

Numerische Mathematik I

Prof. Dr. Christof Büskens

AG Optimierung & Optimale Steuerung
Zentrum für Technomathematik
Universität Bremen
28334 Bremen, Germany

Vorlesungsskript
Sommersemester 2004

(Unkorrigierte Fassung)

Vorwort

Die vorliegende Ausarbeitung entstand während meiner Tätigkeit am Zentrum für Technomathematik der Universität Bremen. Sie entstand im Rahmen einer Vorlesung, die ich im Sommersemester 2004 gehalten habe. An dieser Stelle möchte ich mich bei allen Teilnehmerinnen und Teilnehmern für ihr reges Interesse und ihre aktive Mitarbeit bedanken.

Bremen, Juli 2004

CHRISTOF BÜSKENS

Inhaltsverzeichnis

Inhaltsverzeichnis	5
1 Einleitung	9
1.1 Einführung	9
1.2 Literatur	11
1.3 Ein kurzer geschichtlicher Rückblick	11
1.4 Was ist Numerik?	13
1.5 Motivationsbeispiel	13
1.6 Vorläufiges Fazit	15
2 Fehleranalyse	17
2.1 Maschinenzahlen	17
2.2 Maschinenzahlen auf der Zahlengerade	19
2.3 Rundung	20
2.4 Gleitpunkt-Arithmetik	23
2.5 Fehlerfortpflanzung, Kondition	24
2.6 Algorithmen	26
3 Lineare Gleichungssysteme	29
3.1 Einführung und Aufgabenstellung	29
3.2 <i>LR</i> -Zerlegung und Gauß-Elimination	31
3.2.1 Idee der Gauß-Elimination/ <i>LR</i> -Zerlegung	31
3.2.2 Frobeniusmatrizen	32

3.2.3	Gauß–Elimination/LR–Zerlegung ohne Pivoting	34
3.2.4	Permutationsmatrizen	36
3.2.5	Gauß–Elimination/LR–Zerlegung mit Pivoting	37
3.2.6	Aufwandsbestimmung	41
3.2.7	Algorithmus	42
3.3	Matrizen mit speziellen Eigenschaften	43
3.3.1	Diagonaldominante Matrizen: Diagonalstrategie	43
3.3.2	Positiv definite Matrizen: Cholesky–Verfahren	45
3.3.3	Bandmatrizen: Bandausnutzende Verfahren	49
3.4	Fehleranalyse und Fehlerbehandlung	51
3.4.1	Fehlerabschätzungen	51
3.4.2	Skalierung	54
3.4.3	Iterative Nachverbesserung	55
3.5	Die QR–Zerlegung einer Matrix, das Verfahren von Householder	56
3.5.1	Einleitung und Motivation	56
3.5.2	Householdermatrizen	57
3.5.3	QR–Zerlegung/Verfahren von Householder	59
3.5.4	Erweiterungen	61
3.6	Lineare Ausgleichsrechnung, diskrete Approximation	62
3.6.1	Normalgleichung	62
3.6.2	Numerische Lösung	65
3.6.3	Diskrete Approximation	66
4	Nichtlineare Gleichungen und Gleichungssysteme	69
4.1	Einführung und Aufgabenstellung	69
4.2	Grundlagen	70
4.2.1	Fixpunkte	70
4.2.2	Konvergenz	71
4.3	Nichtlineare Gleichungen	72
4.3.1	Bisektionsverfahren	72
4.3.2	Newton–Verfahren	74
4.3.3	Sekanten–Verfahren	75
4.4	Konvergenz von Iterationsverfahren	77
4.4.1	Kontraktion	77
4.4.2	Fixpunktsatz von Banach	79
4.4.3	Konvergenzsätze	81
4.4.4	Konvergenz des Newton–Verfahrens	82
4.5	Das Newton–Verfahren im \mathbb{R}^n	83
4.5.1	Herleitung des Newton–Verfahrens	83

4.5.2	Praktische Realisierung	85
4.5.3	Newton–Kantorovich	86
4.5.4	Erweiterungen	88
4.5.4.1	Approximation von $f'(x)$ durch Differenzen . . .	88
4.5.4.2	λ -Strategie, Modifiziertes Newton-Verfahren . . .	89
5	Interpolation	91
5.1	Einführung und Aufgabenstellung	91
5.2	Polynominterpolation	92
5.2.1	Existenz und Eindeutigkeit der Polynominterpolation . . .	92
5.2.2	Interpolationsformel von Lagrange	93
5.2.3	Der Algorithmus von Aitken und Neville	94
5.2.3.1	Rekursionsformel von Aitken	94
5.2.3.2	Variante von Neville	94
5.2.4	Die Newton'sche Interpolationsformel, Dividierte Differenzen	95
5.2.5	Interpolationsfehler	98
5.2.6	Konvergenz	99
5.3	Trigonometrische Interpolation	100
5.3.1	Diskrete Fouriertransformation	100
5.3.2	Trigonometrische Interpolation	102
5.3.3	Schnelle Fourier–Transformation (FFT)	103
5.3.4	Anwendungen	105
5.4	Spline–Interpolation	105
5.4.1	Polynom–Splines	105
5.4.2	Kubische Splines	108
5.4.2.1	Einführung und Aufgabenstellung	108
5.4.2.2	Existenz und Eindeutigkeit	109
5.4.2.3	Geometrische und mechanische Interpretation . .	111
5.4.2.4	Die Berechnung von Spline-Funktionen	112
5.4.2.5	Konvergenzeigenschaften	115
5.5	Numerische Differentiation	118
6	Integration	121
6.1	Einführung und Aufgabenstellung	121
6.2	Newton–Cotes–Formeln	121
6.3	Zusammengesetzte Newton–Cotes–Formeln	124
6.3.1	Zusammengesetzte Trapezregel	125
6.3.2	Verfeinerung der zusammengesetzten Trapezregel	126

6.4	Die Gaußsche Integrationsmethode	127
6.4.1	Orthogonalpolynome	127
6.4.2	Gaußintegration	129
6.5	Integration und Extrapolation	132
6.5.1	Euler-Maclaurin'sche Summenformel	132
6.5.2	Anwendung der Extrapolation auf die Integration	133
6.5.3	Integrationsfehler	135

Literaturverzeichnis	139
-----------------------------	------------

Kapitel 1

Einleitung

1.1 Einführung

Gegenstand der numerischen Mathematik (oder einfach *Numerik*) oder auch praktischen Mathematik ist die näherungsweise Lösung mathematischer Probleme durch Zahlenwerte. Die Lösungsberechnung erfolgt dabei durch einen *Algorithmus*, d.h. durch eine Folge von elementaren Anweisungen und Rechenoperationen, die sich auf einem Computer ausführen lassen. Ein solcher Algorithmus stützt sich oft auf Ergebnisse der reinen Mathematik und reflektiert mathematische Eigenschaften des Problems. Die zu behandelnden Probleme stammen oft aus den Ingenieur- und Naturwissenschaften.

Beispiel 1.1. *Als ein erstes praktisches Beispiel sei der Landeanflug eines Verkehrsflugzeuges bei Scherwinden benannt, bei dem es zu 2–3 Unfällen pro Jahr kommt (bereits > 500 Tote), vgl. Abbildung 1.1.*

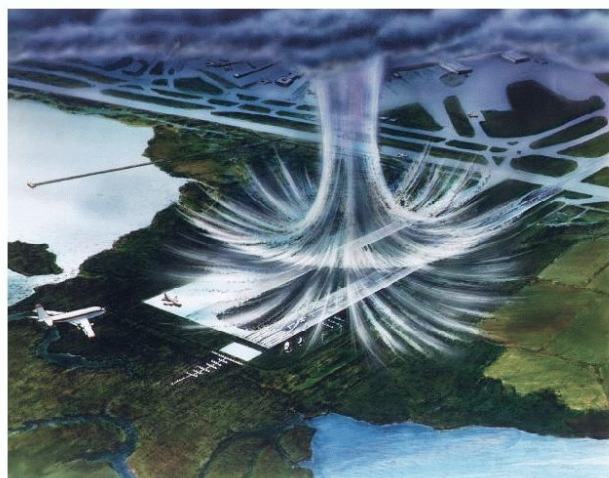


Abbildung 1.1: Scherwinde beim Landeanflug.

Aufgrund der Fallwinde wäre eine sichere Vorgehensweise, den Landeanflug abzurechnen, was aber ist hierzu die sicherste Vorgehensweise? Ein sehr sicherer Weg ist die während des Durchfluges durch den Scherwind angenommene minimale Höhe zu maximieren, vgl. Abbildung 1.2; wie aber kann das erreicht werden?

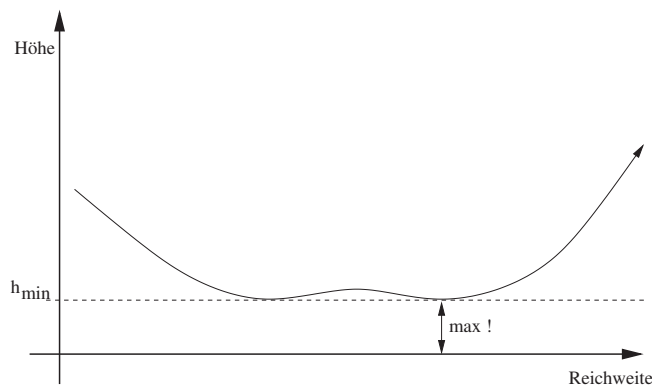


Abbildung 1.2: Maximierung der minimalen Höhe.

Da die physikalischen Vorgänge sehr gut bekannt sind kann zunächst ein sehr realitätsnahes mathematisches Modell erstellt werden.

Die Mathematik kommt dann intensiv bei der Lösung des Problems zur Anwendung. Hierzu muß zunächst eine theoretische Aufarbeitung der zu verwendenden Lösungsmechanismen vorgenommen werden, bzw. neu entwickelt werden. Für unser Beispiel greifen wir auf die sogenannte Variationsrechnung bzw. Optimale Steuerung zurück. Die hierzu angebotenen Lösungsmethoden sind jedoch nicht mehr analytisch auf unser Flugmodell anwendbar und wir werden eine numerische Lösung auf einem Computer bemühen müssen. Zur Anwendung kommen numerische Verfahren für Differentialgleichungen oder lineare und nichtlineare Gleichungssystemlöser.

Ziel der Veranstaltung ist die Einführung in verschiedenen Gebiete der numerischen Mathematik, wie z.B.:

- Lineare Gleichungssysteme,
- Interpolation,
- numerische Integration,
- nichtlineare Gleichungssysteme,
- Numerik der Differentialgleichungen.

Klassischerweise werden in den Vorlesungen Numerische Mathematik 1 und Numerische Mathematik 2 einfache Vorkenntnisse vermittelt, während in den Vorlesungen Numerische Mathematik 3, 4 (die häufig auch anders genannt werden) Spezialisierungen, Vertiefungen und Erweiterungen behandelt werden. Wichtig für alle Vorlesungen zur Numerik sind immer hinreichende Programmierkenntnisse!

1.2 Literatur

In der Numerik gibt es eine Fülle ausgezeichnete Bücher, die die verschiedenen angesprochenen Thematiken umfangreich beleuchten und darüberhinaus ergänzenden Stoff vermitteln. Stellvertretend für andere seien die nachfolgenden Bücher erwähnt:

- Deuffhard/Hohmann: Numerische Mathematik I, Verlag Walter de Gruyter
- Hämmerlin/Hoffmann: Numerische Mathematik, Springer Verlag
- Schwarz: Numerische Mathematik, Teubner Verlag
- Stoer: Numerische Mathematik I, Springer Verlag
- Stoer, Bulirsch: Numerische Mathematik II, Springer Verlag
- Werner: Numerische Mathematik, Vieweg Verlag
- u.v.a.

Es sei erneut erwähnt, dass es sich bei allen Büchern um ausgezeichnete Zusammenstellungen zur Numerik handelt. Das hier zusammengestellte Skript orientiert sich an mehreren Büchern und es ist daher keines im besonderen Maße hervorzuheben.

1.3 Ein kurzer geschichtlicher Rückblick

Ausgangspunkt für numerische Fragestellungen war eine Belebung der Mathematik durch konkrete Fragestellungen aus den Anwendungen. Nicht nur die Existenz, sondern auch die Bestimmung der Lösung, z.B. wie bei der Vorhersage von Himmelserscheinungen, traten in das Zentrum mathematischer Fragestellungen.

Ein aus historischer Sicht vorläufiger Höhepunkt der Numerik im weitesten Sinne wurde von Leonhard Euler (1707/Basel–1781/Petersburg) geschaffen. Euler

untersuchte günstige Verteilungen von Masten auf Segelschiffen. Für diese Arbeiten erhielt er den Preis der Pariser Akademie der Wissenschaften im Alter von nur 20 Jahren; und dies bevor er je den Ozean sah.

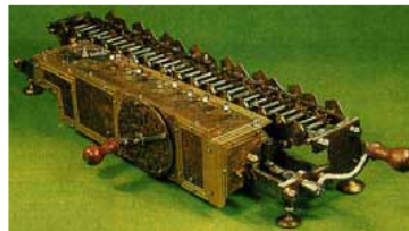
In diese Epoche fällt auch der erste Entwurf einer Rechenmaschine (1672), von Gottfried Wilhelm Leibniz, die er bereits ein Jahr später der Royal Society in London vorführte und die alle vier Grundrechenarten bewältigte.



Gottfried Wilhelm Leibniz:

Erfinder: Differential- und Integralrechnung

Infinitesimalrechnung, Rechenmaschine (1672)



Genie: Psychologe, Ingenieur, Mediziner

Nachlass: 75.000 Seiten

Voltaire: *Leibniz hat Könige gelehrt, Weise erleuchtet, weiser als sie, wusste er um den Zweifel.*

Abbildung 1.3: Gottfried Wilhelm Leibniz und die erste Rechenmaschine.

Die Zeit für die Numerik war jedoch noch nicht reif und kritisch betrachtet, könnte man sagen, dass die Numerik über viele Jahrzehnte hinweg nicht der Durchbruch gelang. Die angewendeten Beweistechniken waren bis ca. 1900 motiviert durch die praktischen Anwendungen/Lösungen meist konstruktiv, doch aus numerischer Sicht nicht brauchbar. Als Folge ist daher (nicht überraschend) zu verzeichnen, dass einer rein logischen Vorgehensweise der Vorzug gegeben wurde. Der berühmte Mathematiker Jacobi äußerte sogar:

'Die Mathematik dient einzig und alleine der Ehre des menschlichen Geistes.'

Heute wissen wir, dass diese Aussage nicht richtig ist!

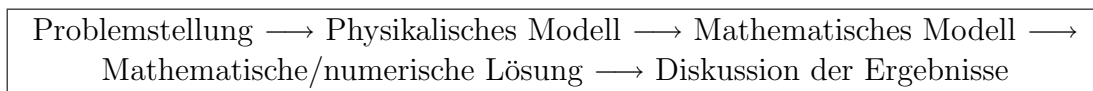
Der eigentliche Aufstieg/Durchbruch der numerischen Mathematik gelang dann mit dem Aufkommen (moderner) Rechenanlagen. Während der Anstieg der Rechengeschwindigkeit bis ca. 1940 um lediglich den Faktor 10 (durch Tricks) gelang, liegt er heute bei 10^{15} (Stand 2004) oder höher.

1.4 Was ist Numerik?

Wir erinnern an das Scherwindbeispiel und stellen einige Dinge fest: Die Anwendung mathematischer Lösungsmethoden auf realistische Aufgabenstellungen der Praxis erfordert fast immer den Einsatz eines Rechners. Die Anforderungen an die Numerik sind dabei vielschichtig:

- Entwicklung von Verfahren zur Konstruktion von Lösungen, meist *Näherungslösungen* mathematischer Aufgabenstellungen
- Effiziente Implementierung auf Rechenanlagen
- Auswahl geeigneter Verfahren
- Aussagen über Güte der Approximation

In diesem Zusammenhang ist die Kette



wichtig, die i.A. mehrfach durchlaufen werden muß.

