

A. NEUMAIER

**Rundungsfehleranalyse einiger Verfahren zur Summation endlicher Summen**

Der bei einer Summation auftretende Rundungsfehler kann als Maß für die Güte des verwendeten Verfahrens gelten. Im folgenden werden für mehrere Summierungsverfahren, unter anderem für das übliche und das Kahan-Babuška-Verfahren, a-priori-Schranken für diese Rundungsfehler angegeben und miteinander verglichen.

The rounding-error arising during summation can be interpreted as a measure for the quality of the procedure used. In the following, a-priori-bounds for this rounding-error are used to compare several summation procedures, e.g. the common procedure and the method of Kahan-Babuška.

Возникающие при суммировании ошибки округления могут рассматриваться как мера качества для применяемого метода. В последующем приведены и сравнены между собой для многих методов суммирования, в том числе для "обычного" метода и для метода Кahan-Бабушка, априорные пределы для этих ошибок округления.

**1. Aufgabenstellung**

Zur Berechnung einer endlichen Summe

$$s := \sum_{m=1}^n a_m \quad (1)$$

können verschiedene Verfahren benutzt werden. Bei Durchführung der Rechnung im Körper der reellen Zahlen liefert jedes Verfahren das Ergebnis  $s$ ; wird die Rechnung aber auf einer Rechenmaschine durchgeführt, so ergeben sich im allgemeinen verschiedene Ergebnisse, da die auftretenden Rundungsfehler vom Verfahren abhängen.

Die Güte einer Summationsvorschrift kann z. B. durch das Verhältnis von Genauigkeit des berechneten Wertes zum Arbeitsaufwand angegeben werden. Gesucht sind daher Schranken für die Abweichung  $\tilde{s} - s$  des nach einem bestimmten Verfahren berechneten Näherungswertes  $\tilde{s}$  für  $s$ .

**2. Summationsvorschriften**

Es werden die folgenden vier Summationsvorschriften untersucht<sup>1)</sup>:

**I. Normalverfahren**

Man setzt

$$\begin{aligned} s_0 &:= 0 \\ s_m &:= s_{m-1} + a_m \quad \text{für } m = 1, \dots, n \\ s &:= s_n, \end{aligned} \quad (I,1)$$

d. h., die Summe (1) wird von links nach rechts abgearbeitet, z. B.

$$s = ((a_1 + a_2) + a_3) + a_4 \quad (n = 4)$$

Der Rechenaufwand beträgt  $n - 1$  Additionen (Addition von 0 wird vernachlässigt).

**II. Verbessertes Normalverfahren**

$2^t$  sei die kleinste Potenz von 2 mit  $2^t \geq n$ . Setzt man

$$\begin{aligned} s_{m0} &:= a_m \quad \text{für } m = 1, \dots, n \\ s_{m0} &:= 0 \quad \text{für } m = n + 1, \dots, 2^t, \end{aligned}$$

so berechnet auch die Vorschrift

$$\begin{aligned} s_{mk} &:= s_{2m-1, k-1} + s_{2m, k-1} \quad \text{für } 1 \leq k \leq t, \quad 0 \leq m \leq 2^{t-k} - 1 \\ s &:= s_{1t} \end{aligned} \quad (II,1)$$

die Summe  $s$ , aber in anderer Klammerung, z. B.

$$s = ((a_1 + a_2) + (a_3 + a_4)) + ((a_5 + a_6) + (a_7 + a_8)) \quad (t = 3, n = 8).$$

Wieder sind  $n - 1$  Additionen auszuführen.

<sup>1)</sup> BABUŠKA [1] behandelte die Verfahren I, II, III unter den Bezeichnungen  $L_1 A$ ,  $L_1 B$ ,  $L_1 C$ .

## III. KAHAN-BABUŠKA-Verfahren

Man setzt [1], [2], [3]

$$\begin{aligned} s_0 &:= 0, & s_m &:= a_m + s_{m-1} & (m = 1, \dots, n) \\ w_0 &:= 0, & w_m &:= w_{m-1} + (a_m + (s_{m-1} - s_m)) & (m = 1, \dots, n) \\ s &:= s_n + w_n. \end{aligned} \quad (\text{III,1})$$

Die  $w_m$  sind bei reeller Rechnung Null; wird auf einer Maschine gerechnet, so sammeln sich die Rundungsfehler in diesem Term an. Durch die Korrektur  $+ w_n$  werden die Rundungsfehler etwa kompensiert. Es ist also eine größere Genauigkeit zu erwarten. Der Rechenaufwand beträgt hier  $4n - 3$  Additionen/Subtraktionen, ist also etwa viermal so groß wie in Verfahren I und II.

## IV. Verbessertes KAHAN-BABUŠKA-Verfahren

Man setzt

$$\begin{aligned} s_0 &:= 0, & s_m &:= a_m + s_{m-1} & (m = 1, \dots, n) \\ w_0 &:= 0, & w_m &:= w_{m-1} + (a_m + (s_{m-1} - s_m)), & \text{wenn } |a_m| \leq |s_{m-1}| \\ & & w_m &:= w_{m-1} + (s_{m-1} - (a_m - s_m)), & \text{wenn } |a_m| > |s_{m-1}| \\ s &:= s_n + w_n. \end{aligned} \quad (\text{IV,1})$$

Das Verfahren unterscheidet sich von Verfahren III nur für  $|a_m| > |s_{m-1}|$ . (Z. B. ist  $\text{IV} = \text{III}$ , wenn  $(a_m)$  eine monoton fallende Folge positiver reeller Zahlen ist, da  $|a_m| > |s_{m-1}|$  genau für  $m = 1$ .) Hier sind außer  $4n - 3$  Additionen/Subtraktionen noch  $n$  Abfragen auszuführen. Rechnet man für eine Abfrage denselben Aufwand wie für eine Subtraktion, so ist der Aufwand für Verfahren IV etwa fünfmal so groß wie für Verfahren I oder II, und 1/4mal größer als in Verfahren III.

## 3. Rundungsfehler beim Normalverfahren (I und II)

Nach WILKINSON [4] genügt die Maschinennäherung  $a \tilde{+} b$  einer Summe  $a + b$  für viele Rechenmaschinen der Formel

$$a \tilde{+} b = (a + b)(1 + \delta), \quad |\delta| \leq \varepsilon. \quad (2)$$

Wenn die Maschine mit  $L$  Gleitkomma-Dualziffern rechnet, ist speziell  $\varepsilon = (1 + c) \cdot 2^{-L}$ ; dabei ist  $c$  eine nichtnegative Maschinenkonstante (bei optimaler Rundung ist  $c = 0$ ).

Auf Grund dieser Formel erhält WILKINSON für Verfahren I folgendes Ergebnis ( $\tilde{s}$  ist der berechnete Wert für  $s$ ):

$$\tilde{s} = \sum_{m=1}^n a_m (1 + \eta_m) = s + \sum_{m=1}^n a_m \eta_m \quad (\text{I,2})$$

$$1 + \eta_m = \prod_{i=m}^n (1 + \varepsilon_i), \quad |\varepsilon_i| \leq \varepsilon \quad \text{für } i > 1, \quad \varepsilon_1 = 0, \quad (\text{I,3})$$

wobei die Größen  $\eta_m$  durch

$$\begin{aligned} |\eta_{n-t+1}| &\leq (1 + \varepsilon)^t - 1 & (1 \leq t \leq n - 1) \\ |\eta_1| &\leq (1 + \varepsilon)^{n-1} - 1 \end{aligned} \quad (\text{I,4})$$

beschränkt sind; diese Schranken werden für  $\varepsilon_i = \varepsilon$  (für  $i > 1$ ) angenommen.

Für Verfahren II gilt: ist  $\tilde{s}_{mk}$  der berechnete Wert für  $s_{mk}$ , so ist nach (2) und (II,1):

$$\begin{aligned} \tilde{s}_{mk} &= (\tilde{s}_{2m-1, k-1} + \tilde{s}_{2m, k-1})(1 + \varepsilon_{mk}) & \text{für } 1 \leq k \leq t, \\ 0 \leq m &\leq 2^{t-k} - 1, & |\varepsilon_{mk}| \leq \varepsilon. \end{aligned}$$

Durch vollständige Induktion erhält man daraus leicht das Ergebnis:

$$\tilde{s}_{mk} = \sum_{i=2^k(m-1)+1}^{2^k m} s_{i0} (1 + \delta_i^{(k)}), \quad 1 + \delta_i^{(k)} = \prod_{L=1}^k \left(1 + \varepsilon_{1 + \left[\frac{i-1}{2^L}\right], L}\right).$$

Für den berechneten Wert  $\tilde{s} = \tilde{s}_{1t}$  der Summe  $s$  gilt also (unter Berücksichtigung von  $s_{i0} = 0$  für  $i > n$ ):

$$\tilde{s} = \sum_{m=1}^n a_m (1 + \eta_m) = s + \sum_{m=1}^n a_m \eta_m \quad (\text{II,2})$$

$$1 + \eta_m = \prod_{i=1}^t \left(1 + \varepsilon_{1 + \left[\frac{m-1}{2^i}\right], i}\right), \quad |\varepsilon_{ik}| \leq \varepsilon, \quad (\text{II,3})$$

woraus man sofort die Schranke

$$|\eta_i| \leq (1 + \varepsilon)^t - 1 \quad (1 \leq i \leq n) \quad (\text{II,4})$$

erhält, die für  $\varepsilon_{ik} = \varepsilon$  angenommen wird.

