

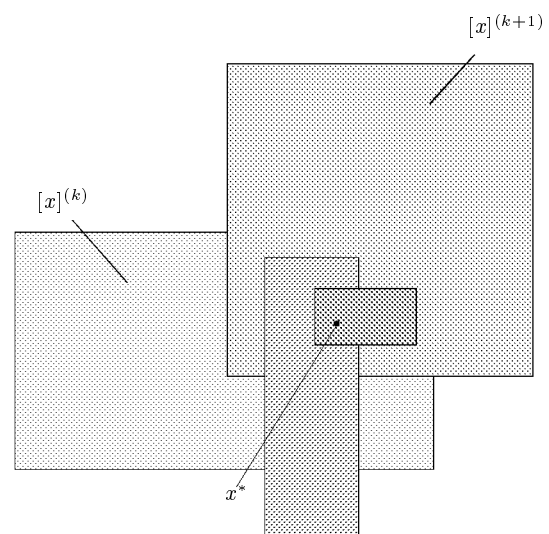
**I**nstitut für  
**A**ngewandte  
**M**athematik

Universität Karlsruhe (TH)  
D-76128 Karlsruhe

## Rückwärtsmethode zur automatischen Berechnung von worst-case Fehlerschranken

Michael Bräuer und Walter Krämer

**F**orschungsschwerpunkt  
**C**omputerarithmetik,  
**I**ntervallrechnung und  
**N**umerische Algorithmen mit  
**E**rgebnisverifikation



Bericht 3/1999

## Impressum

Herausgeber:	Institut für Angewandte Mathematik Lehrstuhl Prof. Dr. Ulrich Kulisch Universität Karlsruhe (TH) D-76128 Karlsruhe
Redaktion:	Dr. Gerd Bohlender

## Internet-Zugriff

Die Berichte sind in elektronischer Form erhältlich über

`ftp://iamk4515.mathematik.uni-karlsruhe.de|`  
im Verzeichnis: `/pub/documents/reports|`

oder über die World Wide Web Seiten des Instituts

`http://www.uni-karlsruhe.de/|iam|`

## Autoren-Kontaktadresse

Rückfragen zum Inhalt dieses Berichts bitte an

Michael Bräuer  
Walter Krämer  
Institut für Angewandte Mathematik  
Universität Karlsruhe (TH)  
D-76128 Karlsruhe  
E-Mail: [walter.kraemer@math.uni-karlsruhe.de](mailto:walter.kraemer@math.uni-karlsruhe.de)

# Rückwärtsmethode zur automatischen Berechnung von worst-case Fehlerschranken

Michael Bräuer und Walter Krämer

## Zusammenfassung

### **Rückwärtsmethode zur automatischen Berechnung von worst-case Fehlerschranken:**

Die insbesondere im Zusammenhang mit der Berechnung von Gradienten häufig verwendete schnelle Rückwärtsmethode wird eingesetzt, um worst-case Fehlerschranken für Gleitkommaalgorithmen automatisch und a priori zu bestimmen. Die numerischen Ergebnisse dieses Verfahrens sowie dessen Handhabung und die benötigten Laufzeiten werden mit den entsprechenden Größen eines Vorwärtsfehlerkalküls verglichen.

## Abstract

### **Computaton of Worst Case Error Bounds in Reverse Mode:**

A priori worst case error bounds for floating point algorithms are computed automatically using the so called reverse mode.

## 1 Einleitung

Bei der Verwendung numerischer Verfahren auf dem Rechner ist nicht nur das Laufzeitverhalten ein entscheidendes Einsatzkriterium. Ebenso wichtig ist das Verhalten dieser Methoden bezüglich der Verstärkung auftretender Rundungsfehler und möglicher Störungen in den Eingangsdaten. Deshalb ist man an der Angabe einer gleichmäßigen Fehlerschranke interessiert, welche eine obere Schranke für die maximale Abweichung eines auf der Maschine errechneten Wertes vom wahren (in  $\mathbb{R}$  berechneten und allgemein unbekanntem) Resultat für beliebige Argumente eines vorgegebenen Bereichs darstellt. Hat man einen Algorithmus in Gestalt einer faktorisierbaren Funktion vorliegen, so läßt sich eine a-priori Fehlerschranke für diese Funktion berechnen, falls die Größe der Eingangsdaten und deren Eingangsfehler bekannt sind.

Im folgenden Abschnitt 2 wird nach einigen Bemerkungen zur Rechnerarithmetik an Begriffe wie *faktorisierebare Funktion*, *Codeliste*, *Darstellungsgraph* erinnert. Diese werden bei der Darstellung des Prinzips der sogenannten Rückwärtsmethode bzw. dessen Realisierung auf einer Rechanlage benötigt.

Abschnitt 3 erläutert das allgemeine Prinzip des Rückwärtsrechnens. Im zentralen Abschnitt 4 wird dieses Prinzip zur Berechnung von worst-case Fehlerschranken herangezogen. Dabei werden Vergleiche bzgl. numerischer Ergebnisse und Handhabbarkeit der Rückwärtsmethode mit einem Vorwärtskalkül angestellt.

Im Anhang sind beispielhaft der detaillierte Algorithmus zur Berechnung der arctan-Funktion und die Daten für die Fehlerabschätzung dieses Algorithmus angegeben.

## 2 Grundlegende Definitionen

### 2.1 Rechnerarithmetik

Bedingt durch die Endlichkeit des Speichers kann auf der Maschine nur mit einer endlichen Teilmenge  $S$  der reellen Zahlen gerechnet werden. Typischerweise bezeichnet  $S$  die Menge der sogenannten Gleitkommazahlen (vgl. [25]).

**Definition 2.1 (Gleitkommasystem)** *Unter einem Gleitkommasystem  $S = S(b, n, e_1, e_2) \subseteq \mathbb{R}$  versteht man die Menge aller auf einer Rechanlage darstellbaren normalisierten Gleitkommazahlen  $x$  mit den folgenden Eigenschaften:*

1.  $x = m \cdot b^e$
2.  $m := \pm 0.d_1d_2 \dots d_n$
3.  $1 \leq d_1 \leq b - 1, 0 \leq d_i \leq b - 1$  ( $2 \leq i \leq n$ )
4.  $e_1 \leq e \leq e_2, e_1 \leq 0$  und  $1 \leq e_2$ , wobei  $e, e_1, e_2$  ganze Zahlen sind.

*Dabei heißt  $m$  die Mantisse mit Länge  $n \in \mathbb{N}$ ,  $b \in \mathbb{N}$  ( $b > 1$ ) die Basis und  $e$  der Exponent der Gleitkommazahl  $x$ . Ferner existiert eine Darstellung der 0, welche darin besteht, daß die Mantisse den Wert 0 und der Exponent den Wert  $e_1$  hat<sup>1</sup>.*

**Bemerkung 2.1** *Mit  $x_{\min}$  wird die kleinste und mit  $x_{\max}$  die größte positive normalisierte Gleitkommazahl bezeichnet. Der Unterlaufbereich ist die Menge  $\{x \in \mathbb{R} \mid |x| \in (0, x_{\min})\}$ . Man definiert  $\mathbb{R}_B := \{x \in \mathbb{R} \mid -x_{\max} \leq x \leq x_{\max}\}$ . Gelegentlich werden auch denormalisierte Gleitkommazahlen betrachtet. Diese Zahlen sind dadurch gekennzeichnet, daß die Forderung  $1 \leq d_1$  in obiger Definition abgeschwächt und durch  $0 \leq d_1$  ersetzt wird.*

Reelle Zahlen müssen bei Berechnungen auf der Maschine durch Gleitkommazahlen ersetzt werden. Dazu sind sie durch eine geeignete Abbildungsvorschrift in die Menge der Gleitkommazahlen zu projizieren.

**Definition 2.2 (Rundung)** *Eine Abbildung  $\square : \mathbb{R}_B \rightarrow S$  heißt Rundung, falls für alle  $x \in S$  die Gleichung*

$$\square(x) = x$$

*erfüllt ist.*

---

<sup>1</sup>Man definiert  $0 := +0.0 \dots 0 \cdot b^{e_1}$ .

Im allgemeinen erhält man beim Verknüpfen zweier Gleitkommazahlen durch Operationen aus  $\{+, -, \cdot, /\}$  keine Gleitkommazahl. Vielmehr müssen die (exakten) Ergebnisse derartiger Operationen wieder ins Gleitkommasystem gerundet werden:

**Definition 2.3 (Gleitkommaoperation)** Sei  $S$  ein Gleitkommasystem mit Rundung  $\square : \mathbb{R}_B \rightarrow S$ . Für  $x, y \in S$  mit  $x \circ y \in \mathbb{R}_B$  wird die zu  $\circ \in \{+, -, \cdot, /\}$  gehörende Gleitkommaoperation  $\boxtimes$  durch

$$x \boxtimes y := \square(x \circ y)$$

definiert, falls kein Unterlauf entsteht.

Auf dem Rechner stehen verschiedene Gleitkommazahlensysteme zur Verfügung. Die wichtigste Rolle spielt das IEEE-double-Format  $S = S(2, 53, -1022, 1023)^2$ . Es ist durch den IEEE-Standard 754 genormt. Jede Gleitkommazahl  $x \in S$  läßt sich hier durch

$$x = (-1)^{s(x)} \cdot m(x) \cdot 2^{e(x)} \quad (1)$$

mit  $s(x) \in \{0, 1\}$ ,  $m(x) = d_0.d_1 \dots d_{52}$  und  $-1022 \leq e(x) \leq 1023$  darstellen. Ist  $d_0 = 1$ , so heißt  $x$  normalisierte, ansonsten denormalisierte Zahl. Die kleinste darstellbare normalisierte positive Zahl wird hier mit `MinReal`, die größte mit `MaxReal` bezeichnet. Der Bereich  $N$  der normalisierten Gleitkommazahlen ist durch

$$N := [-\text{MaxReal}, -\text{MinReal}] \cup [\text{MinReal}, \text{MaxReal}]$$

erklärt. Der Unterlaufbereich ist die Menge aller Gleitkommazahlen aus

$$U := (-\text{MinReal}, \text{MinReal}).$$

Laut Standard müssen verschiedene Rundungen (Rundung zur nächstgelegenen Gleitkommazahl und gerichtete Rundungen) zur Verfügung stehen. Für Verknüpfungen  $\circ \in \{+, -, \cdot, /\}$  werden an die zugehörigen Gleitkommaoperationen  $\boxtimes$  mit  $x, y \in S$  die folgenden Anforderungen gestellt:

$$\left. \begin{array}{ll} \left| \frac{x \circ y - x \boxtimes y}{x \circ y} \right| \leq \varepsilon & , \text{ falls } x \circ y \in N \\ |x \circ y - x \boxtimes y| \leq \text{MinReal} & , \text{ falls } x \circ y \in U \\ x \boxtimes y = x \circ y & , \text{ falls } x \circ y \in S. \end{array} \right\} \quad (2)$$

Die Größe  $\varepsilon$  entspricht dabei dem Maschinenepsilon und beträgt je nach Rundungsmodus  $\varepsilon := 2^{-53}$  (Rundung zur nächstgelegenen Gleitkommazahl) oder  $\varepsilon := 2^{-52}$  (gerichtete Rundung nach unten bzw. nach oben). Für Funktionen  $f$ , z.B. `exp`, `sin`, etc., wird eine Genauigkeitsforderung gemäß (2) nicht verlangt. In dieser Arbeit wird jedoch davon ausgegangen, daß für solche Funktionen Konstanten  $\varepsilon(f)$  und  $u(f)$  bekannt sind, welche die folgenden Eigenschaften erfüllen:

$$\left. \begin{array}{ll} |f(x) - \tilde{f}(x)| \leq \varepsilon(f) \cdot |f(x)| & , \text{ für } x \in S \text{ mit } f(x) \in N \\ |f(x) - \tilde{f}(x)| \leq u(f) & , \text{ für } x \in S \text{ mit } f(x) \in U. \end{array} \right\} \quad (3)$$

(vgl. [5], [12], [20]). Dabei ist  $\tilde{f} : S \rightarrow S$  eine zu  $f$  gehörende Gleitkommanäherung.

