1.4996 \dots \boxminus \boxminus 1.4973 \dots$$

$$= 1.500 \boxminus 1.497 = 3.000 \cdot 10^{-3},$$

$$\boxminus \ln(x \boxminus y) = \boxminus \ln(\boxminus 1.0022 \dots) = \boxminus \ln(1.002) = \boxminus 0.0019980 \dots = 1.998 \cdot 10^{-3},$$

$$(x \boxminus y) \boxminus y = 0.01 \boxminus 4.47 = \boxminus 0.0022371 \dots = 2.237 \cdot 10^{-3}.$$

$$(\text{Exakt: } \boxminus (\ln 4.48 - \ln 4.47) = 2.235 \cdot 10^{-3}.)$$

$$(c) \cos x - \cos y = 2 \sin \frac{x + y}{2} \sin \frac{y - x}{2}.$$

$$\text{Z.B. } x = 4.48, y = 4.47:$$

$$(\boxminus \cos x) \boxminus (\boxminus \cos y) = \boxminus (-0.23030 \dots) \boxminus \boxminus (-0.24002 \dots)$$

$$= -0.2303 \boxminus (-0.2400) = 9.700 \cdot 10^{-3},$$

$$2 \boxminus (\boxminus \sin \frac{x \boxplus y}{2}) \boxminus (\boxminus \sin \frac{y \boxminus x}{2}) = 2 \boxminus (\boxminus \sin(4.475)) \boxminus (\boxminus \sin(-0.005))$$

$$= 2 \boxminus (\boxminus (-0.97195 \dots)) \boxminus (\boxminus (-0.0049999 \dots))$$

$$= 2 \boxminus 0.9720 \boxminus 0.005 = 9.720 \cdot 10^{-3}.$$

8.3 Die Kondition eines mathematischen Problems

Rundungsfehler und ihre Fortpflanzung sind nicht die einzige Fehlerquelle numerischer Berechnungen. Häufig enthalten bereits die Eingangsdaten einer mathematischen Aufgabenstellung Fehler, beispielsweise Messfehler oder Fehler, die auf Vereinfachungen in der Modellierung zurückgehen. Manchmal haben vereinfachende Annahmen keine praktische Auswirkung auf das Ergebnis (z.B. beim freien Fall eines Tennisballs aus 1,50 m Höhe darf man den Luftwiderstand, die Corioliskraft, die Abhängigkeit der Beschleunigung von der augenblicklichen Position sowie relativistische Effekte getrost vernachlässigen), ein anderes Mal führen sie zu unerklärlichen Widersprüchen zwischen Berechnungen und Beobachtungen.

Ein historisches Beispiel dafür liefert die Periheldrehung der Merkurbahn. Zu Beginn des 20. Jahrhunderts lagen präzise astronomische Messungen der Bahnkurve des Merkurs vor, die den Berechnungen auf Basis der Newton'schen Himmelsmechanik widersprachen. Zwischen der beobachteten jährlichen Präzession des Perihels von 5.75 Bogensekunden und dem theoretisch berechneten Wert von 5.32 Bogensekunden klappte eine Lücke von 0.43 Bogensekunden, was immerhin 7.5% Abweichung entspricht. Um die Beobachtung in Einklang mit der Theorie zu bringen (!), suchte man lange Zeit vergeblich nach unbekanntem Himmelskörpern, deren Gravitationskraft die Abweichung erklären sollte. Aufgelöst wurde der Widerspruch erst 1916 durch Einsteins Veröffentlichung seiner allgemeinen Relativitätstheorie, welche die vorliegenden Messungen bestätigte und die Unvollständigkeit der Newton'schen Himmelsmechanik offenbarte.

Sind gegenüber Beobachtungen abweichende Berechnungen auf fehlerhafte Modellierung zurückzuführen, ist dies mit den Mitteln der Numerischen Mathematik nicht zu klären. Manchmal führen aber trotz korrekter mathematischer Modellierung kleine Messfehler in Eingangsgrößen zu unverhältnismäßig großen Fehlern im berechneten Ergebnis. An dieser Stelle ist die Numerik gefordert, die Auswirkungen von Störungen in den Eingangsdaten eines Algorithmus abzuschätzen und gegebenenfalls den Algorithmus so abzuändern, dass Fehlerverstärkung in der Berechnung vermieden wird.

Gegeben ist die folgende Situation:

$$\text{Eingangsdaten} \xrightarrow{f} \text{Ergebnis.}$$

Gesucht ist die Abhängigkeit des Ergebnisses von kleinen Störungen in den Eingangsdaten.

8.3.1 Fehlerfortpflanzung

Wir untersuchen die Fehlerfortpflanzung nur für hinreichend glatte Funktionen. Ist f unstetig, können auch kleinste Fehler in den Eingangsdaten große Fehler im Ergebnis verursachen. Solche schlecht gestellten Probleme betrachten wir hier nicht.