## 4. Rundungsfehler beim Kahan-Babuška-Verfahren (III)

Die berechneten Näherungswerte für  $s_m, w_m, s$  seien bzw.  $\tilde{s}_m, \tilde{w}_m, \tilde{s}$ . Aus (III,1) folgt jetzt mit (2):

$$\begin{aligned}\tilde{s}_0 &= \tilde{w}_0 = 0 \\ \tilde{s}_m &= (\tilde{s}_{m-1} + a_m)(1 + \varepsilon_m) & m = 1, \dots, n \\ \tilde{w}_m &= (1 + \alpha_m)(\tilde{w}_{m-1} + (1 + \beta_m)(a_m + (1 + \gamma_m)(\tilde{s}_{m-1} - \tilde{s}_m))) & m = 1, \dots, n\end{aligned}\quad (\text{III,2})$$

mit  $\alpha_m, \beta_m, \gamma_m, \varepsilon_m$ , die betragsmäßig  $\leq \varepsilon$  sind. Da man annehmen kann, daß  $z + 0 = z$  exakt berechnet wird, kann man voraussetzen:

$$\alpha_1 = \beta_1 = \gamma_1 = \varepsilon_1 = 0. \quad (\text{III,2a})$$

Nun ist  $q_m = \sum_{r=1}^m \left( \psi_r \prod_{i=r}^m \varphi_i \right)$  die Lösung der Differenzgleichung  $q_m = \varphi_m (q_{m-1} + \psi_m)$ , mit der Anfangsbedingung  $q_0 = 0$ , wie man durch Einsetzen feststellt.

Mit den Abkürzungen

$$\prod_{i=r}^s (1 + \varepsilon_i) = b_{rs}, \quad (1 + \beta_r) \prod_{i=r}^n (1 + \alpha_i) = c_r \quad (\text{III,3})$$

ist daher

$$\tilde{s}_m = \sum_{r=1}^m a_r b_{rm} \quad (\text{III,4})$$

$$\tilde{w}_n = \sum_{r=1}^n c_1 (a_r + (1 + \gamma_r)(\tilde{s}_{r-1} - \tilde{s}_r)) \quad (\text{III,5})$$

oder, unter Verwendung von (III,4):

$$\tilde{w}_n = \sum_{r=1}^n a_r \left\{ c_r (1 - (1 + \gamma_r)(1 + \varepsilon_r)) + \sum_{s=r+1}^n c_s (1 + \gamma_s)(b_{r,s-1} - b_{rs}) \right\}. \quad (\text{III,6})$$

Also ist

$$\tilde{s}_n + \tilde{w}_n = \sum_{r=1}^n a_r (1 + \delta_r) \quad (\text{III,7})$$

mit

$$\delta_r = 1 + b_{rn} - c_r \gamma_r (1 + \varepsilon_r) - c_r \varepsilon_r + \sum_{s=r+1}^n c_s (b_{r,s-1} - b_{rs}) (1 + \gamma_s). \quad (\text{III,8})$$

Durch Umformen erhält man

$$\delta_r = -\varepsilon_r (c_r - 1) - c_r \gamma_r (1 + \varepsilon_r) + \sum_{s=r+1}^n (b_{r,s-1} - b_{rs}) (c_s (\gamma_s + 1) - 1)$$

(wegen  $b_{rn} = b_{rr} - \sum_{s=r+1}^n (b_{r,s-1} - b_{rs})$  und  $b_{rr} = 1 + \varepsilon_r$ ) und weiter

$$-\delta_r = \varepsilon_r (c_r - 1) + c_r \gamma_r (1 + \varepsilon_r) + \sum_{s=r+1}^n \varepsilon_s b_{r,s-1} (c_s \gamma_s + c_s - 1)$$

(wegen  $b_{rs} = (1 + \varepsilon_s) b_{r,s-1}$ ), also

$$-\delta_r = \sum_{s=r}^n \varepsilon_s b_{r,s-1} (c_s \gamma_s + c_s - 1) + \gamma_r c_r. \quad (\text{III,9})$$

Aus (III,3) folgt für  $r > 1$ :

$$\begin{aligned}|b_{rs}| &\leq (1 + \varepsilon)^{s-r+1}, & |c_s| &\leq (1 + \varepsilon)^{n-s+2}, \\ |c_s - 1| &\leq (1 + \varepsilon)^{n-s+2} - 1\end{aligned}\quad (\text{III,10a})$$

und für  $r = 1$  wegen (III,2a):

$$|b_{1s}| \leq (1 + \varepsilon)^{s-1} \quad (\text{III,10b})$$

Man erhält damit aus (III,9) die Abschätzung

$$\begin{aligned}|\delta_r| &\leq (1 + \varepsilon)^{n-r+2} \varepsilon + \sum_{s=r}^n \varepsilon (1 + \varepsilon)^{s-r} (\varepsilon (1 + \varepsilon)^{n-s+2} + (1 + \varepsilon)^{n-s+2} - 1) \\ &= \varepsilon (1 + \varepsilon)^{n-r+2} + (n - r + 1) \varepsilon (1 + \varepsilon)^{n-r+3} - (1 + \varepsilon)^{n-r+1} + 1\end{aligned}$$

für  $r > 1$ , also

$$|\delta_{n-t+1}| \leq (\varepsilon(1+\varepsilon)^2 t + \varepsilon(1+\varepsilon) - 1)(1+\varepsilon)^t + 1 \quad (1 \leq t \leq n-1) \quad (\text{III,11a})$$

und

$$\begin{aligned} |\delta_1| &\leq \sum_{s=2}^n \varepsilon(1+\varepsilon)^{s-2} (\varepsilon(1+\varepsilon)^{n-s+2} + (1+\varepsilon)^{n-s+2} - 1) \\ &= (n-1)\varepsilon(1+\varepsilon)^{n+1} - (1+\varepsilon)^{n-1} + 1, \end{aligned}$$

also

$$|\delta_1| \leq (\varepsilon(1+\varepsilon)^2(n-1) - 1)(1+\varepsilon)^{n-1} + 1. \quad (\text{III,11b})$$

Die Schranken in (III,11) werden für  $\alpha_i = \beta_i = \gamma_i = \varepsilon_i = \varepsilon$  ( $i > 1$ ) erreicht, können also i. a. nicht mehr verbessert werden.

Nun ist

$$\tilde{s} = (\tilde{s}_n + \tilde{w}_n)(1+\delta) = \left( s + \sum_{i=1}^n \alpha_i \delta_i \right) (1+\delta),$$

also

$$\tilde{s} = s + \delta s + \sum_{i=1}^n a_i \eta_i \quad (\text{III,12})$$

mit

$$\eta_i = -(\delta+1) \left( \gamma_r c_r + \sum_{s=r}^n \varepsilon_s b_{r,s-1} (c_s(\gamma_s+1) - 1) \right) \quad (\text{III,13})$$

$$|\eta_{n-t+1}| \leq (\varepsilon(1+\varepsilon)^2 t + \varepsilon(1+\varepsilon) - 1)(1+\varepsilon)^{t+1} + 1 + \varepsilon \quad (1 \leq t \leq n-1)$$

$$|\eta_1| \leq (\varepsilon(1+\varepsilon)^2(n-1) - 1)(1+\varepsilon)^n + 1 + \varepsilon \quad (\text{III,14})$$

$$|\delta| \leq \varepsilon.$$

### 5. Forderungen an eine Rechenmaschine. Ein Satz über Rundungsfehler

Zunächst werden einige Forderungen an die zu verwendende Rechenmaschine gestellt, die in der Praxis mehr oder weniger genau erfüllt sind. Zur Vorbereitung der Abschätzung für die Rundungsfehler beim Verfahren IV wird dann ein Satz über die Berechnung von  $B - (A + B)$  für  $|A| \leq |B|$  in einer solchen Rechenmaschine bewiesen.

Über die Struktur der Maschinenzahlen wird folgende Annahme gemacht:

Voraussetzung 1: Maschinenzahlen sind genau die Zahl 0 und die Zahlen  $x = \pm q \cdot Z^{-\alpha}$  mit  $Z^{L-1} \leq q < Z^L$ ,  $q \in \mathbb{N}$ ,  $\alpha \in \mathcal{G}$  \*); dabei sind  $Z > 1$  und  $L$  positive, ganzzahlige Maschinenkonstanten. In der Praxis können natürlich nur endlich viele  $\alpha$  zugelassen werden. Alles Folgende bleibt jedoch richtig, wenn alle auftretenden Summen oder Differenzen eines dieser  $\alpha$  als Exponent besitzen.

Hilfssatz 1: Gilt Vor. 1, sind  $A$  und  $B$  Maschinenzahlen, gilt  $A \geq B \geq 0$  und ist

$$B = q Z^{-\alpha} \quad (q \in \mathbb{N}, \quad Z^{L-1} \leq q < Z^L, \quad \alpha \in \mathcal{G}),$$

so ist  $A = r Z^{-\alpha}$  mit einem  $r \in \mathbb{N}$ ,  $r \geq q$ .

Beweis: Es ist  $A \geq B \geq Z^{L-\alpha-1}$ . Da  $A$  eine Maschinenzahl ist und da  $A > 0$ , ist  $A = r_0 Z^{-\beta}$  mit  $r_0 \in \mathbb{N}$ ,  $Z^{L-1} \leq r_0 < Z^L$ . Also ist  $Z^{L-\beta} > r_0 Z^{-\beta} = A \geq Z^{L-\alpha-1}$  oder  $L-\beta > L-\alpha-1$ . Das bedeutet aber, daß  $-\beta = \gamma - \alpha$  mit  $\gamma > 0$ . Mit  $r = r_0 Z^\gamma$  ( $r \in \mathbb{N}$ ) ist dann  $A = r Z^{-\alpha}$ . Aus  $A \geq B$  folgt schließlich  $r \geq q$ , w. z. z. w.