Aus praktischer Sicht hat man hierbei insbesondere das schwierige Problem, dass mathematische Modell möglichst gut an die Realität anzupassen:

$$\text{Modell} \approx \text{Realität}$$

In dieser Vorlesung werden wir uns genau mit den oben genannten Punkten beschäftigen.

1.5 Motivationsbeispiel

Wir wollen nachfolgend ein konkretes Problem betrachten und gleichzeitig auf eine besondere Problematik aufmerksam machen.

Beispiel 1.2. *Es soll das Integral*

$$I_n = \int_0^1 \frac{x^n}{x+5} dx, \quad n \in \mathbb{N} \cup \{0\} = \mathbb{N}_0$$

für $i = 0, 1, 2, \dots, 20$ berechnet werden. Wir stellen fest:

- *Elementare Integration versagt*
(dennoch: analytische Methoden stets zuerst versuchen!)

- numerische Quadraturverfahren nicht geeignet, da zu aufwendig für das spezielle Problem (→ Auswahl geeigneter Verfahren)
- Lösung: Kombination von analytischer Vorarbeit und numerischer Durchführung

Für die Zahlen I_n kann schnell eine Rekursionsvorschrift angegeben werden:

$$I_0 = \int_0^1 \frac{dx}{x+5} = [\ln|x+5|]_0^1 = \ln \frac{6}{5} \approx 0.182321556\dots$$
$$I_1 = \int_0^1 \frac{x}{x+5} dx = \int_0^1 \left(\frac{x+5}{x+5} - \frac{5}{x+5} \right) dx = 1 - 5 I_0$$

Allgemeiner erhalten wir für $n \rightarrow n+1$:

$$I_{n+1} = \int_0^1 \frac{x^n}{x+5} (x+5-5) dx = \int_0^1 x^n dx - 5 \int_0^1 \frac{x^n}{x+5} dx$$

(1.1)
$$I_{n+1} = \frac{1}{n+1} - 5I_n$$

Ausgehend von $I_0 = \ln \frac{6}{5}$ können wir somit theoretisch alle Werte I_n berechnen. In der rechentechnischen Realisierung erhalten wir jedoch bereits nach wenigen Schritten unbrauchbare Ergebnisse. Mit einem Taschenrechner ergibt sich dann etwa (taschenrechnerspezifisch) folgendes Bild:

$I_1 = 0.08839\dots$	$I_{11} = 0.01377\dots$
$I_2 = 0.05803\dots$	$I_{12} = 0.01445\dots$ (Widerspruch zur Monotonie)
\vdots	\vdots
$I_{10} = 0.01542\dots$	$I_{14} = 0.04814\dots$
	$I_{15} = -0.17404\dots$ (Widerspruch zum Vorzeichen)

Wir wollen die Güte der berechneten Lösung etwas genauer analysieren und stellen fest:

1. Für $x \in [0, 1]$: $\frac{x^n}{x+5} \geq 0 \implies \forall n : I_n \geq 0$.

2. Für $x \in]0, 1[$: $x^{n+1} < x^n \implies \frac{x^{n+1}}{x+5} < \frac{x^n}{x+5} \implies I_{n+1} < I_n$.

Somit ist (I_n) , $n \in \mathbb{N}_0$ streng monoton fallend und wegen $I_n \geq 0$ nach unten beschränkt, also konvergent. Wegen

$$I_n < \int_0^1 x^{n-1} dx = \frac{1}{n}, \quad \forall n \geq 1$$

folgt

$$\lim_{n \rightarrow \infty} I_n = 0,$$

was sich leider mit unseren numerischen Erfahrungen nicht deckt.

Wir wagen einen weiteren Versuch und gehen das Problem von hinten an (Rückwärtsrekursion):

Es ist

$$I_{10} + 5I_9 = \frac{1}{10} \implies I_{10} = \frac{1}{10} - 5I_9 < \frac{1}{10} - 5I_{10} \implies I_{10} < \frac{1}{60}$$

$$I_9 = \frac{1}{50} - \frac{1}{5}I_{10} > \frac{1}{50} - \frac{1}{5}I_9 \implies I_9 > \frac{1}{60}$$

Werten wir nun die Rückwärtsrekursion

$$(1.2) \quad I_{n-1} = \frac{1}{5n} - \frac{1}{5}I_n$$

aus, so erhalten wir:

$$I_9 = \frac{1}{60}$$

$$I_8 = 0.01888\dots$$

⋮

$$I_3 = 0.04313\dots$$

$$I_2 = 0.05803\dots$$

$$I_1 = 0.08839\dots$$

$$I_0 = 0.182321556\dots \text{ alle Stellen richtig!}$$

Bei der Vorwärtsberechnung in (1.1) wird ein Fehler, den wir in I_n z.B. durch Rundung erhalten haben, mit dem Faktor 5 multipliziert und geht so verstärkt in I_{n+1} ein. In der Rückwärtsberechnung (1.2) hingegen, reduziert sich der Fehler um den Faktor $\frac{1}{5}$, so dass die Genauigkeit der Lösung mit jedem weiteren Schritt wächst.

Eine detailliertere Klärung der Situation werden wir später angeben.

1.6 Vorläufiges Fazit

Viele Probleme der Mathematik lassen sich nicht analytisch lösen (es gibt keine explizite Darstellung der Lösung) oder nur sehr schwer lösen (z.B. zu komplex), während eine Lösungsberechnung mit numerischen Verfahren jedoch häufig möglich ist.

Bei numerischen Verfahren können weitere Fehler auftreten, die bei der Analyse der berechneten Lösung zu beachten sind. Hierbei bedeutet Fehler nicht, daß

man etwas falsch gemacht hat. Vielmehr können unvermeidbare Abweichungen vom *exakten*, d.h. realitätsgenauen Ergebnis, auftreten.

Beispiel 1.3. *Wir kommen zurück auf unser Scherwindproblem, die Lösung des Problems kann nur noch numerisch berechnet werden, da eine analytische Lösung nicht existiert. Ihr Lösungsansatz führt auf ein Anfangswertproblem mit einer gewöhnlichen Differentialgleichung. Es können verschiedenen Fehler auftreten:*

Modellfehler: *Die Modellierung ist ungenau, z.B. sind die Windbedingungen nicht beliebig genau modellierbar.*

Datenfehler: *Parameter des DGL-Systems oder Anfangswerte sind nur ungenau angebar.*

Verfahrensfehler: *Das numerische Verfahren zur Lösung der DGL berechnet nur eine genäherte Lösung.*

Rundungsfehler: *Der Computer kann nicht mit beliebig vielen Nachkommastellen rechnen (ein Computer hat nur endlich vielen Speicher), z.B. wird π in der Regel mit nur 16 oder 32 Nachkommastellen berücksichtigt.*

Verfahrensfehler und Rundungsfehler sind Fehler, die aufgrund des Lösungsansatzes durch Numerik auftreten, sie sind innerhalb dieser Vorlesung genauer zu analysieren.

Bemerkung 1.4. *Numerik ist somit nicht nur die Entwicklung von Algorithmen oder Verfahren. Auch die Analyse und die Effizienz der Verfahren sind wesentliche Bestandteile.*

Kapitel 2

Fehleranalyse

Wie bereits festgestellt, können bei der Anwendung mathematischer/numerischer Methoden Fehler z.B. bei der Modellbildung, bei den Eingabeparametern, bei der Approximation oder durch Rundung auftreten. Die beiden letztgenannten wollen wir genauer untersuchen. Hierzu ist es erforderlich die Struktur der Zahlendarstellung auf einem Computer zu untersuchen.

2.1 Maschinenzahlen

Die mit einer bestimmten Codierung darstellbaren Zahlen bezeichnen wir als Menge der *Maschinenzahlen*. Gebräuchlichste Codierungsform ist dabei die sogenannte Gleitpunkt-Darstellung: Die Zahl

$$x = V_M(d_1p^{-1} + d_2p^{-2} + \dots + d_l p^{-l}) \cdot p^E \text{ mit } E = V_E(e_1p^{n-1} + e_2p^{n-2} + \dots + e_{n-1}p + e_n)$$

wird codiert durch

$$\underbrace{V_M d_1 d_2 \dots d_l}_{\text{Mantisse}} \mid \underbrace{V_E e_1 e_2 \dots e_n}_{\text{Exponent}}$$

Dabei ist

$$p \in \mathbb{N}, p > 1, p \text{ fest}$$

Basis

$$d_\lambda \in \{0, 1, \dots, p-1\}, \quad \lambda = 1, \dots, l$$

Ziffern der Mantisse

$$e_v \in \{0, 1, \dots, p-1\}, \quad v = 1, \dots, n$$

Ziffern des Exponenten

$$V_M, V_E \in \{+, -\}$$

Vorzeichen der Mantisse bzw. des Exponenten

Beispiel 2.1. $p = 10, l = 4, n = 3$

$$x = 4711 \quad \rightarrow \quad +4711 \mid +004 \quad \rightarrow \quad 0.4711 \cdot 10^4$$

$$x = -17.5 \quad \rightarrow \quad -1750 \mid +002 \quad \rightarrow \quad -0.1750 \cdot 10^2$$

$$x = 0.008008 \quad \rightarrow \quad +8008 \mid -002 \quad \rightarrow \quad 0.8008 \cdot 10^{-2}$$

Bemerkung 2.2. Fordert man $d_1 \neq 0$ für $x \neq 0$ (normalisierte Gleitpunkt-Darstellung), dann ist für jede Maschinenzahl $x \neq 0$ die Gleitpunktdarstellung eindeutig; lediglich beim Exponenten 0 bleibt V_E unbestimmt.

Definition 2.3 (Gleitpunktzahl). Sei

$$\begin{array}{ll} p > 1, p \in \mathbb{N} & \text{Basis; z. B.: } p = 2, 10, 16 \\ D, E \in \mathbb{Z} & \text{Mantisse, Exponent} \\ l > 0. & \end{array}$$

Dann lautet die Menge der l -stelligen, normalisierten Gleitpunktzahlen zur Basis p :

$$\mathbb{G} = \{Dp^{E-l} \mid D = 0 \vee p^{l-1} \leq |D| < p^l\}$$

Beispiel 2.4. 4-stellige, normalisierte Zahlen zur Basis 10:

$$\mathbb{G} = \{D \cdot 10^{E-4} \mid D = 0 \vee 10^3 \leq |D| < 10^4\}$$

Die übliche Schreibweise lautet:

$$Dp^{E-l} = \pm \underbrace{d_1 d_2 \dots}_{|E| \text{ Stellen vor Komma}}, \underbrace{\dots d_{l-1} d_l}_{l - |E| \text{ Stellen nach Komma}} 0 \dots 0$$

Bezeichnungen:

$\mathbb{M}(p, l, n)$ oder \mathbb{M} : Menge der Maschinenzahlen

Offensichtlich ist diese Menge endlich und somit gilt: $\mathbb{M} \stackrel{\subset}{\neq} \mathbb{R}$.

Bemerkung 2.5.

1. Die Wahl der Basis p wird durch die Rechnerkonstruktion bestimmt.

$$\begin{array}{ll} \text{Dualsystem:} & p = 2 \quad \text{mit Ziffern } 0, 1 \\ \text{Hexadezimalsystem:} & p = 16 \quad \text{mit Ziffern } 0, 1, \dots, 9, A, \dots, F \end{array}$$

2. Der Exponent ist in der Praxis durch den Speicherplatz eingeschränkt.

Beispiel 2.6. Typische Situation auf PC's für PASCAL, FORTRAN: real^*4 :

1 Real-Zahl benötigt 4 Byte Speicherplatz:

1 Byte für Vorzeichen, Exponent

3 Byte für Mantisse

Es stellt sich als nächstes die Frage, welche Zahlen sich überhaupt damit darstellen lassen (1 Byte = 8 Bit, 1 Bit ist entweder 0 oder 1).

i) Stellenzahl:

3 Byte = 24 Bit für die Mantisse, also 24-stellige Dualzahlen:

$$2^{24} = 10^l \implies l = \log_{10} 2^{24} = 24 \log_{10} 2 \approx 24 \cdot 0.3010 \approx 7.2 \dots$$

d.h. 7-stellige Mantisse im 10er System

ii) Exponentenbereich

2 Bit für die beiden Vorzeichen, d.h. 6 Bit für den Exponenten: $\alpha \leq E \leq \beta$
mit $\beta = -\alpha = 63 = \underline{1} \cdot 2^5 + \underline{1} \cdot 2^4 + \underline{1} \cdot 2^3 + \underline{1} \cdot 2^2 + \underline{1} \cdot 2^1 + \underline{1} \cdot 2^0$

Allgemein erkennen wir, dass es E_1 und E_2 gibt mit

$$\mathbb{M} = \{g \in \mathbb{G} \mid E_1 \leq E \leq E_2\}$$

Offensichtlich gibt es eine kleinste positive Maschinenzahl x_{\min} und eine größte Maschinenzahl x_{\max} .

Beispiel 2.7. Aufgrund der Fülle verschiedenener Parametersetzungen gab es 1983 einen Standardisierungsversuch (IEEE) für $p = 2$:

einfache Genauigkeit (32 Bits): $l = 23$; 8 Bit für E
doppelte Genauigkeit (64 Bits): $l = 52$; 11 Bit für E
Register (80 Bits) $l = 64$; 15 Bit für E

Bemerkung 2.8. Integer-Zahlen werden in vergleichbarer Weise codiert.

2.2 Maschinenzahlen auf der Zahlengerade

Wie festgestellt, existieren in \mathbb{M} eine kleinste positive Maschinenzahl x_{\min} und eine größte Maschinenzahl x_{\max} . Die Maschinenzahlen dazwischen sind jedoch nicht gleichmäßig verteilt:

Verteilung auf der Zahlengeraden:

a) Innerhalb jeden Intervalls $[p^{k-1}, p^k)$ liegen die Maschinenzahlen in gleichen Abständen:

$$a \cdot p^k \text{ mit } a = p^{-l} = \underbrace{.00 \dots 01}_l \cdot p^0$$

hierbei bezeichnet a die Einheit der letzten Mantissenstelle beim Exponenten 0.

b) Die Maschinenzahlen sind nicht auf $\mathbb{R} \cap [x_{\min}, x_{\max}]$ gleichabständig verteilt. Der relative Abstand

$$\frac{x_{i+1} - x_i}{x_i}, \quad x_i \neq 0$$

zweier aufeinanderfolgender Maschinenzahlen variiert höchstens um einen Faktor ρ .

Beispiel 2.9. $M(2, 3, 2)$

Verfügbare positive Mantissen	Verfügbare Exponenten
+100 = 1/2	00 = 0
+101 = 5/8	± 01 = ± 1
+110 = 3/4	± 10 = ± 2
+111 = 7/8	± 11 = ± 3

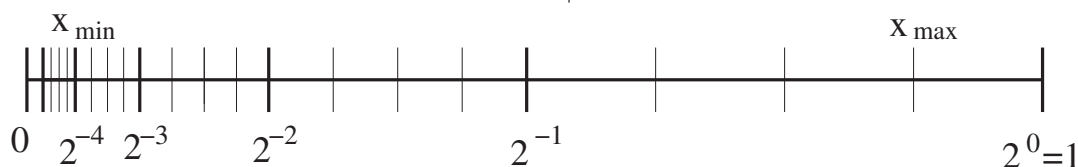


Abbildung 2.1: Beispiel zur Verteilung der Maschinenzahlen auf dem Computer.

2.3 Rundung

Definition 2.10 (Rundung). Eine korrekte Rundung ist die Abbildung

$$rd : \mathbb{R} \rightarrow \mathbb{M},$$

die jedem $r \in \mathbb{R}$ das nächstgelegene $x \in \mathbb{M}$ zuordnet:

$$|rd(r) - r| \leq |x - r| \quad \forall x \in \mathbb{M}.$$

Bemerkung 2.11. Die Definition ist bis auf Ausnahmen eindeutig; dort muß eine Zusatzbedingung zur eindeutigen Behandlung eines Umschlagpunktes \bar{x} angegeben werden.