Für den durch eine kleine Störung Δx verursachten relativen Fehler von $f(x)$ gilt:

1. Gegeben sei eine stetig differenzierbare, reellwertige Funktion $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}$. Für $y = f(x)$, $y + \Delta y = f(x + \Delta x)$ gilt dann:

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

für ein ξ zwischen x und $x + \Delta x$. Für den relativen Fehler des Ergebnisses folgt:

$$\frac{\Delta y}{y} = \frac{f'(\xi)\Delta x}{f(x)} = \frac{x f'(\xi)}{f(x)} \frac{\Delta x}{x} \approx \frac{x f'(x)}{f(x)} \frac{\Delta x}{x}.$$

2. Ist $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ stetig differenzierbar, gilt für $y_i = f_i(x_1, x_2, \dots, x_n)$, $i = 1, 2, \dots, m$, analog

$$\Delta y_i = f_i(x + \Delta x) - f_i(x) \approx \sum_{j=1}^n \Delta x_j \frac{\partial f_i}{\partial x_j}(x),$$

was auf die Abschätzung

$$|\Delta y_i| \lesssim \sum_{j=1}^n |\Delta x_j| \left| \frac{\partial f_i}{\partial x_j}(x) \right| \leq \left(\max_{k=1}^n |\Delta x_k| \right) \sum_{j=1}^n \left| \frac{\partial f_i}{\partial x_j}(x) \right| \leq \|\Delta x\|_\infty \left\| \left(\frac{\partial f_i}{\partial x_j}(x) \right) \right\|_\infty$$

führt, aus welcher

$$\frac{\|\Delta y\|_\infty}{\|y\|_\infty} \lesssim \frac{\|x\|_\infty \left\| \left(\frac{\partial f_i}{\partial x_j}(x) \right) \right\|_\infty}{\|f(x)\|_\infty} \frac{\|\Delta x\|_\infty}{\|x\|_\infty}$$

folgt. Dabei haben wir die Zeilensummennorm in naheliegender Weise für nicht quadratische Matrizen verallgemeinert.

Definition 8.6

1. Gegeben sei eine stetig differenzierbare, reellwertige Funktion $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}$. Dann heißt

$$\kappa_f(x) := (\text{cond } f)(x) := \left| \frac{x f'(x)}{f(x)} \right|$$

(relative) *Kondition* von f an der Stelle x .

2. Ist $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ stetig differenzierbar, dann heißt

$$\kappa_f(x) := (\text{cond } f)(x) := \frac{\|x\|_\infty \left\| \left(\frac{\partial f_i}{\partial x_j}(x) \right) \right\|_\infty}{\|f(x)\|_\infty}$$

(relative) *Kondition* von f an der Stelle x .

Die *Kondition* ist eine Eigenschaft des Problems (der Funktion), nicht eines Lösungsverfahrens (eines Algorithmus).

Beispiel 8.7

1. Addition zweier Zahlen: $f(x) := x + a$, $a \in \mathbb{R} \setminus \{0\}$ fest:

$$\kappa_f(x) = \left| \frac{x \cdot 1}{x + a} \right| = \left| \frac{x}{x + a} \right|.$$

Für $x = 0$ ist die *Kondition* Null (Null addiert man gewissermaßen fehlerfrei), für $x \approx -a$ ist die *Kondition* groß (Auslöschung).

2. Multiplikation zweier Zahlen: $f(x) := ax$, $a \in \mathbb{R}$ fest:

$$\kappa_f(x) = \left| \frac{ax}{ax} \right| = 1.$$

Die Multiplikation ist auf ganz \mathbb{R} gut konditioniert.

8.3.2 Stabilität

Bei der Stabilitätsanalyse will man für einen gegebenen numerischen Algorithmus abschätzen, wie sich Fehler in den Eingangsdaten sowie Rundungsfehler im Verlauf der Rechnung auf das Ergebnis auswirken. Wir gehen von der folgenden Situation aus: Die Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$ wird auf dem Rechner durch den Algorithmus $\tilde{f} : \mathbb{R} \rightarrow \mathcal{R}$ realisiert. Datenfehler (durch Messung oder Rundung der Eingangsdaten) und einzelne Rundungsfehler in Gleitpunktoperationen seien durch die relative Fehlerschranke eps wie in (8.3) beschränkt.

Definition 8.8 Der Algorithmus \tilde{f} heißt auf D vorwärts stabil, falls es eine nicht zu große Konstante C_V gibt, so dass

$$\left| \frac{\tilde{f}(x) - f(x)}{f(x)} \right| \leq C_V \cdot \kappa_f(x) \cdot \text{eps}$$

für alle $x \in D$ gilt.

Der Algorithmus \tilde{f} ist genau dann vorwärts stabil, wenn die Kondition des Problems f durch \tilde{f} höchstens moderat vergrößert wird. Genaue Ergebnisse erhält man in einer numerischen Berechnung, wenn das Problem gut konditioniert und der verwendete Algorithmus stabil ist. Ein instabiler Algorithmus kann auch für ein gut konditioniertes Problem versagen. Schlecht konditionierte Probleme lassen sich auch mit stabilen Algorithmen nicht oder nicht zufriedenstellend lösen. Dann hilft nur die Umformulierung des Problems.