Hilfssatz 2: Gilt Vor. 1, ist  $A$  Maschinenzahl,  $Z^{L-\alpha} > A \geq Z^{L-1-\alpha}$ ,  $\alpha \in \mathcal{G}$ , so ist

$$A = r Z^{-\alpha}, \quad r \in \mathbb{N}, \quad Z^{L-1} \leq r < Z^L.$$

Beweis: Wegen  $Z^{L-1-\alpha} = q Z^{-\alpha}$  mit  $q = Z^{L-1}$  folgt aus Hilfssatz 1:  $A = r Z^{-\alpha}$  mit  $r \in \mathbb{N}$ ,  $Z^{L-1} \leq r$ . Wegen  $A < Z^{L-\alpha}$  gilt dann  $r < Z^L$ , w. z. z. w.

Für die Maschinenoperationen  $\tilde{+}$  (Näherung für  $+$ ) und  $\tilde{-}$  (Näherung für  $-$ ) seien folgende Voraussetzungen erfüllt:

Voraussetzung 2: Sind  $A$  und  $B$  Maschinenzahlen, gilt

$$(V 0) \quad A \geq B \geq 0$$

und existiert ein  $\alpha \in \mathcal{G}$  mit

$$(V 1) \quad A Z^\alpha \in \mathcal{G}, \quad B Z^\alpha \in \mathcal{G}$$

\*)  $\mathcal{G}$ : Menge der ganzen Zahlen.

$$(V 2) \quad B < Z^{L-\alpha}$$

$$(V 3) \quad A - B \leq Z^{L-\alpha}$$

so ist  $A \overset{\sim}{-} B = A - B$ .

Voraussetzung 3: Sind  $A$  und  $B$  Maschinenzahlen, gilt

$$A \geq B \geq 0, \quad B \leq Z^{L-\alpha},$$

so ist

$$(U 1) \quad A \leq A \overset{\sim}{+} B \leq A + Z^{L-\alpha}$$

$$(U 2) \quad A - Z^{L-\alpha} \leq A \overset{\sim}{-} B \leq A$$

$$(U 3) \quad 0 \leq A \overset{\sim}{-} B.$$

Bemerkung: Bei optimaler Rundung<sup>2)</sup> sind unter Annahme von Vor. 1 die Voraussetzungen 2 und 3 erfüllt.

Satz 1: Es mögen Vor. 1–3 erfüllt sein.  $X, Y$  seien Maschinenzahlen mit  $X \geq Y > 0$ . Dann gilt

$$a) \quad (X \overset{\sim}{+} Y) \overset{\sim}{-} X = (X \overset{\sim}{+} Y) - X$$

$$b) \quad X \overset{\sim}{-} (X \overset{\sim}{-} Y) = X - (X \overset{\sim}{-} Y).$$

Beweis: Nach Voraussetzung 1 ist

$$X = q Z^{-\gamma}, \quad Z^{L-1} \leq q < Z^L, \quad (q \in \mathbb{N}, \gamma \in \mathbb{G}),$$

also gilt (3)

$$Y < Z^{L-\gamma}, \quad Z^{L-\gamma-1} \leq X < Z^{L-\gamma}.$$

(4)

a) Setze  $A := X \overset{\sim}{+} Y$ ,  $B := X$ ,  $\alpha := \gamma$ .

Nach (U 1) ist  $A \geq X$ . Daher ist nach Hilfssatz 1:

$$A = r Z^{-\alpha}, \quad r \geq q, \quad (r \in \mathbb{N})$$

$A \geq X = B \geq 0$  liefert (V 0),  $A Z^\alpha = r \in \mathbb{G}$ ,  $B Z^\alpha = q \in \mathbb{G}$  liefert (V 1),  $B = X < Z^{L-\alpha}$  liefert (V 2).

Wegen  $Y < Z^{L-\alpha}$  ist nach Vor. 3 (U 1):

$$A - B = (X \overset{\sim}{+} Y) - X \leq (X + Z^{L-\alpha}) - X = Z^{L-\alpha}.$$

Also gilt auch (V 3) und damit ist nach Vor. 2:  $(X \overset{\sim}{+} Y) \overset{\sim}{-} X = (X \overset{\sim}{+} Y) - X$ .

b) Setze  $A := X$ ,  $B := X \overset{\sim}{-} Y$ .

Nach (U 2), (U 3) gilt

$$0 \leq B \leq A = X < Z^{L-\gamma}$$

(5)

(V 0) ist daher erfüllt.

Es werden folgende Fälle unterschieden:

$$\text{Fall } \alpha: Z^{L-1-\gamma} \leq Y < Z^{L-\gamma}$$

$$\text{Fall } \beta: Z^{L-1-\gamma} \leq B < Z^{L-\gamma}, \quad Y < Z^{L-1-\gamma}$$

$$\text{Fall } \gamma: Z^{L-2-\gamma} \leq B < Z^{L-1-\gamma}, \quad Y < Z^{L-1-\gamma}$$

$$\text{Fall } \delta: Z^{L-2-\gamma} \leq Y < Z^{L-1-\gamma}, \quad B < Z^{L-2-\gamma}.$$

Die Möglichkeit  $Y < Z^{L-2-\gamma}$ ,  $B < Z^{L-2-\gamma}$  kann nicht eintreten, da daraus (im Widerspruch zu (4)) nach (U 2):  $X - Z^{L-2-\gamma} \leq B < Z^{L-2-\gamma}$ , also  $X < 2 Z^{L-2-\gamma} \leq Z^{L-1-\gamma}$  folgen würde.

Fall  $\alpha$ : Setze  $\alpha := \gamma$ . Dann ist  $Z^{L-1-\gamma} \leq Y < Z^{L-\alpha}$ , nach Hilfssatz 2 also

$$Y = r Z^{-\alpha}, \quad Z^{L-1} \leq r < Z^L, \quad (r \in \mathbb{N}).$$

(6)

Nach Voraussetzung ist  $X \geq Y > 0$ . Weiter ist  $X Z^\alpha = q \in \mathbb{G}$ ,  $Y Z^\alpha = r \in \mathbb{N}$ ,  $Y < Z^{L-\alpha}$ ,  $X - Y \leq X < Z^{L-\alpha}$ . Vor. 2 läßt sich daher auf  $X$  und  $Y$  anwenden und liefert

$$B = X \overset{\sim}{-} Y = X - Y = (q - r) Z^{-\alpha}.$$

Nun folgt (V 1) aus  $A Z^\alpha = q \in \mathbb{G}$ ,  $B Z^\alpha = q - r \in \mathbb{G}$ , (V 2) aus  $B \leq A = X < Z^{L-\alpha}$ , (V 3) aus

$$A - B = q Z^{-\alpha} - (q - r) Z^{-\alpha} = r Z^{-\alpha} \leq Z^{L-\alpha} \quad \text{nach (6)}.$$

Fall  $\beta$ : Setze  $\alpha := \gamma$ . Dann ist  $Z^{L-1-\alpha} \leq B < Z^{L-\alpha}$ , nach Hilfssatz 2 also

$$B = r Z^{-\alpha}, \quad Z^{L-1} \leq r < Z^L, \quad (r \in \mathbb{N}).$$

<sup>2)</sup>  $x \overset{\sim}{\pm} y =$  nächste bei  $x \pm y$  gelegene Maschinenzahl; liegt  $x \pm y$  genau in der Mitte zwischen den beiden nächsten Maschinenzahlen, so ist entweder stets  $x \overset{\sim}{\pm} y =$  die betragsmäßig größere der beiden Maschinenzahlen oder stets  $x \overset{\sim}{\pm} y =$  die betragsmäßig kleinere der beiden Maschinenzahlen.

(V 1) folgt aus  $A Z^\alpha = q \in \mathcal{G}$ ,  $B Z^\alpha = r \in \mathcal{G}$ , (V 2) aus  $B \leq A = X < Z^{L-\alpha}$ , (V 3) aus  $A - B \leq A = X < Z^{L-\alpha}$ .

Fall  $\gamma$ : Setze  $\alpha := \gamma + 1$ . Dann ist  $Z^{L-1-\alpha} \leq B < Z^{L-\alpha}$ , nach Hilfssatz 2 also

$$B = r Z^{-\alpha}, \quad Z^{L-1} \leq r < Z^L, \quad (r \in \mathbb{N}). \quad (7)$$

(V 1) folgt aus  $A Z^\alpha = q Z \in \mathcal{G}$ ,  $B Z^\alpha = r \in \mathcal{G}$ , (V 2) aus (7), (V 3) aus

$$A - B = X - (X \sim Y) \leq X - (X - Z^{L-\alpha}) = Z^{L-\alpha}$$

wegen (U 2) und  $Y \leq Z^{L-\alpha}$ .

