

FEHLERRECHNUNG

Prof. Dr. Jörg Hertling, Techn. Univ. Wien

I. EINFÜHRUNG

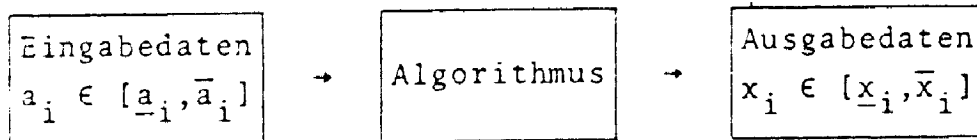
Die zunehmende Verwendung von elektronischen Rechnern birgt die Gefahr, numerischen Ergebnissen kritiklos zu vertrauen. Offensichtlich ist aber ein numerisches Ergebnis ohne eine Vorstellung über dessen Genauigkeit sinnlos.

Die Fehlerquellen können in drei Klassen aufgeteilt werden

- A) Datenfehler.
- B) Rundungsfehler.
- C) Verfahrensfehler.

Im allgemeinen wird ein numerisches Ergebnis von all diesen Fehlerquellen beeinflusst.

ad A) In realistischen Problemen haben Eingabedaten a_i im allgemeinen auf Grund ihrer Herkunft (etwa durch Messung) nur eine beschränkte (absolute oder relative) Genauigkeit, man spricht vom "Datenfehler". Tatsächlich wäre es vernünftiger, davon auszugehen, daß die Eingabedaten a_i jeweils in einem Intervall $[\underline{a}_i, \bar{a}_i]$ liegen. Naturgemäß kann man dann nach der Rechnung auch von den Ausgabedaten x_i nur erwarten, daß sie in einem Intervall $[\underline{x}_i, \bar{x}_i]$ liegen. Diesen Effekt bezeichnet man als "Datenfehlereffekt".



Obwohl es Maschinen gibt, die tatsächlich in jedem Schritt mit Intervallen rechnen, ist dies i.a. zu aufwendig.

Grundsätzlich muß man aber zwei Klassen von Problemen unterscheiden:

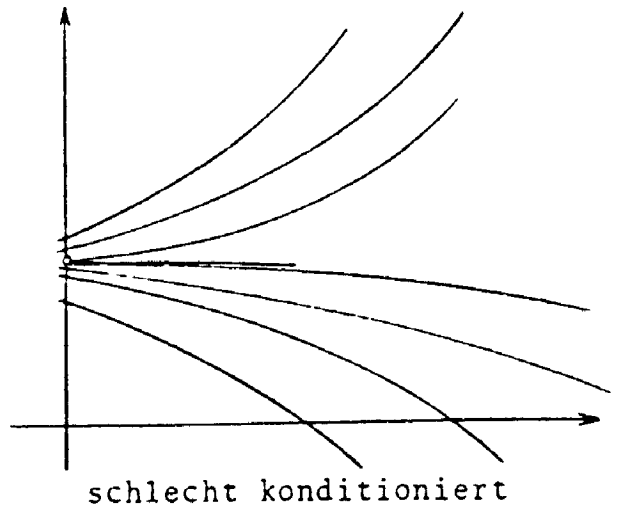
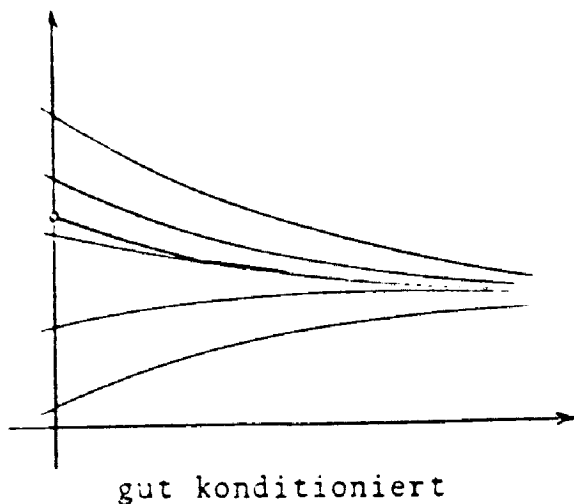
- 1) Probleme, bei denen eine geringe Änderung der Eingabedaten auch eine geringe Änderung des Ergebnisses bewirkt, diese Probleme nennt man "gut konditioniert" oder "stabil".
- 2) Probleme, bei denen eine geringe Änderung von Eingabedaten sehr große Änderungen des Ergebnisses bewirken, diese Probleme nennt man "schlecht konditioniert" oder "instabil".

Allgemein nennt man die Empfindlichkeit eines Ergebnisses bezüglich der Änderung eines Datenelements "Kondition".

Dafür sollen nun einige Beispiele gegeben werden.

Bsp.: Stellt man ein Gleichungssystem von zwei Gleichungen in zwei Unbekannten graphisch dar, so ist das Problem schlecht konditioniert, wenn die beiden Geraden einen "schleifenden" Schnitt aufweisen; eine leichte Änderung des Anstiegs einer der beiden Geraden kann dann nämlich eine sehr große Verschiebung des Schnittpunktes (d.h. der Lösung) bewirken. Ist der Schnitt nicht schleifend, so ist das Problem gut konditioniert.

Bsp.: Löst man ein Anfangswertproblem für eine gewöhnliche Differentialgleichung erster Ordnung, numerisch (etwa mit dem Eulerverfahren), so ist dieses Problem bei einem "zusammenlaufenden" Richtungsfeld gut konditioniert; bei einem "auseinanderlaufenden" Richtungsfeld ist das Problem schlecht konditioniert, weil man dann schon nach wenigen Schritten des Verfahrens auf eine andere Lösungskurve gerät als jene, die dem Anfangswert entspricht.