Beispiel 8.9 Es sei $f(x) := \sqrt{x+1} - \sqrt{x}$, $x \geq 0$. Aus

$$f'(x) = \frac{1}{2\sqrt{x+1}} - \frac{1}{2\sqrt{x}} = -\frac{\sqrt{x+1} - \sqrt{x}}{2\sqrt{x+1}\sqrt{x}}$$

folgt

$$\kappa_f(x) = \left| \frac{x f'(x)}{f(x)} \right| = \frac{1}{2} \frac{\sqrt{x}}{\sqrt{x+1}} < \frac{1}{2}.$$

Die Funktion ist gut konditioniert und sollte sich daher genau auswerten lassen.

Es sei nun $x \in \mathcal{R}$ und $\tilde{f}_1(x)$ sei der folgende Algorithmus:

$$\begin{aligned} z_1 &= \square(x+1), \\ z_2 &= \square\sqrt{z_1}, \\ z_3 &= \square\sqrt{x}, \\ z_4 &= z_2 \ominus z_3, \\ \tilde{f}_1 &= z_4. \end{aligned}$$

Dann gilt:

$$\begin{aligned} z_1 &= (x+1)(1+\varepsilon_1), \\ z_2 &= \sqrt{z_1}(1+\varepsilon_2) = \sqrt{(x+1)(1+\varepsilon_1)}(1+\varepsilon_2) \approx \sqrt{(x+1)}(1+\frac{1}{2}\varepsilon_1 + \varepsilon_2), \\ z_3 &= \sqrt{x}(1+\varepsilon_3), \\ z_4 &= (z_2 - z_3)(1+\varepsilon_4) = (\sqrt{x+1} - \sqrt{x} + \sqrt{x+1}(\frac{1}{2}\varepsilon_1 + \varepsilon_2) - \varepsilon_3\sqrt{x})(1+\varepsilon_4) \\ &\approx \sqrt{x+1} - \sqrt{x} + \sqrt{x+1}(\frac{1}{2}\varepsilon_1 + \varepsilon_2 + \varepsilon_4) - \sqrt{x}(\varepsilon_3 + \varepsilon_4) = \tilde{f}_1(x). \end{aligned}$$

Die Fehlerterme ε_i können negativ sein. Der absolute Gesamtfehler ist durch

$$\frac{3}{2}\sqrt{x+1}\text{ eps} + \sqrt{x}\text{ eps} + (\sqrt{x+1} - \sqrt{x})\text{ eps}$$

bestmöglich abgeschätzt. Also gilt

$$\left| \frac{\tilde{f}_1(x) - f(x)}{f(x)} \right| \leq \left(\underbrace{\frac{\frac{3}{2}\sqrt{x+1} + \sqrt{x}}{\sqrt{x+1} - \sqrt{x}} + 1}_{\approx 5x \text{ für } x \rightarrow \infty} \right) \text{ eps.}$$

Der Algorithmus \tilde{f}_1 ist nicht vorwärts stabil.

Alternativ betrachten wir den Algorithmus $\tilde{f}_2(x)$ gegeben durch

$$z_1 = \square(x+1),$$

$$z_2 = \square\sqrt{z_1},$$

$$z_3 = \square\sqrt{x},$$

$$z_4 = z_2 \boxplus z_3,$$

$$z_5 = 1 \boxtimes z_4,$$

$$\tilde{f}_1 = z_5,$$

welcher auf der Darstellung

$$f(x) = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

beruht. Analog zu oben gilt die Fehlerfortpflanzung

$$z_2 \approx \sqrt{x+1} (1 + \frac{1}{2}\varepsilon_1 + \varepsilon_2)$$

$$z_3 = \sqrt{x} (1 + \varepsilon_3)$$

$$z_4 = (z_2 + z_3)(1 + \varepsilon_4)$$

$$\approx \sqrt{x+1} + \sqrt{x} + \sqrt{x+1} (\frac{1}{2}\varepsilon_1 + \varepsilon_2 + \varepsilon_4) + \sqrt{x} (\varepsilon_3 + \varepsilon_4)$$

$$= (\sqrt{x+1} + \sqrt{x}) \left(1 + \frac{\sqrt{x+1}}{\sqrt{x+1} + \sqrt{x}} (\frac{1}{2}\varepsilon_1 + \varepsilon_2) + \frac{\sqrt{x}}{\sqrt{x+1} + \sqrt{x}} \varepsilon_3 + \varepsilon_4 \right)$$

Der Satz von Taylor liefert die Entwicklung

$$\frac{1}{z_4} \approx \frac{1}{\sqrt{x+1} + \sqrt{x}} \left(1 - \frac{\sqrt{x+1}}{\sqrt{x+1} + \sqrt{x}} (\frac{1}{2}\varepsilon_1 + \varepsilon_2) - \frac{\sqrt{x}}{\sqrt{x+1} + \sqrt{x}} \varepsilon_3 - \varepsilon_4 \right),$$

welche für große Werte von x mit $\sqrt{x+1} \approx \sqrt{x}$ in

$$\frac{1}{z_4} \approx \frac{1}{\sqrt{x+1} + \sqrt{x}} \left(1 - \frac{1}{4}\varepsilon_1 - \frac{1}{2}\varepsilon_2 - \frac{1}{2}\varepsilon_3 - \varepsilon_4 \right)$$