<sup>2</sup>Im folgenden sind Gleitkommazahlen immer als Maschinenzahlen im IEEE-double-Format anzusehen.

## 2.2 Faktorisierbarkeit von Funktionen

Bei der Lösung mathematischer Probleme durch eine Rechenanlage verwendet man Programmiersprachen wie z. B. Fortran, C oder C++. Dazu werden die benötigten Funktionen (Algorithmen) in einer dieser Sprachen kodiert und dann durch einen Compiler oder Interpreter in ein ausführbares Programm umgewandelt. Grundsätzlich können nur solche Funktionen ausgewertet werden, die mit Sprachkomponenten beschreibbar sind, welche die verwendete Programmiersprache bereitstellt, also z. B. durch Aufrufe von elementaren Operationen, Standardfunktionen, Bedingungsanweisungen und Schleifen. In diesem Zusammenhang definiert man eine eingeschränkte Klasse solcher Funktionen, nämlich die Klasse der faktorisierten Funktionen. Doch zunächst zwei weitere Definitionen.

**Definition 2.4 (elementare Operationen  $\mathcal{OP}$ )** Die Menge aller unären und binären reellen elementaren Operationen  $\{-, +, -, \cdot, /\}$  wird mit  $\mathcal{OP}$  bezeichnet.

**Definition 2.5 (Funktionenmenge  $\mathcal{F}$ )** Das Symbol  $\mathcal{F}$  bezeichnet die Menge der reellen Standardfunktionen  $\{\text{sqr}(\cdot), \text{sqrt}(\cdot), \text{sqrt}(\cdot, \cdot), \text{power}(\cdot, \cdot), \sin(\cdot), \cos(\cdot), \tan(\cdot), \cot(\cdot), \arcsin(\cdot), \arccos(\cdot), \arctan(\cdot), \text{arccot}(\cdot), \exp(\cdot), \ln(\cdot), \sinh(\cdot), \cosh(\cdot), \tanh(\cdot), \coth(\cdot), \text{arsinh}(\cdot), \text{arcosh}(\cdot), \text{artanh}(\cdot), \text{arcoth}(\cdot)\}$ .

Mit Hilfe dieser Definitionen wird nun die Klasse der faktorisierten Funktionen erklärt. Bei solchen Funktionen sind die Operanden nur durch Operationen aus  $\mathcal{OP}$  und  $\mathcal{F}$  zusammengesetzt. Die folgende Definition lehnt sich eng an die von Kedem [18] an.

**Definition 2.6 (Faktorisierte Funktion)** Sei  $\pi_i^n : \mathbb{R}^n \rightarrow \mathbb{R}$  die Projektion eines reellen  $n$ -Tupels auf seine  $i$ -te Komponente, d.h.  $\pi_i^n(x_1, x_2, \dots, x_n) = x_i$ . Eine Funktion  $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$  heißt faktorisiert, falls es eine endliche Folge von Funktionen  $f_1, \dots, f_k$  gibt, für die gilt:

- (1)  $f_j : D \rightarrow \mathbb{R}$  für alle  $j \in \{1, \dots, k\}$ ,
- (2)  $f = f_k$ ,
- (3)  $f_1 = \pi_1^n, f_2 = \pi_2^n, \dots, f_n = \pi_n^n$  und
- (4) für  $n < j \leq k$ :  $f_j = g(f_l)$  mit  $g \in \mathcal{F}$  und  $l < j$  oder  $f_j = f_l \circ f_r$  mit  $\circ \in \mathcal{OP}^3$  und  $l, r < j$  oder  $f_j$  ist eine reelle Konstante.

Ein Beispiel soll diese Definition verdeutlichen.

**Beispiel 2.1** Sei  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  mit  $f(x_1, x_2, x_3) = 2 \cdot \sin(x_1 \cdot x_3) + e^{x_2^2}$  gegeben. Dann kann  $f$  mittels

$$\begin{aligned} f_1 &= x_1 = \pi_1^3(x_1, x_2, x_3) \\ f_2 &= x_2 = \pi_2^3(x_1, x_2, x_3) \\ f_3 &= x_3 = \pi_3^3(x_1, x_2, x_3) \end{aligned}$$

---

<sup>3</sup>Entspricht  $\circ$  dem unären  $-$ , so wird immer davon ausgegangen, daß nur ein Operand in die Rechnung eingeht.

$$\begin{aligned}
f_4 &= f_1 \cdot f_3 \\
f_5 &= \sin(f_4) \\
f_6 &= 2 \\
f_7 &= f_6 \cdot f_5 \\
f_8 &= \text{sqr}(f_2) \\
f_9 &= \exp(f_8) \\
f &= f_{10} = f_7 + f_9
\end{aligned}$$

zerlegt werden.

**Bemerkung 2.2** Eine solche Darstellung einer zerlegbaren Funktion wird *Codeliste* genannt. Die Funktionen  $f_1, \dots, f_n$  werden als *Unabhängige* oder *Veränderliche* von  $f$  bezeichnet. Die Ergebnisfunktion  $f_k$  heißt *Abhängige* von  $f$ . Alle anderen Funktionen der Funktionsfolge  $f_1, \dots, f_k$  heißen *Zwischenergebnisse* der Codeliste. Eine Codeliste ist im allgemeinen nicht eindeutig bestimmt. Für die Funktion des obigen Beispiels ist

$$\begin{aligned}
f_1 &= x_1 \\
f_2 &= x_2 \\
f_3 &= x_3 \\
f_4 &= f_1 \cdot f_3 \\
f_5 &= 2 \\
f_6 &= \sin(f_4) \\
f_7 &= \text{sqr}(f_2) \\
f_8 &= \exp(f_7) \\
f_9 &= f_5 \cdot f_6 \\
f &= f_{10} = f_9 + f_8
\end{aligned}$$

eine weitere Darstellung. Die Auswertungsreihenfolge auf dem Rechner wird u.a. durch die Grammatik der benutzten Programmiersprache festgelegt. Vorsicht ist bei optimierenden Compilern geboten. Hier kann eine gewünschte Auswertungsreihenfolge durch entsprechende Klammerung erzwungen werden.

**Bemerkung 2.3** Funktionen der Form  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  heißen *faktorisierbar*, falls alle Komponentenfunktionen  $g_j$  ( $1 \leq j \leq m$ ) von  $g$  faktorisierbar sind.

**Bemerkung 2.4** Funktionen, die Intervalle als Argumente haben, können in gleicher Art und Weise wie in Definition 2.6 als *faktorisierbar* erklärt werden. Dann sind die elementaren Operationen und Funktionen als *Intervalloperationen* bzw. *Intervallfunktionen* aufzufassen und werden ebenfalls mit  $\mathcal{OP}$  bzw.  $\mathcal{F}$  bezeichnet.

## 2.3 Darstellungsgraph

Äquivalent zur Beschreibung faktorisierbarer Funktionen durch eine Codeliste ist die visuelle Darstellung mittels eines Graphen. Ein solcher Darstellungsgraph veranschaulicht die Codeliste in einem gerichteten, zyklensfreien Graphen<sup>4</sup>. Abbildung 1 (S.

<sup>4</sup>Eine genaue Definition ist z.B. in [27] zu finden.

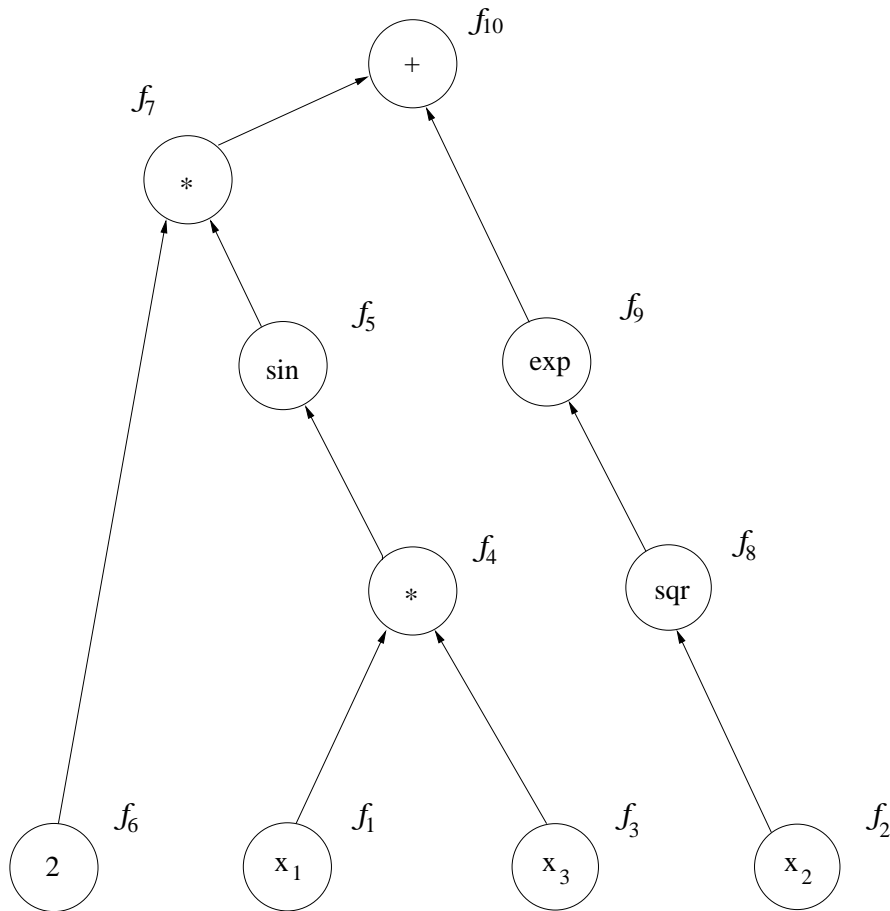


Abbildung 1: Darstellungsgraph der Funktion aus Bsp. 3.1.1

8) zeigt den Darstellungsgraphen der Funktion  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  mit  $f(x_1, x_2, x_3) = 2 \cdot \sin(x_1 \cdot x_3) + e^{x_2^2}$  aus Beispiel 2.1. Wie man ihr entnehmen kann, ist jeder Operation aus  $\mathcal{F}$  oder  $\mathcal{OP}$  ein Knoten zugeordnet. Jeder dieser Knoten hat entweder einen oder zwei Vorgänger, abhängig davon, ob dieser Knoten eine unäre oder binäre Operation repräsentiert. Nur die Knoten der in der Funktion vorkommenden Konstanten und der Eingangsgrößen, welche die Unabhängigen darstellen, haben keine Vorgänger. Der Knoten ohne Nachfolger wird als Ergebnisknoten (Wurzel) bezeichnet. Seine Operation liefert den Wert von  $f$ . Offensichtlich gibt es zu jeder faktorisierten Funktion einen Darstellungsgraphen und umgekehrt.

Im allgemeinen treten in Computerprogrammen Funktionen nicht nur in Form zerlegbarer Funktionen auf. Wie man Definition 2.6 entnehmen kann, dürfen Bedingungenweisungen und Schleifen<sup>5</sup> innerhalb faktorisierbarer Funktionen nicht vorkommen. Auf Verallgemeinerungen und dabei auftretende zusätzliche Probleme wird in [6] näher eingegangen.

<sup>5</sup>Damit sind Schleifen gemeint, bei denen die Anzahl der Schleifendurchläufe erst zur Laufzeit bekannt ist.