Bsp.: Die Ritzsche Achsenkonstruktion ist verschieden konditioniert, je nachdem, ob die konjugierten Durchmesser schon nahezu orthogonal sind, oder nicht.

ad B) Der Algorithmus kann durchgeführt werden auf

a) Analogrechnern:

Dabei werden mathematische Größen und Algorithmen durch "kontinuierliche" physikalische Größen (z.B. Längen oder Spannungen) und physikalische Vorgänge simuliert. Beispiele dafür sind etwa der Rechenschieber und das Planimeter. Auf Analogrechnern ist die Genauigkeit durch die physikalische Meßgenauigkeit begrenzt.

b) Digitalrechnern:

Dabei werden mathematische Größen und Algorithmen durch "diskrete" physikalische Zustände und Verknüpfungen simuliert, die den arithmetischen Verknüpfungen entsprechen. Beispiele dafür sind der Abacus und die üblichen elektronischen Taschenrechner. Auf Digitalrechnern ist die Genauigkeit durch die Zahldarstellung begrenzt. Es treten daher *Rundungsfehler* auf. Natürlich können sich bei schlecht konditionierten Problemen auch die Rundungsfehler verheerend auswirken.

Unsere weiteren Betrachtungen gelten ausschließlich Digitalrechnern.

ad C) Selbst bei rundungsfehlerfreier Rechnung liefern viele Algorithmen nach endlich vielen Schritten die "Lösung" zwar beliebig genau, aber gewöhnlich nicht exakt, man spricht vom *Abbrechfehler*. Es sind keine infinitesimalen Operationen und Grenzprozesse möglich. Differentialquotienten müssen i.a. durch Differenzenquotienten ersetzt werden. Den hierbei auftretenden Fehler bezeichnet man als *Diskretisierungsfehler*. Fehler treten auf, wenn man, wie etwa bei der Interpolation oder Approximation eine Funktion durch eine "Ersatzfunktion" aus einer bestimmten Funktionenklasse ersetzt. Ähnliches gilt für die numerische Integration (etwa bei der Trapez- oder Simpsonregel). Solche Fehler können oft durch sogenannte "Glattheitsvoraussetzungen" abgeschätzt werden.

Naturgemäß kann über Verfahrensfehler wenig Allgemeines gesagt werden, sie müssen bei jedem Algorithmus für sich untersucht werden.

Unsere weiteren Betrachtungen gelten daher den Rundungs- und Verfahrensfehlern. Dazu definieren wir

Def.:

absoluter Fehler := berechneter Wert - idealer Wert

relativer Fehler := $\frac{\text{absoluter Fehler}}{\text{idealer Wert}}$

Entscheidend ist natürlich stets der relative Fehler.

II. ZAHLDARSTELLUNG AUF DIGITALRECHNERN, RUNDUNGSFEHLERANALYSE

Meist wird auf Digitalrechnern eine Zahldarstellung im Dualsystem verwendet:

$$x = \pm (\alpha_n 2^n + \alpha_{n-1} 2^{n-1} + \dots + \alpha_0 2^0 + \alpha_{-1} 2^{-1} + \alpha_{-2} 2^{-2} + \dots)$$

$$\alpha_i = 0 \text{ oder } 1 .$$

Die Anzahl der Dualstellen (bzw. Dezimalstellen) ist festgelegt durch die sogenannte *Wortlänge*.

Bei der sogenannten *Festpunktdarstellung* (oder *Festkommadarstellung*) ist eine Anzahl von Stellen vor dem Komma fest reserviert und eine Anzahl von Stellen nach dem Komma. Solche Maschinen finden bei kommerziellen Rechnungen Verwendung und sind für wissenschaftliche Rechnungen ungeeignet.

Bei der sogenannten *Gleitpunktdarstellung* (oder *Gleitkommadarstellung*) wird die Lage des Kommas durch einen Exponenten angegeben, d.h.

$$x = a \cdot 2^b \quad (\text{bzw. } x = a \cdot 10^b)$$

mit $|a| < 1$, b ganz. a heißt *Mantisse*, b heißt *Exponent* und die Voraussetzung $|a| < 1$ bedeutet, daß der Exponent so zu wählen ist, daß die Mantisse erst nach dem Komma Ziffern aufweist, die ungleich Null sind.

Auf Rechenanlagen ist nun eine feste Anzahl t von Stellen für die Mantisse reserviert und eine feste Anzahl e von Stellen für den Exponenten. Zusätzlich sind noch zwei Stellen für die Vorzeichen von Mantisse und Exponent reserviert. Damit ist aber die Darstellung einer Zahl noch nicht eindeutig, denn man kann die Mantisse ja nach dem Komma noch mit einer "beliebigen" Anzahl von Nullen beginnen lassen und diese "Verschie-

bung" der Mantisse durch den Exponenten kompensieren. Man wählt daher die sogenannte *Normalisierte Gleitpunktdarstellung*, d.h. der Exponent wird so gewählt, daß die erste Ziffer der Mantisse von Null verschieden ist. Das bedeutet im Dualsystem $|a| \geq 2^{-1}$, im Dezimalsystem $|a| \geq 10^{-1}$ und allgemein $|a| \geq p^{-1}$, wobei p die Basis des Zahlensystems ist. Eine Zahl in binärer normalisierter Gleitkommadarstellung hat somit folgendes Aussehen:

$$x = \pm (\alpha_{-1} 2^{-1} + \dots + \alpha_{-t} 2^{-t}) \cdot 2^{\pm (\beta_0 2^0 + \dots + \beta_{e-1} 2^{e-1})},$$

$$\begin{aligned} \alpha_{-1} &= 1 & \alpha_{-i} &= 0 \text{ oder } 1 \text{ für } i = 2, 3, \dots, t \\ \beta_i &= 0 \text{ oder } 1 \text{ für } i = 0, 1, \dots, e-1 \end{aligned}$$

Offensichtlich kann die Zahl 0 nicht in normalisierter Gleitkommadarstellung dargestellt werden und besitzt daher eine Sonderstellung.