übergeht. Hieraus folgt zunächst

$$z_5 = \frac{1}{z_4} (1 + \varepsilon_5) \approx f(x) \left(1 - \frac{1}{4}\varepsilon_1 - \frac{1}{2}\varepsilon_2 - \frac{1}{2}\varepsilon_3 - \varepsilon_4 + \varepsilon_5 \right)$$

und schließlich

$$\left| \frac{\tilde{f}_2(x) - f(x)}{f(x)} \right| \leq \left(\frac{1}{4} + \frac{1}{2} + \frac{1}{2} + 1 + 1 \right) \text{ eps} = \frac{13}{4} \text{ eps.}$$

Somit ist der Algorithmus \tilde{f}_2 vorwärts stabil (mit $C_V = \frac{13}{2}$).

Häufig ist die Vorwärtsanalyse eines Algorithmus zu aufwändig. Eventuell ist dann zumindest die von Wilkinson entwickelte Rückwärtsanalyse durchführbar. Bei dieser interpretiert man den vom Algorithmus \tilde{f} berechneten Näherungswert $\tilde{f}(x)$ als exakten Funktionswert von f an der gestörten Stelle \tilde{x} :

$$\tilde{f}(x) = f(\tilde{x}).$$

Definition 8.10 Der Algorithmus \tilde{f} heißt auf D rückwärts stabil, falls es eine nicht zu große Konstante C_R gibt, so dass zu jedem $x \in D$ ein $\tilde{x} \in D$ mit

$$\tilde{f}(x) = f(\tilde{x}), \quad \left| \frac{\tilde{x} - x}{x} \right| \leq C_R \cdot \text{eps}$$

existiert.

Bemerkung 8.11 Ist der Algorithmus \tilde{f} rückwärts stabil, dann gilt

$$\left| \frac{\tilde{f}(x) - f(x)}{f(x)} \right| = \left| \frac{f(\tilde{x}) - f(x)}{f(x)} \right| \approx \kappa_f(x) \left| \frac{\tilde{x} - x}{x} \right| \leq C_R \cdot \kappa_f(x) \cdot \text{eps}.$$

Jeder rückwärts stabile Algorithmus ist also auch vorwärts stabil. Die Umkehrung gilt im Allgemeinen nicht.

Beispiel 8.12 (Rückwärtsanalyse der quadratischen Gleichung)

Wir betrachten wieder das in Beispiel 8.4 behandelte Nullstellenproblem:

$$x^2 - 2px + q = 0, \quad p^2 \gg q.$$

Zur Berechnung der betragskleineren Nullstelle

$$x_2 := f(p, q) := p - \sqrt{p^2 - q}$$

hatten wir in Beispiel 8.4 einen Algorithmus mit der Fehlerfortpflanzung

$$\begin{aligned} z_1 &= p^2(1 + \varepsilon_1), \\ z_2 &= (z_1 - q)(1 + \varepsilon_2), \\ z_3 &= \sqrt{z_2}(1 + \varepsilon_3), \\ z_4 &= (p - z_3)(1 + \varepsilon_4), \\ \tilde{f}_1 &= z_4 \end{aligned}$$

benutzt. Die Rückwärtsanalyse liefert für diesen Algorithmus

$$\begin{aligned} & \left(p - \sqrt{(p^2(1 + \varepsilon_1) - q)(1 + \varepsilon_2)(1 + \varepsilon_3)} \right) (1 + \varepsilon_4) \\ &= p(1 + \varepsilon_4) - \sqrt{(p(1 + \varepsilon_4))^2 \underbrace{(1 + \varepsilon_1)(1 + \varepsilon_2)(1 + \varepsilon_3)^2}_{=: 1 + \varepsilon_5} - q \underbrace{(1 + \varepsilon_2)(1 + \varepsilon_3)^2(1 + \varepsilon_4)^2}_{=: 1 + \varepsilon_6}} \\ &= p(1 + \varepsilon_4) - \sqrt{(p(1 + \varepsilon_4))^2 - q \underbrace{\left(1 + \varepsilon_6 - \varepsilon_5(1 + \varepsilon_4)^2 \frac{p^2}{q} \right)}_{=: 1 + \varepsilon_7}} \\ &= p(1 + \varepsilon_4) - \sqrt{(p(1 + \varepsilon_4))^2 - q(1 + \varepsilon_7)} = f(p(1 + \varepsilon_4), q(1 + \varepsilon_7)) \end{aligned}$$

mit

$$|\varepsilon_4| \leq \text{eps},$$

$$|\varepsilon_5| \approx |\varepsilon_1 + \varepsilon_2 + 2\varepsilon_3| \leq 4 \text{eps},$$

$$|\varepsilon_6| \approx |\varepsilon_2 + 2\varepsilon_3 + 2\varepsilon_4| \leq 5 \text{eps},$$

$$|\varepsilon_7| \approx \left| \varepsilon_6 - \varepsilon_5 \frac{p^2}{q} \right| \lesssim \left(5 + 4 \frac{p^2}{q} \right) \text{eps}.$$

Nach Voraussetzung ist $4 \frac{p^2}{q} \gg 1$, also ist der Algorithmus nicht rückwärts stabil.