### 3 Allgemeines Prinzip der Rückwärtsrechnung

#### 3.1 Ein einführendes Beispiel

Das Prinzip der Rückwärtsrechnung wird zunächst anhand der laufzeiteffizienten Berechnung des Gradienten einer Funktion veranschaulicht (Literatur zum Thema *Automatische Differentiation* [4, 8, 7, 9].) Dazu sei die in  $\mathbb{R}^3$  beliebig oft differenzierbare und faktorisierte Funktion  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  mit  $f(x_1, x_2, x_3) = x_3 + \sin(x_1 \cdot x_2)$  gegeben.  $f$  läßt sich zum Beispiel mittels der Codeliste

$$f_1 = x_1 \quad (4)$$

$$f_2 = x_2 \quad (5)$$

$$f_3 = x_3 \quad (6)$$

$$f_4 = f_1 \cdot f_2 \quad (7)$$

$$f_5 = \sin(f_4) \quad (8)$$

$$f_6 = f_3 + f_5. \quad (9)$$

beschreiben. Differenziert man nun (4) bis (9) nach  $x_1$ ,  $x_2$  und  $x_3$ , so erhält man das Gleichungssystem

$$\left. \begin{array}{lll} \frac{\partial f_1}{\partial x_1} = 1 & \frac{\partial f_1}{\partial x_2} = 0 & \frac{\partial f_1}{\partial x_3} = 0 \\ \frac{\partial f_2}{\partial x_1} = 0 & \frac{\partial f_2}{\partial x_2} = 1 & \frac{\partial f_2}{\partial x_3} = 0 \\ \frac{\partial f_3}{\partial x_1} = 0 & \frac{\partial f_3}{\partial x_2} = 0 & \frac{\partial f_3}{\partial x_3} = 1 \\ \frac{\partial f_4}{\partial x_1} = \frac{\partial f_1}{\partial x_1} f_2 + \frac{\partial f_2}{\partial x_1} f_1 & \frac{\partial f_4}{\partial x_2} = \frac{\partial f_1}{\partial x_2} f_2 + \frac{\partial f_2}{\partial x_2} f_1 & \frac{\partial f_4}{\partial x_3} = \frac{\partial f_1}{\partial x_3} f_2 + \frac{\partial f_2}{\partial x_3} f_1 \\ \frac{\partial f_5}{\partial x_1} = \frac{\partial f_4}{\partial x_1} \cos(f_4) & \frac{\partial f_5}{\partial x_2} = \frac{\partial f_4}{\partial x_2} \cos(f_4) & \frac{\partial f_5}{\partial x_3} = \frac{\partial f_4}{\partial x_3} \cos(f_4) \\ \frac{\partial f_6}{\partial x_1} = \frac{\partial f_3}{\partial x_1} + \frac{\partial f_5}{\partial x_1} & \frac{\partial f_6}{\partial x_2} = \frac{\partial f_3}{\partial x_2} + \frac{\partial f_5}{\partial x_2} & \frac{\partial f_6}{\partial x_3} = \frac{\partial f_3}{\partial x_3} + \frac{\partial f_5}{\partial x_3} \end{array} \right\} \quad (10)$$

durch Anwendung der elementaren Differentiationsregeln. (10) ist in Matrixform gleichbedeutend mit

$$\underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ -f_2 & -f_1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -\cos(f_4) & 1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 1 \end{pmatrix}}_{=: I - \alpha =: A} \cdot \underbrace{\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} \\ \frac{\partial f_4}{\partial x_1} & \frac{\partial f_4}{\partial x_2} & \frac{\partial f_4}{\partial x_3} \\ \frac{\partial f_5}{\partial x_1} & \frac{\partial f_5}{\partial x_2} & \frac{\partial f_5}{\partial x_3} \\ \frac{\partial f_6}{\partial x_1} & \frac{\partial f_6}{\partial x_2} & \frac{\partial f_6}{\partial x_3} \end{pmatrix}}_{=: H = (h^1, h^2, h^3)} = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}}_{=: Z = (z^1, z^2, z^3)},$$

wobei

$$\alpha := \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ f_2 & f_1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \cos(f_4) & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{pmatrix}$$

eine untere Dreiecksmatrix darstellt, deren Diagonalelemente Null sind. Die Berechnung des Gradienten von  $f$  an einer bestimmten Stelle  $x_0 \in \mathbb{R}^3$  entspricht somit dem Lösen der drei linearen Gleichungssysteme

$$A h^1 = z^1 \quad (11)$$

$$A h^2 = z^2 \quad (12)$$

$$A h^3 = z^3. \quad (13)$$

Diese Systeme können durch Vorwärtseinsetzen gelöst werden. Zum einen gilt für  $1 \leq l \leq 3$ :

$$h_i^l = z_i^l = \delta_{li} \quad (1 \leq i \leq 3) \text{ und} \quad (14)$$

$$h_i^l = \sum_{j=1}^{i-1} \alpha_{ij} h_j^l + z_i^l \quad (15)$$

$$= \sum_{j=1}^{i-1} \alpha_{ij} h_j^l \quad (4 \leq i \leq 6). \quad (16)$$

Zum anderen ist  $\det(I - \alpha) \neq 0$ . Für  $D = (d_{ij})_{1 \leq i, j \leq 6} := (I - \alpha)^{-1}$  ist wegen (11)-(13)  $D z^l = h^l$  ( $1 \leq l \leq 3$ ). Es gilt:

$$h_i^l = \sum_{j=1}^6 d_{ij} z_j^l \quad (1 \leq i \leq 6)$$

und speziell

$$\frac{\partial f_6}{\partial x_l} = h_6^l = \sum_{j=1}^6 d_{6j} z_j^l = d_{6j} \cdot \delta_{jl}, \quad (17)$$

also

$$\text{grad}(f)(x_0) = (d_{61}, d_{62}, d_{63}).$$

Die Größen  $d_{ij}$  ( $1 \leq i, j \leq 6$ ) lassen sich wegen

$$\begin{aligned} D(I - \alpha) &= I \\ D &= I + D\alpha \end{aligned} \quad (18)$$

wieder rekursiv berechnen. (18) ist gleichbedeutend mit

$$d_{ij} = \begin{cases} 0 & , \quad i < j \\ 1 & , \quad i = j \\ \sum_{m=j+1}^i d_{im} \alpha_{mj} & , \quad j = i - 1, \dots, 1 \end{cases}$$

und speziell für  $i = 6$  kann man die Größen  $d_{6j}$  ( $j = 5, \dots, 1$ ) beginnend mit  $d_{66} := 1$  und  $d_{6j} := 0$  ( $1 \leq j \leq 5$ ) rekursiv ermitteln.

### 3.2 Der allgemeine Fall

Hat man einen Algorithmus gegeben, bei dem die Funktionswerte der Codeliste als Lösung eines linearen Gleichungssystem darstellbar sind und es sich bei der Systemmatrix um eine untere Dreiecksmatrix handelt, können die Aussagen des letzten Abschnitts generalisiert werden. Die folgenden Aussagen gehen auf [8] zurück. Gegeben sei das System  $(\Gamma - \alpha) h = z$ , welches folgende Gestalt besitzt:

$$\begin{pmatrix} \frac{1}{\gamma_1} & 0 & 0 & \cdots & 0 & 0 \\ -\alpha_{21} & \frac{1}{\gamma_2} & 0 & \cdots & 0 & 0 \\ -\alpha_{31} & -\alpha_{32} & \frac{1}{\gamma_3} & \cdots & 0 & 0 \\ \vdots & & \ddots & \ddots & \vdots & \vdots \\ \vdots & & & \ddots & \frac{1}{\gamma_{s-1}} & 0 \\ -\alpha_{s1} & \cdots & \cdots & \cdots & -\alpha_{ss-1} & \frac{1}{\gamma_s} \end{pmatrix} \cdot \begin{pmatrix} h_1 \\ h_2 \\ h_3 \\ \vdots \\ h_{s-1} \\ h_s \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_{s-1} \\ z_s \end{pmatrix} \quad (19)$$

Es ist  $\Gamma := (\delta_{ij} \cdot \gamma_i^{-1})_{1 \leq i, j \leq s} \in \mathbb{R}^{s \times s}$ ,  $\gamma_i \neq 0$  für  $1 \leq i \leq s$  und  $\alpha := (\alpha_{ij})_{1 \leq i, j \leq s}$  mit  $\alpha \in \mathbb{R}^{s \times s}$ , wobei  $\alpha_{ij} = 0$  für  $i \leq j$ .  $h$  und  $z$  sind Vektoren des  $\mathbb{R}^s$  und  $h$  kann wieder wegen der speziellen Bauart der Matrix  $\Gamma - \alpha$  durch Vorwärtsauffösen analog zu (14) und (15) bestimmt werden. Es gelten wegen (19) die Gleichungen:

$$\begin{aligned} h_1 &= \gamma_1 z_1 \\ h_i &= \gamma_i \left( \sum_{j=1}^{i-1} \alpha_{ij} h_j + z_i \right) \quad (1 < i \leq s). \end{aligned}$$

Wie im speziellen Fall der Gradientenberechnung kann man (19) wieder rückwärts auflösen. Sei dazu  $D = (d_{ij})_{1 \leq i, j \leq s} := (\Gamma - \alpha)^{-1}$ . Aufgrund von  $\det(\Gamma - \alpha) = \det(\Gamma) \neq 0$  existieren  $D$  und  $\Gamma^{-1}$ . Es gilt:

$$\begin{aligned} D(\Gamma - \alpha) &= I \\ D\Gamma &= I + D\alpha \\ D &= (I + D\alpha)\Gamma^{-1}. \end{aligned}$$

Zusammenfassend kann der folgende Satz formuliert werden (entnommen aus [8]):

**Satz 3.1** Seien  $\Gamma, \alpha, D \in \mathbb{R}^{s \times s}$  ( $s \in \mathbb{N}$ ) wie oben definiert. Mit  $z \in \mathbb{R}^s$  gilt für die Lösung  $h \in \mathbb{R}^s$  des linearen Gleichungssystems  $(\Gamma - \alpha) h = z$ :

$$\begin{aligned} h_1 &= \gamma_1 z_1 \\ h_i &= \gamma_i \left( \sum_{j=1}^{i-1} \alpha_{ij} h_j + z_i \right) \quad (1 < i \leq s) \\ &= \sum_{j=1}^s d_{ij} z_j, \end{aligned} \quad (20)$$

wobei die Berechnung der Elemente der Matrix  $D$  rekursiv durch

$$d_{ij} = 0 \quad (i < j),$$

$$d_{ij} = \gamma_j \quad (i = j), \quad (21)$$

$$d_{ij} = \gamma_j \sum_{m=j+1}^i d_{im} \alpha_{mj} \quad (j = i-1, i-2, \dots, 1) \quad (22)$$

erfolgen kann.

Hat man einmal alle  $d_{ij}$  ( $1 \leq j \leq i \leq s$ ) berechnet, so kann man mit (20) das Gleichungssystem (19) für verschiedene rechte Seiten  $z$  auswerten. Dies wurde bei der Berechnung des Gradienten im letzten Abschnitt benutzt (vgl. (17)) und begründet die Überlegenheit der Rückwärtsmethode gegenüber der Vorwärtsmethode hinsichtlich des Aufwandes zur Berechnung des Gradienten.

Setzt man in (21) und (22) speziell  $i = s$ , so läßt sich das folgende Korollar ableiten.

**Korollar 3.1 (allgemeine Rückwärtsmethode)** Für  $s \in \mathbb{N}$  sei  $\gamma_i, z_i \in \mathbb{R}$  ( $1 \leq i \leq s$ ) und  $\alpha_{ij} \in \mathbb{R}$  ( $1 \leq i, j \leq s$ ) mit  $\alpha_{ij} = 0$  ( $i \leq j$ ),  $\gamma_i \neq 0$  ( $1 \leq i \leq s$ ) gegeben. Ist dann die (endliche) Folge  $h_1, \dots, h_s$  durch

$$h_1 := \gamma_1 z_1 \quad (23)$$

$$h_i := \gamma_i \left( \sum_{j=1}^{i-1} \alpha_{ij} h_j + z_i \right) \quad (1 < i \leq s) \quad (24)$$

definiert, so gilt:

$$\begin{aligned} h_s &= \sum_{m=1}^s d_m z_m \text{ mit} \\ d_s &= \gamma_s \text{ und} \\ d_j &= \gamma_j \sum_{m=j+1}^s d_m \alpha_{mj} \quad (j = s-1, \dots, 1). \end{aligned}$$

Beweis: Siehe [8].  $\square$

Mit Hilfe des Korollars hat man für Algorithmen, welche durch (23) und (24) beschreibbar sind, die Möglichkeit, die Berechnung durch Rückwärtseinsetzen durchzuführen. Als Anwendungen dieses Korollars wird im nächsten Abschnitt ein Verfahren zur automatischen Fehlerabschätzung von Gleitkommaalgorithmen vorgestellt. In [6] wird die Methodik auch zur schnellen Berechnung von Intervallsteigungen (siehe auch [28, 29]) herangezogen.