Ausführung der korrekten Rundung:

Bei der technischen Realisierung der korrekten Rundung muß r nicht exakt bekannt sein, z.B. $r = \sqrt{2}, \pi$. Offenbar reicht es, die $(l + 1)$ ste Mantisse zu kennen:¹

$$r = \pm .d_1 d_2 \dots d_l d_{l+1} \dots p^E \quad (\text{normalisiert})$$

$$r' := \begin{cases} .d_1 d_2 \dots d_{l-1} d_l & \text{falls } 0 \leq d_{l+1} < p/2 \\ .d_1 d_2 \dots d_{l-1} (d_l + 1) & \text{falls } p/2 \leq d_{l+1}, \end{cases} \quad (\text{Aufrunden in } \bar{x})$$

←←←←←
(evtl. Übertrag)

¹Eine Realisierung der korrekten Rundung läßt sich jedoch erst mit zwei Schutzziffern realisieren.

$$rd(r) := \text{sign}(r) \cdot r' \cdot p^E$$

Beispiel 2.12. $l = 4, p = 10$

$$rd(0.142842_{10^2}) = 0.1428_{10^2}$$

$$rd(0.14285_{10^0}) = 0.1429_{10^0}$$

$$rd(0.14997_{10^{-1}}) = 0.1500_{10^{-1}}$$

Schranken für den absoluten und den relativen Fehler:

Sei $r = \pm .d_1 d_2 \dots d_l d_{l+1} \dots p^E$ ($d_1 \neq 0$).

$$(2.1) \quad |rd(r) - r| \leq \frac{a}{2} p^E = \frac{1}{2} p^{E-l} \quad (a = p^{-l}) \quad \underline{\text{absoluter Fehler}}$$

$$(2.2) \quad \frac{|rd(r) - r|}{|r|} \leq \frac{p}{2} p^{-l} \quad \underline{\text{relativer Fehler}}$$

$$|r| \geq 0.1 p^E$$

Wir erhalten

Satz 2.13. Für alle $r \in \mathbb{R} \cap ([-x_{\max}, -x_{\min}] \cup [x_{\min}, x_{\max}])$ gibt es ein $\varepsilon \in \mathbb{R}$, $|\varepsilon| \leq \text{eps}$,

$$(2.3) \quad \text{eps} := \frac{p}{2} p^{-l}, \quad (\text{relative Maschinengenauigkeit}),$$

so dass

$$(2.4) \quad rd(r) = r(1 + \varepsilon)$$

gilt.

Es gibt einige Besonderheiten zu berücksichtigen:

- Zwischenresultate können IM verlassen: z. B. bei $c = \sqrt{a^2 + b^2}$
- Exponentenunterlauf: Für $r \approx 0$, d.h. $r \in] -x_{\min}, x_{\min}[$ gilt:

$$rd(r) := 0.$$

Für $r \in] -x_{\min}, x_{\min}[\setminus \{0\}$ gilt: Der relative Fehler ist stets 1, während der absolute Fehler klein, nämlich $< x_{\min}$ ist. In diesem Fall sollte eine Warnung "underflow" gesetzt werden; eventuell können massive Probleme auftreten, z.B. bei $c = \frac{1}{r}$.

- Exponentenüberlauf: Für $|r| > x_{\max}$ gilt:

$$rd(r) := \text{sign}(r) \cdot x_{\max},$$

gleichzeitig sollte eine Warnung "exponential overflow" gesetzt werden; eventuell können auch hier massive Probleme auftreten, z.B. bei $c = r^2$, $r \rightarrow \infty$. Hier ist der relative Fehler durch 1 beschränkt, während der absolute Fehler beliebig groß werden kann.

Bemerkung 2.14. Die Fehlerschranken (2.1) und (2.2) gelten nur solange kein Exponentenunter- bzw. überlauf vorkommt.

Bemerkung 2.15. Es gibt verschiedene Vereinbarungen an Umschlagpunkte, wie z.B. Aufrunden, Abrunden oder Abschneiden. Stets muß die Rundungsvorschrift jedoch eine idempotente ($c \in \mathbb{M}$, $\Rightarrow rd(c) = c$) und monotone ($c_1 \leq c_2 \Rightarrow rd(c_1) \leq rd(c_2)$) Abbildung sein.

Sondercodierungen: Für $0, +\infty, -\infty, \infty, \text{NaN}$ (Not-A-Number;) hat man oft Sondercodierungen mit speziellen Rechenregeln, z. B.: $x = 0, \frac{1}{x} = \text{NaN}; \infty - \infty = \text{NaN}; \dots$

Beispiel 2.16. Wir wollen nachfolgend die Rekursionsformel aus Beispiel 1.2 genauer untersuchen:

$$I_n + 5I_{n-1} = \frac{1}{n}, \quad I_0 = \ln \frac{6}{5} \notin \mathbb{M}$$

Auswirkungen des Eingabefehlers in I_0 :

$$\hat{I}_0 = I_0 + \Delta I_0$$

$$\Rightarrow I_n + \Delta I_n = \hat{I}_n = \frac{1}{n} - 5\hat{I}_{n-1} = \frac{1}{n} - 5(I_{n-1} + \Delta I_{n-1}) = I_n - 5\Delta I_{n-1}$$

$$\Rightarrow \Delta I_n = -5\Delta I_{n-1} = (-5)^2 \Delta I_{n-2} = \dots = (-5)^n \Delta I_0$$

Ist beispielsweise $\Delta I_0 = 10^{-10}$, dann ist bereits $\Delta I_{15} = (-5)^{15} 10^{-10} \approx -3.05 \dots$. Hierbei ist $(-5)^i$ oszillierend und stark anwachsend.

Umgekehrt tritt bei der Rückwärtsrekursion Fehlerdämpfung auf:

$$\Delta I_0 = \frac{1}{(-5)^n} \Delta I_n$$

Mit beispielsweise $\Delta I_{15} = 1$ folgt

$$\Rightarrow \Delta I_0 = -3.27 \dots \cdot 10^{-11}.$$

2.4 Gleitpunkt-Arithmetik

Sind $x, y \in \mathbb{M}$, so braucht $x \cdot y$ mit $\cdot \in \{+, -, \times, /\}$ nicht aus \mathbb{M} zu sein.

Beispiel 2.17. $\mathbb{M}(10, 5, 2)$, mit $x = .25684_{10^1}$, und $y = .32791_{10^{-2}}$:

$$\left. \begin{aligned} x + y &= .25716791_{10^1} \\ x \times y &= .842204044_{10^{-2}} \\ x/y &= .7832637004 \dots_{10^3} \end{aligned} \right\} \notin \mathbb{M}$$

Gleitpunkt-Operationen $\oplus, \ominus, \otimes, \oslash$:

Die korrekte Rundung lautet

$$x \odot y := rd(x \cdot y), \cdot \in \{+, -, \times, /\},$$

wobei wegen (2.4) gilt:

$$\begin{aligned} x \oplus y &= (x + y)(1 + \alpha) \\ x \ominus y &= (x - y)(1 + \beta) \\ x \otimes y &= (x \times y)(1 + \gamma) \\ x \oslash y &= (x/y)(1 + \delta) \end{aligned}$$

mit $|\alpha|, |\beta|, |\gamma|, |\delta| \leq eps$ (Maschinengenauigkeit).

Bemerkung 2.18. *Es ist anzumerken, dass $\alpha, \beta, \gamma, \delta$ von x und y abhängen, nicht jedoch von ihrer Schranke (eps ist a priori bekannt). $\alpha, \beta, \gamma, \delta$ sind die relativen Fehler der Gleitpunkt-Operationen.*

Bemerkung 2.19. *Nur die Kommutativität der Addition und der Multiplikation bleiben auch bei der Gleitpunkt-Arithmetik erhalten. Assoziativ- und Distributivgesetze gelten nicht mehr!*

Beispiel 2.20. $\mathbb{M}(10, 8, 2)$, $a = .23371258_{10^{-4}}$, $b = .33678429_{10^2}$, $c = -.33677811_{10^2}$.

$$\begin{aligned} a \oplus (b \oplus c) &= a \oplus .61800000_{10^{-3}} &&= .64137126_{10^{-3}} \\ (a \oplus b) \oplus c &= .33678452_{10^2} \oplus c &&= .64100000_{10^{-3}} \\ a + b + c &= &&= .641371258_{10^{-3}} \end{aligned}$$

Bemerkung 2.21. *Die korrekt rundende Gleitpunkt-Arithmetik kann technisch realisiert werden, wenn das Rechenwerk über mindestens zwei Stellen mehr verfügt als die Mantissenlänge l . Es gibt jedoch nur wenige Anlagen, die korrekt runden!*

2.5 Fehlerfortpflanzung, Kondition

Wir hatten festgestellt, dass bei der Berechnung mit dem Computer verschiedene Fehler auftreten können. Betrachten wir für $D \subset \mathbb{R}^n$ und $f : D \rightarrow \mathbb{R}^m$ das Problem der Berechnung

$$y = f(x), \quad x \in D,$$

so bezeichnen wir mit x die *Eingabedaten*, während y *Ausgabe-* oder *Resultatdaten* genannt werden. Die Genauigkeit der Berechnung von $y = f(x)$ wird durch die *Fehlertypen*

- Fehler in den Eingabedaten,
- Abbrechfehler oder Diskretisierungsfehler,
- Rundungsfehler während der Rechnungen.

begrenzt. Anstelle einer exakten Rechnung $y = f(x)$ wird man daher eine Approximation $\tilde{f}(\tilde{x})$, mit

\tilde{x} : Approximation für x , z.B. durch $\tilde{x} = rd(x)$

\tilde{f} : Approximation für f

berechnen, vgl. Abbildung 2.2.

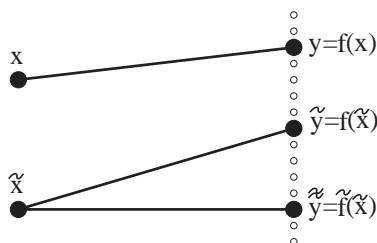


Abbildung 2.2: Fehlereinflüsse.

Wir befassen uns zunächst damit, wie sich Fehler in x auf das Ergebnis $y = f(x)$ auswirken. Sei \tilde{x} eine Näherung von x und sei

$\Delta x = \tilde{x} - x$: der *absolute Fehler*,

$\frac{\tilde{x}_i - x_i}{x_i}$: der *relative Fehler*, ($i = 1, \dots, n$),

$\tilde{y} = f(\tilde{x})$: der Näherungswert für $y = f(x)$.

Die Funktion f sei eine C^1 -Funktion. Die Taylor-Entwicklung erster Ordnung liefert für den absoluten Fehler $\Delta y = \tilde{y} - y$ die Approximation

$$(2.5) \quad \Delta y_i = f_i(x + \Delta x) - f_i(x) \approx \frac{\partial f_i(x)}{\partial x} \Delta x = \sum_{j=1}^n \frac{\partial f_i(x)}{\partial x_j} \Delta x_j, \quad i = 1, \dots, m.$$

Für den relativen Fehler erhält man dann

$$(2.6) \quad \frac{\Delta y_i}{y_i} \approx \sum_{j=1}^n \left(\frac{\partial f_i(x)}{\partial x_j} \frac{x_j}{f_i(x)} \right) \cdot \left(\frac{\Delta x_j}{x_j} \right), \quad x_j, y_i \neq 0.$$

Definition 2.22 (Kondition).

1. Die Zahlen

$$(2.7) \quad k_{ij}(x) = \left| \frac{\partial f_i(x)}{\partial x_j} \frac{x_j}{f_i(x)} \right|$$

heißen Verstärkungsfaktoren bzw. (relative) Konditionszahlen.

2. Das Problem 'Berechne $y = f(x)$ ' heißt gut konditioniert, falls alle $k_{ij}(x)$ die Größenordnung 1 haben. Andernfalls heißt das Problem schlecht konditioniert.

Zuerst untersuchen wir damit die *arithmetischen Operationen* $+$, $-$, $*$, $/$.

Multiplikation: $y = f(x_1, x_2) = x_1 * x_2$

Es gilt $k_{11}(x) = k_{12}(x) = 1$: gutartig.

Division: $y = f(x_1, x_2) = x_1/x_2$

Es gilt $k_{11}(x) = k_{12}(x) = 1$: gutartig.

Addition, Subtraktion: $y = f(x_1, x_2) = x_1 + x_2$

Es gilt

$$k_{11}(x) = \left| \frac{x_1}{x_1 + x_2} \right|, \quad k_{12}(x) = \left| \frac{x_2}{x_1 + x_2} \right|$$

Das Problem ist schlecht konditioniert, falls $x_1 \approx -x_2$. Daher ist die Subtraktion nahezu gleichgroßer Zahlen mit gleichen Vorzeichen schlecht konditioniert.

Dieses Phänomen heißt Auslöschung.

Beispiel 2.23.

$$1.31 - 1.25 = 0.06$$

$$1.32 - 1.24 = 0.08 \quad (\text{Störrechnung})$$

Es gilt:

$$x = (1.31, -1.25)$$

$$y = x_1 + x_2 = 0.06$$

$$\Delta x = (0.01, 0.01)$$

$$\left| \frac{\Delta x_i}{x_i} \right| \leq 0.008, \quad \text{d.h. relativer Eingabefehler ca. 0.8\%}$$

$$k_{1,i}(x) \leq 22, \quad i = 1, 2$$

Der relative Fehler im Ergebnis ist ca. 40 mal (Summe) größer als der relative Fehler in den Daten.

Wurzel: $y = f(x_1) = \sqrt{x_1}$, $x_1 > 0$
Es gilt $k(x) = 1/2$: gutartig. (Übung)

Bemerkung 2.24. Bei einigen Problemen kann die Auslöschung durch geeignete Umformulierung vermieden werden, vgl. die nachfolgenden Beispiele.

2.6 Algorithmen

Ein *Algorithmus* zur Berechnung der Lösung $y = f(x)$ eines Problems ist eine Sequenz von endlich vielen 'elementaren Operationen' (+, -, *, /, $\cos(x)$, \sqrt{x} , ...). Es gibt i.A. mehrere Anordnungen der Rechenschritte, welche zum gleichen Ergebnis $y = f(x)$ führen. In jedem Rechenschritt fallen Rundungsfehler an. Dabei kann der Fall auftreten, daß bei der Lösung eines an sich gut konditionierten Problems eine ungünstige Anordnung der Rechenschritte zum Aufschaukeln der Rundungsfehler führt. Der zugehörige Algorithmus ist numerisch instabil.

Beispiel 2.25. Gesucht ist die betragskleinere Lösung von $x^2 + 2px - q = 0$ mit $p \gg q$.

Die exakte Lösung ist gegeben durch

$$(2.8) \quad y = f(p, q) = -p + \sqrt{p^2 + q}$$

Für die Konditionszahlen gilt: (Übung)

$$(2.9) \quad k_p(p, q) < 1, \quad k_q(p, q) < 1, \quad (\text{wegen } p \gg q)$$

also ist die Aufgabe gut konditioniert für $q > 0$.

(Wäre jedoch etwa $q \approx -p^2$, wäre das Problem schlecht konditioniert.)

Beispiel 2.26. Wir betrachten erneut die Aufgabenstellung aus Beispiel 2.25 und untersuchen zwei Algorithmen:

Algorithmus 1: $y = f(p, q) = -p + \sqrt{p^2 + q}$

$$s := p * p$$

$$t := s + q$$

$$u := \sqrt{t}$$

$$y := -p + u$$

Für $p \gg q$ tritt in $y := -p + u$ Auslöschung auf, der Algorithmus ist in diesem Fall schlecht konditioniert, obgleich das Problem gut konditioniert ist.