Aus alledem ergeben sich zunächst folgende Folgerungen:

- 1) Zur Durchführung eines Algorithmus stehen nur endlich viele Zahlen einer bestimmten Bauart zur Verfügung, sogenannte "Maschinenzahlen".
- 2) Es gibt keine beliebig großen und beliebig kleinen Zahlen und es gibt keine beliebig benachbarten Zahlen.
- 3) Die Körperaxiome sind verletzt. Summe, Differenz, Produkt und Quotient zweier Maschinenzahlen sind i.a. keine Maschinenzahl; Assoziativgesetze und das Distributivgesetz sind verletzt.

Will man eine Zahl, die nicht im Bereich der Maschinenzahlen liegt, durch eine Maschinenzahl darstellen, muß somit ab der Mantissenstelle $\alpha_{-(t+1)}$... gerundet oder abgeschnitten werden.

Der zur Verfügung stehende Bereich von Maschinenzahlen ist somit durch die Basis p , die Länge der Mantisse t und die Länge des Exponenten e eindeutig festgelegt und wird mit $M(p, t, e)$ bezeichnet. Ist zusätzlich noch festgelegt, ob nach der letzten Stelle der Mantisse abgeschnitten oder gerundet werden soll, so spricht man von einer *Maschinenarithmetik*. Eine Maschinenarithmetik mit runden bezeichnen wir mit $[M(p, t, e), \tilde{rd}]$, eine Maschinenarithmetik mit abschneiden bezeichnen wir mit $[M(p, t, e), \hat{rd}]$.

Berücksichtigt man den Exponenten nicht, so kann der absolute Fehler beim Runden der Mantisse somit abgeschätzt werden durch

$$|\text{rd}(a) - a| \leq \frac{1}{2} p^{-t} .$$

Der relative Fehler beim Runden auf Zahlen in normalisierter Gleitkommadarstellung kann (unabhängig vom Exponenten, der sich wegekürzt!) somit abgeschätzt werden durch

$$\left| \frac{\text{rd}(x) - x}{x} \right| \leq \frac{1}{2} p^{-t+1} =: \text{eps} .$$

"eps" nennt man die *Maschinengenauigkeit*, und sie kann auch als die kleinste positive Maschinenzahl g definiert werden, für die gilt

$$\text{rd}(1 + g) > 1 .$$

Wird somit eine Zahl x auf eine Maschinenzahl gerundet, so gilt

$$\text{rd}(x) = x(1 + \epsilon) \quad \text{mit} \quad |\epsilon| \leq \text{eps} ,$$

wobei ϵ der relative Fehler ist, der beim Runden entsteht.

Bsp.:

Für $[M(2,27,7), \tilde{r}d]$ ist $\text{eps} = \frac{1}{2} \cdot 2^{-26} \approx 7,45 \cdot 10^{-9}$

Für $[M(2,48,10), \tilde{r}d]$ ist $\text{eps} = \frac{1}{2} \cdot 2^{-47} \approx 3,55 \cdot 10^{-15}$

Offensichtlich ist es sinnlos, auf Maschinen mit solchen Maschinengenauigkeiten kleinere relative Fehler zu verlangen. Zwei Beispiele sollen nun binäre Maschinenarithmetiken erläutern.

Bsp.: Welche Zahlen sind in der (unrealistisch einfachen) Maschinenarithmetik $[M(2,2,1), \tilde{r}d]$ exakt darstellbar?

Als Mantisse kommen in normalisierter Gleitkommadarstellung nur ± 0.10 und ± 0.11 in Frage, als Exponent nur ± 1 oder 0 . Symmetrisch zum Nullpunkt sind daher 12 Zahlen exakt darstellbar:

$\pm 0,25, \pm 0.375, \pm 0.5, \pm 0.75, \pm 1, \pm 1.5$.

Zusätzlich muß noch die Zahl 0 dargestellt werden.

Bsp.: Man untersuche die Maschinenarithmetik $[M(2,t,e), \tilde{r}d]$.

Die größte Maschinenzahl ist

$$(2^{-1} + \dots + 2^{-t}) \cdot 2^{2^0 + \dots + 2^{e-1}} = (1 - 2^{-t}) 2^{2^e - 1}$$

Die zweitgrößte Maschinenzahl ist

$$(2^{-1} + \dots + 2^{-t+1}) \cdot 2^{2^0 + \dots + 2^{e-1}} = (1 - 2^{-t+1}) 2^{2^e - 1}$$

Die kleinste positive Maschinenzahl ist

$$2^{-1} \cdot 2^{-(2^0 + \dots + 2^{e-1})} = 2^{-2^e}$$

Die zweitkleinste positive Maschinenzahl ist

$$(2^{-1} + 2^{-t}) 2^{-(2^0 + \dots + 2^{e-1})} = (2^{-1} + 2^{-t}) 2^{1-2^e}.$$

Die halbe Differenz der beiden größten Maschinenzahlen ergibt (im Bereich der Maschinenzahlen) den größten absoluten Fehler beim Runden, nämlich 2^{2^e-2-t} . Die halbe Differenz der beiden kleinsten Maschinenzahlen ergibt (im Bereich der Maschinenzahlen) den kleinsten absoluten Fehler beim Runden, nämlich 2^{-t-2^e} . Während also der absolute Fehler enorm variiert, gilt für den relativen Fehler (im Bereich der Maschinenzahlen), wie schon besprochen, die obere Schranke $\epsilon_s = 2^{-t}$. Schließlich läßt sich auch die Anzahl der Maschinenzahlen (beiderlei Vorzeichens und ohne Null) leicht berechnen; sie ist gegeben durch $2^t(2^{e+1} - 1)$.