## 4 Automatische Fehlerkontrolle

### 4.1 Problemstellung

Hofschuster und Krämer leiten in [12] ein Fehlerkalkül her, um maximale Fehler von auf der Maschine implementierten Standardfunktionen abzuschätzen. Dieses Verfahren

beruht auf dem Vorwärtsdurchlaufen des Darstellungsgraphen. Iri verwendet mittels Rückwärtsrechnung gewonnene Fehlerschranken, um optimale Abbruchkriterien für Iterationsverfahren zu ermitteln (vgl. [26]). Weiterführende Literatur zur automatischen Fehlerkontrolle findet man in [1], [2], [14], [15], [21], [22], [23] und [24].

In diesem Kapitel soll ein Verfahren zur automatischen Berechnung einer Gesamtfehlerschranke einer faktorisierten Funktion durch Rückwärtsrechnung vorgestellt werden. Um Korollar 3.1 benutzen zu können, wird zunächst in Analogie zu [8] ein Verfahren hergeleitet, welches den Gesamtfehler durch Vorwärtsrechnung berechnet. Aus diesem soll anschließend das gewünschte Verfahren abgeleitet werden. Desweiteren wird im letzten Abschnitt dieses Kapitels der gewonnene Algorithmus mit einem bereits implementierten Verfahren (siehe dazu [1], [2]) zur sicheren Abschätzung einer Gesamtfehlerschranke einer faktorisierten Funktion anhand einiger numerischer Beispiele verglichen.

## 4.2 Fehlerkontrolle durch Vorwärtsrechnung

Für die Herleitung sind einige Definitionen notwendig. Es wird ähnlich wie in [8] vorgegangen.

**Definition 4.1** Mit  $\widetilde{\mathcal{OP}}$  und  $\widetilde{\mathcal{F}}$  werden die zu  $\mathcal{OP}$  und  $\mathcal{F}$  entsprechenden Operationen auf der Maschine bezeichnet. Dabei entspricht  $\widetilde{\mathcal{OP}}$  der in Definition 2.3 vorgestellten Menge der Gleitkommaoperationen (vgl. S. 5).

**Definition 4.2 (Datenfehler)** Seien  $\widetilde{f}_l \in S$  und  $\widetilde{f}_r \in S$  Näherungen von  $f_l$  und  $f_r$ . Dann ist der absolute Datenfehler  $\Delta_{D,f}$  für die elementaren Operationen und Funktionen durch

$$\begin{aligned}\Delta_{D,f} &:= (\widetilde{f}_l \circ \widetilde{f}_r) - (f_l \circ f_r), \quad f = f_l \circ f_r \\ \Delta_{D,f} &:= g(\widetilde{f}_l) - g(f_l), \quad f = g(f_l)\end{aligned}$$

erklärt.

**Definition 4.3 (Rundungsfehler)** Gegeben seien die Gleitkommanäherungen  $\widetilde{f}_l, \widetilde{f}_r$  von  $f_l$  und  $f_r$ . Für  $f = f_l \circ f_r$  ist der absolute Rundungsfehler  $\Delta_{R,f}$  durch

$$\Delta_{R,f} := (\widetilde{f}_l \boxtimes \widetilde{f}_r) - (\widetilde{f}_l \circ \widetilde{f}_r), \quad \circ \in \mathcal{OP}, \quad \boxtimes \in \widetilde{\mathcal{OP}}$$

bestimmt. Ist  $f = g(f_l)$ , so definiert man  $\Delta_{R,f}$  mittels

$$\Delta_{R,f} := \widetilde{g}(\widetilde{f}_l) - g(\widetilde{f}_l), \quad g \in \mathcal{F}, \quad \widetilde{g} \in \widetilde{\mathcal{F}}.$$

**Definition 4.4 (absoluter Gesamtfehler)** Sei die faktorisierte Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  mit einer Codeliste  $f_1, \dots, f_k$  und deren Näherung  $\widetilde{f}$  mit Codeliste  $\widetilde{f}_1, \dots, \widetilde{f}_k$  gegeben. Weiterhin sei

$$\begin{aligned}f_j = x_j &\quad \text{bzw.} \quad \widetilde{f}_j = \widetilde{x}_j \quad \text{für} \quad 1 \leq j \leq n \quad \text{und o.B.d.A.} \\ f_j = c_j &\quad \text{bzw.} \quad \widetilde{f}_j = \widetilde{c}_j \quad \text{für} \quad n+1 \leq j \leq n+m < k\end{aligned}$$

gesetzt. Dann ist der Gesamtfehler  $\Delta_{G,f_j}$  für  $1 \leq j \leq k$  wie folgt definiert:

$$\Delta_{G,f_j} := \begin{cases} \widetilde{f}_j - f_j & , \quad 1 \leq j \leq n+m \\ \widetilde{f}_j(\widetilde{x}_1, \dots, \widetilde{x}_n) - f_j(x_1, \dots, x_n) & , \quad n+m+1 \leq j \leq k. \end{cases}$$

**Bemerkung 4.1** Aus dem Betrag des absoluten Gesamtfehlers ergibt sich der Betrag des relativen Gesamtfehlers  $\Delta_{G,f_j}^{rel}$  unter den Voraussetzungen von 4.4 durch

$$|\Delta_{G,f_j}^{rel}| := \frac{|\Delta_{G,f_j}|}{|f_j|}, \quad f_j \neq 0, \quad 1 \leq j \leq k.$$

Die beiden folgenden Lemmata und der darauffolgende Satz beantworten die Frage, wie der Datenfehler und der Gesamtfehler bei elementaren Operationen und Funktionen aus den Datenfehlern, Rundungsfehlern bzw. den Gesamtfehlern der Operanden entstehen.

**Lemma 4.1** Für  $n + m + 1 \leq j \leq k$  gilt mit den Voraussetzungen aus Definition 4.4:

$$\Delta_{G,f_j} = \Delta_{D,f_j} + \Delta_{R,f_j}.$$

Beweis: Durch Anwendung der Definitionen folgt die Behauptung.  $\square$

**Lemma 4.2** Seien die Näherungen  $\tilde{f}_l \in S$  und  $\tilde{f}_r \in S$  von  $f_l$  und  $f_r$  und bekannte Fehlerschranken  $\Delta_{G,f_l} = \tilde{f}_l - f_l$  und  $\Delta_{G,f_r} = \tilde{f}_r - f_r$  gegeben. Dann gilt für den Datenfehler  $\Delta_{D,f}$  bei Anwendung einer elementaren Operation:

$$\Delta_{D,f} = -\Delta_{G,f_l}, \quad \text{falls } f = -f_l \quad (25)$$

$$\Delta_{D,f} = \Delta_{G,f_l} + \Delta_{G,f_r}, \quad \text{falls } f = f_l + f_r \quad (26)$$

$$\Delta_{D,f} = \Delta_{G,f_l} - \Delta_{G,f_r}, \quad \text{falls } f = f_l - f_r \quad (27)$$

$$\Delta_{D,f} = \tilde{f}_r \Delta_{G,f_l} + f_l \Delta_{G,f_r}, \quad \text{falls } f = f_l \cdot f_r \quad (28)$$

$$\Delta_{D,f} = (\Delta_{G,f_l} - (f_l/f_r) \Delta_{G,f_r})/\tilde{f}_r, \quad \text{falls } f = f_l/f_r, \quad (\tilde{f}_r, f_r \neq 0). \quad (29)$$

Für den Datenfehler bei Anwendung elementarer Funktionen  $g \in \mathcal{F}$  gilt:

$$\Delta_{D,f} = g'(\xi) \Delta_{G,f_l}, \quad \text{falls } f = g(f_l). \quad (30)$$

Dabei wird vorausgesetzt, daß  $g'(\xi)$  für alle  $\xi \in f_l \sqcup \tilde{f}_l$  existiert.

Beweis:

Die Formeln (25)-(30) lassen sich durch Anwendung der Definitionen zeigen. In (30) benötigt man zusätzlich den Mittelwertsatz der Differentialrechnung. Hier wird nur der Beweis von (28) durchgeführt. Wegen Definition 4.2 gilt:

$$\begin{aligned} \Delta_{D,f} &= \tilde{f}_l \tilde{f}_r - f_l f_r \\ &= \tilde{f}_l \tilde{f}_r - \tilde{f}_r f_l + \tilde{f}_r f_l - f_l f_r \\ &= \tilde{f}_r (\tilde{f}_l - f_l) + (\tilde{f}_r - f_r) f_l \\ &= \tilde{f}_r \Delta_{G,f_l} + f_l \Delta_{G,f_r}. \end{aligned}$$

Somit folgt die Behauptung.  $\square$

Aus den beiden Lemmata und Definition 4.4 ergibt sich

**Satz 4.1** *Gelten die Voraussetzungen von Definition 4.4 und Lemma 4.2, so erhält man für den Gesamtfehler der  $j$ -ten Operation der Codeliste:*

$$\Delta_{G,f_j} = \widetilde{f}_j - f_j = \Delta_{R,f_j}, \text{ falls } 1 \leq j \leq n + m \quad (31)$$

$$\Delta_{G,f_j} = -\Delta_{G,f_l} + \Delta_{R,f_j}, \text{ falls } f_j = -f_l \quad (32)$$

$$\Delta_{G,f_j} = \Delta_{G,f_l} + \Delta_{G,f_r} + \Delta_{R,f_j}, \text{ falls } f_j = f_l + f_r \quad (33)$$

$$\Delta_{G,f_j} = \Delta_{G,f_l} - \Delta_{G,f_r} + \Delta_{R,f_j}, \text{ falls } f_j = f_l - f_r \quad (34)$$

$$\Delta_{G,f_j} = \widetilde{f}_r \Delta_{G,f_l} + f_l \Delta_{G,f_r} + \Delta_{R,f_j}, \text{ falls } f_j = f_l \cdot f_r \quad (35)$$

$$\Delta_{G,f_j} = (\Delta_{G,f_l} - (f_l/f_r) \Delta_{G,f_r})/\widetilde{f}_r + \Delta_{R,f_j}, \text{ falls } f_j = f_l/f_r \quad (36)$$

$$\Delta_{G,f_j} = g'(\xi) \Delta_{G,f_l} + \Delta_{R,f_j}, \text{ falls } f_j = g(f_l), \xi \in f_l \sqcup \widetilde{f}_l. \quad (37)$$

### 4.3 Fehlerkontrolle durch Rückwärtsrechnung

Die Gleichungen (31)-(37) bilden ein lineares Gleichungssystem mit einer unteren Dreiecksmatrix als Systemmatrix, wie in Abschnitt 3.2 (S. 11f) beschrieben. Die Größen  $\Delta_{G,f_j}$  und  $\Delta_{R,f_j}$  ( $1 \leq j \leq k$ ) entsprechen den Vektoren  $h$  bzw.  $z$  des dort beschriebenen Systems. Somit kann der Gesamtfehler  $\Delta_{G,f_k}$  durch Rückwärtsrechnung bestimmt werden, da die Voraussetzungen von Korollar 3.1 (S. 12) erfüllt sind. Es gilt:

$$\Delta_{G,f} = \Delta_{G,f_k} = \sum_{j=1}^k d_j \Delta_{R,f_j}, \quad (38)$$

wobei die Größen  $d_j$  nach Korollar 3.1 durch die folgende Akkumulation, beginnend mit  $d_k := 1$  und  $d_j := 0$  ( $1 \leq j \leq k - 1$ ), berechnet werden:

$$d_l := d_l - d_j, \text{ falls } f_j = -f_l \quad (39)$$

$$d_l := d_l + d_j, \quad d_r := d_r + d_j, \text{ falls } f_j = f_l + f_r \quad (40)$$

$$d_l := d_l + d_j, \quad d_r := d_r - d_j, \text{ falls } f_j = f_l - f_r \quad (41)$$

$$d_l := d_l + \widetilde{f}_r d_j, \quad d_r := d_r + f_l d_j, \text{ falls } f_j = f_l \cdot f_r \quad (42)$$

$$d_l := d_l + d_j/\widetilde{f}_r, \quad d_r := d_r - f_j d_j/\widetilde{f}_r, \text{ falls } f_j = f_l/f_r \quad (43)$$

$$d_l := d_l + g'(\xi) d_j, \text{ falls } f_j = g(f_l), \xi \in f_l \sqcup \widetilde{f}_l \quad (44)$$