Algorithmus 2: $y = f(p, q) = \frac{q}{p + \sqrt{p^2 + q}}$

$$s := p * p$$

$$t := s + q$$

$$u := \sqrt{t}$$

$$v := p + u$$

$$y := q/v$$

Algorithmus 2 ist für $p \gg q$ gutartig.

Zahlenwerte: $p = 6.0002$, $q = 0.01$ und einer Mantissenlänge von 5 erhalten wir

Algorithmus 1: 0.0008

Algorithmus 2: 0.00083326

exakte Lösung: 0.00083325 (gerundet auf Mantissenlänge)

(Übung: Nachrechnen)

Kapitel 3

Lineare Gleichungssysteme

3.1 Einführung und Aufgabenstellung

Algorithmen zur Lösung linearer Gleichungssysteme bilden die Basis für viele Anwendungen der Numerik.

Aufgabenstellung: Sei A eine (m, n) -Matrix und sei $b \in \mathbb{R}^m$. Gesucht ist ein Vektor $x \in \mathbb{R}^n$, welcher das lineare Gleichungssystem (LGS)

$$(3.1) \quad Ax = b$$

löst.

In der Numerik unterscheidet man in *direkte Methoden* zur Lösung von $Ax = b$, bei der eine Lösung x in *endlich* vielen Schritten berechnet wird und *indirekte Methoden*, bei denen eine Näherungslösung x von $Ax = b$ iterativ bestimmt wird. In diesem Kapitel werden wir uns mit den direkten Methoden auseinandersetzen, die indirekten Methoden sind dann Bestandteil der Numerik II.

Bemerkung 3.1. *Abhängig vom "Aussehen" der Koeffizientenmatrix A unterscheidet man in*

<i>kleine</i>	- <i>große</i>	<i>Systeme</i>
<i>symmetrische</i>	- <i>nichtsymmetrische</i>	<i>Matrizen</i>
<i>mit</i>	- <i>ohne</i>	<i>Bandstruktur</i>
<i>schwach</i>	- <i>vollbesetzte</i>	<i>Matrizen</i>

Danach richtet sich auch die Auswahl der Verfahren.

Beispiel 3.2. *Wir betrachten das Gleichstromnetzwerk in Abbildung 3.1. Nach den KIRCHHOFFSCHEN Gesetzen müssen sich zunächst an allen Knoten*

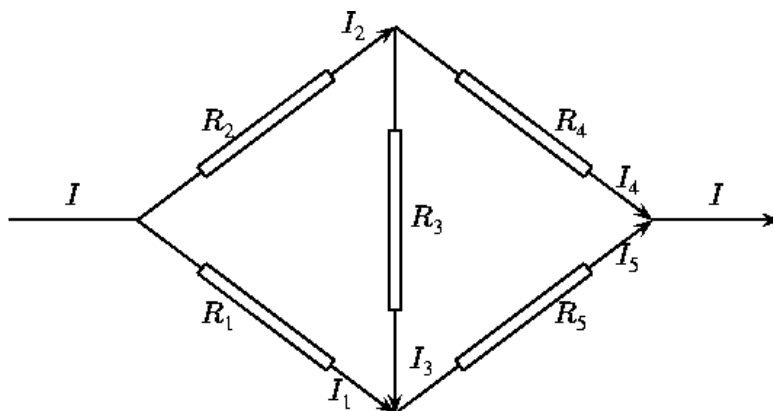


Abbildung 3.1: Gleichstromnetzwerk.

die eingehenden und die ausgehenden Ströme zu Null ergänzen. Für Anfang und Ende erhalten wir

$$I_1 + I_2 = I = I_4 + I_5,$$

für Oben und Unten

$$I_2 = I_3 + I_4 \quad \text{und} \quad I_1 + I_3 = I_5.$$

Darüberhinaus müssen sich die Spannungen in den beiden Dreiecken zu Null summieren, nach dem OHmschen Gesetz ($U = R \cdot I$) führt dies für die bekannten Widerstände R_i auf die beiden Gleichungen

$$R_2 I_2 + R_3 I_3 - R_1 I_1 = 0 \quad \text{und} \quad R_3 I_3 + R_5 I_5 - R_4 I_4 = 0$$

Wir erhalten somit nach Umsortierung ein lineares Gleichungssystem der Form $Ax = b$:

$$\begin{array}{ccccccc} I_1 & + & I_2 & & & & = I \\ & & & & I_4 & + & I_5 = I \\ & & I_2 & - & I_3 & - & I_4 = 0 \\ I_1 & & & + & I_3 & & - I_5 = 0 \\ -R_1 I_1 & + & R_2 I_2 & + & R_3 I_3 & & = 0 \\ & & & & -R_3 I_3 & + & R_4 I_4 - R_5 I_5 = 0 \end{array}$$

Während man das letzte Beispiel sicher noch von Hand lösen kann, werden größere Probleme am Computer gelöst, vgl. Abbildung 3.2.

Bei der Lösung linearer Gleichungssysteme sind verschiedene Fälle möglich:

1. $m = n$: $\text{rang}(A) = n$, d.h. $Ax = b$ ist eindeutig lösbar. Da A invertierbar ist folgt $x = A^{-1}b$. Für numerische Rechnungen ist diese Darstellung jedoch nicht geeignet: auch die Cramersche Regel ist für $n \geq 3$ numerisch nicht brauchbar.

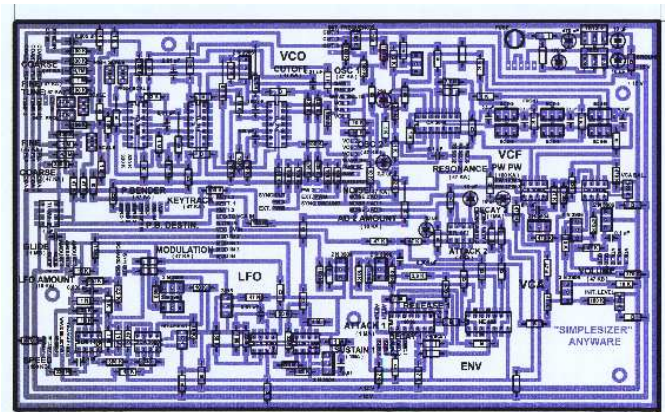


Abbildung 3.2: Eine eher kleine Platine.

2. $m > n$: Das LGS $Ax = b$ heißt *überbestimmt* und hat im allgemeinen keine Lösung. Stattdessen wird ein Ersatzproblem gelöst, vgl. Lineare Ausgleichsrechnung.
3. $m < n$: Das LGS $Ax = b$ heißt *unterbestimmt*. Wenn eine Lösung existiert, dann hat der Lösungsraum die Dimension $n - \text{rang}(A)$. Anwendungen findet man etwa in der Linearen Optimierung.

3.2 LR-Zerlegung und Gauß-Elimination

3.2.1 Idee der Gauß-Elimination/LR-Zerlegung

Sei $A = (a_{i,k})$ eine (n, n) -Matrix und $b \in \mathbb{R}^n$. Zu lösen sei das LGS

$$(3.2) \quad Ax = b$$

Das *Gauß'sche-Eliminationsverfahren* zur Lösung von LGS haben sie bereits im Rahmen ihres bisherigen Studiums kennengelernt. Es ist ein recht anschauliches Verfahren, das sich zudem leicht implementieren läßt.

Wir betrachten zunächst den vereinfachten Fall, daß die Matrix A in oberer Dreiecksform vorliegt, d.h.

$$(3.3) \quad A = R = \begin{pmatrix} r_{1,1} & r_{1,2} & \cdots & r_{1,n} \\ 0 & r_{2,2} & \cdots & r_{2,n} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & r_{n,n} \end{pmatrix}$$

Man spricht dann von einem *gestaffelten Gleichungssystem*, der Grund ist leicht ersichtlich. In diesem Fall kann für $r_{i,i} \neq 0$ leicht eine Lösung von $Rx = c$ (hier

3.2.3 Gauß–Elimination/LR–Zerlegung ohne Pivoting

Sei die Matrix

$$A := \begin{pmatrix} a_{1,1}^{(1)} & \cdots & a_{1,n}^{(1)} \\ \vdots & & \vdots \\ a_{n,1}^{(1)} & \cdots & a_{n,n}^{(1)} \end{pmatrix}$$

gegeben.

1. Schritt: Sei $a_{1,1}^{(1)} \neq 0$

$$L_1 A = \begin{pmatrix} a_{1,1}^{(1)} & \cdots & \cdots & a_{1,n}^{(1)} \\ 0 & a_{2,2}^{(2)} & \cdots & a_{2,n}^{(2)} \\ \vdots & \vdots & & \vdots \\ 0 & a_{n,2}^{(2)} & \cdots & a_{n,n}^{(2)} \end{pmatrix} \quad (\text{vgl. (3.6) mit } j = 1)$$

mit

$$l_{i1} := \frac{a_{i1}^{(1)}}{a_{11}^{(1)}}, \quad i = 2, \dots, n,$$

$$a_{ik}^{(2)} = a_{ik}^{(1)} - l_{i1} a_{1k}^{(1)}, \quad i, k = 2, \dots, n.$$

In Worten: Subtrahiere von der i -ten Zeile der Matrix A das l_{i1} -fache der 1. Zeile, $i = 2, \dots, n$.

Allgemein sei nun unsere Ausgangssituation vor dem j -ten Schritt ($j \geq 2$) bekannt.

Ausgangsmatrix vor dem j -ten Schritt ($j \geq 2$):

$$L_{j-1} \dots L_1 A = \begin{pmatrix} a_{11}^{(1)} & & & a_{1n}^{(1)} \\ & \ddots & & \vdots \\ & & a_{jj}^{(j)} & \cdots & a_{jn}^{(j)} \\ 0 & & \vdots & & \vdots \\ & & a_{nj}^{(j)} & \cdots & a_{nn}^{(j)} \end{pmatrix}.$$

Mit diesem Wissen können wir den j -ten Schritt in Angriff nehmen:

j-ter Schritt ($j \geq 2$): Sei $a_{jj}^{(j)} \neq 0$

$$L_j L_{j-1} \dots L_1 A = \begin{pmatrix} a_{1,1}^{(1)} & \dots & \dots & \dots & \dots & a_{1,n}^{(1)} \\ & \ddots & & & & \vdots \\ & & a_{j,j}^{(j)} & \dots & \dots & a_{j,n}^{(j)} \\ & & 0 & a_{j+1,j+1}^{(j+1)} & \dots & a_{j+1,n}^{(j+1)} \\ & & \vdots & \vdots & & \vdots \\ & & 0 & a_{n,j+1}^{(j+1)} & \dots & a_{n,n}^{(j+1)} \end{pmatrix}$$

mit

$$\begin{aligned} l_{ij} &= \frac{a_{ij}^{(j)}}{a_{jj}^{(j)}}, & i &= j+1, \dots, n, \\ a_{ik}^{(j+1)} &= a_{ik}^{(j)} - l_{ij} a_{jk}^{(j)}, & i, k &= j+1, \dots, n. \end{aligned}$$

Nach $n-1$ Schritten erhalten wir dann das gewünschte Resultat:

$$(3.9) \quad L_{n-1} \dots L_1 A = \begin{pmatrix} a_{11}^{(1)} & \dots & \dots & a_{1n}^{(1)} \\ & \ddots & & \vdots \\ & & \ddots & \vdots \\ & & & a_{nn}^{(n)} \\ 0 & & & \end{pmatrix} =: R = (r_{ik})$$

mit $r_{ii} = a_{ii}^{(i)} \neq 0$.

Wendet man die Matrizen L_j direkt auf die erweiterte Matrix (A, b) an, so ergibt sich

$$L_{n-1} \dots L_1 (A, b) = (R, c).$$

Das LGS $Ax = b$ ist dann äquivalent zu $Rx = c$ und kann gemäß (3.5) gelöst werden (Gauß-Elimination).

Aus (3.9) folgt mittels der Formeln (3.7)(3.8) die LR-Zerlegung der Matrix A :

$$(3.10) \quad A = L_1^{-1} \dots L_{n-1}^{-1} R =: LR$$

$$L = \begin{pmatrix} 1 & & & & 0 \\ & l_{2,1} & \ddots & & \\ & \vdots & \ddots & \ddots & \\ & \vdots & & \ddots & \ddots \\ l_{n,1} & \cdots & \cdots & l_{n,n-1} & 1 \end{pmatrix} \quad \text{linke Dreiecksmatrix}$$

Bei gegebener LR -Zerlegung $A = LR$ ist das LGS $Ax = b$ äquivalent zu den beiden leicht auflösbaren LGS

$$Lc = b, \quad Rx = c.$$

Insbesondere folgt noch aus (3.10)

$$\det(A) = \det(L) \det(R) = \prod_{j=1}^n r_{jj}.$$

Bei unseren bisherigen Überlegungen hatten wir stets $a_{jj}^{(j)} \neq 0$ voraussetzen müssen und es stellt sich die Frage, wann dies gesichert anzunehmen ist.

Problem: Wann gilt $a_{jj}^{(j)} \neq 0$?

Satz 3.3. Sei A eine (n, n) -Matrix, deren Hauptabschnittsmatrizen A_j regulär sind. Dann gibt es eine eindeutige Zerlegung

$$A = LR,$$

L linke Dreiecksmatrix mit $l_{jj} = 1, j = 1, \dots, n$,

R reguläre rechte Dreiecksmatrix.

Beweis: Vgl. Satz 3.9

3.2.4 Permutationsmatrizen

Zur Behandlung des Falles $a_{jj}^{(j)} = 0$ für ein j benötigen wir sogenannte *Permutationsmatrizen*. Hierzu sei

$$e_i = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i \quad \text{der } i\text{-te kanonische Einheitsvektor}$$

Eine Matrix P heißt Permutationsmatrix, wenn eine Permutation (i_1, \dots, i_n) von $(1, \dots, n)$ existiert mit

$$P = \begin{pmatrix} e_{i_1}^T \\ \vdots \\ e_{i_n}^T \end{pmatrix}.$$

Insbesondere haben Permutationsmatrizen die Eigenschaften: $P^2 = I$, also $P^{-1} = P$.

3.2.5 Gauß-Elimination/LR-Zerlegung mit Pivoting

Wir gehen davon aus, dass wir vielleicht bereits einige Schritte zur Gauß-Elimination bzw. LR-Zerlegung durchgeführt haben und befinden uns im j -ten $j \geq 1$ Schritt.

j -ter Schritt ($j \geq 1$): Die Ausgangsmatrix sei

$$(3.11) \quad A^{(j)} := \begin{pmatrix} a_{11}^{(1)} & & a_{1n}^{(1)} \\ & \ddots & \vdots \\ & & a_{jj}^{(j)} \cdots a_{jn}^{(j)} \\ 0 & & \vdots \\ & & a_{nj}^{(j)} \cdots a_{nn}^{(j)} \end{pmatrix}, \quad A^{(1)} := A$$

Wir führen eine Spaltenpivot-Suche durch: Wähle eine Zeile r mit

$$|a_{rj}^{(j)}| = \max_{i \geq j} |a_{ij}^{(j)}|.$$

Wir haben verschiedenen Fälle zu unterscheiden:

1. Fall: $a_{rj}^{(j)} = 0$: A ist singulär, setze $L_j = I$.

2. Fall: $a_{rj}^{(j)} \neq 0$: Vertausche die j -te Zeile mit der r -ten Zeile in $A^{(j)}$. Dies entspricht

Beispiel 3.7.

$$\begin{pmatrix} 3 & 1 & 6 \\ 2 & 1 & 3 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 2 \\ 7 \\ 4 \end{pmatrix}$$

Die Pivot-Elemente in den erweiterten Matrizen werden durch einen Unterstrich markiert.