Soll eine Zahl berechnet werden, die betragsgrößer als die größte Maschinenzahl ist, so tritt ein sogenannter *Exponentenüberlauf* auf, der i.a. als Fehler gemeldet wird. Soll eine Zahl berechnet werden, die betragskleiner als die betragskleinste Maschinenzahl ist, tritt also ein *Exponentenunterlauf* auf, so wird die Zahl meist gleich Null gesetzt. Der relative Fehler kann dann größer als ϵ_s sein, da ja für die Zahl Null keine normalisierte Gleitkommadarstellung existiert!

Für arithmetische Verknüpfungen bzw. Funktionsauswertungen gilt also, wenn man sie als Gleitpunktoperationen betrachtet folgendes:

$$\left. \begin{aligned} \text{rd}(x + y) &=: x +^* y = (x + y)(1 + \epsilon_1) \\ \text{rd}(x - y) &=: x -^* y = (x - y)(1 + \epsilon_2) \\ \text{rd}(x \cdot y) &=: x \cdot^* y = (x \cdot y)(1 + \epsilon_3) \\ \text{rd}(x / y) &=: x /^* y = (x / y)(1 + \epsilon_4) \\ \text{rd}(\sqrt{x}) &=: \sqrt{x}^* = \sqrt{x}(1 + \epsilon_5) \\ \text{rd}(\sin x) &=: \sin x^* = (\sin x)(1 + \epsilon_6) \end{aligned} \right\} |\epsilon_i| \leq \epsilon_s$$

u. s. w.

Werden zwei Zahlen mit verschiedenen Exponenten in einer Maschinarithmetik addiert, so muß der Exponent der kleineren Zahl dem Exponenten der größeren Zahl angeglichen werden.

Bsp.: [M[10,6,3],rd] :

$$\begin{array}{rcl}
 0.999727 \cdot 10^2 & \text{Angleichung} & 0.999727 \cdot 10^2 \\
 + 0.854621 \cdot 10^{-1} & \text{d. Exponenten} & + 0.000855 \cdot 10^2 \\
 \hline
 & \longrightarrow & \hline
 & & 1.000582 \cdot 10^2 \\
 \\
 \text{Normalisieren} & & \\
 \text{d. Ergebnisses} & \longrightarrow & 0.100058 \cdot 10^3
 \end{array}$$

Man beachte, daß es bei der Subtraktion von nahezu gleich großen Zahlen zu einer Verringerung der gültigen Stellen der Mantisse kommt, man nennt diesen Effekt *Auslöschung*.

Bsp.: [M[10,7,2],rd] :

$$\begin{array}{rcl}
 0.3141592 \cdot 10^4 & & \\
 - 0.3141333 \cdot 10^4 & \text{normalisieren} & \\
 \hline
 0.0000259 \cdot 10^4 & \longrightarrow & 0.2590000 \\
 & & \underbrace{\hspace{2cm}} \\
 & & \text{ungültige Stellen}
 \end{array}$$

Man beachte auch, daß die Analyse der Rundungsfehler von der Reihenfolge abhängt, in der die arithmetischen Operationen durchgeführt werden. Diese Reihenfolge kann durch Klammern, durch einen geordneten Baum oder durch Angabe des Algorithmus spezifiziert werden.

Bsp.: $\frac{a^2 - \sin b}{c}$

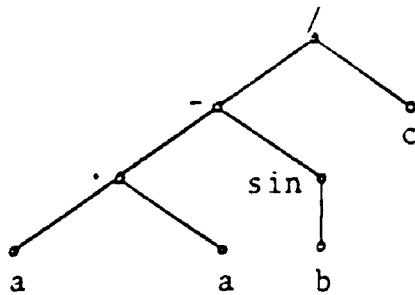
Angabe der Reihenfolge der Operationen durch Klammern:

$$((a \cdot a) - \sin b) / c$$

bzw. in Maschinarithmetik

$$((a \cdot a) - \sin b) / c$$

Angabe der Reihenfolge der Operationen durch einen geordneten Baum:



Angabe der Reihenfolge der Operationen durch Angabe des Algorithmus

$$\begin{aligned} n_1 &= a \cdot a \\ n_2 &= \sin b \\ n_3 &= n_1 - n_2 \\ n_4 &= n_3 / c \end{aligned}$$

bzw. in Maschinearithmetik

$$\begin{aligned} \tilde{n}_1 &= a \cdot a(1 + \epsilon_1) \\ \tilde{n}_2 &= \sin b(1 + \epsilon_2) \\ \tilde{n}_3 &= (\tilde{n}_1 - \tilde{n}_2)(1 + \epsilon_3) \\ \tilde{n}_4 &= (\tilde{n}_3 / c)(1 + \epsilon_4) \end{aligned}$$

Entwickelt man in erster Ordnung in den ϵ_i , so muß das Ergebnis folgende Gestalt haben:

$$\text{exakte Formel} \cdot (1 + \text{rel. Fehler})$$

wobei der relative Fehler linear in den ϵ_i sein muß.

Rundungsfehleranalyse

$$\begin{aligned}\tilde{n}_4 &= ((a \cdot a(1 + \epsilon_1) - \sin b(1 + \epsilon_2))(1 + \epsilon_3)/c)(1 + \epsilon_4) \\ &= (a^2 - \sin b) \left(1 + \frac{a^2}{a^2 - \sin b} \epsilon_1 - \frac{\sin b}{a^2 - \sin b} \epsilon_2\right) (1 + \epsilon_3)(1 + \epsilon_4)/c \\ &\approx \frac{a^2 - \sin b}{c} \left(1 + \frac{a^2}{a^2 - \sin b} \epsilon_1 - \frac{\sin b}{a^2 - \sin b} \epsilon_2 + \epsilon_3 + \epsilon_4\right) \\ & \qquad \qquad \qquad |\epsilon_i| \leq \text{eps}\end{aligned}$$

Setzen wir etwa die Arithmetik $[M(10,6,3), \tilde{rd}]$ voraus, so ist

$$\text{eps} = \frac{1}{2} 10^{-3} .$$

Der relative Fehler kann somit abgeschätzt werden durch

$$\left| \frac{a^2}{a^2 - \sin b} \right| \text{eps} + \left| \frac{\sin b}{a^2 - \sin b} \right| \text{eps} + 2\text{eps} .$$

Kritisch ist somit der Fall, wo

$$a^2 \approx \sin b .$$

Man beachte, daß bei einer Funktionsauswertung der entstehende Rundungsfehler unabhängig von der Funktion ist.