Fall  $\delta$ : Setze  $\alpha := \gamma + 1$ . Dann ist  $Z^{L-1-\alpha} \leq Y \leq Z^{L-\alpha}$ , nach Hilfssatz 2 also

$$Y = r Z^{-\alpha}, \quad Z^{L-1} \leq r \leq Z^L, \quad (r \in \mathbb{N}) \quad (8)$$

Nach (U 2) gilt:  $X - Z^{L-\alpha} \leq B < Z^{L-\alpha-1}$ , also

$$X - Y \leq Z^{L-\alpha} + Z^{L-\alpha-1} - Z^{L-\alpha-1} = Z^{L-\alpha}.$$

Wegen  $X Z^\alpha = q Z \in \mathcal{G}$ ,  $Y Z^\alpha = r \in \mathcal{G}$ ,  $Y < Z^{L-\alpha}$  läßt sich Vor. 2 anwenden und man erhält:

$$B = X \sim Y = X - Y = (Zq - r) Z^{-\alpha}. \quad (9)$$

(V 1) folgt jetzt aus  $A Z^\alpha = Zq \in \mathcal{G}$ ,  $B Z^\alpha = Zq - r \in \mathcal{G}$ , (V 2) aus (8), (9), (V 3) aus  $A - B = X - (X - Y) = Y < Z^{L-\alpha}$ . Also sind in allen vier Fällen die Bedingungen (V 1), (V 2), (V 3) erfüllt. Wegen (5) ist daher Vor. 2 anwendbar und liefert  $X \sim (X \sim Y) = X - (X \sim Y)$ , w. z. z. w.

Um aus diesem Satz Aussagen über  $B \sim (A \tilde{+} B)$  zu gewinnen, werden folgende weitere Voraussetzungen getroffen (die bei optimaler Rundung wieder erfüllt sind):

Voraussetzung 4: Für alle Maschinenzahlen  $A$  gilt:

- (a)  $A \tilde{+} 0 = 0 \tilde{+} A = A$
- (b)  $A \sim A = 0$
- (c)  $0 \sim A = -A$ .

Voraussetzung 5: Für alle positiven Maschinenzahlen  $X, Y$  gilt:

- (d)  $Y \tilde{+} X = X \tilde{+} Y$
- (e)  $Y \sim X = -(X \sim Y)$
- (f)  $(-X) \tilde{+} (-Y) = -(Y \tilde{+} X)$
- (g)  $(-X) \sim (-Y) = -(X \sim Y)$
- (h)  $X \tilde{+} (-Y) = X \sim Y$ .

Satz 2: Unter Vor. 1–5 gilt:

Sind  $A, B$  Maschinenzahlen mit  $|A| \leq |B|$ , so ist

$$B \sim (A \tilde{+} B) = B - (A \tilde{+} B).$$

Beweis: Ist  $A = 0$ , so ist  $B \sim (A \tilde{+} B) = B \sim B = 0 = B - B = B - (A \tilde{+} B)$  nach Vor. 4.

Ist  $B = 0$ , so ist  $B \sim (A \tilde{+} B) = B \sim A = 0 \sim A = -A = 0 - A = B - A = B - (A \tilde{+} B)$  nach Vor. 4.

Ist  $A > 0, B > 0$ , so ist mit  $A = Y, B = X$ :

$$\begin{aligned} B \sim (A \tilde{+} B) &\stackrel{(e)}{=} -((Y \tilde{+} X) - X) \stackrel{(d)}{=} -((X \tilde{+} Y) \sim X) \stackrel{\text{Satz 1}}{=} -((X \tilde{+} Y) - X) \\ &= X - (X \tilde{+} Y) \stackrel{(d)}{=} B - (A \tilde{+} B). \end{aligned}$$

Ist  $A > 0, B < 0$ , so ist mit  $A = Y, B = -X$ :

$$\begin{aligned} B \sim (A \tilde{+} B) &\stackrel{(h)}{=} (-X) \sim (Y \sim X) \stackrel{(e)}{=} (-X) \sim -(X \sim Y) \stackrel{(g)}{=} -(X \sim (X \sim Y)) \\ &\stackrel{\text{Satz 1}}{=} -(X - (X \sim Y)) \stackrel{(e)}{=} (-X) - (Y \sim X) \stackrel{(h)}{=} (-X) - (Y \tilde{+} (-X)) = B - (A \tilde{+} B). \end{aligned}$$

Ist  $A < 0, B > 0$ , so ist mit  $A = -Y, B = X$ :

$$B \sim (A \tilde{+} B) \stackrel{(d)}{=} X \sim (X \tilde{+} (-Y)) \stackrel{(h)}{=} X \sim (X \sim Y) \stackrel{\text{Satz 1}}{=} X - (X \sim Y) \stackrel{(h)}{=} X - (X \tilde{+} (-Y)) \stackrel{(d)}{=} B - (A \tilde{+} B).$$

Ist schließlich  $A < 0, B < 0$ , so ist mit  $A = -Y, B = -X$ :

$$\begin{aligned} B \sim (A \tilde{+} B) &\stackrel{(f)}{=} (-X) \sim -(X \tilde{+} Y) \stackrel{(g)}{=} -(X \sim (X \tilde{+} Y)) \stackrel{(e)}{=} (X \tilde{+} Y) \sim X \stackrel{\text{Satz 1}}{=} (X \tilde{+} Y) - X \\ &= -X - (-X \tilde{+} Y) \stackrel{(f)}{=} B - (A \tilde{+} B). \end{aligned}$$

Damit ist Satz 2 bewiesen.

### 6. Das verbesserte Kahan-Babuška-Verfahren (IV)

Um einfache Abschätzungen zu gewinnen, sei angenommen, daß die verwendete Maschine die Voraussetzungen 1 bis 5 erfüllt (also z. B. optimale Rundung besitzt). Für andere Maschinen kann man die erhaltene Abschätzung als näherungsweise gültig ansehen. Die Bezeichnung der in der Maschine berechneten Werte sei wie in Nummer 4. Dann ist für  $|a_m| \leq |\tilde{s}_{m-1}|$  unter Verwendung von (2):

$$\begin{aligned} a_m \tilde{+} (\tilde{s}_{m-1} \tilde{-} \tilde{s}_m) &= a_m \tilde{+} (\tilde{s}_{m-1} \tilde{-} (a_m \tilde{+} \tilde{s}_{m-1})) \stackrel{\text{Satz 2}}{=} a_m \tilde{+} (\tilde{s}_{m-1} - (a_m \tilde{+} \tilde{s}_{m-1})) \\ &= a_m \tilde{+} (\tilde{s}_{m-1} - \tilde{s}_m) = (1 + \beta_m) (a_m + \tilde{s}_{m-1} - \tilde{s}_m) \end{aligned}$$

und für  $|a_m| > |\tilde{s}_{m-1}|$ :

$$\begin{aligned} \tilde{s}_{m-1} \tilde{+} (a_m \tilde{-} \tilde{s}_m) &= \tilde{s}_{m-1} \tilde{+} (a_m \tilde{-} (\tilde{s}_{m-1} \tilde{+} a_m)) \stackrel{\text{Satz 2}}{=} \tilde{s}_{m-1} \tilde{+} (a_m - (\tilde{s}_{m-1} \tilde{+} a_m)) \\ &= \tilde{s}_{m-1} \tilde{+} (a_m - \tilde{s}_m) = (1 + \beta_m) (\tilde{s}_{m-1} + a_m - \tilde{s}_m). \end{aligned}$$

Daher ist in beiden Fällen:

$$\tilde{s}_m = (1 + \varepsilon_m) (\tilde{s}_{m-1} + a_m) \tag{IV,2}$$

$$\tilde{w}_m = (1 + \alpha_m) (\tilde{w}_{m-1} + (1 + \beta_m) (a_m + \tilde{s}_{m-1} - \tilde{s}_m)). \tag{IV,3}$$

Das sind aber gerade die Formeln (III,2) für  $\gamma_m = 0$  ( $m \geq 1$ ). (Man kann also eine Verbesserung gegenüber Verfahren III erwarten. Dies ist der Grund für die Einführung dieses Verfahrens.) Daher folgen aus (IV,2, 3) die folgenden Formeln ((III,7, 9, 10) für  $\gamma_m = 0$ ):

$$s' := \tilde{s}_n + \tilde{w}_n = \sum_{r=1}^n a_r (1 + \delta_r) \tag{IV,4}$$

$$\delta_r = - \sum_{s=r}^n \varepsilon_s b_{r,s-1} (c_s - 1) \tag{IV,5}$$

$$|b_{rs}| \leq (1 + \varepsilon)^{s-r+1}, \quad |c_s - 1| \leq (1 + \varepsilon)^{n-s+2} - 1, \quad |b_{1s}| \leq (1 + \varepsilon)^{s-1}. \tag{IV,6}$$

Daraus erhält man für  $r > 1$ :

$$|\delta_r| \leq \sum_{s=r}^n \varepsilon (1 + \varepsilon)^{s-r} ((1 + \varepsilon)^{n+2-s} - 1) = \varepsilon (1 + \varepsilon)^{n-r+2} (n - r + 1) - (1 + \varepsilon)^{n-r+1} + 1,$$

also

$$|\delta_{n-t+1}| \leq (\varepsilon (1 + \varepsilon) t - 1) (1 + \varepsilon)^t + 1 \quad (1 \leq t \leq n - 1) \tag{IV,7a}$$

und für  $r = 1$ :

$$|\delta_1| \leq \sum_{s=2}^n \varepsilon (1 + \varepsilon)^{s-2} ((1 + \varepsilon)^{n-s+2} - 1) = \varepsilon (n - 1) (1 + \varepsilon)^n - (1 + \varepsilon)^{n-1} + 1,$$

also

$$|\delta_1| \leq (\varepsilon (1 + \varepsilon) (n - 1) - 1) (1 + \varepsilon)^{n-1} + 1. \tag{IV,7b}$$

Nun ist  $\tilde{s} = (\tilde{s}_n + \tilde{w}_n) (1 + \delta) = s' (1 + \delta) = \left( s + \sum_{i=1}^n a_i \delta_i \right) (1 + \delta)$ , also ist

$$\tilde{s} = s + \delta s + \sum_{i=1}^n a_i \eta_i \tag{IV,8}$$

$$\eta_r = - (1 + \delta) \sum_{s=r}^n \varepsilon_s b_{r,s-1} (c_s - 1) \tag{IV,9}$$

$$|\eta_{n-t+1}| \leq (\varepsilon (1 + \varepsilon) t - 1) (1 + \varepsilon)^{t+1} + 1 + \varepsilon \quad (1 \leq t \leq n - 1)$$

$$|\eta_1| \leq (\varepsilon (1 + \varepsilon) (n - 1) - 1) (1 + \varepsilon)^n + 1 + \varepsilon \tag{IV,10}$$

$$|\delta| \leq \varepsilon.$$

Wie man aus (IV,9) sieht, werden auch die Schranken in (IV,10) angenommen, und zwar für  $\alpha_m = \beta_m = \varepsilon_m = \varepsilon$  ( $m > 1$ ).