1. Schritt:

$$\left(\begin{array}{ccc|c} \underline{3} & 1 & 6 & 2 \\ 2 & 1 & 3 & 7 \\ 1 & 1 & 1 & 4 \end{array} \right)$$

Anwendung von L_1 :

$$\left(\begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ \boxed{2/3} & 1/3 & -1 & 17/3 \\ \boxed{1/3} & \underline{2/3} & -1 & 10/3 \end{array} \right)$$

2. Schritt:

Vertausche Zeile 2 und 3

$$\left(\begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ \boxed{1/3} & \underline{2/3} & -1 & 10/3 \\ \boxed{2/3} & 1/3 & -1 & 17/3 \end{array} \right)$$

Anwendung von L_2 :

$$\left(\begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ \boxed{1/3} & 2/3 & -1 & 10/3 \\ \boxed{2/3} & \boxed{1/2} & -1/2 & 4 \end{array} \right)$$

$$\Rightarrow L = \begin{pmatrix} 1 & 0 & 0 \\ 1/3 & 1 & 0 \\ 2/3 & 1/2 & 1 \end{pmatrix}, \quad R = \begin{pmatrix} 3 & 1 & 6 \\ 0 & 2/3 & -1 \\ 0 & 0 & -1/2 \end{pmatrix}$$

$$\Rightarrow P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}, \quad PA = \begin{pmatrix} 3 & 1 & 6 \\ 1 & 1 & 1 \\ 2 & 1 & 3 \end{pmatrix}$$

Somit gilt $PA = LR$ Wir erhalten das gestaffelte Gleichungssystem $Rx = c$:

$$\begin{pmatrix} 3 & 1 & 6 \\ 0 & 2/3 & -1 \\ 0 & 0 & -1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 2 \\ 10/3 \\ 4 \end{pmatrix}, \quad \begin{array}{l} x_1 = 19 \\ x_2 = -7 \\ x_3 = -8 \end{array}$$

Das Gauß-Verfahren zur Lösung von $Ax = b$ untergliedert sich somit in drei wesentliche Schritte:

Gauß-Elimination:

1. $PA = LR$
 $p = (p_1, \dots, p_n)$ Permutationsvektor

2. $Lc = Pb$
 Vorwärtseinsetzen: $i = 1, \dots, n$:

$$c_i = b_{p_i} - \sum_{k=1}^{i-1} l_{ik}c_k$$

3. $Rx = c$
 Rückwärtseinsetzen: $i = n, n-1, \dots, 1$:

$$x_i = \frac{1}{r_{ii}} \left(c_i - \sum_{k=i+1}^n r_{ik}x_k \right)$$

3.2.6 Aufwandsbestimmung

Ein wichtiger Aspekt bei der Analyse numerischer Verfahren ist es zu untersuchen, wie lange diese Verfahren in der Regel benötigen, um zum gewünschten Ergebnis zu gelangen. Da sich die Rechenzeiten von Computer zu Computer unterscheiden, orientiert man sich nicht an der *Rechenzeit*, sondern an der *Anzahl der Rechenoperationen*, die ein Algorithmus benötigt.

Das vorgestellte Gauß-Verfahren liefert nach endlich vielen Schritten ein Ergebnis, wobei die Anzahl der elementaren Rechenoperationen von der Dimension n der Matrix A abhängt. Multiplikationen und Divisionen sind sogenannte *wesentliche Rechenoperationen*. Die Auswertung einer wesentlichen Rechenoperation war im Allgemeinen noch vor einigen Jahren deutlich 'teurer' als eine Addition oder Subtraktion (rechnerintern wird nicht in Addition und Subtraktion unterschieden). Die Unterschiede verschmelzen jedoch mehr und mehr mit moderenen Rechnerarchitekturen.

Zur *Aufwandsbestimmung* zählen wir die Rechenoperationen einfach ab. Zuvor erinnern wir uns an:

$$\sum_{i=1}^n i = \frac{1}{2}(n+1)n$$

und

$$\sum_{i=1}^n i^2 = \frac{1}{6}n(n+1)(2n+1).$$

Anzahl der Operationen: (ohne Additionen)

$$\begin{aligned} 1. \quad PA &= LR \\ &[(n-1) + (n-1)^2] + [(n-2) + (n-2)^2] + \dots + [1 + 1^2] \\ &= \sum_{j=1}^{n-1} [(n-j) + (n-j)^2] = \frac{1}{2}n(n-1) + \frac{1}{6}n(n-1)(2n-1) \\ &= \frac{1}{3}(n^3 - n) \end{aligned}$$

$$\begin{aligned} 2. \quad Lc &= Pb \\ 1 + 2 + \dots + (n-1) &= \frac{1}{2}(n^2 - n) \end{aligned}$$

$$\begin{aligned} 3. \quad Rx &= c \\ 1 + 2 + \dots + n &= \frac{1}{2}(n^2 + n) \end{aligned}$$

Gesamtaufwand: $\frac{1}{3}n^3 + n^2 - \frac{1}{3}n$ Multiplikationen. (Additionen: $\frac{1}{3}n^3 + \frac{1}{2}n^2 - \frac{5}{6}n$)

Bemerkung 3.8. *Der Aufwand und damit die Rechenzeit steigt mit der dritten Potenz der Zahl der Unbekannten an: $\mathcal{O}(n^3)$.*

3.2.7 Algorithmus

Wir formulieren abschließend den Algorithmus.

Programm: $PA = LR$

für $j = 1, \dots, n$:

$$p_j = j$$

für $j = 1, \dots, n-1$:

Pivotsuche:

$$\max = |a_{jj}|, \quad r = j$$

für $i = j+1, \dots, n$:

falls $|a_{ij}| > \max$:

$$\max = |a_{ij}|, \quad r = i$$

falls $\max = 0$: STOP A singular

Zeilentausch:

falls $r > j$:

für $k = 1, \dots, n$:

$$h_r = a_{jk}, \quad a_{jk} = a_{rk}, \quad a_{rk} = h_r$$

$$h_i = p_j, \quad p_j = p_r, \quad p_r = h_i$$

Transformation:

für $i = j+1, \dots, n$:

$$a_{ij} = a_{ij}/a_{jj}$$

für $k = j+1, \dots, n$:

$$a_{ik} = a_{ik} - a_{ij}a_{jk}$$

3.3 Matrizen mit speziellen Eigenschaften

Besitzen Matrizen spezielle Eigenschaften, so kann es sich lohnen diese Eigenschaften gewinnbringend bei der Implementierung zu berücksichtigen.

3.3.1 Diagonaldominante Matrizen: Diagonalstrategie

Zunächst geben wir Bedingungen an, die die Durchführung der Gauß-Elimination ohne Pivotsuche ermöglichen (Diagonalstrategie).

Satz 3.9. Sei A eine (n, n) -Matrix, deren Hauptabschnittsmatrizen A_j regulär sind. Dann gibt es eine eindeutige Zerlegung

$$A = LR$$

L : linke Dreiecksmatrix mit $l_{jj} = 1$, $j = 1, \dots, n$,

R : reguläre rechte Dreiecksmatrix.

Beweis: Der Beweis wird durch Induktion über n geführt.

IA: Für $n = 1$ ist die Beh. trivial.

IV: Die Beh. sei richtig für $n - 1$.

IS: Für eine (n, n) -Matrix ist die folgende Zerlegung zu zeigen.

$$A = \left(\begin{array}{c|c} A_{n-1} & c \\ \hline a^T & a_{nn} \end{array} \right) = \left(\begin{array}{c|c} L_{n-1} & 0 \\ \hline l^T & 1 \end{array} \right) \left(\begin{array}{c|c} R_{n-1} & r \\ \hline 0 & r_{nn} \end{array} \right).$$

Nach der Induktionsvoraussetzung gibt es eine Zerlegung

$$A_{n-1} = L_{n-1}R_{n-1}.$$

Für die gesuchten $l, r \in \mathbb{R}^{n-1}$, $r_{nn} \in \mathbb{R}$ erhält man die Gleichungen

$$(3.13) \quad c = L_{n-1}r$$

$$(3.14) \quad l^T R_{n-1} = a^T \quad \Rightarrow \quad R_{n-1}^T l = a$$

$$(3.15) \quad l^T r + r_{nn} = a_{nn}.$$

Diese Gleichungen sind eindeutig auflösbar, da nach Voraussetzung L_{n-1} , R_{n-1} regulär sind. \diamond

Mit

$$D = \text{diag}(r_{jj}) = \text{diag}(a_{jj}^{(j)})$$

erhält man somit die Zerlegung

$$\boxed{A=LDR}, \quad l_{jj} = 1, \quad r_{jj} = 1.$$

Die Regularität der Hauptabschnittsmatrizen von A kann mit einer einfachen Bedingung für die Elemente a_{ij} von A nachgeprüft werden.

Definition 3.10 (Diagonaldominanz). Die Matrix A heißt diagonaldominant, wenn

$$|a_{ii}| > \sum_{\substack{k=1 \\ k \neq i}}^n |a_{ik}|, \quad (i = 1, \dots, n).$$

Satz 3.11. Bei einer diagonaldominanten Matrix A sind alle Hauptabschnittsmatrizen regulär, also existiert die LR-Zerlegung $A = LR$.

Beweis: Für die j -te Abschnittsmatrix A_j gelte

$$A_j x = 0 \quad \text{für ein } x \in \mathbb{R}^j.$$

Zu zeigen ist dann $x = 0$. Wäre

$$|x_r| = \max_{1 \leq i \leq j} |x_i| > 0,$$

so betrachten wir die r -te Gleichung

$$\sum_{i=1}^j a_{ri} x_i = 0.$$

Zusammen mit

$$|a_{rr}| > \sum_{\substack{k=1 \\ k \neq r}}^j |a_{rk}|$$

ergibt sich hieraus ein Widerspruch:

$$\begin{aligned} |a_{rr}| |x_r| &= \left| \sum_{\substack{k=1 \\ k \neq r}}^j a_{rk} x_k \right| \\ &\leq \sum_{k \neq r} |a_{rk}| |x_k| \\ &\leq \sum_{k \neq r} |a_{rk}| |x_r| < |a_{rr}| |x_r|. \end{aligned}$$

◇

Beispiel 3.12. Die bei der Berechnung von Spline-Funktionen (vgl. Kapitel zur Interpolation) auftretende tridiagonale Matrix

$$\begin{pmatrix} 4 & 1 & & 0 \\ 1 & 4 & 1 & \\ & \ddots & \ddots & \ddots \\ & & 1 & 4 & 1 \\ 0 & & & 1 & 4 \end{pmatrix}$$

ist diagonal dominant und damit LR-zerlegbar.

Bemerkung 3.13. *Spezielle Matrizen, die das Kriterium in Satz 3.9 erfüllen, sind die positiv definiten Matrizen (vgl. nächsten Abschnitt).*

3.3.2 Positiv definite Matrizen: Cholesky-Verfahren

Definition 3.14. *Eine (n, n) -Matrix A heißt symmetrisch, falls $A = A^T$ gilt.*

Definition 3.15. *Eine symmetrische (n, n) -Matrix A heißt positiv definit, falls*

$$(3.16) \quad x^T A x > 0, \quad \text{für alle } x \in \mathbb{R}^n, x \neq 0$$

gilt.

Die positive Definitheit scheint sehr einschränkend zu sein, dennoch ist sie in vielen Anwendungen erfüllt.

Bemerkung 3.16. *Für positiv definite Matrizen kann eine LR-Zerlegung ohne Pivoting durchgeführt werden.*

Satz 3.17. *Sei A positiv definit.*

1. *Alle Hauptabschnitt-Matrizen von A sind positiv definit und regulär. Insbesondere ist A regulär.*
2. *Es gibt genau eine linke Dreiecksmatrix L mit $l_{ii} > 0, i = 1, \dots, n$, so dass gilt*

$$\boxed{A = LL^T}$$

(Beachte: $l_{ii} = 1$ wird nicht gefordert)

Beweis:

zu 1: Übung

zu 2: Nach Satz 3.9 gibt es genau eine Zerlegung

$$\begin{aligned} A &= UV \\ U &= (u_{ik}) : \text{ linke Dreiecksmatrix, } u_{ii} = 1, \\ V &= (v_{ik}) : \text{ reguläre rechte Dreiecksmatrix} \end{aligned}$$

Sei

$$D = \begin{pmatrix} v_{11} & & 0 \\ & \ddots & \\ 0 & & v_{nn} \end{pmatrix}, \quad v_{ii} \neq 0.$$

Setze

$$\begin{aligned} R &= D^{-1}V : \text{ rechte Dreiecksmatrix, } r_{ii} = 1 \\ \Rightarrow A &= UDR, \quad A = A^T = R^T D^T U^T = R^T D U^T. \end{aligned}$$

Wegen der Eindeutigkeit der Zerlegung folgt:

$$R^T = U, \text{ d.h. } A = UDU^T = R^T DR.$$

Behauptung: D ist positiv definit, d.h. $v_{ii} > 0$.

Für alle $x \neq 0$ gilt:

$$\begin{aligned} 0 &< x^T Ax = x^T R^T DRx = (Rx)^T DRx \\ \Rightarrow 0 &< y^T Dy \quad \text{für alle } y \neq 0, \text{ da } R \text{ regulär,} \\ \Rightarrow D &\quad \text{positiv definit.} \end{aligned}$$

Mit

$$D^{1/2} := \begin{pmatrix} \sqrt{v_{11}} & & 0 \\ & \ddots & \\ 0 & & \sqrt{v_{nn}} \end{pmatrix}, \quad L := UD^{1/2}$$

gilt

$$\boxed{A = LL^T}$$

◇

Bemerkung 3.18. Ist L eine linke untere Dreiecksmatrix, so ist L^T eine rechte obere Dreiecksmatrix, d.h. für positiv definite Matrizen existiert eine LR-Zerlegung mit $R = L^T$. (Achtung: Hier sind die Diagonalelemente von L nicht normiert.)

Bemerkung 3.19. Offensichtlich reicht es aufgrund von Satz 3.9 für eine Cholesky-Zerlegung $A = LL^T$ die Matrix L zu bestimmen.

Leider ist das Cholesky-Verfahren nicht so anschaulich wie die Gauß-Elimination. Zur Bestimmung der Komponenten von L geht man induktiv Spaltenweise vor: Sei $L = (l_{i,j})$ die linke untere (n, n) -Dreiecksmatrix mit $A = LL^T$, die nach Satz 3.9 existiert und eindeutig ist. In Komponentenschreibweise ergibt sich

$$(3.17) \quad A = \begin{pmatrix} a_{1,1} & \cdots & a_{1,n} \\ \vdots & & \vdots \\ a_{n,1} & \cdots & a_{n,n} \end{pmatrix} = \begin{pmatrix} l_{1,1} & & 0 \\ \vdots & \ddots & \\ l_{n,1} & \cdots & l_{n,n} \end{pmatrix} \begin{pmatrix} l_{1,1} & \cdots & l_{n,1} \\ & \ddots & \vdots \\ 0 & & l_{n,n} \end{pmatrix}.$$

Offensichtlich gilt

$$(3.18) \quad a_{1,1} = l_{1,1} \cdot l_{1,1}, \quad \text{also } l_{1,1} = \sqrt{a_{1,1}},$$

d.h. $l_{1,1}$ läßt sich einfach berechnen. Ebenso gilt

$$(3.19) \quad a_{i,1} = l_{i,1} \cdot l_{1,1}, \quad \text{also} \quad l_{i,1} = \frac{a_{i,1}}{l_{1,1}}, \quad i = 2, \dots, n,$$

womit die erste Spalte von L bekannt ist (Dieses war der Induktionsanfang). Seien also die $l_{i,j}$, für $j \leq k-1$ bekannt (Induktionsvoraussetzung). Wir möchten als nächstes die Elemente der k -ten Spalte berechnen. Aus (3.17) ergibt sich

$$(3.20) \quad a_{k,k} = l_{k,1}^2 + \dots + l_{k,k}^2,$$

und somit aufgrund der Eindeutigkeit von L

$$(3.21) \quad l_{k,k} = \sqrt{a_{k,k} - \sum_{j=1}^{k-1} l_{k,j}^2}.$$

Ebenso ergibt sich aus (3.17)

$$(3.22) \quad a_{i,k} = \sum_{j=1}^k l_{i,j} l_{k,j}$$

und damit

$$(3.23) \quad l_{i,k} = \frac{1}{l_{k,k}} \left(a_{i,k} - \sum_{j=1}^{k-1} l_{i,j} l_{k,j} \right), \quad i \geq k+1.$$

Auf diesem Wege können wir die vollständige Matrix L bestimmen. Der Zusammenhang mit dem Ausgangsproblem ist durch

$$Ax = LL^T x = Lc = b$$

gegeben. Eine Abschätzung des Aufwandes ergibt, daß außer n Quadratwurzeln noch

$$(3.24) \quad \boxed{\frac{1}{6}n^3 + \mathcal{O}(n^2)}$$

Rechenoperationen durchgeführt werden müssen.