Bsp.: Die Zahlen

$$a = 0,1732051$$

$$b = 0,3141592 \cdot 10^4$$

$$c = -0,3141333 \cdot 10^4$$

sollen in der Arithmetik $[M(10,7,2), \tilde{rd}]$ addiert werden und zwar nach folgenden Algorithmen:

$$\alpha) (a + b) + c$$

$$\beta) a + (b + c)$$

$$\begin{array}{r} \alpha) \quad 0,3141592 \cdot 10^4 \\ \quad 0,0000173 \cdot 10^4 \\ \hline \quad 0,3141765 \cdot 10^4 \\ -0,3141333 \cdot 10^4 \\ \hline \quad 0,0000432 \cdot 10^4 \end{array}$$

Das normalisierte Ergebnis ist somit
 $(a + *b) + *c = 0,4320000$.

$$\begin{array}{r} \beta) \quad 0,3141592 \cdot 10^4 \\ -0,3141333 \cdot 10^4 \\ \hline \quad 0,0000259 \cdot 10^4 \\ \\ \quad 0,1732051 \\ \quad 0,2590000 \\ \hline \quad 0,4322051 \end{array}$$

Das normalisierte Ergebnis ist somit
 $a + * (b + *c) = 0,4322051$.

Offensichtlich ist das Assoziativitätsgesetz verletzt, das Ergebnis im Fall β) ist wesentlich genauer als im Fall α), wie sich auch durch Rundungsfehleranalyse zeigen läßt.

Rundungsfehleranalyse für α):

$$\begin{aligned} (a + *b) + *c &= ((a + b)(1 + \epsilon_1) + c)(1 + \epsilon_2) \\ &\approx (a + b + c) \left(1 + \frac{a+b}{a+b+c} \epsilon_1 + \epsilon_2 \right) \end{aligned}$$

$$\frac{a+b}{a+b+c} \approx \frac{0,31 \cdot 10^4}{0,43} \approx 0,72 \cdot 10^4$$

Für diese Maschinendarithmetik ist $\epsilon_{ps} = \frac{1}{2} 10^{-6}$ und somit kann der relative Fehler abgeschätzt werden durch

$$0,72 \cdot 10^4 \cdot \frac{1}{2} \cdot 10^{-6} + \frac{1}{2} \cdot 10^{-6} \approx 0,36 \cdot 10^{-2}$$

und der absolute Fehler durch

$$0,36 \cdot 10^{-2} \cdot 0,43 \approx 0,0015 .$$

Daraus folgt, daß nicht einmal die dritte Stelle nach dem Komma garantiert werden kann.

Rundungsfehleranalyse für β):

$$\begin{aligned} (a + (b + c)) &= (a + (b + c)(1 + \varepsilon_1))(1 + \varepsilon_2) \\ &\approx (a + b + c) \left(1 + \frac{b+c}{a+b+c} \varepsilon_1 + \varepsilon_2 \right) \end{aligned}$$

$$\frac{b+c}{a+b+c} \approx \frac{0,26}{0,43} \approx 0,6 .$$

Somit kann man den relativen Fehler abschätzen durch

$$0,6 \cdot \frac{1}{2} 10^{-6} + \frac{1}{2} 10^{-6} = 0,8 \cdot 10^{-6}$$

und den absoluten Fehler durch

$$0,8 \cdot 10^{-6} \cdot 0,43 \approx 0,34 \cdot 10^{-6} .$$

Eine Abweichung vom "exakten" Resultat kann also höchstens in der letzten Stelle des Ergebnisses erwartet werden.

Bsp.: Nun soll $a^2 - b^2$ berechnet werden und zwar nach folgenden Algorithmen

α) $(a \cdot a) - (b \cdot b)$

β) $(a + b) \cdot (a - b)$

Rundungsfehleranalyse für α):

$$\begin{aligned} ((a \cdot a)(1 + \varepsilon_1) - (b \cdot b)(1 + \varepsilon_2))(1 + \varepsilon_3) \\ \approx (a^2 - b^2) \left(1 + \frac{a^2}{a^2 - b^2} \varepsilon_1 - \frac{b^2}{a^2 - b^2} \varepsilon_2 + \varepsilon_3 \right) \end{aligned}$$

Der Algorithmus ist somit "gefährlich", falls $a^2 \approx b^2$ (Auslöschung!). Der gesamte absolute Rundungsfehlereinfluß kann abgeschätzt werden durch

$$|a^2 \varepsilon_1 - b^2 \varepsilon_2 + (a^2 - b^2) \varepsilon_3| \leq (a^2 + b^2 + |a^2 - b^2|) \text{eps} .$$

Rundungsfehleranalyse für β):

$$((a + b)(1 + \varepsilon_1) \cdot (a - b)(1 + \varepsilon_2))(1 + \varepsilon_3)$$

$$\approx (a^2 - b^2)(1 + \varepsilon_1 + \varepsilon_2 + \varepsilon_3)$$

Der gesamte absolute Rundungsfehlereinfluß kann somit abgeschätzt werden durch

$$|(a^2 - b^2)(\varepsilon_1 + \varepsilon_2 + \varepsilon_3)| \leq 3\text{eps} \cdot |a^2 - b^2| .$$

Vergleicht man nun die beidem Rundungsfehler, so erkennt man, daß für $\frac{1}{3} \leq \left| \frac{a}{b} \right|^2 \leq 3$ Algorithmus β) "stabiler" ist, andernfalls aber Algorithmus α).