### 7. Vereinfachung der Abschätzungen für kleine $\varepsilon$

Es ist

$$\varepsilon - \frac{\varepsilon^2}{2} \leq r \leq \varepsilon \quad \text{für} \quad r = \ln(1 + \varepsilon). \tag{10}$$

Daher ist für Verfahren I nach (I,4):

$$\begin{aligned} |\eta_{n-t+2}| &\leq (1 + \varepsilon)^{t-1} - 1 = \varepsilon(t-1) + r^2 t^2 \frac{e^{rt} - rt - 1}{(rt)^2} - t(\varepsilon - r) - \varepsilon((1 + \varepsilon)^{t-1} - 1) \\ &\leq \varepsilon(t-1) + \varepsilon^2 t^2 \frac{e^k - k - 1}{k^2} \quad \text{für } \varepsilon t \leq k, \end{aligned}$$

ebenso

$$|\eta_1| \leq (1 + \varepsilon)^{n-1} - 1 \leq \varepsilon(n-1) + \varepsilon^2 n^2 \frac{e^k - k - 1}{k^2} \quad \text{für } \varepsilon n \leq k,$$

also

$$|\eta_i| \leq \varepsilon(t-1) + \varepsilon^2 t^2 \frac{e^k - k - 1}{k^2} \quad (\varepsilon n \leq k) \quad (\text{I,5})$$

$$t = n - i + 2 \quad \text{für } i > 1, \quad t = n \quad \text{für } i = 1.$$

Für Verfahren II gilt (mit  $u = t - 1$ ):

$$\begin{aligned} |\eta_i| &\leq (1 + \varepsilon)^{u+1} - 1 = \varepsilon(u+1) + r^2 u^2 \frac{e^{ru} - ru - 1}{(ru)^2} + \varepsilon r u \frac{e^{ru} - 1}{ru} - u(\varepsilon - r), \\ |\eta_i| &\leq \varepsilon t + \varepsilon^2 \left( (t-1)^2 \frac{e^k - k - 1}{k^2} + (t-1) \frac{e^k - 1}{k} \right) \quad (\varepsilon(t-1) \leq k) \quad (\text{II,5}) \end{aligned}$$

$$t \leq \log_2 n + 1.$$

Für Verfahren III ist nach (III,14):

$$\begin{aligned} |\eta_{n-t+2}| &\leq (\varepsilon(1 + \varepsilon)^2 t - \varepsilon^2(1 + \varepsilon) - 1)(1 + \varepsilon)^t + 1 + \varepsilon \\ &= \varepsilon + r^2 t^2 \frac{e^{rt}(rt-1) + 1}{(rt)^2} + \varepsilon r t \frac{5}{2} e^{rt} + \varepsilon^2 e^{rt} (\dots) \\ &\leq \varepsilon + \varepsilon^2 \left( t^2 \frac{e^k(k-1) + 1}{k^2} + t \frac{5}{2} e^k \right) \quad \left( \varepsilon t \leq k \leq \frac{4}{9}, \varepsilon \leq \frac{1}{2} \right), \end{aligned}$$

denn  $\frac{e^x(x-1) + 1}{x^2} = \sum_{i=0}^{\infty} \frac{i+1}{(i+2)!} x^i$  ist für  $x > 0$  monoton wachsend, und aus (10) erhält man durch eine einfache Abschätzung  $(\dots) \leq 0$  für  $\varepsilon t \leq \frac{4}{9}, \varepsilon \leq \frac{1}{2}$ . Aus (III,14) sieht man, daß  $|\eta_1|$  durch die Schranke von  $|\eta_2|$  beschränkt wird. Daher gilt:

$$|\eta_i| \leq \varepsilon + \varepsilon^2 \left( t^2 \frac{e^k(k-1) + 1}{k^2} + t \frac{5}{2} e^k \right) \quad \left( \varepsilon t \leq k \leq \frac{4}{9}, \varepsilon \leq \frac{1}{2} \right) \quad (\text{III,15})$$

$$t = n - i + 2 \quad \text{für } i > 1, \quad t = n \quad \text{für } i = 1.$$

Für Verfahren IV erhält man aus (IV,10):

$$\begin{aligned} |\eta_{n-t+2}| &\leq (\varepsilon(1 + \varepsilon)(t-1) - 1)(1 + \varepsilon)^t + 1 + \varepsilon \\ &= r^2 t^2 \frac{e^{rt}(rt-1) + 1}{(rt)^2} + \varepsilon r t \frac{e^{rt}(3rt-2) + 2}{2rt} + \varepsilon^2 e^{rt} (\dots) \\ &\leq \varepsilon^2 \left( t^2 \frac{e^k(k-1) + 1}{k^2} + t \frac{e^k(3k-2) + 2}{2k} \right) \quad \left( \varepsilon t \leq k \leq 1, \varepsilon \leq \frac{1}{2} \right), \end{aligned}$$

da  $\frac{e^x(3x-2) + 2}{x^2} = \sum_{i=0}^{\infty} \frac{3i+1}{(i+1)!} x^i$  für  $x > 0$  monoton wächst und da  $(\dots) \leq 0$  für  $\varepsilon t \leq 1, \varepsilon \leq \frac{1}{2}$ .

Wieder wird  $|\eta_1|$  durch die Schranke von  $|\eta_2|$  beschränkt, also:

$$|\eta_i| \leq \varepsilon^2 \left( t^2 \frac{e^k(k-1) + 1}{k^2} + t \frac{e^k(3k-2) + 2}{2k} \right) \quad \left( \varepsilon t \leq k \leq 1, \varepsilon \leq \frac{1}{2} \right), \quad (\text{IV,11}) \text{ al}$$

$$t = n - i + 2 \quad \text{für } i > 1, \quad t = n \quad \text{für } i = 1.$$

Nun gilt für alle vier Verfahren eine Gleichung der Form

$$\tilde{s} = s + \delta s + \sum_{i=1}^n a_i \eta_i \quad (11)$$



( $|\delta| \leq \varepsilon$  für Verfahren III, IV;  $\delta = 0$  für Verfahren I, II). Daraus lassen sich zwei einfache Formeln für den absoluten Fehler  $s - \tilde{s}$  ableiten. Denn aus (II) folgt:

$$|s - \tilde{s}| \leq |\delta| |s| + \sum_{i=1}^n |a_i| |\eta_i|;$$

also

$$|s - \tilde{s}| \leq |\delta| |s| + \left( \text{Max}_{1 \leq i \leq n} |\eta_i| \right) \left( \sum_{i=1}^n |a_i| \right) \quad (12)$$

$$|s - \tilde{s}| \leq |\delta| |s| + \left( \sum_{i=1}^n |\eta_i| \right) \left( \text{Max}_{1 \leq i \leq n} |a_i| \right). \quad (13)$$

Formel (12) ist nützlich, wenn alle  $a_i$  dasselbe Vorzeichen haben, Formel (13), wenn alle  $a_i$  von derselben Größenordnung sind. Auf diesen Formeln baut die Diskussion der Verfahren in diesem und dem nächsten Abschnitt auf.