Bemerkung 3.20. Auch wenn Gauß- und Cholesky-Verfahren beide die Ordnung $\mathcal{O}(n^3)$ besitzen, so ist das Cholesky-Verfahren für große n dennoch etwa doppelt so schnell wie das Gauß-Verfahren, man vergleiche die jeweiligen Faktoren vor n^3 .

Die Lösung des LGS $Ax = b$ nach der Methode von Cholesky erfolgt in den drei Schritten

$$(3.25) \quad \begin{array}{l} 1. \quad A = LL^T \quad : \quad \text{Cholesky-Zerlegung} \\ 2. \quad Lc = b \quad : \quad \text{Vorwärtseinsetzen} \\ 3. \quad L^T x = c \quad : \quad \text{Rückwärtseinsetzen} \end{array}$$

Bei positiv definiten Matrizen A sind die Hauptdiagonalelemente $a_{ii} = e_i^T A e_i > 0$ d.h. positiv. Darüberhinaus kann man leicht zeigen, dass diagonal-dominante Matrizen (vgl. Definition 3.10) mit $a_{ii} > 0$, d.h.

$$a_{ii} > \sum_{k \neq i} |a_{ik}| \quad (i = 1, \dots, n),$$

positiv definit sind.

Für eine positiv definite Matrix A ist die Reduktion der quadratischen Form $x^T A x$ auf eine Summe von Quadraten (im Körper der reellen Zahlen) möglich:

$$\begin{aligned} x^T A x &= x^T L L^T x = (L^T x)^T (L^T x) \\ &= \sum_{j=1}^n \left(\sum_{k=j}^n l_{kj} x_k \right)^2. \end{aligned}$$

Zusätzlich ergibt sich für die Hauptabschnittsmatrizen (Hauptmenoren):

$$\det A = \prod_{j=1}^n l_{jj}^2 = \prod_{j=1}^n a_{jj}^{(j)} > 0, \quad \det \begin{pmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{pmatrix} = \prod_{j=1}^k l_{jj}^2 > 0.$$

Folgerung 3.21. Eine symmetrische Matrix A ist genau dann positiv definit, wenn

$$\det \begin{pmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{pmatrix} > 0 \quad \text{für } k = 1, \dots, n.$$

Beispiel 3.22. Die bei der Diskretisierung von Randwertproblemen für Differentialgleichungen auftretende Matrix

$$A_n = \left(\begin{array}{cccc} 2 & -1 & & 0 \\ -1 & 2 & -1 & \\ & \ddots & \ddots & \ddots \\ & & -1 & 2 & -1 \\ 0 & & & -1 & 2 \end{array} \right) \Bigg\} n$$

Besonders einfach zu behandelnde LGS erhält man für tridiagonale Matrizen A der Bandbreite $m = 1$. Zu lösen sei das LGS $Ax = d$ mit

$$A = \begin{pmatrix} a_1 & b_1 & & & \\ c_2 & a_2 & b_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-1} & a_{n-1} & b_{n-1} \\ & & & c_n & a_n \end{pmatrix}, \quad d = \begin{pmatrix} d_1 \\ \vdots \\ d_n \end{pmatrix}$$

Es existiere die LR -Zerlegung $A = LR$. Nach Satz 3.25 sind L und R bidiagonal und können in der Form angesetzt werden

$$L = \begin{pmatrix} 1 & & & & \\ l_2 & 1 & & & \\ & l_3 & 1 & & \\ & & \ddots & \ddots & \\ & & & l_n & 1 \end{pmatrix}, \quad R = \begin{pmatrix} m_1 & r_1 & & & \\ & m_2 & r_2 & & \\ & & & \ddots & \\ & & & & m_{n-1} & r_{n-1} \\ & & & & & m_n \end{pmatrix}.$$

Die Ausmultiplikation $A = LR$ führt auf die Beziehung $r_i = b_i$, $i = 1, \dots, n$ und den folgenden Algorithmus zur Lösung von $Ax = d$:

$$(3.26) \quad \begin{array}{l} \underline{A = LR} \\ m_1 = a_1 \\ \text{für } i = 2, \dots, n : \\ \quad l_i = c_i / m_{i-1} \\ \quad m_i = a_i - l_i \cdot b_{i-1} \\ \underline{Ly = d} \\ y_1 = d_1 \\ \text{für } i = 2, \dots, n : \\ \quad y_i = d_i - l_i \cdot y_{i-1} \\ \underline{Rx = y} \\ x_n = y_n / m_n \\ \text{für } i = n-1, n-2, \dots, 1 : \\ \quad x_i = (y_i - b_i \cdot x_{i+1}) / m_i \end{array}$$

3.4 Fehleranalyse und Fehlerbehandlung

3.4.1 Fehlerabschätzungen

Wie in der Einleitung ausgeführt können Computer nicht alle reellen Zahlen darstellen, daher werden die meisten Zahlen intern gerundet. Als Konsequenz ergeben sich Rundungsfehler. Selbst wenn Eingabedaten und das Ergebnis eines Algorithmus frei von Rundungsfehlern wären, können Zwischenergebnisse gerundet worden sein. Aus diesem Grund wird in der Regel nicht die Lösung x des Gleichungssystems $Ax = b$ berechnet, sondern die Lösung \tilde{x} eines 'benachbarten' oder 'gestörten' Gleichungssystems

$$(3.27) \quad (A + \Delta A)\tilde{x} = b + \Delta b$$

Δb : Fehler im Vektor b (Residuum),

ΔA : Fehler in der Matrix A ,

$\Delta x := \tilde{x} - x$: Fehler der Näherungslösung

Um die nachfolgende Analyse durchzuführen benötigen wir den Begriff der *zugeordneten Matrixnorm*. Wir erinnern zunächst an verschiedene Normen für Vektoren $x \in \mathbb{R}^n$. In dieser Vorlesung verwenden wir üblicherweise

$$\|x\|_2 = \sqrt{x^T x} = \sqrt{\sum_{i=1}^n x_i^2}, \quad (\text{euklidische Norm oder 2-Norm})$$

welche wir meistens einfach mit $\|\cdot\|$ bezeichnen. Weitere Normen sind die

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \quad (1\text{-Norm})$$

oder die

$$\|x\|_\infty = \max_{i=1, \dots, n} |x_i|, \quad (\text{Maximumsnorm oder } \infty\text{-Norm}).$$

Für alle Vektornormen kann man eine zugeordnete Matrixnorm definieren.

Definition 3.27. Sei A eine (n, n) -Matrix und $\|\cdot\|_p$ eine Vektornorm im \mathbb{R}^n . Die Zahl

$$\|A\|_p := \max_{\|x\|_p=1} \|Ax\|_p$$

heißt die der Vektor-Norm $\|\cdot\|_p$ zugeordnete Matrixnorm.

Bemerkung 3.28. Wir bezeichnen nachfolgend die Vektornormen und die ihnen zugeordneten Matrixnormen mit dem gleichen Symbol.

Sei nachfolgend $A = (a_{i,j})$.

Beispiel 3.29.

$$\|A\|_1 := \max_{j=1,\dots,n} \sum_{i=1}^n |a_{i,j}|, \quad (\text{Spaltensummennorm})$$

Beispiel 3.30.

$$\|A\|_\infty := \max_{i=1,\dots,n} \sum_{j=1}^n |a_{i,j}|, \quad (\text{Zeilensummennorm})$$

Beispiel 3.31.

$$\|A\|_2 := \sqrt{\rho(A^T A)}, \quad (\text{Spektralnorm})$$

wobei $\rho(B)$ den Betrag des betragsgrößten Eigenwert einer symmetrischen Matrix B bezeichnet.

Definition 3.32. Als Kondition von A bzgl. einer Matrixnorm $\|\cdot\|_p$ bezeichnen wir die Zahl

$$\text{cond}_p(A) := \|A\|_p \|A^{-1}\|_p.$$

Satz 3.33 (Fehleranalyse). Sei x die eindeutige Lösung von $Ax = b$, und ΔA , Δb Störungen von A , b mit

$$(3.28) \quad q = \text{cond}(A) \frac{\|\Delta A\|}{\|A\|} < 1.$$

Dann ist auch das gestörte System

$$(3.29) \quad (A + \Delta A)(x + \Delta x) = b + \Delta b$$

eindeutig lösbar und es gilt

$$(3.30) \quad \boxed{\frac{\|\Delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - q} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right)}$$

Beweis: Sei $x + \Delta x$ Lösung von (3.29). Nach Ausmultiplizieren ergibt sich

$$Ax + A\Delta x + \Delta Ax + \Delta A\Delta x - b - \Delta b = 0$$

und weiter wegen $Ax - b = 0$

$$A\Delta x = \Delta b - \Delta Ax - \Delta A\Delta x$$

und somit

$$\Delta x = -A^{-1}(-\Delta b + \Delta Ax + \Delta A \Delta x).$$

Für *verträgliche Matrixnormen* folgt hieraus (Dreiecksungleichung)

$$\begin{aligned} \|\Delta x\| &\leq \|A^{-1}\| \cdot \|\Delta b + \Delta Ax + \Delta A \Delta x\| \\ &\leq \|A^{-1}\| (\|\Delta b\| + \|\Delta A\| \cdot \|x\| + \|\Delta A\| \cdot \|\Delta x\|) \end{aligned}$$

und weiter

$$(1 - \|A^{-1}\| \cdot \|\Delta A\|) \|\Delta x\| \leq \|A^{-1}\| (\|\Delta b\| + \|\Delta A\| \cdot \|x\|)$$

Aufgrund der Voraussetzung (3.28) mit $q < 1$ folgt

$$\|\Delta x\| \leq \frac{\|A^{-1}\|}{1 - q} (\|\Delta b\| + \|\Delta A\| \cdot \|x\|)$$

als absolutem Fehler. Es ergibt sich weiter

$$(3.31) \quad \frac{\|\Delta x\|}{\|x\|} \leq \frac{\|A^{-1}\|}{1 - q} \left(\frac{\|\Delta b\|}{\|x\|} + \|\Delta A\| \right).$$

Aus $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$ folgt

$$\|x\| \geq \frac{\|b\|}{\|A\|}$$

und somit aus (3.31)

$$\frac{\|\Delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - q} \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right)$$

und hieraus die Behauptung. ◇

Die Zahl $\text{cond}(A)$ hat also die Bedeutung eines Verstärkungsfaktors und mißt die Empfindlichkeit der Lösung x gegenüber Störungen in A und b . Das LGS $Ax = b$ heißt schlecht konditioniert, wenn $\text{cond}(A) \gg 1$.

Beispiel 3.34. *Auswirkung schlechter Kondition:*

$$\begin{aligned} A &= \begin{pmatrix} 1 & 1 \\ 1 & 0.99 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad x = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ \Delta A &= \begin{pmatrix} 0.01 & 0.01 \\ 0 & 0 \end{pmatrix}, \quad \Delta b = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad x + \Delta x = \begin{pmatrix} 200/101 \\ -100/101 \end{pmatrix} \end{aligned}$$

Obwohl der Fehler in A bei 1% liegt, haben $x, x + \Delta x$ nichts mehr miteinander zu tun.

Erklärung:

$$\|A\|_\infty = 2, \quad A^{-1} = \begin{pmatrix} -99 & 100 \\ 100 & -100 \end{pmatrix}, \quad \|A^{-1}\|_\infty = 200, \quad \text{cond}_\infty(A) = 400!$$

Geometrisch: Die Zeilenvektoren a_1, a_2 von A haben beinahe die gleiche Richtung.

3.4.2 Skalierung

Die Kondition eines Problems kann ggf. durch Skalierung der Matrix A verbessert werden. Unter Skalierung versteht man den Übergang

$$A \rightarrow DA, \quad D = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{pmatrix}, \quad d_i \neq 0,$$

d.h. die i -te Zeile von A wird mit d_i multipliziert. Die optimale Wahl einer Diagonalmatrix D , welche $\text{cond}(DA)$ möglichst klein macht, erhält man durch den folgenden Satz (ohne Beweis):

Satz 3.35 (Van der Sluis). Für $A = (a_{ik})$ sei

$$\sum_{k=1}^n |a_{ik}| = 1, \quad i = 1, \dots, n \quad (\text{insbesondere } \|A\|_\infty = 1).$$

Dann gilt für jede Diagonalmatrix D mit $\det D \neq 0$

$$\text{cond}_\infty(DA) \geq \text{cond}_\infty(A).$$

Folgerung 3.36. Für eine beliebige reguläre Matrix $A = (a_{ik})$ ist mit der Skalierung

$$D = \text{diag}(d_i), \quad d_i := \left(\sum_{k=1}^n |a_{ik}| \right)^{-1}$$

die Kondition $\text{cond}_\infty(DA)$ möglichst klein.

3.4.3 Iterative Nachverbesserung

Unabhängig von einer schlechten Kondition der Matrix A liefern numerische Verfahren zur Lösung linearer Gleichungssysteme nicht die exakte Lösung. Wir erhalten dann lediglich eine Näherungslösung, die unseren Anforderungen aber möglicherweise nicht hinreichend gerecht wird. Mit einem kleinen Trick läßt sich die berechnete Näherungslösung aber dennoch weiter verbessern:

x sei exakte Lösung von $Ax = b$

\tilde{x} sei irgendeine Näherungslösung, z.B. aus Gauß-Algorithmus.

Verbesserung in drei Schritten:

1. Berechne $r := b - A\tilde{x}$ "Residuum"
2. Bestimme Δx aus $A\Delta x = r$ "Korrektur"
3. Berechne $x' := \tilde{x} + \Delta x$

Begründung:

$$\begin{aligned} x' &= \tilde{x} + \Delta x = \tilde{x} + A^{-1}r = \tilde{x} + A^{-1}(b - A\tilde{x}) \\ &= \tilde{x} + x - \tilde{x} = x \end{aligned}$$

In der praktischen Anwendung/Implementierung bewirken Rundungsfehler, dass i.A. $x' \neq x$ ist. Die Verbesserung kann wiederholt werden, solange Δx nur mit einer Stelle korrekt berechnet wird. In diesem Fall ist x' besser als \tilde{x} .

Bei der algorithmischen Durchführung sind einige Dinge zu beachten:

- Wurde \tilde{x} durch den Gauß-Algorithmus gewonnen, so erfüllt \tilde{x} das Gleichungssystem meist sehr gut, d. h.

$$b \approx A\tilde{x} \quad \Rightarrow \quad r = b - A\tilde{x} \quad \text{auslöschungsgefährdet!}$$

Um dem entgegenzuwirken berechnet man im Schritt 1. das Residuum mit doppelter Genauigkeit.

- Bei der Auflösung von $A \cdot \Delta x = r$ benutze man die bereits berechnete *LR-Zerlegung*.
- Rundungsfehler im 3. Schritt begrenzen i.A. die erreichbare Genauigkeit.