**Bemerkung 4.2** *Beim Auftreten des unären Operators – wird auch auf der Maschine nicht gerundet. Es wird lediglich das Vorzeichenflag neu gesetzt. Der entsprechende Summand kann in (38) weggelassen werden.*

**Bemerkung 4.3** *Die Größen  $d_j$  in (38) bestimmen die Größe des Gesamtfehlers. Sie werden als Verstärkungsfaktoren bezeichnet.*

Da die reellen Werte  $f_j$  ( $1 \leq j \leq k$ ) und  $\xi$  nicht bekannt sind, kann durch Einsatz der Intervallrechnung (Intervallauswertung von (39)-(44) und der Funktion  $f$ ) eine obere Schranke für den Rundungsfehler hergeleitet werden. Dazu sei  $F_j, D_j \in IS$  mit

$f_j, \widetilde{f}_j \in F_j$  und  $d_j \in D_j$  ( $1 \leq j \leq k$ ) auf der Maschine berechnet. Dann gilt aufgrund von (38) für den Betrag des Gesamtfehlers die Ungleichung:

$$|\Delta_{G,f}| \leq \sum_{j=1}^k |D_j| |\Delta_{R,f_j}|. \quad (45)$$

Handelt es sich bei  $f_j$  um eine elementare Operation, so erhält man für den Betrag des Rundungsfehlers durch Anwendung von (2) (S. 5) die Abschätzung:

$$|\Delta_{R,f_j}| \leq \begin{cases} \varepsilon |F_j| & , F_j \subseteq N \\ \text{MinReal} & , F_j \subseteq U \\ \varepsilon |F_j| + \text{MinReal} & , \text{sonst.} \end{cases} \quad (46)$$

Führt  $f_j$  eine elementare Funktion aus, so gilt wegen der Annahme (3) (S. 5) die folgende Ungleichung:

$$|\Delta_{R,f_j}| \leq \begin{cases} \varepsilon(f_j) |F_j| & , F_j \subseteq N \\ u(f_j) & , F_j \subseteq U \\ \varepsilon(f_j) |F_j| + u(f_j) & , \text{sonst.} \end{cases} \quad (47)$$

Mit den Ungleichungen (45), (46) und (47) läßt sich eine obere Schranke für den Betrag des entstehenden Gesamtfehlers berechnen.

Wenn in der zugrundeliegenden Funktion binäre Operationen aus  $\mathcal{OP}$  auftreten, so ist es möglich, zusätzliche Aussagen über die Größe der bei diesen Operationen auftretenden Rundungsfehler zu treffen. In [1] werden hinreichende Kriterien dafür angegeben, unter welchen Umständen diese Operationen rundungsfehlerfrei verlaufen. So erhält man zum Beispiel für die Addition zweier Gleitkommazahlen folgende Aussage:

**Satz 4.2** *Sei die Anzahl der abschließenden Nullen der Mantisse einer Gleitkommazahl  $z$  durch die Funktion  $z_{\text{trail}}(z)$  berechenbar. Wenn mindestens eine der Bedingungen*

1.  $-2 \leq \frac{x}{y} \leq -\frac{1}{2}$
2.  $e(x+y) \leq \min(e(x) + z_{\text{trail}}(x), e(y) + z_{\text{trail}}(y))$
3.  $x = 0 \vee y = 0$

*erfüllt ist und bei der Berechnung kein Unterlauf entsteht, so verläuft die Addition zweier Gleitkommazahlen  $x, y \in S$  auf der Maschine rundungsfehlerfrei, d.h. es gilt  $|(x \boxplus y) - (x + y)| = 0$ .*

Beweis: Siehe [1].  $\square$

Ebenso können Aussagen für die verbleibenden elementaren Operationen formuliert werden (vgl. dazu [1]).



## 4.4 Implementierung

Durch die zur Klasse `rivall` (siehe [6]) gehörenden `friend`-Funktionen

```
real reverse_error(TAPE& wt, int y_c);
void reverse_error_j(TAPE& wt, int& j,
                    intvector& rnd_free_op,
                    rvector& delta);

real reverse_error(TAPE& wt, int y_c, rvector& stoerung);
void reverse_error_j(TAPE& wt, int& j,
                    intvector& rnd_free_op,
                    rvector& delta,
                    rvector& stoerung);
```

wurde (45) für alle Operationen aus  $\mathcal{OP}$  in `reverse.cpp` implementiert. Um eine garantierte Oberschranke des Betrags des absoluten Gesamtfehlers auf der Maschine zu erhalten, müssen gerichtete Rundungen verwendet werden. Nachdem mit der Funktion

```
void forward_sweep(TAPE& working_tape,
                  rivall (*f)(TAPE& working_tape, rivall_vec& x),
                  ivector& d_vec,
                  interval& f_val)
```

alle Funktionswerteinschließungen  $F_j$  ( $1 \leq j \leq k$ ) berechnet wurden, kann anschließend der Betrag des absoluten Gesamtfehlers mit einer der beiden Funktionen `reverse_error(...)` bestimmt werden. Bei der zweiten dieser Funktionen ist es möglich, Eingangsstörungen direkt zu übergeben. Die Abschätzung des Rundungsfehlers der entsprechenden Operation erfolgt mittels (46)<sup>6</sup> unter Beachtung möglicher rundungsfehlerfrei verlaufender Operationen. Ob eine der binären Operationen aus  $\widetilde{\mathcal{OP}}$  exakt verläuft, wird durch die Funktionen

```
int check_rnd_free(real& delta_l, real& delta_r,
                  interval& f_l, interval& f_r,
                  int _op );
```

und

```
int check_rnd_free_U(real& delta_l, real& delta_r,
                    interval& f_l, interval& f_r,
                    int _op );
```

festgestellt. Diese sind maschinenabhängig und entsprechen den jeweiligen Funktionen aus der Implementierung des Kalküls in [1, 2].

Sind Eingangsgrößen oder Konstanten exakt darstellbar, so benutzt man, nachdem ihre Einschließungswerte den Variablen vom Typ `rivall` übergeben worden sind, die `friend`-Funktionen

---

<sup>6</sup>Falls auch elementare Funktionen in der Implementierung berücksichtigt werden sollen, so kann dies durch geeignete Fallunterscheidungen bei der Berechnung des Rundungsfehlers geschehen. Dann muß (47) benutzt werden.

```
void exactVar(rivall& s)
void exactConst(rivall& s).
```

Dabei wird an der Stelle, an der die Variable im Tape erscheint, eine Marke gesetzt, welche der reverse-sweep dazu benutzt, diese Variable als fehlerfreie Größe zu identifizieren. Analog ist vorzugehen, falls bekannt ist, daß eine eingehende Größe oder Konstante nach der nächstgelegenen Maschinenzahl gerundet wird. Es sind dann die beiden `friend`-Funktionen

```
void rnd_n_Var(rivall& s)
void rnd_n_Const(rivall& s)
```

zu verwenden.

## 4.5 Vergleich

An einigen einfachen Testfunktionen wird nun das Verhalten des durch obige Implementierung realisierten Algorithmus (Rückwärtsmethode (R)) mit der Implementierung des Fehlerkalküls von Hofschuster und Krämer (im folgenden immer als Vorwärtsmethode (V) bezeichnet) verglichen.

**Beispiel 4.1 (FormelAuswertung)** *Gegeben sei die Formel*

$$y = a^2 + b^2 - c^2 + d^2 - e^2 \quad (48)$$

mit den Argumenten  $a = 320000.0$ ,  $b = 19.0/32768.0$ ,  $c = 39.0/65536.0$ ,  $d = 240000.0$  und  $e = 400000.0$ . Diese sind so gewählt, daß bei Hintereinanderausführung der Operationen gemäß (48) starke Auslöschung auftritt. Man kann (48) auf verschiedene Weise auswerten, wie z.B.

$$\begin{aligned} y_1 &= (((a^2 + b^2) - c^2) + d^2) - e^2 \\ y_2 &= ((a^2 + b^2) - c^2) + (d^2 - e^2) \\ y_3 &= a^2 + ((b^2 - c^2) + (d^2 - e^2)) \\ y_4 &= (b^2 - c^2) + (d^2 - e^2 + a^2). \end{aligned}$$

Je nach Auswertungsreihenfolge ergeben sich verschiedene Fehlerschranken für den Betrag des absoluten Fehlers, aus denen man erkennen kann, welche Reihenfolge optimal ist.

Bei der intervallmäßigen Auswertung treten bei Anwendung von  $y_1$ ,  $y_2$  und  $y_3$  starke Überschätzungen auf:

$$\begin{aligned} y_1 &\in [-3.051758E - 005, 3.051758E - 005] \\ y_2 &\in [-1.525879E - 005, 1.525879E - 005] \\ y_3 &\in [-1.525879E - 005, 0.000000E + 000] \\ y_4 &\in [-1.792796E - 008, -1.792795E - 008] \end{aligned}$$

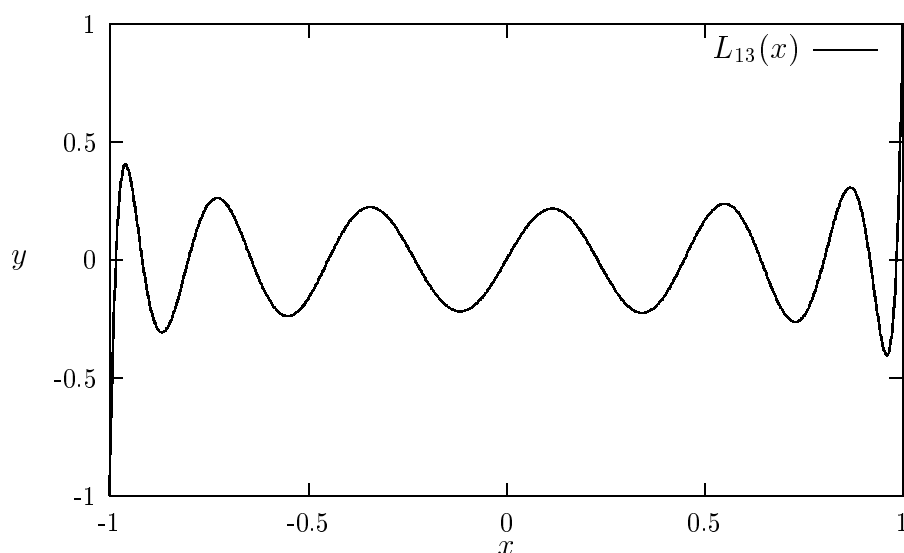
Wie Tabelle 1 zeigt, ist der absolute Fehler auch dort am kleinsten, wo die Intervallüberschätzung am geringsten ausfällt. Ebenso ist zu erkennen, daß die Fehlerschranken, die durch Rückwärtsrechnung bestimmt wurden, den durch Vorwärtsrechnung ermittelten entsprechen.

Auswertung	abs. Fehler (V)	abs. Fehler (R)
$y_1$	8.100187E-005	8.100187E-005
$y_2$	4.547474E-005	4.547474E-005
$y_3$	2.273737E-005	2.273737E-005
$y_4$	4.598607E-022	4.598607E-022

Tabelle 1: Absolute Fehler bei verschiedener Auswertungsreihenfolge in Bsp. 4.1

Das nächste Beispiel beschäftigt sich mit der Auswertung von Polynomen. Dabei ist zu erwarten, daß sich für verschiedene Auswertungsarten unterschiedliche Fehlerschranken ergeben, da zum einen die Anzahl der auszuführenden Operationen verschieden ist und zum anderen die Auswertungsreihenfolge der Operationen differiert.