8.4 Beispiele

Die Stabilitätsanalyse eines Algorithmus ist in der Regel schwierig. Zum Abschluss dieses Kapitels illustrieren wir das Vorgehen anhand von Beispielen.

8.4.1 Addition Es sei (vgl. Beispiel 8.7)

$$f(x) := x + a, \quad a \in \mathbb{R} \setminus \{0\} \text{ fest},$$

$$\tilde{f}(x) = \square(x + a) = (x + a)(1 + \varepsilon), \quad |\varepsilon| \leq \text{eps}.$$

Die Vorwärtsanalyse liefert

$$\left| \frac{\tilde{f}(x) - f(x)}{f(x)} \right| = \left| \frac{(x + a)(1 + \varepsilon) - (x + a)}{x + a} \right| = |\varepsilon| = \left| \frac{x + a}{x} \right| \kappa_f(x) |\varepsilon| \leq \left| 1 + \frac{a}{x} \right| \kappa_f(x) \text{eps}.$$

Die Addition ist vorwärts stabil, sofern $|x| \ll |a|$ gilt.

Im Fall $|x| \ll |a|$ ist die Addition auch rückwärts stabil, denn es gilt:

$$\tilde{f}(x) = (x + a)(1 + \varepsilon) = \underbrace{x \left(1 + \varepsilon + \frac{\varepsilon a}{x} \right)}_{\tilde{x}} + a.$$

Für $|x| \ll |a|$ ist $\frac{\varepsilon a}{x}$ groß, die Addition also nicht rückwärts stabil. Man kann dies auch durch die folgende Überlegung leicht einsehen: In \mathcal{S} gilt für $a = 1$, $|x| \leq \frac{1}{2} \cdot 10^{-3} (= \text{eps})$

$$\square(1 + x) = 1 = 1 + 0 =: 1 + \tilde{x}.$$

$\tilde{x} = 0$ ist eindeutig bestimmt. Will man die Bedingung

$$\left| \frac{\tilde{x} - x}{x} \right| = 1 \stackrel{!}{\leq} C_R \text{eps}$$

erfüllen, muss $C_R \geq \frac{1}{\text{eps}}$ gelten.

8.4.2 Multiplikation Es sei

$$f(x) := ax, \quad a \in \mathbb{R} \text{ fest,}$$

$$\tilde{f}(x) = \square(ax) = ax(1 + \varepsilon), \quad |\varepsilon| \leq \text{eps.}$$

Vorwärtsanalyse: Mit der Kondition aus Beispiel 8.7 gilt:

$$\left| \frac{\tilde{f}(x) - f(x)}{f(x)} \right| = |\varepsilon| = \kappa_f(x) |\varepsilon| \leq \kappa_f(x) \text{eps.}$$

Die Multiplikation ist somit vorwärts stabil.

Rückwärtsanalyse:

$$\tilde{f}(x) = (ax)(1 + \varepsilon) = a \underbrace{(x(1 + \varepsilon))}_{\tilde{x}}.$$

Die Multiplikation ist auch rückwärts stabil.

8.4.3 Logarithmus Es sei

$$f(x) := \ln x, \quad \tilde{f}(x) = \square \ln x = \ln x (1 + \varepsilon), \quad |\varepsilon| \leq \text{eps.}$$

Falls x nicht nahe bei 1 liegt, ist f gut konditioniert:

$$\kappa_f(x) = \left| \frac{xf'(x)}{f(x)} \right| = \left| \frac{1}{\ln x} \right|.$$

Vorwärtsanalyse:

$$\left| \frac{\tilde{f}(x) - f(x)}{f(x)} \right| = |\varepsilon| \leq |\ln x| \kappa_f(x) \text{eps.}$$

Der Algorithmus zur Berechnung der Logarithmusfunktion ist vorwärts stabil, wenn $|\ln x|$ nicht zu groß ist (z.B. für $x \in [e^{-10}, e^{10}]$).

Rückwärtsanalyse:

$$\tilde{f}(x) = (1 + \varepsilon) \ln x = \ln(x^{1+\varepsilon}) =: \ln \tilde{x}$$

(\tilde{x} ist eindeutig bestimmt). Es gilt:

$$\left| \frac{\tilde{x} - x}{x} \right| = \left| \frac{xe^{\varepsilon \ln x} - x}{x} \right| \approx \left| \frac{x(1 + \varepsilon \ln x) - x}{x} \right| = |\varepsilon| |\ln x|,$$

so dass der Algorithmus zur Berechnung der Logarithmusfunktion auch rückwärts stabil ist, wenn $|\ln x|$ nicht zu groß ist.