Beispiel 3.37. *Sechs Dezimalstellen; unterstrichene Stellen sind falsch*

$$\begin{bmatrix} 0.566012 & 0.765456 \\ 0.389953 & 0.527611 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.395102 \\ 0.272744 \end{bmatrix}$$

$$\text{exakt: } x_1 = -2.20227459\dots$$

$$x_2 = 2.14462470\dots$$

$$\text{Gauß } \tilde{x}_1 = -2.19453$$

$$\tilde{x}_2 = 2.13889$$

1. Nachverbesserung 2. Nachverbesserung

$$x'_1 = -2.20226 \quad x''_1 = -2.20227$$

$$x'_2 = 2.14461 \quad x''_2 = 2.14462$$

Bemerkung 3.38. *Meist reicht nur ein Schritt, um das Resultat deutlich zu verbessern.*

3.5 Die QR-Zerlegung einer Matrix, das Verfahren von Householder

3.5.1 Einleitung und Motivation

Sei A eine (n, n) -Matrix (reell, nicht notwendig regulär).

Bei der LR-Zerlegung (ohne Pivotsuche) hatten wir das Ergebnis:

$$\boxed{A = LR}$$

L : linke Dreiecksmatrix

R : rechte Dreiecksmatrix.

Bei der QR-Zerlegung suchen wir hingegen eine Zerlegung der Form:

$$\boxed{A = QR}$$

Q : orthogonal, d. h. $Q^T Q = I$,

R : rechte Dreiecksmatrix.

Motivation zur QR-Zerlegung:

Zur Lösung des LGS $Ax = b$ erzeugt man bei der LR-Zerlegung und Gauß-Elimination eine Sequenz

$$(A, b) = (A^{(1)}, b^{(1)}) \rightarrow \dots \rightarrow (A^{(j)}, b^{(j)}) \rightarrow \dots \rightarrow (A^{(n)}, b^{(n)}) = (R, c)$$

$$(A^{(j+1)}, b^{(j+1)}) = L_j(A^{(j)}, b^{(j)}).$$

Sei $\varepsilon^{(j)}$ der Rundungsfehler bei der Berechnung von $(A^{(j)}, b^{(j)})$. Für irgendeine Vektornorm $\|x\|$ gilt nach Satz 3.33 die Abschätzung

$$\frac{\|\Delta x\|}{\|x\|} \leq \sum_{j=1}^n \varepsilon^{(j)} \operatorname{cond}(A^{(j)}).$$

Die Gauß-Elimination ist daher nicht gutartig, falls

$$\operatorname{cond}(A^{(j)}) \gg \operatorname{cond}(A^{(1)}) = \operatorname{cond}(A).$$

Idee: Wähle Matrix Q_j mit Übergang

$$(A^{(j+1)}, b^{(j+1)}) = Q_j(A^{(j)}, b^{(j)}), \quad \operatorname{cond}(A^{(j+1)}) = \operatorname{cond}(A^{(j)}).$$

Dazu beschränken wir uns auf die euklidische Norm

$$\|x\| = \|x\|_2 = (x^T x)^{1/2}, \quad \|A\| = \|A\|_2$$

und notieren eine später zu benutzende Hilfsaussage:

Hilfssatz 3.39. *Sei Q orthogonal, dann gilt:*

(i) $\|Q\|_2 = 1$

(ii) $\|QA\|_2 = \|A\|_2$ für alle A

(iii) Wenn A regulär ist, gilt $\operatorname{cond}_2(QA) = \operatorname{cond}_2(A)$.

Beweis: Übung

3.5.2 Householdermatrizen

Sei $w \in \mathbb{R}^n$ mit $w^T w = 1$ und sei die *Householdermatrix* Q definiert durch

$$Q := I - 2ww^T, \quad ww^T = (w_i \cdot w_k).$$

Dann hat die so konstruierte Matrix Q folgende Eigenschaften:

Q ist symmetrisch:

$$Q^T = I - 2(ww^T)^T = I - 2ww^T = Q$$

Q ist orthogonal wegen $w^T w = 1$:

$$\begin{aligned} Q^T Q &= (I - 2ww^T)(I - 2ww^T) \\ &= I - 2ww^T - 2ww^T + 4ww^T ww^T = I \end{aligned}$$

Für $x \in \mathbb{R}^n$ bedeutet

$$Qx = (I - 2ww^T)x = x - 2(w^T x)w$$

eine Spiegelung an der Hyperebene

$$H = \{z \in \mathbb{R}^n \mid w^T z = 0\} :$$

$x = y + z$ mit $w^T z = 0$, (y aus dem orthogonalen Komplement)

$$= \alpha w + z$$

$$\Rightarrow w^T x = \alpha w^T w + w^T z = \alpha$$

$$\Rightarrow Qx = x - 2(w^T x)w = x - 2\alpha w = \alpha w + z - 2\alpha w = -\alpha w + z = -y + z$$

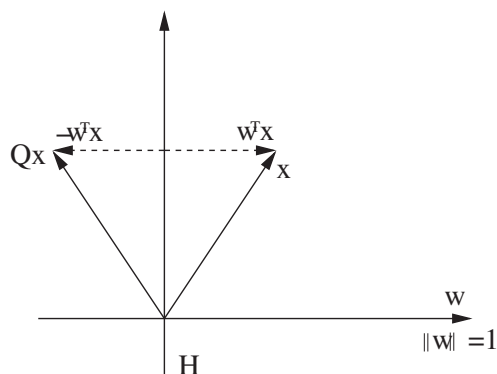


Abbildung 3.3: Spiegelung an Hyperebene.

Problem: Sei $x = (x_1, \dots, x_n)^T \neq 0$ vorgegeben. Bestimme $w \in \mathbb{R}^n$, $w^T w = 1$, mit

$$Qx = k e_1, \quad k \in \mathbb{R}.$$

In diesem Fall ist Q eine spezielle Spiegelung an einer Hyperebene, vgl. die nachfolgende Abbildung:

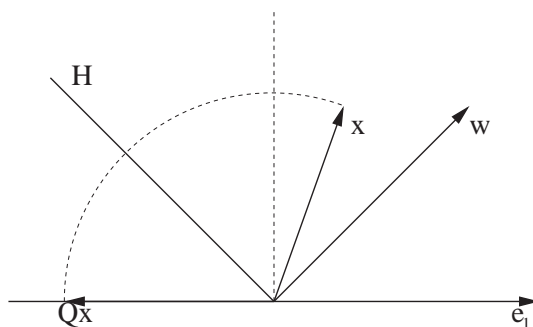


Abbildung 3.4: Spiegelung.

Analytische Berechnung von Q : (Für $Qx = ke_1$)

$$\begin{aligned} \Rightarrow |k| &= \|Qx\| = \|x\|, & k &= \pm\|x\|. \\ Qx &= (I - 2ww^T)x = x - 2w(w^Tx) = ke_1 \\ \Rightarrow w &= \frac{x - ke_1}{2(w^Tx)} = c(x - ke_1) & (w \text{ ist Vielfaches vom Vektor } x - ke_1) \\ \stackrel{\|w\|=1}{\Rightarrow} w &= \frac{x - ke_1}{\|x - ke_1\|} \end{aligned}$$

An dieser Stelle ist lediglich das Vorzeichen von $k = \pm\|x\|$ noch unbekannt. Aus Stabilitätsgründen (Vermeidung von Auslöschung) wählen wir k in geeigneter Weise. Es ist

$$\|x - ke_1\| = ((x_1 - k)^2 + x_2^2 + \dots + x_n^2)^{1/2}.$$

Keine Auslöschung tritt auf für

$$\begin{aligned} k &= -\text{sign}(x_1)\|x\|, & (x_1 - k)^2 &= (|x_1| + \|x\|)^2. \\ \Rightarrow \|x - ke_1\|^2 &= \|x\|^2 + 2\|x\||x_1| + \|x\|^2 = 2\|x\|(\|x\| + |x_1|) \end{aligned}$$

Insgesamt erhalten wir

$$(3.32) \quad \begin{array}{l} Q = I - 2ww^T = I - 2\frac{(x - ke_1)(x - ke_1)^T}{\|x - ke_1\|^2} \\ \quad = I - \beta uu^T \\ k = -\text{sign}(x_1)\|x\|, \quad \beta = \frac{1}{\|x\|(|x_1| + \|x\|)} \\ u := x - ke_1 = \begin{pmatrix} \text{sign}(x_1)(|x_1| + \|x\|) \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \end{array}$$

Householder-Transformation

3.5.3 QR-Zerlegung/Verfahren von Householder

Zur Zerlegung der Matrix A bilden wir die Sequenz

$$\begin{aligned} A &= A^{(1)} \rightarrow A^{(2)} \rightarrow \dots \rightarrow A^{(n)} = R, \\ A^{(j+1)} &= Q_j A^{(j)}, \quad Q_j \text{ orthogonal.} \end{aligned}$$

j-ter Schritt ($j \geq 1$): Sei

$$A^{(j)} = \left(\begin{array}{ccc|ccc} * & \cdots & * & * & \cdots & * \\ & & \vdots & \vdots & & \vdots \\ & & \vdots & \vdots & & \vdots \\ 0 & & * & * & \cdots & * \\ \hline & & & a_{jj}^{(j)} & \cdots & a_{jn}^{(j)} \\ & & & \vdots & & \vdots \\ & & & a_{nj}^{(j)} & \cdots & a_{nn}^{(j)} \end{array} \right) \left. \begin{array}{l} \\ \\ \\ \\ \\ \\ \end{array} \right\} \begin{array}{l} j-1 \\ \\ \\ n-j+1 \end{array}$$

$$x := (a_{jj}^{(j)}, \dots, a_{nj}^{(j)})^T \in \mathbb{R}^{n-j+1}.$$

1. Fall: $x = 0$: A ist singulär (Beweis!), setzt $Q_j = I$
2. Fall: $x \neq 0$: Bestimme nach (3.32) die orthogonale $(n-j+1, n-j+1)$ -Matrix \tilde{Q}_j mit

$$\tilde{Q}_j \begin{pmatrix} a_{jj}^{(j)} \\ \vdots \\ a_{nj}^{(j)} \end{pmatrix} = k \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{n+1-j}.$$

Setzen wir nun jeweils

$$Q_j = \begin{pmatrix} I_{j-1} & 0 \\ 0 & \tilde{Q}_j \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad \text{orthogonal, symmetrisch}$$

so erhalten wir nach n Schritten

$$(3.33) \quad R := A^{(n)} = Q_{n-1} Q_{n-2} \cdots Q_1 A.$$

Definieren wir die orthogonale Matrix

$$Q := (Q_{n-1} \cdots Q_1)^{-1} = Q_1 \cdots Q_{n-1}, \quad (\text{da } Q_j \text{ orthogonal, symmetrisch}) \\ \Rightarrow A = QR$$

Satz 3.40 (QR-Zerlegung). *Zu jeder (n, n) -Matrix A existiert eine orthogonale (n, n) -Matrix Q und eine rechte Dreiecksmatrix R mit*

$$A = QR.$$

Ist A regulär, so ist R regulär.

Bei einer regulären Matrix A bildet man zur Lösung des LGS $Ax = b$ analog zu (3.33) den Ausdruck

$$c := b^{(n)} = Q_{n-1} \cdots Q_1 b$$

und löst dann das gestaffelte LGS $Rx = c$.

Anzahl der Operationen: $\approx \frac{2}{3}n^3$ Jedoch sind keine zusätzlichen Permutationsmatrizen notwendig.

3.5.4 Erweiterungen

Die QR -Zerlegung kann unmittelbar auf nichtquadratische (m, n) -Matrizen A ($m > n$) erweitert werden. Hier bildet man eine Sequenz

$$A^{(j+1)} = Q_j A^{(j)} \quad (j \geq 1), \quad A^{(1)} = A$$

Q_j orthogonale (m, m) -Matrix.

Wegen $m > n$ erhält man nach n Schritten

$$(3.34) \quad \begin{array}{l} A^{(n+1)} = \tilde{Q}A = \underbrace{\begin{pmatrix} R \\ 0 \end{pmatrix}}_n \begin{array}{l} \}n \\ \}m-n \end{array} \\ \tilde{Q} = Q_n \cdots Q_1 \quad \text{orthogonal,} \\ R = \begin{pmatrix} r_{11} & \cdots & r_{1n} \\ & \ddots & \vdots \\ 0 & & r_{nn} \end{pmatrix} \quad \text{obere Dreiecksmatrix} \end{array}$$

Praktische Durchführung mit (3.32)

$$A^{(j)} = \underbrace{\begin{pmatrix} * & * \\ 0 & \tilde{A}^{(j)} \end{pmatrix}}_n \begin{array}{l} \}j-1 \\ \}m-j+1 \end{array}, \quad A^{(1)} = A$$

$$Q_j = \underbrace{\begin{pmatrix} I_{j-1} & 0 \\ 0 & \tilde{Q}_j \end{pmatrix}}_m \begin{array}{l} \}j-1 \\ \}m-j+1 \end{array}$$

$$\tilde{Q}_j = I - \beta_j u_j u_j^T, \quad j = 1, \dots, n,$$

wobei nach (3.32) gilt:

$$\begin{aligned}x_j &= (a_{jj}^{(j)}, \dots, a_{mj}^{(j)})^T \in \mathbb{R}^{m-j+1}, \\k_j &= -\text{sign}((x_j)_1) \|x_j\|, \\\beta_j &= \frac{1}{\|x_j\|(|x_{j1}| + \|x_j\|)}, \\u_j &= x_j - k_j e_j. \\\tilde{Q}_j \tilde{A}^{(j)} &= \tilde{A}^{(j)} - u_j s_j^T \\s_j^T &= \beta_j u_j^T \tilde{A}^{(j)}.\end{aligned}$$

$$\text{d.h. } (u_j s_j^T)_{i,k} = a_{ij} \beta_j \sum_{l=j}^m a_{lj} a_{lk}$$

Programm QR(A,d)

Die u_j stehen spaltenweise im linken Teil von A , $R \setminus \text{diag}(R)$ steht im rechten Teil von A , $\text{diag}(R)$ steht auf $d = (d_1, \dots, d_n)$.

für $j = 1, \dots, n$:

$$xnorm = \left(\sum_{i=j}^m a_{ij}^2 \right)^{1/2}$$

falls $xnorm=0$: *STOP*

$$d_j = -\text{sign}(a_{jj}) \cdot xnorm$$

$$beta = 1/(xnorm(|a_{jj}| + xnorm))$$

$$a_{jj} = a_{jj} - d_j$$

für $k = j + 1, \dots, n$:

$$s = beta \cdot \sum_{i=j}^m a_{ij} a_{ik}$$

für $i = j, \dots, m$:

$$a_{ik} = a_{ik} - a_{ij} \cdot s.$$

3.6 Lineare Ausgleichsrechnung, diskrete Approximation

3.6.1 Normalgleichung

Ausgleichsrechnungen sind für viele praktische Zwecke besonders wichtig.

Beispiel 3.41. Wir haben bei einem Experiment für die Eingabewerte t_1, \dots, t_m Messwerte s_1, \dots, s_m erhalten. Aufgrund theoretischer Überlegungen (etwa physikalische Gesetze) kennt man eine Funktion $f(t)$, für die $f(t_i) = s_i$ gelten soll.

Die Funktion f hängt aber in der Regel von unbekanntem Parametern x_1, \dots, x_n ab; wir schreiben $f(t; x)$ für $x = (x_1, \dots, x_n)^T$ um dies zu betonen. Beispielsweise könnte f durch eine Parabel

$$(3.35) \quad f(t; x) = x_1 + x_2 t + x_3 t^2$$

gegeben sein. In der Regel hat man mehr Messwerte als Parameter ($m > n$) und es liegen Messfehler vor, so dass der naheliegende Versuch durch Lösen des m dimensionalen (i.A. nichtlinearen) Gleichungssystems

$$(3.36) \quad \begin{aligned} f(t_1; x) &= s_1 \\ &\vdots \\ f(t_m; x) &= s_m \end{aligned}$$

den n dimensionalen Lösungsvektor zu bestimmen scheitern muß.