**Beispiel 4.2 (Polynomauswertung)** *Es wird die Auswertung des in Abbildung 2 (S. 19) dargestellten Legendre-Polynoms  $L_{13}(x)$  vom Grad 13 auf dem Intervall  $[-1, 1]$  betrachtet. Die Koeffizienten lassen sich rekursiv berechnen (vgl. dazu [32]). Für  $L_{13}$*

Abbildung 2: Legendre-Polynom  $L_{13}(x)$  auf dem Bereich  $[-1, 1]$ 

ergeben sich diese zu

$$\begin{aligned}
 a_1 &= 3003/1024 & a_3 &= -90090/1024 & a_5 &= 765765/1024 \\
 a_7 &= -2771340/1024 & a_9 &= 4849845/1024 & a_{11} &= -4056234/1024 \\
 a_{13} &= 1300075/1024
 \end{aligned}$$

und  $a_i = 0$  für  $i = 0, 2, \dots, 12$ . Gesucht sei wieder eine gleichmäßig gültige Fehlerschranke auf dem vorgegebenen Intervall  $[-1, 1]$ .

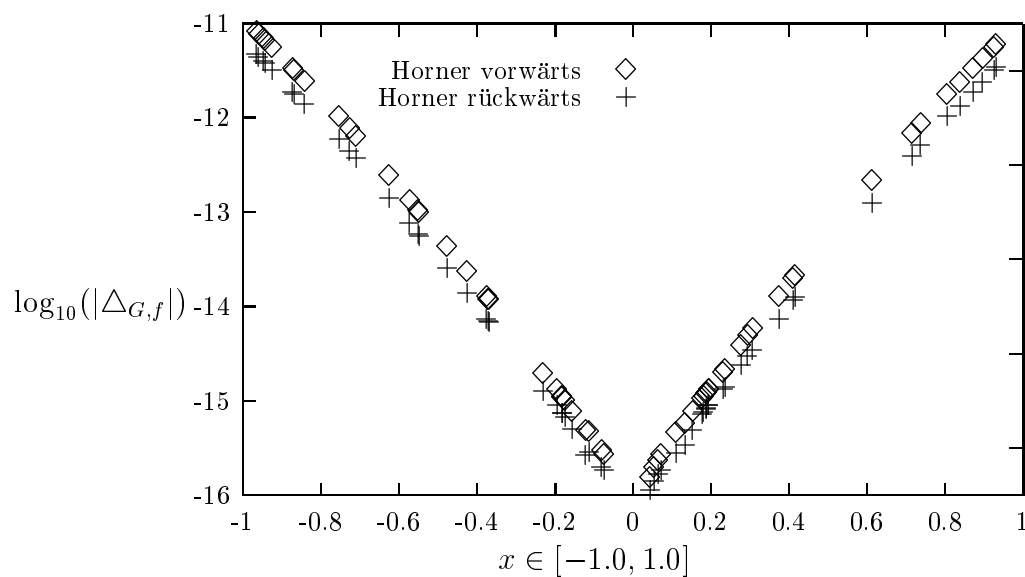


Abbildung 3: Absolute Fehler bei Auswertung durch das Hornerschema in Bsp. 4.2 mit zufällig erzeugten Argumenten aus  $[-1, 1]$

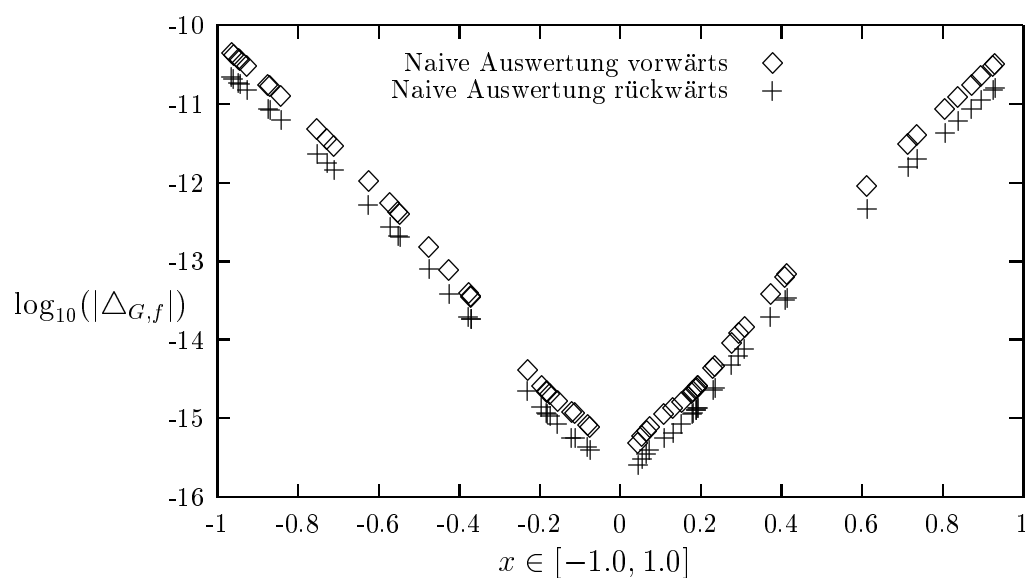


Abbildung 4: Absolute Fehler bei naiver Auswertung in Bsp. 4.2 mit zufällig erzeugten Argumenten aus  $[-1, 1]$

Die Auswertung erfolgte sowohl durch naive Auswertung des Polynoms und Auswertung durch das Hornerschema im Vorwärts- als auch im Rückwärtsmodus. Mit zufällig erzeugten Punktintervallen wurden zunächst Unterschranken für die von beiden Auswertungsarten berechneten Fehlerschranken bestimmt. In den Abbildungen 3 und 4 (S. 20) ist der prinzipielle Verlauf der Fehlerfunktion zu sehen. Es ist zu erkennen, daß die Fehler stark anwachsen, sobald der Betrag der Argumente gegen Eins geht. Das Alternieren der Koeffizientenvorzeichen macht sich hier verstärkt bemerkbar.

Für die Berechnung des Gesamtfehlers wurde das zugrundeliegende Intervall in 5 gleichgroße Teilbereiche  $B_1, \dots, B_5$  mit

$$\begin{aligned} B_1 &:= [-1.0, -0.6] & B_2 &:= [-0.6, -0.2] \\ B_3 &:= [-0.2, 0.2] & B_4 &:= [0.2, 0.6] \\ B_5 &:= [0.6, 1.0] \end{aligned}$$

unterteilt<sup>7</sup>, die dann wiederum in  $n$  Intervalle aufgeteilt wurden. Aus den Abbildungen 3 und 4 (S. 20) sowie den Ergebnistabellen 2 bis 5 (S. 22f) ist zu erkennen, daß die Auswertung durch das Hornerschema sowohl für  $n = 1$  als auch für  $n = 100$  der naiven Polynomauswertung überlegen ist. Mit der Rückwärtsmethode werden in diesem Beispiel bessere absolute Fehlerschranken erzielt als mit der Vorwärtsmethode.

Daß die durch Rückwärtsrechnung bestimmten Fehlerschranken nicht immer besser sein müssen, zeigt das nächste Beispiel.

**Beispiel 4.3 (Fehlerabschätzung arctan-Funktion)** Eine Möglichkeit zur Berechnung der arctan-Funktion auf der Maschine stellt der Algorithmus `q_atan` aus [13] (vgl. auch S. 26) dar. Der Verlauf des relativen Fehlers für diesen Algorithmus bei sinnvoller Unterteilung des Argumentbereichs ( $|x| > \text{q_atnt}$ , vgl. S. 26) ist in Abbildung 5 dargestellt<sup>8</sup>. Es fällt auf, daß beide Fehlerkurven prinzipiell gleich verlaufen, die Rückwärtsmethode in den Bereichen  $B_2$  bis  $B_6$  jedoch etwas größere Fehlerschranken liefert. Eine mögliche Ursache dafür ist die hier verwendete Methode zur Berücksichtigung des Approximationsfehlers (vgl. Abb. 6, S. 24), welcher bei der Polynomapproximation in Schritt 3b des Algorithmus `q_atan` auftritt. Dieser Approximationsfehler kann nicht wie bei der Vorwärtsmethode auf die bis dahin berechnete Fehlerschranke addiert werden. Vielmehr muß er als Störung in der weiteren Rechnung angesehen werden. Genauer: Es wird der Gesamtfehler bis einschließlich der Polynomapproximation berechnet, dann wird die Berechnung abgebrochen. Anschließend gehen Approximationsfehler und der bis dahin berechnete Gesamtfehler als Störung in eine erneute Rückwärtsrechnung ein<sup>9</sup>. Es muß beachtet werden, daß der dabei durchgeführte forward-sweep mit dem um die Störung aufgeblähten Funktionswert gestartet werden muß, um korrekte Einschließungen zu erhalten. Folglich können sich in den weiteren Berechnungen die Einschließungen  $F_j$  der in der Codeliste vorkommenden Funktionswerte  $f_j$  vergrößern und sich deshalb die Fehlerschranken verschlechtern. Durch weitere Unterteilung konnte dieser Effekt teilweise beseitigt werden.

<sup>7</sup>Auf der Maschine müssen Einschließungen von  $B_1$  bis  $B_5$  benutzt werden, da die Zahlen  $-0.6$ ,  $-0.2$ ,  $0.2$  und  $0.6$  nicht exakt darstellbar sind.

<sup>8</sup>Die dabei zugrundeliegenden Werte sind im Abschnitt 6.2 (S. 27f) ausführlich angegeben.

<sup>9</sup>Dazu wird die im letzten Abschnitt vorgestellte Funktion `real reverse_error(TAPE& wt, int y_c, rvector& stoerung)` benutzt.

Bereich	abs. Fehler (V)	abs. Fehler (R)
$B_1$	2.334474E-011	1.768227E-011
$B_2$	3.812549E-013	2.846412E-013
$B_3$	3.466742E-015	2.830968E-015
$B_4$	3.812549E-013	2.846412E-013
$B_5$	2.334474E-011	1.768227E-011

Tabelle 2: Absolute Fehler bei Auswertung durch das Hornerschema in Bsp. 4.2,  $n = 1$ 

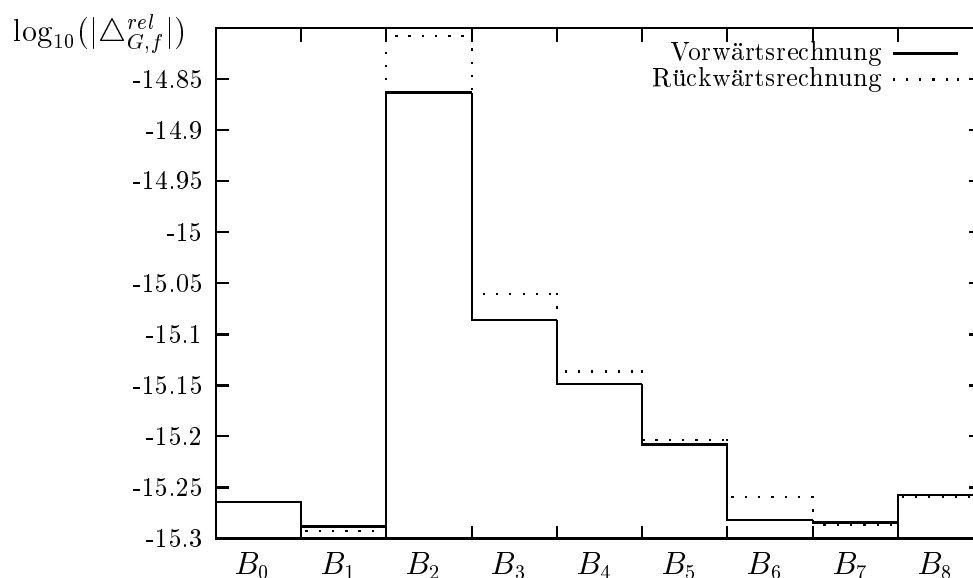
Bereich	abs. Fehler (V)	abs. Fehler (R)
$B_1$	6.720361E-011	5.296701E-011
$B_2$	9.064231E-013	7.393192E-013
$B_3$	6.142146E-015	6.011912E-015
$B_4$	9.064231E-013	7.393192E-013
$B_5$	6.720361E-011	5.296701E-011