8.4.4 Produkte Es sei

$$f(x_1, x_2, \dots, x_n) := \prod_{i=1}^n x_i.$$

Der Algorithmus $\tilde{f}(x_1, x_2, \dots, x_n)$ zur Berechnung des Produkts sei definiert durch

$$z_1 = x_1 x_2 (1 + \varepsilon_1),$$

$$z_i = z_{i-1} x_{i+1} (1 + \varepsilon_i), \quad i = 2, 3, \dots, n-1,$$

$$\tilde{f} = z_{n-1}.$$

Rückwärtsanalyse:

$$\tilde{f}(x) = \prod_{i=1}^n x_i \prod_{i=1}^{n-1} (1 + \varepsilon_i) = x_n \prod_{i=1}^{n-1} (x_i (1 + \varepsilon_i)).$$

Die Berechnung von Produkten ist daher rückwärts stabil.

8.4.5 Summen Es sei

$$f(x_1, x_2, \dots, x_n) := \sum_{i=1}^n x_i.$$

Der Algorithmus $\tilde{f}(x_1, x_2, \dots, x_n)$ zur Berechnung der Summe sei definiert durch

$$\begin{aligned} z_1 &= (x_1 + x_2)(1 + \varepsilon_2), \\ z_i &= (z_{i-1} + x_{i+1})(1 + \varepsilon_{i+1}), \quad i = 2, 3, \dots, n-1, \\ \tilde{f} &= z_{n-1}. \end{aligned}$$

Dann gilt:

$$\tilde{f}(x) = x_1 \prod_{i=2}^n (1 + \varepsilon_i) + \sum_{j=2}^n x_j \prod_{i=j}^n (1 + \varepsilon_i)$$

Wegen

$$\prod_{i=j}^n (1 + \varepsilon_i) \approx 1 + \sum_{i=j}^n \varepsilon_i$$

gilt

$$\left| \tilde{f}(x) - f(x) \right| \lesssim (n-1) \text{eps} |x_1| + \sum_{j=2}^n (n+1-j) \text{eps} |x_j|.$$

Die Fehlerschranke für den absoluten Fehler wird am kleinsten, wenn die Zahlen nach aufsteigenden Beträgen summiert werden.

8.4.6 Lineare Gleichungssysteme Es sei $A \in \mathbb{R}^{n \times n}$ eine feste reguläre Matrix. Für variables $b \in \mathbb{R}^n$ betrachten wir die Funktion

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad x = f(b) := A^{-1}b.$$

Mit der Funktionalmatrix

$$\frac{\partial f}{\partial b} = A^{-1}$$

und $Ax = b$ erhalten wir die Kondition von f :

$$\kappa_f(b) = \frac{\|b\| \|A^{-1}\|}{\|A^{-1}b\|} = \frac{\|Ax\| \|A^{-1}\|}{\|x\|}.$$

Ihr maximaler Wert ist

$$\sup_{b \neq 0} \kappa_f(b) = \sup_{x \neq 0} \frac{\|Ax\| \|A^{-1}\|}{\|x\|} = \|A\| \|A^{-1}\| = \kappa(A).$$

8.4.7 Rekursionsformel für Integrale Für $n \in \mathbb{N}$ sei

$$I_n := \int_1^e (\ln x)^n dx.$$

Dieses Integral kann mit Hilfe einer Rekursionsformel geschlossen gelöst werden. Mit partieller Integration folgt:

$$\begin{aligned} n = 1: \quad I_1 &= \int_1^e \ln x dx = (x \ln x - x) \Big|_1^e = 1. \\ n > 1: \quad I_n &= x(\ln x)^n \Big|_1^e - \int_1^e x n (\ln x)^{n-1} \frac{1}{x} dx = e - n I_{n-1}. \end{aligned}$$

Aus der Rekursionsformel

$$I_n = e - n I_{n-1} \tag{8.4}$$

erhält man wegen $\ln x \geq 0$ für $x \geq 1$ die triviale Grobabschätzung

$$0 \leq I_n \leq e.$$

Benutzt man die genauere Abschätzung

$$\ln x \geq \frac{x-1}{e-1} \quad \text{für } 1 \leq x \leq e,$$

dann folgt

$$I_n \geq \int_1^e \left(\frac{x-1}{e-1} \right)^n dx = \frac{(x-1)^{n+1}}{(n+1)(e-1)^n} \Big|_1^e = \frac{e-1}{n+1}.$$

Mit Hilfe von (8.4) erhält man

$$I_n \leq e - n \frac{e-1}{n} = 1,$$

also insgesamt die Schranken

$$\frac{e-1}{n+1} \leq I_n \leq 1.$$

Berechnet man I_n für einige Werte von n , lässt sich vermuten, dass

$$I_n = e + (-1)^n n! \left\{ e \sum_{j=0}^{n-1} \frac{(-1)^j}{j!} - I_1 \right\} \tag{8.5}$$

gilt. Mit vollständiger Induktion weist man die angegebene Darstellung leicht nach.