Im folgenden wird  $\varepsilon n \leq k = 1/3$  vorausgesetzt, was in der Praxis wohl meist erfüllt ist. Dann gilt:

$$\frac{e^k - 1}{k} < \frac{6}{5}, \quad \frac{e^k (k - 1) + 1}{k^2} < \frac{3}{4}, \quad \frac{e^k - k - 1}{k^2} < \frac{3}{5}, \quad \frac{5}{2} e^k < \frac{7}{2}, \quad \frac{e^k (3k - 2) + 2}{2k} < 1. \quad (14)$$

Für die einzelnen Verfahren erhält man aus (12), (13), damit folgende Formeln:

Verfahren I: Aus (I,5) erhält man mit (14):

$$\text{Max}_{1 \leq i \leq n} |\eta_i| \leq \varepsilon (n - 1) + \frac{3}{5} \varepsilon^2 n^2$$

$$\sum_{i=1}^n |\eta_i| \leq \varepsilon (n - 1) + \frac{3}{5} \varepsilon^2 n^2 + \sum_{t=2}^n \left( \varepsilon (t - 1) + \frac{3}{5} \varepsilon^2 t^2 \right),$$

$$\begin{aligned} \sum_{i=1}^n |\eta_i| &\leq \varepsilon \frac{n^2 + n - 2}{2} + \varepsilon^2 \frac{2n^3 + 9n^2 + n - 6}{10}, \\ &\leq \varepsilon \frac{n^2 + n - 2}{2} + \varepsilon^2 \frac{n^3 + 5n^2}{5} \end{aligned}$$

also

$$|s - \tilde{s}| \leq \left( \varepsilon (n - 1) + \frac{2}{5} \varepsilon^2 n^2 \right) \sum_{i=1}^n |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right) \quad (I,6)$$

$$|s - \tilde{s}| \leq \left( \varepsilon \frac{n^2 + n - 2}{2} + \varepsilon^2 \frac{n^3 + 5n^2}{5} \right) \text{Max}_{1 \leq i \leq n} |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right) \quad (I,7)$$

Verfahren II: Aus (II,2), (II,5), (14) folgt unmittelbar:

$$|s - \tilde{s}| \leq \left( \varepsilon (1 + \log_2 n) + \frac{3}{5} \varepsilon^2 \log_2 n (\log_2 n + 2) \right) \sum_{i=1}^n |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right), \quad (II,6)$$

$$|s - \tilde{s}| \leq \left( \varepsilon (1 + \log_2 n) + \frac{3}{5} \varepsilon^2 \log_2 n (\log_2 n + 2) \right) n \text{Max}_{1 \leq i \leq n} |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right). \quad (II,7)$$

Verfahren III: Nach (III,15) ist

$$\text{Max}_{1 \leq i \leq n} |\eta_i| \leq \varepsilon + \varepsilon^2 \left( \frac{3}{4} n^2 + \frac{7}{2} n \right)$$

$$\begin{aligned} \sum_{i=1}^n |\eta_i| &\leq \varepsilon + \varepsilon^2 \left( \frac{3}{4} n^2 + \frac{7}{2} n \right) + \sum_{t=2}^n \left( \varepsilon + \varepsilon^2 \left( \frac{3}{4} t^2 + \frac{7}{2} t \right) \right) \\ &= \varepsilon n + \varepsilon^2 \left( \frac{1}{4} n^3 + \frac{23}{8} n^2 + \frac{43}{8} n - \frac{17}{4} \right) \leq \varepsilon n + \varepsilon^2 \left( \frac{1}{4} n^3 + 3n^2 + 4n \right), \end{aligned}$$

also

$$|s - \tilde{s}| \leq \varepsilon |s| + \left( \varepsilon + \varepsilon^2 \left( \frac{3}{4} n^2 + \frac{7}{2} n \right) \right) \sum_{i=1}^n |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right) \quad (III,16)$$

$$|s - \tilde{s}| \leq \varepsilon |s| + \left( \varepsilon n + \varepsilon^2 \left( \frac{1}{4} n^3 + 3n^2 + 4n \right) \right) \text{Max}_{1 \leq i \leq n} |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right). \quad (III,17)$$

Verfahren IV: Aus (IV,11) folgt

$$\begin{aligned} \text{Max}_{1 \leq i \leq n} |\eta_i| &\leq \varepsilon^2 \left( \frac{3}{4} n^2 + n \right) \\ \sum_{i=1}^n |\eta_i| &\leq \varepsilon^2 \left( \frac{3}{4} n^2 + n \right) + \sum_{t=2}^n \varepsilon^2 \left( \frac{3}{4} t^2 + t \right) = \varepsilon^2 \left( \frac{1}{4} n^3 + \frac{17}{8} n^2 + \frac{25}{8} n - \frac{11}{4} \right) \\ &\leq \varepsilon^2 \left( \frac{1}{4} n^3 + \frac{5}{2} n^2 + n \right), \end{aligned}$$

also

$$|s - \tilde{s}| \leq \varepsilon |s| + \varepsilon^2 \left( \frac{3}{4} n^2 + n \right) \sum_{i=1}^n |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right) \quad (\text{IV,12})$$

$$|s - \tilde{s}| \leq \varepsilon |s| + \varepsilon^2 \left( \frac{1}{4} n^3 + \frac{5}{2} n^2 + n \right) \text{Max}_{1 \leq i \leq n} |a_i| \quad \left( \varepsilon n \leq \frac{1}{3} \right). \quad (\text{IV,13})$$

Die Ergebnisse ((I,6, 7), (II,6, 7), (III,16, 17), (IV,12, 13)) sind in Tabelle I zusammengefaßt.

Bemerkung: Der  $\varepsilon$ -Anteil der Verfahrensfehler von I, II, III wurde — in etwas anderer Form — schon von BABUŠKA [3] angegeben.

Tabelle I. Fehlerkoeffizienten der Summierungsverfahren

$$\left\{ \begin{aligned} |s - \tilde{s}| &\leq A \varepsilon |s| + (B \varepsilon + C \varepsilon^2) \sum_{i=1}^n |a_i| \\ |s - \tilde{s}| &\leq A \varepsilon |s| + (D \varepsilon + E \varepsilon^2) \text{Max}_{1 \leq i \leq n} |a_i| \end{aligned} \right\} \text{für } \varepsilon n \leq \frac{1}{3} \quad \left\{ \begin{array}{l} (15) \\ (16) \end{array} \right.$$

Verf. Nr.	A	B(n)	C(n)	D(n)	E(n)
I	0	$n - 1$	$0,6 n^2$	$0,5 n^2 + 0,5 n - 1$	$0,2 n^3 + n^2$
II	0	$1 + \log_2 n$	$0,6 (\log_2 n)^2 + 1,2 \log_2 n$	$n (1 + \log_2 n)$	$0,6 n \log_2 n$
III	1	1	$0,75 n^2 + 3,5 n$	$n$	$0,25 n^3 + 3 n^2 + 4 n$
IV	1	0	$0,75 n^2 + n$	0	$0,25 n^3 + 2,5 n^2 + n$

## 8. Diskussion der Ergebnisse

8.1. Zunächst sollen alle  $a_i$  dasselbe Vorzeichen haben. Dann ist mit der Abkürzung

$$s_{\text{abs}} := \sum_{i=1}^n |a_i| \quad (17)$$

$s_{\text{abs}} = |s|$  und der relative Fehler ist nach (15)

$$\left| \frac{s - \tilde{s}}{s} \right| \leq (A + B) \varepsilon + C \varepsilon^2. \quad (18)$$

Solange  $n < \varepsilon^{-1/2}$  ist, spielt der Term  $C \varepsilon^2$  in (18) gegenüber  $(A + B) \varepsilon$  keine Rolle. Das erste Verfahren hat dann einen maximalen Fehler, der proportional zur Zahl  $n$  der Summanden wächst. Im zweiten Verfahren hat man trotz gleichen Rechenaufwands einen maximalen Fehler, der nur mit dem Logarithmus von  $n$  wächst (ist z. B.  $n = 10^6$ , so ist der Verfahren-I-Fehler  $10^6 \varepsilon$ , der Verfahren-II-Fehler nur  $20 \varepsilon$ !). Wenn daher die Zahl der Summanden von vornherein bekannt ist, empfiehlt es sich, anstelle von Verfahren I stets Verfahren II zu benutzen. (Bei der Summierung einer unendlichen Reihe ist das natürlich nicht möglich.)

Beim KAHAN-BABUŠKA-Verfahren (III) ist der maximale Fehler fast konstant. Mit etwa  $2 \varepsilon$  liegt er in derselben Größenordnung wie der Fehler bei einer Normalverfahrenrechnung mit doppelter Genauigkeit ( $\varepsilon^2$  statt  $\varepsilon$ ), bei der das Ergebnis auf einfache Genauigkeit gerundet wird. Da aber der Rechenaufwand für doppelte Genauigkeit nur doppelt so groß ist wie der für einfache, lohnt sich Verfahren II nur dann, wenn auf der Maschine keine Möglichkeit, mit doppelter Genauigkeit zu rechnen, besteht. Dann ist das Verfahren ein guter Ersatz dafür (das gilt nur für  $s_{\text{abs}} \approx s$ ; vgl. unten!). Die Fehlerschranke von Verfahren IV hat ebenfalls die Größenordnung  $\varepsilon$ . Die geringfügige Verbesserung gegenüber III lohnt den höheren Aufwand zur Berechnung der Summe jedoch nicht.

Ist  $\varepsilon^{-1} \leq n^2 \ll \varepsilon^{-1} \log_2(\varepsilon^{-1})$ , so wachsen die Schranken in (18) zwar quadratisch in  $n$  bzw.  $\log_2 n$ , aber das Obengesagte bleibt qualitativ richtig. Für  $n^2 \approx \varepsilon^{-1} \log_2(\varepsilon^{-1})$  sind Verfahren II und III etwa gleich genau; für  $n^2 \geq 2 \varepsilon^{-1} \log_2(\varepsilon^{-1})$  ist schließlich Verfahren II das genaueste Verfahren (genauer als III und IV), da

$$0,6 (\log_2 n)^2 + 1,2 \log_2 n \leq n \quad \text{für alle } n,$$

$$\log_2 n \leq \frac{3}{4} n^2 \varepsilon \quad \text{für } n^2 \geq 2 \varepsilon^{-1} \log_2(\varepsilon^{-1})$$

( $\frac{\log_2 t}{t}$  ist für  $t \geq 3$  monoton fallend und  $\leq \frac{3}{2} \varepsilon$  für  $t = 2 \varepsilon^{-1} \log_2(\varepsilon^{-1})$ , also ist  $\log_2 n = \frac{1}{2} n^2 \frac{\log_2 n^2}{n^2} \leq \frac{3}{4}$  für  $n \geq 2 \varepsilon^{-1} \log_2(\varepsilon^{-1})$ ).