Bsp.: Nun soll die Wurzel $y = -p + \sqrt{p^2 + q}$ der quadratischen Gleichung $y^2 + 2py - q = 0$ berechnet werden und zwar nach folgendem Algorithmus:

$$n_1 = p \cdot p$$

$$n_2 = n_1 + q$$

$$n_3 = \sqrt{n_2}$$

$$n_4 = -p + n_3$$

bzw.

$$\tilde{n}_1 = p \cdot p(1 + \varepsilon_1)$$

$$\tilde{n}_2 = (\tilde{n}_1 + q)(1 + \varepsilon_2)$$

$$\tilde{n}_3 = \sqrt{\tilde{n}_2} (1 + \varepsilon_3)$$

$$\tilde{n}_4 = (-p + \tilde{n}_3)(1 + \varepsilon_4)$$

$$\tilde{n}_4 = (-p + \sqrt{(p^2(1 + \varepsilon_1) + q)(1 + \varepsilon_2)(1 + \varepsilon_3)})(1 + \varepsilon_4)$$

Entwickelt man in erster Ordnung in ε_i , so erhält man

$$\tilde{n}_4 \approx y \left(1 - \frac{p}{y} \varepsilon_4 + \frac{\sqrt{p^2+q}}{y} \left(\frac{1}{2} \frac{p^2}{p^2+q} \varepsilon_1 + \frac{1}{2} \varepsilon_2 + \varepsilon_3 + \varepsilon_4 \right) \right)$$

Gefährlich ist offensichtlich jener Fall, wo $p > 0$ und $q \approx 0$ ist, weil dann die Koeffizienten der ε_i sehr groß werden können. Um den Algorithmus für diesen Fall "stabil" zu gestalten, kann man so vorgehen, daß man die zweite Wurzel der quadratischen Gleichung $y_2 = -p - \sqrt{p^2 + q}$ berechnet und dann aus der Relation $y_1 \cdot y_2 = -q$ das gewünschte y_1 .

Bsp.: Man berechne $d = \sqrt{a^2 + b^2 - 2ab \cos \gamma}$ nach folgendem Algorithmus:

$$\begin{aligned} \tilde{n}_1 &= a \cdot a(1 + \varepsilon_1) & \tilde{n}_6 &= \tilde{n}_4 \cdot \tilde{n}_5(1 + \varepsilon_6) \\ \tilde{n}_2 &= b \cdot b(1 + \varepsilon_2) & \tilde{n}_7 &= (\tilde{n}_1 + \tilde{n}_2)(1 + \varepsilon_7) \\ \tilde{n}_3 &= a \cdot b(1 + \varepsilon_3) & \tilde{n}_8 &= (\tilde{n}_7 - \tilde{n}_6)(1 + \varepsilon_8) \\ \tilde{n}_4 &= 2 \cdot \tilde{n}_3(1 + \varepsilon_4) & \tilde{n}_9 &= \sqrt{\tilde{n}_8} (1 + \varepsilon_9) \\ \tilde{n}_5 &= \cos \gamma (1 + \varepsilon_5) \end{aligned}$$

Rundungsfehleranalyse

$$\begin{aligned} \tilde{n}_9 &= [((a^2(1+\varepsilon_1) + b^2(1+\varepsilon_2))(1 + \varepsilon_7) - \\ &\quad - 2ab(1+\varepsilon_3)(1+\varepsilon_4)\cos \gamma (1+\varepsilon_5)(1+\varepsilon_6))(1+\varepsilon_8)]^{\frac{1}{2}}(1+\varepsilon_9) \\ &\approx d \left(1 + \frac{a^2}{2d^2} (\varepsilon_1 + \varepsilon_7 + \varepsilon_8) + \right. \\ &\quad \left. + \frac{b^2}{2d^2} (\varepsilon_2 + \varepsilon_7 + \varepsilon_8) - \frac{ab \cos \gamma}{d^2} (\varepsilon_3 + \varepsilon_4 + \varepsilon_5 + \varepsilon_6 + \varepsilon_8) + \varepsilon_9 \right) \end{aligned}$$

Obwohl unsere Darstellung der Maschinenzahlen die wesentlichen Charakteristika trifft, weicht aber doch die Realisierung auf Maschinen gelegentlich etwas davon ab.

III. FORTPFLANZUNG DER DATENFEHLER

Während bei der Analyse der Rundungsfehler bei jeder arithmetischen Operation und bei Funktionsauswertungen, die durchgeführt werden, ein relativer Fehler zu berücksichtigen ist, wird nun die Auswirkung der relativen Fehler sämtlicher (oder einiger) Eingangsdaten untersucht.

Es seien ϵ_x und ϵ_y die relativen Fehler von x und y ($x, y \neq 0$), d.h.

$$\epsilon_x = \frac{\tilde{x} - x}{x}, \quad \epsilon_y = \frac{\tilde{y} - y}{y}$$

oder

$$\tilde{x} = x(1 + \epsilon_x), \quad \tilde{y} = y(1 + \epsilon_y).$$

Wir untersuchen, wie sich diese relativen Fehler auf den relativen Fehler des Ergebnisses auswirken, wenn x und y durch arithmetische Operationen verknüpft werden, oder wenn Funktionen dieser Eingangsdaten gebildet werden.