Die Summe in (8.5) ist die n -te Partialsumme von e^{-1} . Der Wert in den geschweiften Klammern strebt also gegen Null, wodurch die Berechnungsformel für große Werte von n zum uneigentlichen Grenzwert "0 · ∞" mutiert.

Außerdem gewinnt man aus der Rekursionsformel (8.4) für $n > k \geq 1$ die Darstellung

$$I_n = (-1)^{n-k} \frac{n!}{k!} I_k + c_{n,k} \tag{8.6}$$

mit gewissen Zahlen $c_{n,k} \in \mathbb{R}$. Wir fassen nun I_n in (8.6) als Funktion von I_k auf,

$$I_n = f(I_k), \quad n > k \geq 1,$$

und wollen untersuchen, wie sich Fehler in I_k auf I_n auswirken. Aus

$$|f'(I_k)| = \frac{n!}{k!}$$

erhalten wir

$$\kappa_f(I_k) = \left| \frac{I_k f'(I_k)}{I_n} \right| = \frac{n! I_k}{k! I_n} \geq \frac{(e-1)n!}{(k+1)!}.$$

Für $n \gg k$ ist f beliebig schlecht konditioniert. Die Rekursionsformel (8.4) eignet sich daher nicht zur Berechnung von I_n . Ein Zahlenbeispiel soll dies verdeutlichen. Maple 9.5 liefert (mit der Standardeinstellung für die Genauigkeit) ausgehend von $I_1 := 1.0$ die folgenden Näherungen \tilde{I}_n (gültige Ziffern sind schwarz dargestellt, ungültige blau):

n	\tilde{I}_n
1	1.000000000
2	0.718281828
3	0.563436344
4	0.464536452
5	0.395599568
6	0.344684420
7	0.305490888
8	0.274354724
9	0.249089312
10	0.227388708
11	0.217006040
12	0.114209348
13	1.233560304
14	-14.55156243
15	220.9917182

Der Genauigkeitsverlust ist hier nicht auf mangelnde Stabilität des Berechnungsalgorithmus zurückzuführen, sondern auf die schlechte Kondition der Formel (8.4). Abhilfe ist nur durch eine grundlegende Änderung der Berechnungsweise zu erwarten. Es ist

$$\begin{aligned} I_n &= e + (-1)^n n! \left\{ e \left(e^{-1} - \sum_{j=n}^{\infty} \frac{(-1)^j}{j!} \right) - 1 \right\} \\ &= e \sum_{j=n+1}^{\infty} \frac{(-1)^{n+1+j}}{j!} n! \\ &\approx e \left(\frac{1}{n+1} - \frac{1}{(n+1)(n+2)} + \dots - \frac{1}{(n+1)(n+2)\dots(n+N)} \right) =: J_{n,N} \quad (8.7) \end{aligned}$$

für geeignete Werte von N . Die zuletzt aufgetretene Reihe konvergiert sehr schnell. Mit N Summanden in (8.7) bekommt man bereits für $n = 10$ ungefähr N gültige Dezimalstellen von I_n . Außerdem ist (8.7) keine Rekursionsformel. Man kann $J_{n,N}$ für beliebige Werte von n und N bestimmen, ohne vorhergehende Integrale zu berechnen. Zum Vergleich mit der obigen Tabelle geben wir auch für $J_{n,N}$ einige Zahlenwerte an. Die Vorteile der geänderten Berechnungsformel sind deutlich zu sehen.

n	$J_{n,2}$	$J_{n,5}$	$J_{n,10}$	I_n
5	0.3883259755	0.3956070875	0.3955995478	0.3955995478
10	0.2265234857	0.2280019606	0.2280015154	0.2280015155
15	0.1598989310	0.1604263754	0.1604263089	0.1604263089
20	0.1235582649	0.1238038466	0.1238038307	0.1238038308
50	0.05227465055	0.05229363840	0.05229363829	0.05229363830
100	0.02664982184	0.02665235919	0.02665235919	0.02665235919