8.2. Es sei  $s_{\text{abs}} \approx |s|$ . Dann gilt (18) angenähert. Daher lassen sich alle obigen Aussagen näherungsweise übertragen. Insbesondere ist Verfahren III fast so gut wie Verfahren IV, für kleine  $n$  ist Verfahren III (IV) das beste, für große  $n$  ist Verfahren II am besten geeignet. Verfahren I ist wieder das schlechteste Verfahren von allen.

8.3. Es sei  $s_{\text{abs}} \gg |s|$ . Jetzt kann man das Glied  $A \varepsilon |s|$  vernachlässigen und man erhält

$$|s - \tilde{s}| \leq (B + C \varepsilon) (s_{\text{abs}} \varepsilon). \tag{19}$$

8.3.1. Ist  $s_{\text{abs}}$  so groß, daß  $s_{\text{abs}} \cdot \varepsilon$  mindestens die Größenordnung von  $|s|$  hat (etwa wenn  $s$  fast Null ist), so ist der Näherungswert  $\tilde{s}$  i. a. nur dann brauchbar, wenn  $B = 0$  ist, also wenn  $s$  nach IV berechnet wird. In den anderen Fällen kann eine so starke Auslöschung auftreten, daß nur noch eine oder keine Ziffer von  $\tilde{s}$  mehr richtig ist.

8.3.2. Auch wenn der Extremfall 8.3.1 nicht vorliegt, bestimmt die Größe von  $B$  zunehmend die Fehler-schranke. Für Verfahren I ist der maximale relative Fehler jetzt  $\gg n \varepsilon$ , für große  $n$  ist das Verfahren unbrauchbar. Verfahren II hat jetzt einen Fehler von ungefähr  $\varepsilon_{\text{abs}} \log_2 n s_{\text{abs}}/|s|$ , der in nicht zu ungünstigen Fällen etwa die Größenordnung  $n \varepsilon$  haben kann. Dieses Verfahren ist also noch so gut wie Verfahren I im günstigsten Fall, also bedingt verwendbar. Für das KAHAN-BABUŠKA-Verfahren ist die Schranke ungefähr  $\varepsilon s_{\text{abs}}/|s|$  (für kleine und mittlere  $n$ ), wird also in der Mehrzahl der Fälle (z. B. Summation absolut konvergenter Reihen mit alternierenden Vorzeichen) noch gute Ergebnisse liefern. Das verbesserte KAHAN-BABUŠKA-Verfahren hat jedoch eine Schranke von größenordnungsmäßig  $n^2 \varepsilon^2 \frac{s_{\text{abs}}}{|s|}$  für den relativen,  $n^2 \varepsilon^2 s_{\text{abs}}$  für den absoluten Fehler, liefert also für nicht allzu große  $n$  (z. B. für  $n \ll \varepsilon^{-1/2}$ ) einen Maximalfehler, der in der Größenordnung von  $\varepsilon$  liegt! Hier wird also die Auslöschung fast vollständig kompensiert. In diesem Fall wird sich daher der gegenüber I fünffache Rechenaufwand praktisch immer lohnen (wenn man nicht mit doppelter Genauigkeit rechnen kann).

8.4. Wenn  $s_{\text{abs}}$  nicht sehr groß gegenüber  $|s|$  ist, so ist wieder Verfahren II dem ersten Verfahren vorzuziehen. Ob man Verfahren III (oder IV) verwendet, wird von der benötigten Genauigkeit der Summe abhängen. Verfahren III und IV dürften dann etwa gleichwertig sein, der Vorteil größerer Genauigkeit bei IV wird durch den größeren Rechenaufwand aufgewogen.

Das hier Gesagte schließt natürlich nicht aus, daß in Einzelfällen Verfahren I bessere Ergebnisse liefert als Verfahren II oder III; oder das Verfahren I auch einmal für ein kleines  $n$  besser ist als Verfahren III; Beispiele dafür werden im nächsten Abschnitt gegeben.

### 9. Beispiele

Für die ersten drei Beispiele sei eine  $L$ -ziffrige Gleitkomma-Dualmaschine mit optimaler Rundung<sup>3)</sup> zugrunde gelegt; eine genaue 1 in der  $(L + 1)^{\text{ten}}$  Dualstelle soll abgerundet werden. Nach der Bemerkung in Abschnitt 3 ist daher  $\varepsilon = 2^{-L}$ .

a) Es sei  $n = 7$ ,

$$a_1 = (2^{-1} + 2^{1-L} + 2^{-L}) 2^k,$$

$$a_2 = a_3 = a_4 = a_5 + (2^{-1}) 2^{k-1-L}$$

$$a_6 = (2^{-1}) 2^{k+2}$$

$$a_7 = - (2^{-1} + 2^{-3} + 2^{-L}) 2^{k+2}.$$

Es ist  $s = 0$ . Tabelle 2 gibt die Zwischenwerte für die Summation nach Verfahren I, III und IV wieder, Tabelle 3 die für Verfahren II.

Tabelle 2. Zwischenwerte zu Beispiel a)

$m$	$\tilde{s}_m$	$(\tilde{w}_m \text{III})$	$\tilde{w}_m \text{(IV)}$	Ergebnisse
1	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	0	0	$\tilde{s}_I = 0$
2	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$(2^{-1}) 2^{k-1-L}$	$(2^{-1}) 2^{k-1-L}$	$\tilde{s}_{\text{III}} = (2^{-1}) 2^{k+1-L}$
3	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$(2^{-1}) 2^{k-L}$	$(2^{-1}) 2^{k-L}$	$\tilde{s}_{\text{IV}} = 0$
4	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$(2^{-1} + 2^{-2}) 2^{k-L}$	$(2^{-1} + 2^{-2}) 2^{k-L}$	
5	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$(2^{-1}) 2^{k+1-L}$	$(2^{-1}) 2^{k+1-L}$	
6	$(2^{-1} + 2^{-3} + 2^{-L}) 2^{k+2}$	$(2^{-1}) 2^{k+1-L}$	0	
7	0	$(2^{-1}) 2^{k+1-L}$	0	

<sup>3)</sup> siehe Fußnote 2).

Tabelle 3. Zwischenwerte zu Beispiel a)

$m$	$\tilde{s}_{m0}$	$\tilde{s}_{m1}$	$\tilde{s}_{m2}$	$\tilde{s}_{m3}$
1	$(2^{-1} + 2^{1-L}) 2^k$	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$(2^{-1} + 2^{1-L} + 2^{-L}) 2^k$	$-(2^{-1}) 2^{k+1}$
2	$(2^{-1}) 2^{k-1-L}$	$(2^{-1}) 2^{k-L}$	$-(2^{-1} + 2^{2-L}) 2^k$	
3	$(2^{-1}) 2^{k-1-L}$	$(2^{-1}) 2^{k+2}$		
4	$(2^{-1}) 2^{k-1-L}$	$-(2^{-1} + 2^{-3} + 2^{-L}) 2^{k+2}$		
5	$(2^{-1}) 2^{k-1-L}$			
6	$(2^{-1}) 2^{k+2}$		Ergebnis: $\tilde{s}_{II} = -(2^{-1}) 2^{k+1-L}$	
7	$-(2^{-1} + 2^{-3} + 2^{-L}) 2^{k+2}$			
8	0			

Für die absoluten Fehler erhält man

$$|s - \tilde{s}_I| = |s - \tilde{s}_{IV}| = 0, \quad |s - \tilde{s}_{II}| = |s - \tilde{s}_{III}| = 2^{k-L} = \frac{1}{2} \varepsilon \max_{1 \leq i \leq 7} |a_i|.$$

Die in Tabelle 1 angegebenen Schranken überschätzen den wirklichen Fehler bei Verfahren III um das Doppelte, bei Verfahren II um das Achtfache. Der nach I berechnete Wert ist sogar genau. Trotz der wesentlich schlechteren Schranke in Tab. 1 ist also Verfahren I hier genauer als Verfahren II oder III. Verfahren IV liefert das exakte Ergebnis.

b) Es sei  $n = 4$ ,  $a_1 = 1$ ,  $a_2 = 2^{-2} + 2^{-1-L}$ ,  $a_3 = -1$ ,  $a_4 = -(2^{-2} + 2^{-1-L})$ . Wie man leicht nachprüft, ist jetzt

$$s = 0, \quad \tilde{s}_{II} = \tilde{s}_{IV} = 0, \quad \tilde{s}_I = \tilde{s}_{III} = -2^{-L-1},$$

also Verfahren II genauer als Verfahren I oder III. Verfahren IV gibt wieder das exakte Ergebnis.

c) Sei  $n = 3$ ,  $a_1 = 1$ ,  $a_2 = 2^{L+1} + 2^L$ ,  $a_3 = -(2^{L+1} + 2^L)$ . Jetzt ist

$$s = 1, \quad \tilde{s}_I = \tilde{s}_{II} = 0, \quad \tilde{s}_{III} = \tilde{s}_{IV} = 1,$$

also Verfahren I und II ungenauer als Verfahren III oder IV.