Tabelle 3: Absolute Fehler bei naiver Auswertung in Bsp. 4.2,  $n = 1$ 

Bereich	abs. Fehler (V)	abs. Fehler (R)
$B_1$	1.124859E-011	6.784536E-012
$B_2$	1.832245E-013	1.111242E-013
$B_3$	1.418401E-015	1.005887E-015
$B_4$	1.832245E-013	1.111242E-013
$B_5$	1.124859E-011	6.784536E-012

Tabelle 4: Absolute Fehler bei Auswertung durch das Hornerschema in Bsp. 4.2,  $n = 100$

Bereich	abs. Fehler (V)	abs. Fehler (R)
$B_1$	6.160155E-011	3.263988E-011
$B_2$	7.509274E-013	4.020721E-013
$B_3$	2.766735E-015	1.571294E-015
$B_4$	7.509274E-013	4.020721E-013
$B_5$	6.160155E-011	3.263988E-011

Tabelle 5: Absolute Fehler bei naiver Auswertung in Bsp. 4.2,  $n = 100$ Abbildung 5: Relative Fehler bei Auswertung durch Vorwärts- und Rückwärtsrechnung von `q_atan`

Um eine möglichst kleine Gesamtfehlerschranke für den Algorithmus `q_atan` zu erhalten, können die Ergebnisse beider Berechnungsmethoden benutzt werden. Man bildet das Maximum der Minima der durch beide Verfahren bestimmten Fehlerschranken auf allen Teilbereichen. Dieser Wert ist dann  $1.555012E - 015$ .

**Bemerkung 4.4 (Fehlerabschätzung der exp-Funktion)** Für ein Tabellenverfahren, welches die Exponentialfunktion implementiert, wurde ebenfalls versucht, eine Fehlerschranke für den relativen Fehler zu berechnen. Dabei brach die Rückwärtsrechnung bei großen Argumenten ab. Die Ursache liegt darin, daß Einschließungen der Verstärkungsfaktoren, die im Bereich von `MaxReal` liegen, in Operationen eingehen, deren Ergebnis einen Überlauf liefert. Für die Berechnung einer Gesamtfehlerschranke ist hier die Rückwärtsmethode nicht geeignet.

Ohne explizit die Laufzeiten zur Berechnung der Fehlerschranken in den Beispielen aufzulisten, ist zu bemerken, daß sich die Rückwärtsrechnung hinsichtlich des Laufzeitverhaltens schlechter verhält als die Vorwärtsrechnung. Größenordnungsmäßig

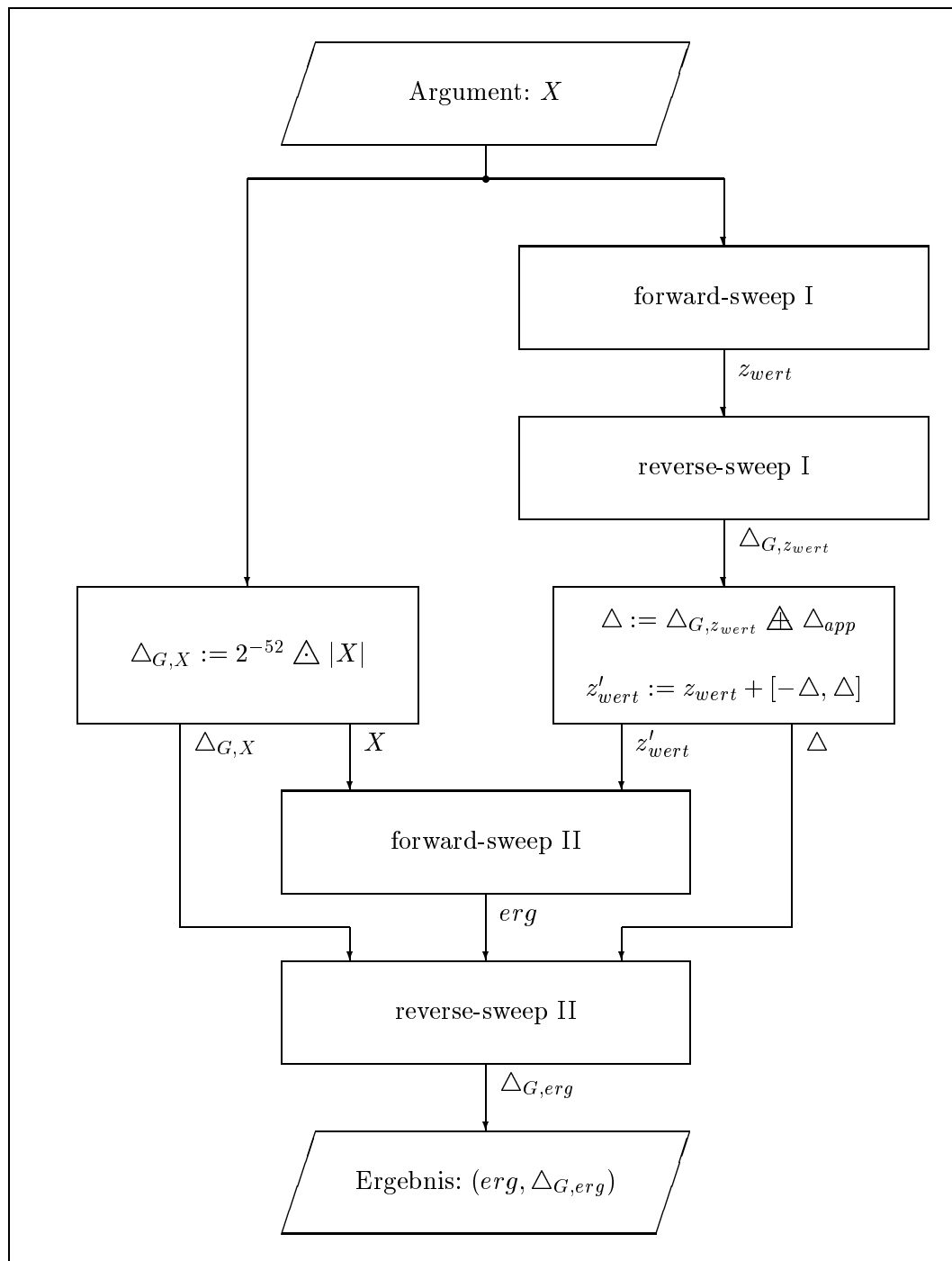


Abbildung 6: Prinzipielle Vorgehensweise bei Einbeziehung des Approximationsfehlers



entspricht die Anzahl der Operationen zur Berechnung der Verstärkungsfaktoren in (45) der Anzahl der Operationen zur Berechnung der absoluten Fehler der einzelnen Teilfunktionen der Codeliste in der Vorwärtsrechnung. Jedoch müssen in der Rückwärtsrechnung ebenfalls alle absoluten Fehler der einzelnen Teilfunktionen der Codeliste berechnet werden, um überprüfen zu können, ob rundungsfehlerfreie Operationen vorliegen. Das erhöht den Aufwand erheblich.

Zwar wurden in der Implementierung [2] noch keine Verzweigungen mit einbezogen, aber man kann auch hier eine Aussage über das Verhalten von Vorwärts- und Rückwärtsrechnung treffen. Treten Verzweigungen im zugrundeliegenden Algorithmus auf, so kann das Maximum der berechneten Fehlerschranken beider Zweige als Fehlerschranke dienen. Erweitert man  $\mathcal{F}$  um die in [16, 17] (siehe auch [6, 28], Kap. 7) vorgestellte Verzweigungsfunktion  $\chi$ , so werden durch (45) die Fehlerschranken beider Zweige addiert. Aufgrund der Ungleichung

$$\max\{x, y\} \leq x + y \quad x, y \in \mathbb{R} \text{ mit } x, y \geq 0$$

ist die mittels Vorwärtsrechnung bestimmte Fehlerschranke der durch Rückwärtsrechnung ermittelten Schranke im allgemeinen überlegen.

Abschließend ist zu bemerken, daß sich die Anwendung der Rückwärtsrechnung im Vergleich zur Vorwärtsrechnung mühsamer gestaltet, das Umschreiben der zu testenden Algorithmen ist *erheblich aufwendiger*. Bei der Vorwärtsrechnung müssen größtenteils nur Datentypen ausgetauscht werden.

## 5 Schlussbemerkung

Es können gleichmäßige Fehlerschranken numerischer Verfahren auf gegebenen Argumentbereichen durch Vorwärts- als auch durch Rückwärtsrechnung automatisch ermittelt werden. Solche Fehlerschranken geben Auskunft über die Qualität des betrachteten Verfahrens hinsichtlich der Verstärkung von Rundungsfehlern. Außerdem ist mit ihnen eine Sensitivitätsanalyse möglich, bei der die Auswirkung von Störungen in den Eingangsparametern auf die Ergebniswerte untersucht werden kann. Auch hier wird die Faktorisierbarkeit der Funktionen vorausgesetzt, für welche die Fehlerschranken ermittelt werden sollen. Durch Anwendung einer speziellen Verzweigungsfunktion ([6], Kap. 7) können aber zusätzlich Bedingungsanweisungen einbezogen werden. Dabei liefert die Vorwärtsmethode im allgemeinen bessere Fehlerschranken als die Rückwärtsrechnung. In der vorliegenden Arbeit wurden Fehlerschranken für Funktionen auf dem Rechner bestimmt, welche aus elementaren Operationen  $\{-, +, \cdot, /\}$  zusammengesetzt sind. Treten in diesen Funktionen Approximationen auf, wie sie zum Beispiel in Verfahren vorkommen, welche Gleitkommanäherungen von Standardfunktionen auf der Maschine bestimmen, so konnte festgestellt werden, daß die Vorwärtsmethode kleinere Fehlerschranken ermittelt als die Rückwärtsrechnung. Ebenfalls zeigt die Vorwärtsmethode in den gerechneten Beispielen ein besseres Laufzeitverhalten, wenn man für alle Operationen aus  $\{+, -, \cdot, /\}$  überprüft, ob diese möglicherweise rundungsfehlerfrei verlaufen.

Da in den XSC-Sprachen zusätzlich das hochgenaue Skalarprodukt als fünfte Grundoperation zur Verfügung steht, sollte in weiterführenden Arbeiten untersucht werden, wie man gleichmäßige Fehlerschranken erhält, wenn die Menge der elementaren Operationen um das Skalarprodukt erweitert wird.

Die vorgestellte Methodik wurde in [6] unter Verwendung der C++ Klassenbibliothek C-XSC [19, 11] in Softwarewerkzeuge umgesetzt. Quellcodeauszüge finden sich im Anhang A 3 in [6].