Index

- 3/8-Regel, 127
- a posteriori-Abschätzung, 29, 41
- a priori-Abschätzung, 29, 41
- Abstiegsverfahren, 67
- Algorithmus, 11
 - determinierter, 12
 - deterministischer, 12
- Ansatzfunktion
 - bei der Methode der kleinsten Quadrate, 119
- Approximation, 98
- Argumentreduktion, 97
- Ausgleichsgerade, 120
- Ausgleichsproblem, 85
- Ausgleichsrechnung
 - nichtlineare, 124
- Auslöschung, 148
- Auswahanweisung, 16
- Banach'scher Fixpunktsatz
 - im \mathbb{R}^n , 41
 - in \mathbb{R} , 29
- Bisektionsverfahren, 37
- cg-Verfahren, 70
 - Konvergenzverhalten, 73
- diagonaldominante Matrix, 52
- Einzelschrittverfahren, 63
- Eliminationsmatrix, 48
- Energienorm, 66
- Extrapolation, 98
- Fehlerfortpflanzung, 151
- Fixpunkt, 27
 - abstoßender, 31
 - anziehender, 31
- Fixpunktiteration, 27
 - im \mathbb{R}^n , 42
- Fixpunktsatz
 - von Banach
 - im \mathbb{R}^n , 41
 - in \mathbb{R} , 29
- Fourier-Reihe, 116
- Fourier-Transformation
 - diskrete, 116
- Gauß-Markov
 - Satz von, 119
- Gauß-Newton-Verfahren, 124
- Gauß-Quadratur, 133
- Gauß-Seidel-Verfahren, 63
- Genauigkeitsgrad
 - einer Quadraturformel, 128
- Gesamtschrittverfahren, 62
- Gewichtsfunktion, 134
- Gleitpunktsystem, 9
- Gleitpunktzahl, 9, 144
 - normalisierte, 9, 144
- Gram-Schmidt-Orthogonalisierung, 77
- Hermite-Interpolation, 107
- Householder-Matrix, 79
- Householder-Transformation, 79
- Interpolation, 98
 - Polynom-Interpolation, 98
 - Spline-Interpolation, 108
 - stückweise lineare, 108
 - trigonometrische, 116
- Interpolationsfehler
 - bei der Polynom-Interpolation, 100, 103
 - bei der Spline-Interpolation, 114
- Interpolationspolynom
 - Darstellung von Lagrange, 98
 - Darstellung von Newton, 99
- iterationsfähige Gestalt, 26
- Iterationsverfahren, 27
- Jacobi-Verfahren, 62
- kleinste Quadrate
 - Methode der, 85, 118
- Kondition
 - einer Matrix, 57
 - relative, 152
- kontrahierende Abbildung, 28
 - im \mathbb{R}^n , 41

- Kontraktion, 28
- Kontraktionskonstante, 28
- Konvergenzordnung
 - eines Iterationsverfahrens, 32
- Krylov-Raum, 71
- Landau-Symbol, 22
- LGS
 - überbestimmtes, 85
 - unterbestimmtes, 87
- LU*-Zerlegung, 49
 - Aufwand der, 52
- Mantisse, 9, 144
- Maple, 15
 - Prozedur, 19
- Maschinengenauigkeit, 147
- Maschinenzahl, 147
- Matrix
 - orthogonale, 76
 - positiv definite, 66
- maxreal, 144
- Merkur
 - Periheldrehung, 151
- Methode der kleinsten Quadrate, 85, 118
- Methode des steilsten Abstiegs, 68
 - Konvergenzverhalten, 69
- Milne-Regel, 127
- mininorm, 144
- Modellbildung, 7
 - Modellfehler, 7
- Newton-Cotes-Formeln, 125
 - abgeschlossene, 126
 - Fehlerabschätzung, 128
 - offene, 126
- Newton-Verfahren
 - für reellwertige Funktionen, 34
 - quadratische Konvergenzordnung, 35
 - vereinfachtes, 36
 - im \mathbb{R}^n , 44
- Norm, 38
 - Matrixnorm, 38
 - Frobenius-Norm, 38
 - induzierte, 39
 - Spaltensummennorm, 40
 - Spektralnorn, 40
 - submultiplikative, 39
 - verträgliche, 39
 - Zeilensummennorm, 40
 - Vektornorm, 38
- Normalgleichungssystem, 85
- PALU*-Zerlegung, 55
- Permutationsmatrix, 53
 - elementare, 53
- Pivotelement, 56
- Polynom-Interpolation, 98
- positiv definite Matrix, 66
- Prozedur, 18
 - globale Variable, 20
 - lokale Variable, 20
 - Parameterliste, 19
- QR*-Zerlegung, 76, 79
 - reduzierte, 79
- Quadraturformel, 126
 - summierte, 129
- rückwärts stabil, 155
- Rechnen
 - numerisches, 8
 - symbolisches, 8
- regula falsi, 38
- Relaxation, 33
- Residuum, 60
- Restglied
 - eines Taylor-Polynoms, 94
- Romberg-Schema, 132
- Rundung, 147
 - zur nächstgelegenen Maschinenzahl, 147
- Rundungsfehler, 10
- Satz von
 - Gauß-Markov, 119
 - Taylor, 93
- Schleife, 15
- Sekantenverfahren, 36
- Selbstabbildung, 27
- Simpson-Regel, 127
- Spline
 - kubischer, 109
 - Minimierungseigenschaft, 113
 - natürlicher, 110
 - not-a-knot Spline, 110
 - periodischer, 110
 - Randbedingungen, 110
 - zur Parameterdarstellung einer Kurve, 115
- Spline-Interpolation, 108
- Splitting-Verfahren, 62
- Stabilität
 - rückwärts stabil, 155
 - vorwärts stabil, 153
- Standardfunktion, 91

Steigung, 99

n -ter Ordnung, 99

 verallgemeinerte, 111

Steigungsschema, 101

Taylor

 Satz von, 93

Taylor-Polynom, 93

Taylor-Reihe, 94

Trapezregel, 127

 summierte, 129

Tschebyscheff-Polynom, 105

Tschebyscheff-Stützstellen, 106

Turing-Maschine, 11

Vorkonditionierung, 74

vorwärts stabil, 153