Meistens werden von allen Verfahren ungenaue Ergebnisse geliefert, z. B. stets, wenn  $s$  keine Maschinenzahl ist. Ist aber  $s$  als Maschinenzahl darstellbar, so ist oft  $\tilde{s}_{IV} = s$ ; vgl. dazu Beispiel d), Tab. 4, 5. ( $n = 7$ ; die  $a_m$  und die Zwischenwerte sind in den Tabellen angegeben.)

Tabelle 4. Zwischenwerte zu Beispiel d)

$m$	$a_m$	$\tilde{s}_m$	$\tilde{w}_m(\text{III})$	$\tilde{w}_m(\text{IV})$	Ergebnisse
1	0.5555 $10^0$	0.5555 $10^0$	0	0	$\tilde{s}_I = 0.778010^1$
2	0.5555 $10^1$	0.6111 $10^1$	-0.1000 $10^{-2}$	-0.5000 $10^{-2}$	$\tilde{s}_{III} = 0.777410^1$
3	0.5555 $10^0$	0.6667 $10^1$	-0.1500 $10^{-2}$	-0.1000 $10^{-2}$	$\tilde{s}_{IV} = 0.777710^1$
4	0.5555 $10^1$	0.7223 $10^1$	-0.2000 $10^{-2}$	-0.1500 $10^{-2}$	
5	0.5555 $10^2$	0.6277 $10^2$	-0.2000 $10^{-2}$	+0.1500 $10^{-2}$	
6	0.5555 $10^0$	0.6333 $10^2$	-0.6500 $10^{-2}$	-0.3000 $10^{-2}$	
7	-0.5555 $10^2$	0.7780 $10^1$	-0.6500 $10^{-2}$	-0.3000 $10^{-2}$	

Tabelle 5. Zwischenwerte zu Beispiel d)

$m$	$\tilde{s}_{m0}$	$\tilde{s}_{m1}$	$\tilde{s}_{m2}$	$\tilde{s}_{m3}$
1	0.5555 $10^0$	0.6111 $10^1$	0.7222 $10^1$	0.7782 $10^1$
2	0.5555 $10^1$	0.1111 $10^1$	0.5600 $10^0$	
3	0.5555 $10^0$	0.5611 $10^2$		
4	0.5555 $10^1$	-0.5555 $10^2$		
5	0.5555 $10^2$			
6	0.5555 $10^0$		Ergebnis: $\tilde{s}_{II} = 0.778210^1$	
7	-0.5555 $10^2$			
8	0			

Die Werte wurden unter Annahme einer Dezimalmaschine mit vier geltenden Ziffern und der üblichen Rundung berechnet. Das genaue Ergebnis ist  $s = 0.777710^1 = \tilde{s}_{IV}$ . Wie zu erwarten war, nimmt jetzt die Genauigkeit der Verfahren in der Reihenfolge I, II, III, IV zu.

An den Tabellen 2 und 4 sieht man gut, woran es liegt, daß Verfahren III i. a. schlechter ist als Verfahren IV. Immer wenn ein großes  $a_m$  zu einem kleinen  $s_{m-1}$  addiert wird, berechnet der Korrektor  $w_m$  des KAHAN-BABUŠKA-Verfahrens nicht mehr, wie geplant, den Rundungsfehler, sondern einen Wert, der mit dem Rundungsfehler nicht mehr viel zu tun hat. Die im verbesserten KAHAN-BABUŠKA-Verfahren eingebaute Abfrage vertauscht hier sozusagen die Rollen von  $a_m$  und  $s_{m-1}$  und berechnet dadurch — nach Satz 2 des fünften Abschnitts — den richtigen Fehler.

NICKEL [3] berechnete auf einer Maschine mit  $\varepsilon = 2^{-30}$  die Summe  $\sum_{i=1}^{2^t} \frac{1}{i}$  für  $0 \leq t \leq 14$  nach Verfahren I und III (= IV in diesem Fall). Für  $s = \sum_{i=1}^{2^{12}} \frac{1}{i}$  z. B. erhält er

$$|s - \tilde{s}_I| \approx 25 \cdot 10^{-8}, \quad |s - \tilde{s}_{IV}| \leq 0.5 \cdot 10^{-8}.$$

Unter Verwendung von  $s < 9$  ergibt sich aus Tabelle 1:

$$|s - \tilde{s}_I| \leq 3500 \cdot 10^{-8}, \quad |s - \tilde{s}_{IV}| \leq 1.7 \cdot 10^{-8}.$$

Der Fehler von Verfahren I wird also durch die Schranke in Tab. 1 um das 140fache, bei Verfahren IV mindestens um das dreifache überschätzt. Jedoch liegt die Tab. 1-Schranke für Verfahren IV immer noch niedriger als der maschinell berechnete Fehler von Verfahren I.

Das letzte Beispiel soll wieder auf einer  $L$ -ziffrigen Dualmaschine gerechnet werden. Diesmal soll die Rundung jedoch durch Abschneiden der  $(L+i)$ -ten Ziffern ( $i \geq 1$ ) geschehen ( $\varepsilon = 2^{1-L}$ ). Nun sei  $n = 2^t$ ,  $2 \leq t \leq L-1$ ,

$$a_1 = 1, \quad a_2 = 2^{-L}, \dots, a_{2^{t+1}} = \dots = a_{2^i+1} = 2^{-L} + 2^{i-2L} \quad (1 \leq i \leq t-1).$$

Man findet leicht:

$$\tilde{s} = 1 + 2^{t-L} - 2^{-L} + \frac{1}{3} 2^{-2L} (2^{2t} - 4),$$

$$\tilde{s}_I = 1,$$

$$\tilde{s}_{II} = 1 + 2^{t-L} - 2^{1-L} \quad \left(t > \frac{L}{2} + 1\right)$$

$$\tilde{s}_{II} = 1 + 2^{t-L} - 2^{1-L} + \frac{1}{3} \cdot 2^{-2L} (2^{2t} - 2^{L+r}) \quad \left(t > \frac{L}{2} + 1; r=2 \text{ für gerade } L, r=1 \text{ für ungerade } L\right),$$

$$\tilde{s}_{III} = 1 + 2^{t-L} - 2^{1-L}.$$

Man sieht unmittelbar, daß Verfahren II für  $n^2 \geq 8 \varepsilon^{-1}$  genauer ist als Verfahren III (hier = Verfahren IV); das ist in guter Übereinstimmung mit der in Abschnitt 8 abgeleiteten Bedingung  $n^2 \geq \varepsilon^{-1} \log_2 \varepsilon^{-1}$ . Daß  $\tilde{s}_{II}$  für kleine  $n$  hier =  $\tilde{s}_{III}$  ist, folgt daraus, daß wegen der speziellen Form der Summanden bei Verfahren II nur die  $\tilde{s}_{1i}$  gerundet sind.

Für die Fehler erhält man

$$|\tilde{s}_I - s| = 2^{t-L} + \frac{1}{3} (2^{t-L})^2 - 2^{-L} - \frac{1}{3} 2^{2-2L} = \frac{1}{2} (n-1) \varepsilon + \left(\frac{1}{12} n^2 - \frac{1}{3}\right) \varepsilon^2,$$

$$|\tilde{s}_{II} - s| = 2^{-L} + \frac{1}{3} 2^{-2L} (2^{2t} - 4) = \frac{1}{2} \varepsilon + \left(\frac{1}{12} n^2 - \frac{1}{3}\right) \varepsilon^2 \quad (\text{für } n^2 < 8 \varepsilon^{-1}),$$

$$|\tilde{s}_{II} - s| = \left(\frac{1}{3} 2^r + 1\right) 2^{-L} - \frac{1}{3} 2^{2-2L} \leq \frac{7}{6} \varepsilon - \frac{1}{3} \varepsilon^2 \quad (\text{für } n^2 \geq 8 \varepsilon^{-1}),$$

$$|\tilde{s}_{III} - s| = 2^{-L} + \frac{1}{3} 2^{-2L} (2^{2t} - 4) = \frac{1}{2} \varepsilon + \left(\frac{1}{12} n^2 - \frac{1}{3}\right) \varepsilon^2.$$

Die Fehlerschranken aus Tab. 1 überschätzen hier den Fehler bei Verfahren I um den Faktor 2, bei Verfahren II um den Faktor  $2 \log_2 \varepsilon^{-1}$ , bei Verfahren III um den Faktor 4 (da  $s \approx 1$ ). Der hohe Fehler (im Vergleich zur Schranke) bei Verfahren I und III und der kleine Fehler bei Verfahren II rühren von der speziellen Form der  $a_i$  her; im allgemeinen werden die Fehler bei Verfahren I und III kleiner, bei Verfahren II größer sein. Trotzdem gibt dieses Beispiel die Unterschiede der benutzten Summierungsverfahren gut wieder, wie sie in Abschnitt 8.1 diskutiert wurden.

### Literatur

- 1 BABUŠKA, I., Numerical Stability in Mathematical Analysis, Information Processing 68, Amsterdam, 11-22 (1969).
- 2 KAHAN, W., Further Remarks on Reducing Truncation Errors, Comm. ACM, 8, 40 (1965).
- 3 NICKEL, K., Das Kahan-Babuškache Summierungsverfahren in Triplex-ALGOL 60, ZAMM 50, 369-373 (1970).
- 4 WILKINSON, J. H., Rounding Errors in Algebraic Processes, London 1963; dt. Übersetzung: Rundungsfehler, Berlin-Heidelberg-New York 1969.

Eingereicht am 6. 3. 1973

Anschrift: ARNOLD NEUMAIER, 75 Karlsruhe 41, Kieselweg 1, BRD