Es gelten die folgenden

Fehlerfortpflanzungsformeln 1. Ordnung:

1. $x \cdot y$: $\epsilon_{xy} \approx \epsilon_x + \epsilon_y$
2. x / y : $\epsilon_{x/y} \approx \epsilon_x - \epsilon_y$
3. $x \pm y$: $\epsilon_{x \pm y} = \frac{x}{x \pm y} \epsilon_x \pm \frac{y}{x \pm y} \epsilon_y$ falls $(x \pm y) \neq 0$
4. x^n : $\epsilon_{x^n} \approx n \epsilon_x$
5. $\sqrt[n]{x}$: $\epsilon_{\sqrt[n]{x}} \approx \frac{1}{n} \epsilon_x$

6. e^x : $\epsilon_{e^x} \approx x \epsilon_x$
 7. $\ln x$: $\epsilon_{\ln x} \approx \frac{1}{\ln x} \epsilon_x$
 8. $\sin x$: $\epsilon_{\sin x} \approx (x \operatorname{ctg} x) \epsilon_x$
 9. $\cos x$: $\epsilon_{\cos x} \approx (-x \operatorname{tg} x) \epsilon_x$
 10. $\operatorname{tg} x$: $\epsilon_{\operatorname{tg} x} \approx \frac{2x}{\sin 2x} \epsilon_x$

Beweis von 2:

$$\begin{aligned} \frac{\tilde{x}}{\tilde{y}} &= \frac{x(1+\epsilon_x)}{y(1+\epsilon_y)} = \frac{x}{y} (1+\epsilon_x)(1-\epsilon_y+\epsilon_y^2-\dots) \\ &= \frac{x}{y} \left(1 + \underbrace{\epsilon_x - \epsilon_y}_{\substack{\text{Glieder höherer Ordnung} \\ \text{(i.a. sehr klein)}}} + \dots \right) \end{aligned}$$

Näherung für den relativen Fehler des Ergebnisses;

Beweis von 5:

$$\begin{aligned} \sqrt[n]{\tilde{x}} &= \sqrt[n]{x(1+\epsilon_x)} = \sqrt[n]{x} \sqrt[n]{1+\epsilon_x} \\ &= \sqrt[n]{x} \left(1 + \frac{1}{n} \epsilon_x + \underbrace{\frac{1}{n} \left(\frac{1}{n} - 1 \right) \frac{1}{2} \epsilon_x^2 + \dots}_{\substack{\text{Glieder höherer Ordnung} \\ \text{(vernachlässigt)}}} \right) \end{aligned}$$

Beweis von 6:

$$e^{\tilde{x}} = e^{x(1+\epsilon_x)} = e^x e^{x\epsilon_x} = e^x \left(1 + \frac{x\epsilon_x}{1!} + \frac{(x\epsilon_x)^2}{2!} + \dots \right) \approx e^x (1 + x\epsilon_x)$$

ad 3: Die Addition ist gefährlich, wenn $x \approx -y$. Die Faktoren $\frac{x}{x+y}$, $\frac{y}{x+y}$ heißen Verstärkerfaktoren oder *Konditionszahlen*.

Allgemein gilt (nach der Taylorformel):

Fehlerfortpflanzungsformel 1. Ordnung in einer Variablen:

$$\begin{aligned} f(\tilde{x}) &= f(x(1 + \epsilon_x)) \approx f(x) + x\epsilon_x f'(x) \\ &= f(x) \left(1 + \underbrace{x \frac{f'(x)}{f(x)} \epsilon_x}_{\epsilon_{f(x)}} \right) = f(x) (1 + \epsilon_{f(x)}) \end{aligned}$$

Fehlerfortpflanzungsformel 1. Ordnung in zwei Variablen:

$$\begin{aligned} f(\tilde{x}, \tilde{y}) &= f(x(1 + \epsilon_x), y(1 + \epsilon_y)) \approx f(x, y) + x\epsilon_x \frac{\partial f}{\partial x} + y\epsilon_y \frac{\partial f}{\partial y} \\ &= f(x, y) \left(1 + \underbrace{\frac{x}{f(x, y)} \frac{\partial f}{\partial x} \epsilon_x + \frac{y}{f(x, y)} \frac{\partial f}{\partial y} \epsilon_y}_{\epsilon_{f(x, y)}} \right) \\ &= f(x, y) (1 + \epsilon_{f(x, y)}) \end{aligned}$$

Nun soll für jene Beispiele, die wir einer Rundungsfehleranalyse unterzogen haben, auch eine Datenfehleranalyse durchgeführt werden.

Bsp.: $\frac{a^2 - \sin b}{c}$

Datenfehleranalyse

$$\begin{aligned} &\frac{a^2(1 + \epsilon_a)^2 - \sin b(1 + (b \operatorname{ctg} b)\epsilon_b)}{c(1 + \epsilon_c)} \approx \\ &\approx \frac{a^2 - \sin b}{c} \left(1 + \frac{2a^2}{a^2 - \sin b} \epsilon_a - \frac{b \cos b}{a^2 - \sin b} \epsilon_b \right) \frac{1}{1 + \epsilon_c} \\ &\approx \frac{a^2 - \sin b}{c} \left(1 + \frac{2a^2}{a^2 - \sin b} \epsilon_a - \frac{b \cos b}{a^2 - \sin b} \epsilon_b - \epsilon_c \right) \end{aligned}$$

Dasselbe Ergebnis erhält man durch direkte Anwendung der analogen Fehlerfortpflanzungsformel 1. Ordnung in drei Variablen. Kritisch ist wieder der Fall $a^2 \approx \sin b$.

Man beachte, daß die Datenfehleranalyse unabhängig von der Reihenfolge der Auswertung ist.