## 6 Anhang: Fehlerabschätzung der arctan-Funktion

### 6.1 Algorithmus `q_atan` aus Beispiel 4.3

Für die Berechnung der Funktion  $\arctan(\cdot)$  auf der Maschine kann der Algorithmus `q_atan` aus [13] benutzt werden. Hier findet man auch eine ausführliche Erklärung dieses Verfahrens und eine Beschreibung aller verwendeter Konstanten. Der Algorithmus `q_atan` ist in [13] durch folgende vier Hauptschritte beschrieben<sup>10</sup>:

1. Abprüfen von Spezialfällen und deren Behandlung
  - (a)  $x$  ist NaN (not a number)  $\implies$  Fehlermeldung: Invalid Argument
  - (b)  $|x| \leq \mathbf{q\_atnt} = 1.807032 \dots \cdot 10^{-8} \implies \mathit{res} := x$
2. Initialisierung und Argumentreduktion
  - (a) Falls  $x < 0$  :  $y := -x$ ,  $\mathit{neg} := \mathbf{true}$   
sonst:  $y := x$ ,  $\mathit{neg} := \mathbf{false}$
  - (b) Falls  $y < 8$  :  $\mathit{vzi} := 1$ ,  $y_m := 0$   
sonst:  $\mathit{vzi} := -1$ ,  $y_m := \frac{\pi}{2} \mid_{IEEE}$ ,  $y := 1 \boxminus y$
  - (c) Bestimme  $i := \max_{k=1,2,\dots,7} \{k \mid b_k \leq y\}$
  - (d) Berechne  $y := (y \boxminus c_i) \boxminus (1.0 \boxplus y \boxtimes c_i)$
3. Approximation:  $y + y \cdot (\sum_{i=0}^5 d_i \cdot y^2)$ , d.h. auf dem Rechner
  - (a)  $y_{sq} := y \boxtimes y$
  - (b)  $z_{wert} := (((d_5 \boxtimes y_{sq} \boxplus d_4) \boxtimes y_{sq} \boxplus d_3) \boxtimes y_{sq} \boxplus d_2) \boxtimes y_{sq} \boxplus d_1) \boxtimes y_{sq} \boxplus d_0$
  - (c)  $\mathit{res} := y \boxtimes z_{wert} \boxplus y$
4. Ergebnisanpassung
  - (a)  $\mathit{res} := \mathit{res} \boxplus a_i$
  - (b)  $\mathit{res} := \mathit{vzi} \cdot \mathit{res} \boxplus y_m$
  - (c) Falls  $\mathit{neg} = \mathbf{true}$  :  $\mathit{res} := -\mathit{res}$

Nun gilt:  $\mathit{res} \approx \arctan(x)$ .

---

<sup>10</sup>Die Bezeichnungen wurden von dort übernommen.

## 6.2 Ergebnisse der Fehlerabschätzung

Die in Abbildung 5 gezeigten Fehlerschranken ergeben sich aus den folgenden Ergebnissen. Die Größe  $n$  gibt dabei die Anzahl nochmaliger Unterteilungen der Bereiche  $B_0 - B_8$  an. Bei Anwendung der Vorwärtsrechnung erhält man:

B0	[1.807031E-008,1.000000E-007]	n: 40	4.944260E-016
B0	[9.999999E-008,1.000000E-006]	n: 40	5.440093E-016
B0	[9.999999E-007,1.000001E-005]	n: 40	5.440093E-016
B0	[1.000000E-005,1.000001E-004]	n: 40	5.440093E-016
B0	[1.000000E-004,1.000001E-003]	n: 40	5.440093E-016
B0	[1.000000E-003,1.000001E-002]	n: 40	5.440114E-016
B0	[1.000000E-002,1.000001E-001]	n: 40	5.442172E-016
B1	[1.000000E-001,1.250000E-001]	n: 2	5.153864E-016
B2	[1.250000E-001,3.913935E-001]	n: 10000	1.369861E-015
B3	[3.913934E-001,7.165921E-001]	n: 1000	8.212020E-016
B4	[7.165920E-001,1.186492E+000]	n: 10	7.105762E-016
B5	[1.186491E+000,2.061722E+000]	n: 10	6.197640E-016
B6	[2.061721E+000,4.860929E+000]	n: 10	5.223889E-016
B7	[4.860928E+000,8.000000E+000]	n: 1	5.197107E-016
B8	[8.000000E+000,1.000000E+001]	n: 10	4.221772E-016
B8	[1.000000E+i,1.000000E+(i+1)]	n: je 10	
		mit i=1(1)307:	gesamt 5.528541E-016
B8	[1.000000E+308,1.797694E+308]	n: 1000	3.330669E-016

Gesamtfehler: 1.369861E-015

Die Rückwärtsrechnung ergibt bei gleicher Unterteilung:

B0	[1.807031E-008,1.000000E-007]	n: 40	4.944260E-016
B0	[9.999999E-008,1.000000E-006]	n: 40	5.440093E-016
B0	[9.999999E-007,1.000001E-005]	n: 40	5.440093E-016
B0	[1.000000E-005,1.000001E-004]	n: 40	5.440093E-016
B0	[1.000000E-004,1.000001E-003]	n: 40	5.440093E-016
B0	[1.000000E-003,1.000001E-002]	n: 40	5.440107E-016
B0	[1.000000E-002,1.000001E-001]	n: 40	5.441492E-016
B1	[1.000000E-001,1.250000E-001]	n: 2	5.097375E-016
B2	[1.250000E-001,3.913935E-001]	n: 10000	1.555012E-015
B3	[3.913934E-001,7.165921E-001]	n: 1000	8.702787E-016
B4	[7.165920E-001,1.186492E+000]	n: 10	7.300869E-016
B5	[1.186491E+000,2.061722E+000]	n: 10	6.256731E-016
B6	[2.061721E+000,4.860929E+000]	n: 10	5.504527E-016
B7	[4.860928E+000,8.000000E+000]	n: 1	5.175578E-016
B8	[8.000000E+000,1.000000E+001]	n: 10	4.204228E-016
B8	[1.000000E+i,1.000000E+(i+1)]	n: je 10	
		mit i=1(1)307:	gesamt 5.510245E-016
B8	[1.000000E+308,1.797694E+308]	n: 1000	3.330669E-016

Gesamtfehler: 1.555012E-015

und bei Verfeinerung der Teilbereiche  $B_4 - B_6$ :

B4 [7.165920E-001,1.186492E+000] n: 100 7.057976E-016  
 B5 [1.186491E+000,2.061722E+000] n: 100 6.046695E-016  
 B6 [2.061721E+000,4.860929E+000] n: 100 5.255126E-016

## Literatur

- [1] Bantle, A., Krämer, W.: *Ein Kalkül für verlässliche absolute und relative Fehlerabschätzungen* Preprint 98/5 des IWRMM, Universität Karlsruhe, 1998
- [2] Bantle, A., Krämer, W.: *Implementierung, Handhabung und Beispielanwendungen eines verlässlichen Vorwärtsfehlerkalküls* Forschungsschwerpunkt Computerarithmetik, Intervallrechnung und Numerische Algorithmen mit Ergebnisverifikation, Bericht 2/1999, Karlsruhe, 1999
- [3] Beck, T.: *Automatisches Differenzieren von Algorithmen* Dissertation, Technische Universität München, 1991
- [4] Berz, M., Bischof, C., Corliss, G. and Griewank, A. (eds.): *Computational Differentiation: Techniques, Applications, and Tools* SIAM, Philadelphia, Penn., 1996
- [5] Braune, K., Krämer, W.: *High-Accuracy Standard Functions for Intervals* in: [30], pp. 341-347
- [6] Bräuer, M.: *Berechnungsmethoden für Ableitungen und Steigungen und deren Realisierung in C-XSC* Diplomarbeit, Universität Karlsruhe, 1999.
- [7] Fischer, H.: *Special Problems in Automatic Differentiation* in: [9], pp. 43-50
- [8] Fischer, H.-C.: *Schnelle Automatische Differentiation, Einschließungsmethoden und Anwendungen* Dissertation, Institut für Angewandte Mathematik, Karlsruhe, 1990
- [9] Griewank, A., Corliss, G. F. (eds.): *Automatic Differentiation of Algorithms: Theory, Implementation, and Application* SIAM, Philadelphia, Penn., 1991
- [10] Griewank, A., Juedes, D. and Utke, Jean: *ADOL-C: A Package for the Automatic Differentiation of Algorithms Written in C++* ACM Trans. Math. Software, 22 (1996), pp. 131-167
- [11] Hammer, R., Hocks, M., Kulisch, U. und Ratz, D.: *C++ Toolbox for Verified Computing I - Basic Numerical Problems* Springer Verlag, Berlin, 1995
- [12] Hofschuster, W. und Krämer, W.: *Ein rechnergestützter Fehlerkalkül mit Anwendung auf ein genaues Tabellenverfahren* Preprint Nr. 96/5, Institut für Wissenschaftliches Rechnen und Mathematische Modellbildung, Universität Karlsruhe, 1996

- [13] Hofschuster, W. und Krämer, W.: *FLLIB, eine schnelle und portable Funktionsbibliothek für reelle Argumente und reelle Intervalle im IEEE-double-Format* Preprint Nr. 98/7, Institut für Wissenschaftliches Rechnen und Mathematische Modellbildung, Universität Karlsruhe, 1998
- [14] Iri, M.: *Simultaneous Computation of Functions, Partial Derivatives and Estimates of Rounding Errors* Japan Journal of Applied Mathematics, 1 (1994), pp. 223-252
- [15] Iri, M., Tsuchiya, T. and Hoshi, M.: *Automatic computation of partial derivatives and rounding error estimates with application to large-scale systems of nonlinear equations* Journal of Computational and Applied Mathematics 24, 1988, pp. 365-392
- [16] Kearfott, R.: *Interval Extensions of Non-Smooth Functions for Global Optimization and Nonlinear System Solvers* Computing 57 (1996), pp. 149-162
- [17] Kearfott, R.: *Automatic Differentiation of Conditional Branches in an Operator Overloading Context* in: [7], pp. 75-81
- [18] Kedem, G.: *Automatic Differentiation of Computer Programs* ACM Trans. Math. Software, 6 (1980), pp. 150-165
- [19] Klatte, R., Kulisch, U., Lawo, C., Rauch, M., Wiethoff, A.: *C-XSC, A C++ Class Library for Extended Scientific Computing* Springer-Verlag Berlin, Heidelberg, New York 1993
- [20] Krämer, W.: *Inverse Standard Functions for Real and Complex Point and Interval Arguments with Dynamic Accuracy* Computing, Suppl. 6, 1988, pp. 185-212
- [21] Krämer, W.: *Sichere und genaue Abschätzung des Approximationsfehlers bei rationalen Approximationen* Forschungsschwerpunkt Computerarithmetik, Intervallrechnung und Numerische Algorithmen mit Ergebnisverifikation, Bericht 3/1996, Karlsruhe, 1996
- [22] Krämer, W.: *Eine Fehlerfaktorarithmetik für zuverlässige a priori Fehlerabschätzungen* Forschungsschwerpunkt Computerarithmetik, Intervallrechnung und Numerische Algorithmen mit Ergebnisverifikation, Bericht 5/1997, 21 Seiten, Karlsruhe, 1997
- [23] Krämer, W.: *Constructive Error Analysis* Journal of Universal Computer Science (JUCS), Vol. 4, No. 2, 1998, pp. 147-163
- [24] Krämer, W.: *A priori Worst Case Error Bounds for Floating-Point Computations* IEEE Transactions on Computers, Vol. 47, No. 7, July 1998
- [25] Kulisch, U.: *Grundlagen des Numerischen Rechnens* Bibliographisches Institut, Mannheim, Wien, Zürich 1976

- [26] Kubota, K., Iri, M.: *Estimates of Roundings Errors with Fast Automatic Differentiation and Interval Analysis* Technical Report METR-88-15, Dep. of Mathematical Engineering and Information Physics, University of Tokyo, October 1988
- [27] Neumann, K., Morlock M.: *Operations research* Carl Hanser Verlag, München, Wien, 1993
- [28] Ratz, D.: *Automatic Slope Computation and its Application in Nonsmooth Global Optimization* Habilitationsschrift, Institut für Angewandte Mathematik, Karlsruhe, 1998
- [29] Rump, S. M.: *Expansion and Estimation of the Range of Nonlinear Functions* Mathematics of Computation, Vol. 65, No. 216, 1996, pp. 1503-1512
- [30] Ruschitzka, M. (editor): *Computer Systems: Performance and Simulation* Elsevier Science Publishers B.V. (North-Holland), IMACS, 1986
- [31] Sander, P., Stucky, W., Herschel, R.: *Automaten, Sprachen, Berechenbarkeit* B.G.Teubner, Stuttgart, 1991
- [32] Schwarz, H. R.: *Numerische Mathematik* B.G.Teubner, Stuttgart, 1993
- [33] Stroustrup, B.: *Die C++ Programmiersprache* Addison-Wesley-Longman, Bonn, 1998