Bsp.: $a + b + c$

Datenfehleranalyse

$$\begin{aligned} & a(1 + \epsilon_a) + b(1 + \epsilon_b) + c(1 + \epsilon_c) = \\ & = (a + b + c) \left(1 + \frac{a}{a + b + c} \epsilon_a + \frac{b}{a + b + c} \epsilon_b + \frac{c}{a + b + c} \epsilon_c \right) \end{aligned}$$

Bsp.: $a^2 - b^2$

Datenfehleranalyse

$$a^2(1 + \epsilon_a)^2 - b^2(1 + \epsilon_b)^2 \approx (a^2 - b^2) \left(1 + \frac{2a^2}{a^2 - b^2} \epsilon_a - \frac{2b^2}{a^2 - b^2} \epsilon_b \right)$$

Bsp.: $y = -p + \sqrt{p^2 + q}$

Datenfehleranalyse

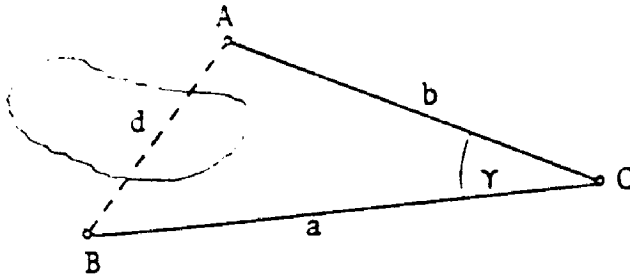
$$\begin{aligned} & -p(1 + \epsilon_p) + \sqrt{p^2(1 + \epsilon_p)^2 + q(1 + \epsilon_q)} \approx \\ & \approx (-p + \sqrt{p^2 + q}) \left(1 - \frac{p}{\sqrt{p^2 + q}} \epsilon_p + \frac{p + \sqrt{p^2 + q}}{2\sqrt{p^2 + q}} \epsilon_q \right) \end{aligned}$$

Bsp.: $d = \sqrt{a^2 + b^2 - 2ab \cos \gamma}$

Datenfehleranalyse

$$\begin{aligned} & \sqrt{a^2(1 + \epsilon_a)^2 + b^2(1 + \epsilon_b)^2 - 2a(1 + \epsilon_a)b(1 + \epsilon_b) \cos \gamma (1 - \gamma \epsilon_\gamma \tan \gamma)} \approx \\ & \approx d \left(1 + \frac{a^2 - ab \cos \gamma}{d^2} \epsilon_a + \frac{b^2 - ab \cos \gamma}{d^2} \epsilon_b + \frac{ab \gamma \sin \gamma}{d^2} \epsilon_\gamma \right) \end{aligned}$$

Diese Aufgabe kann etwa in folgenden Text eingekleidet werden:
 Der Abstand d zwischen den Punkten A und B kann wegen eines Hindernisses nicht gemessen werden. Statt dessen werden die Distanzen a und b zu einem Punkt C und der Winkel γ gemessen.



Es gilt dann $d = \sqrt{a^2 + b^2 - 2ab \cos \gamma}$.

a) Man bestimme die absoluten Konditionszahlen K_a , K_b und K_γ für die Bestimmung von d bezüglich der absoluten Fehler in a, b und γ . Man zeige, daß $|K_a|$ und $|K_b|$ stets kleiner gleich 1 sind.

b) Bei einer Messung ergab sich

$$a = 74,83 \pm 0,02 \text{ m} \quad b = 53,19 \pm 0,02 \text{ m}$$

$$\gamma = 35^\circ 12' \pm 3' .$$

Wie stark können sich diese Meßungenauigkeiten auf den berechneten Abstand d auswirken?

Nach der obigen Datenfehleranalyse gilt

$$\Delta d = \underbrace{\frac{a-b \cos \gamma}{d}}_{K_a} \underbrace{(\Delta a)}_{\Delta a} + \underbrace{\frac{b-a \cos \gamma}{d}}_{K_b} \underbrace{(\Delta b)}_{\Delta b} + \underbrace{\frac{ab \sin \gamma}{d}}_{K_\gamma} \underbrace{(\Delta \gamma)}_{\Delta \gamma}$$

$$|K_a| = \frac{|a-b \cos \gamma|}{\sqrt{a^2+b^2-2ab \cos \gamma}} = \sqrt{\frac{a^2+b^2 \cos^2 \gamma - 2ab \cos \gamma}{a^2+b^2-2ab \cos \gamma}} \leq 1$$

$$|\Delta d| \leq 1 \cdot |\Delta a| + 1 \cdot |\Delta b| + \frac{a_{\max} b_{\max} \sin \gamma_{\max}}{\sqrt{a_{\min}^2 + b_{\min}^2 - 2a_{\max} b_{\max} \cos \gamma_{\min}}} |\Delta \gamma|$$

Durch einsetzen erhält man $|\Delta d| \leq 0,09 \text{ m}$.

IV. SCHLUSSFOLGERUNGEN

Wir haben in unseren vorgehenden Betrachtungen die Analyse der Datenfehler und der Rundungsfehler streng getrennt. Dies sollte aus didaktischen Erwägungen stets getan werden, weil doch die Behandlung dieser beiden Fehlertypen durchaus verschieden ist. Natürlich treten bei der Auswertung jeder Formel i. a. sowohl Daten- wie auch Rundungsfehler auf, aber es ist leicht einzusehen, daß sich der relative Gesamtfehler als Summe des relativen Datenfehlers und des relativen Rundungsfehlers darstellen läßt.

Vorbedingung ist natürlich die Fähigkeit, eine "gestörte" Formel in erster Ordnung zu entwickeln. Mit der im III. Abschnitt angegebenen Formeltabelle ist dies in der höheren Schule durchaus möglich und man braucht dann weder auf die Formel von Taylor noch auf die Reihenentwicklungen der elementaren transzendenten Funktionen zurückzugreifen.

Jedenfalls rücken Analysen dieser Art die Mathematik in die Nähe des "realistischen Problems", sie schärfen den kritischen Verstand und es kommt ihnen ein bedeutender Bildungswert zu.

Weiterführende Literatur:

Stetter, H.J.: Numerik für Informatiker:
Computergerechte numerische Verfahren
R. Oldenbourg Verlag, München 1976.

Stoer, J.: Einführung in die Numerische Mathematik I
Springer Verlag, Berlin, 1972.