In vielen praktischen Fällen ist das Gleichungssystem (3.36) linear, so z.B. auch für die obige Parabel und wir erhalten aus (3.36) ein überbestimmtes (i.A. nicht lösbares) LGS

$$\boxed{Cx = s},$$

mit $C \in \mathbb{R}^{m \times n}$, $m \geq n$ und $s \in \mathbb{R}^m$.

Beispiel 3.42. Für unser Beispiel einer Parabel in (3.35) erhalten wir etwa

$$C = \begin{pmatrix} 1 & t_1 & t_1^2 \\ \vdots & \vdots & \vdots \\ 1 & t_m & t_m^2 \end{pmatrix} \quad \text{und} \quad s = \begin{pmatrix} s_1 \\ \vdots \\ s_m \end{pmatrix}.$$

Anstelle also den (vermutlich vergeblichen Versuch) zu unternehmen, eine exakte Lösung x des Systems $Cx = s$ zu finden, begnügen wir uns damit, ein x zu finden, so daß Cx 'möglichst nahe' bei s liegt. Als Ersatzproblem betrachtet man das Optimierungsproblem

$$(3.37) \quad \min_{x \in \mathbb{R}^n} \|s - Cx\|_2^2$$

Um eine Lösung von (3.37) zu bestimmen, bildet man die erste Ableitung, die in einem Minimum zwangsweise gleich Null sein muß.

$$0 = \frac{\partial \|s - Cx\|_2^2}{\partial x} = \frac{\partial [(s - Cx)^T (s - Cx)]}{\partial x} = 2C^T Cx - 2C^T s.$$

Dies ergibt die *Normalgleichung*

$$(3.38) \quad \boxed{C^T C x = C^T s},$$

die nach Definition von $A := C^T C$ und $b := C^T s$ dem LGS

$$Ax = b$$

entspricht.

Für eine beliebige Matrix C ist die Lösung von (3.38) nicht eindeutig und es gilt:

Satz 3.43. *Das lineare Ausgleichsproblem (3.37) besitzt mindestens eine Lösung x_0 , d.h. $C^T C x_0 = C^T s$.*

Beweisidee:

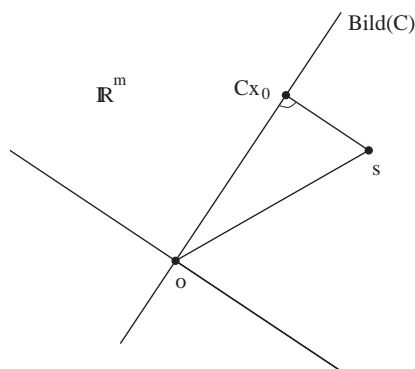


Abbildung 3.5: Lösung im Unterraum.

Nach linearer Algebra gilt die Zerlegung $\mathbb{R}^m = \text{Bild}(C) \oplus \text{Kern}(C^T)$. Daher kann $s \in \mathbb{R}^m$ zerlegt werden in

$$s = y + r, \quad y = Cx_0, \quad C^T r = 0,$$

und es folgt

$$C^T s = C^T y = C^T C x_0,$$

d.h. x_0 ist eine Lösung des linearen Ausgleichsproblems. Wegen $\text{Kern}(C^T C) = \text{Kern}(C)$ prüft man nun leicht nach, dass

$$x_0 + \text{Kern}(C)$$

die Gesamtheit der Lösungen ist:

$$\begin{aligned} C^T C(x_0 + \text{Kern}(C)) &= C^T s \\ \Leftrightarrow C^T C(x_0 + w) &= C^T s \text{ mit } Cw = 0 \\ \Leftrightarrow C^T C x_0 &= C^T s \end{aligned}$$

◇

3.6.2 Numerische Lösung

Von nun an sei $\text{rang}(C) = n < m$: Die Matrix $C^T C$ ist dann positiv definit, und die Normalgleichung $C^T C x = C^T s$ kann im Prinzip mit dem Cholesky-Verfahren gelöst werden. Dieses Verfahren hat jedoch zwei Nachteile:

- (1) $C^T C$ ist schwierig auszurechnen, z. B.:

$$C = \begin{pmatrix} 1 & 1 \\ \varepsilon & 0 \\ 0 & \varepsilon \end{pmatrix}, \quad C^T C = \begin{pmatrix} 1 + \varepsilon^2 & 1 \\ 1 & 1 + \varepsilon^2 \end{pmatrix}.$$

Mit $\varepsilon = \frac{1}{2}\sqrt{\text{eps}}$ ist auf der Maschine

$$\text{rd}(C^T C) = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \text{ singular.}$$

- (2) Die Kondition und damit die Empfindlichkeit gegenüber Störungen in C, y beträgt

$$\text{cond}(C^T C).$$

Beide Nachteile können mit der Householder-Transformation vermieden werden. Nach (3.34) gibt es eine QR-Zerlegung mit

$$QC = \underbrace{\begin{pmatrix} R \\ 0 \end{pmatrix}}_n \left. \begin{array}{l} \} n \\ \} m - n \end{array} \right.$$

R obere reguläre (n, n) -Dreiecksmatrix

Q orthogonale (m, m) -Matrix

Mit

$$Qs = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix}, \quad h_1 \in \mathbb{R}^n, \quad h_2 \in \mathbb{R}^{m-n}$$

berechnet man

$$\begin{aligned} \|s - Cx\|_2^2 &= \|Q(s - Cx)\|_2^2 \\ &= \|h_1 - Rx\|_2^2 + \|h_2\|_2^2. \end{aligned}$$

Dieser Ausdruck wird minimal für $x \in \mathbb{R}^n$ mit $Rx = h_1$. Die Lösung des linearen Ausgleichsproblems ist also

$$(3.39) \quad \boxed{x = R^{-1}h_1, \quad \|s - Cx\|_2 = \|h_2\|_2}.$$

Die Kondition bei Anwendung der Householder-Transformation beträgt i.W. $\text{cond}_2(R)$; (vergleiche Stoer I, §4.8.3).

3.6.3 Diskrete Approximation

Als eine Anwendung der linearen Ausgleichsrechnung betrachten wir die diskrete Approximation.

Zu $n + 1$ Basisfunktionen

$$f_0(t), f_1(t), \dots, f_n(t)$$

und $m \geq n + 2$ Meßpunkten

$$(t_i, s_i), \quad i = 1, \dots, m$$

wird eine Linearkombination

$$f(t) = \sum_{k=0}^n x_k f_k(t)$$

gesucht, die die Werte s_i in den Punkten t_i möglichst gut annähert. Dies führt (vergleichbar zu linearen Ausgleichsrechnung) auf das Optimierungsproblem

$$\min_{x \in \mathbb{R}^{n+1}} \sum_{i=1}^m \left(s_i - \sum_{k=0}^n x_k f_k(t_i) \right)^2$$

also zu einem Problem der Form (3.37) mit

$$C = \begin{pmatrix} f_0(t_1) & \cdots & f_n(t_1) \\ \vdots & & \vdots \\ f_0(t_m) & \cdots & f_n(t_m) \end{pmatrix}$$

und $x = (x_0, \dots, x_n)^T$, $s = (s_1, \dots, s_m)^T$. Für die Basisfunktionen $f_k(t) = t^k$, $k = 0, \dots, n$, ist

$$C = \begin{pmatrix} 1 & t_1 & t_1^2 & \cdots & t_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & t_m & t_m^2 & \cdots & t_m^n \end{pmatrix}$$

die Vandermonde-Matrix, für die gilt:

$$\text{rang}(C) = n + 1, \quad \text{falls } t_i \neq t_j \quad \text{für } i \neq j.$$

Also ist $C^T C$ positiv definit und die Lösung x ist eindeutig bestimmt.

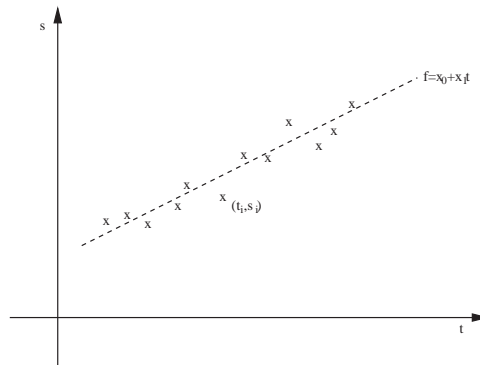


Abbildung 3.6: Ausgleichsgrade.

Beispiel 3.44. Der Fall $n = 1$: (Ausgleichsgerade)

Die Normalgleichungen für das Problem

$$\min_{x_0, x_1} \sum_{i=1}^m (s_i - (x_0 + x_1 t_i))^2$$

lauten

$$\begin{aligned} x_0 m + x_1 \sum_{i=1}^m t_i &= \sum_{i=1}^m s_i \\ x_0 \sum_{i=1}^m t_i + x_1 \sum_{i=1}^m t_i^2 &= \sum_{i=1}^m s_i t_i \end{aligned}$$

und können explizit nach x_0, x_1 aufgelöst werden. In der Statistik spricht man dabei von Regressionsrechnung.

Beispiel 3.45. Der Fall $n = 2$: (Ausgleichsparabel)

Zu den Meßpunkten:

i	1	2	3	4	5	6	7
t_i	0.04	0.32	0.51	0.73	1.03	1.42	1.60
s_i	2.63	1.18	1.16	1.54	2.65	5.41	7.67

erhält man die Ausgleichsparabel (Schwarz, S.288)

$$\begin{aligned} f(t) &= x_0 + x_1 t + x_2 t^2, \\ x_0 &= 2.749198 \\ x_1 &= -5.954657 \\ x_2 &= 5.607247 \end{aligned}$$

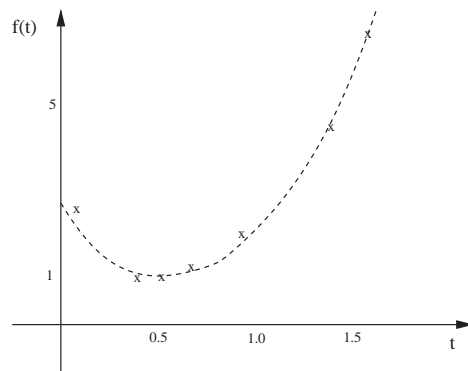


Abbildung 3.7: Ausgleichsparabel.

Kapitel 4

Nichtlineare Gleichungen und Gleichungssysteme

4.1 Einführung und Aufgabenstellung

Die Berechnung von Nullstellen nichtlinearer Gleichungssysteme bildet eine natürliche Erweiterung der LGS aus dem vorhergehenden Kapitel. Nichtlineare Gleichungen und Gleichungssysteme müssen in vielen Anwendungen der Mathematik gelöst werden. Typischerweise werden die Lösungen nichtlinearer Gleichungen über die Nullstellen einer Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ definiert, für die dann ein $x^* \in \mathbb{R}^n$ mit

$$(4.1) \quad f(x^*) = 0$$

gesucht wird.

Beispiel 4.1.

1. *Polynome:*

$$f(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n = 0, \quad x \in \mathbb{R}$$

Z.B. sind Eigenwerte von Matrizen Nullstellen des charakteristischen Polynoms.

2. *Bei der Berechnung der Schwingungen eines Balkens tritt das Problem*

$$f(x) = x - \tan(x) = 0$$

auf. In den Intervallen $((k - 1/2)\pi, (k + 1/2)\pi)$ liegen Nullstellen.

3. *Bei der Lösung nichtlinearer Optimierungsprobleme.*

4. *Gleichgewichtslösungen chemischer Prozesse.*

4.2 Grundlagen

4.2.1 Fixpunkte

Fixpunktgleichungen lassen sich an vielen Stellen leichter analysieren als die Nullstellen nichtlinearer Funktionen, tatsächlich sind aber beide Problemklassen ineinander überföhrbar.

Sei $D \subset \mathbb{R}^n$ und $g : D \rightarrow \mathbb{R}^n$. Gesucht sind die L6sungen $x \in D$ der Gleichung

$$(4.2) \quad x = g(x).$$

Ein Punkt $x^* \in D$ hei6t *Fixpunkt von g* , wenn $x^* = g(x^*)$ gilt.

Bemerkung 4.2. Durch Definition von $f(x) := x - g(x)$ wird eine Fixpunktgleichung in eine Nullstellenberechnung überföhrt. Ist umgekehrt $A(x)$ eine reguläre (n, n) -Matrix (z.B. die Einheitsmatrix), $x \in D$, dann ist die Nullstellenbestimmung

$$f(x) = 0$$

äquivalent zur Fixpunktgleichung

$$x = g(x) := x + A(x)f(x).$$

Für gegebene Startwerte x^0, x^1, \dots, x^s , $s \geq 0$ werden Fixpunkte mit *Iterationsverfahren* der Form

$$(4.3) \quad \boxed{x^{k+1} = \psi(x^k, x^{k-1}, \dots, x^{k-s}), \quad k \geq s}$$

bestimmt.

ψ hei6t *Iterationsfunktion* und hängt von g ab. Oft kann $\psi = g$ und $s = 0$ gewählt werden, so da6 die Iteration dann lautet

$$(4.4) \quad \boxed{x^{k+1} = g(x^k), \quad k = 1, 2, 3, \dots, \quad \text{für gegebenes } x^0 \in D.}$$

Es stellen sich die folgenden Fragen:

1. Wie findet man eine passende Iterationsfunktion?
2. Wie findet man passende Anfangspunkte?
3. Wann konvergiert die Folge gegen einen Fixpunkt?
4. Wie schnell konvergiert die Folge?

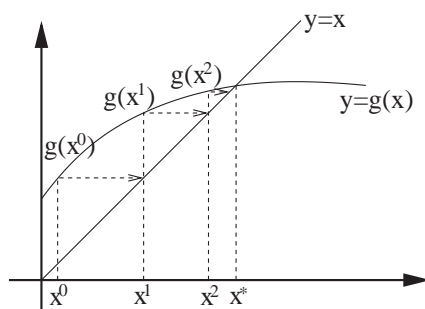


Abbildung 4.1: Graphische Darstellung eines Fixpunktes.

4.2.2 Konvergenz

Der Begriff der Konvergenzordnung erlaubt es, iterative Verfahren auf ihre Konvergenzgeschwindigkeit hin zu untersuchen.

Iterationsverfahren liefern eine Folge $\{x^k\} \subset \mathbb{R}^n$ approximativer Lösungen, die gegen die exakte Lösung x^* konvergieren. Die Konvergenzordnung gibt an, wie schnell der Fehler $\|x^k - x^*\|$ gegen Null konvergiert.

Definition 4.3 (Konvergenzordnung). Sei $\|\cdot\|$ eine Norm für \mathbb{R}^n und sei $\{x^k\} \subset \mathbb{R}^n$ eine aus einem Iterationsverfahren entstandene Folge mit

$$x^* = \lim_{k \rightarrow \infty} x^k.$$

1. Existiert eine Konstante $c \in (0, 1)$, so daß

$$(4.5) \quad \frac{\|x^{k+1} - x^*\|}{\|x^k - x^*\|} \leq c, \quad \forall k = 0, 1, 2, \dots$$

gilt, so heißt $\{x^k\}$ linear konvergent. Das zugehörige Iterationsverfahren wird linear konvergent oder konvergent von der Ordnung 1 genannt.

2. Existieren Konstanten $c > 0$ und $p > 1$, so daß

$$(4.6) \quad \frac{\|x^{k+1} - x^*\|}{\|x^k - x^*\|^p} \leq c, \quad \forall k = 0, 1, 2, \dots$$

gilt, so heißt $\{x^k\}$ konvergent von Grade p . Das zugehörige Iterationsverfahren wird konvergent von der Ordnung p genannt. Im Sonderfall $p = 2$ spricht man auch von quadratischer Konvergenz.

Bemerkung 4.4. Aus der aus (4.6) hergeleiteten Schreibweise

$$\|x^{k+1} - x^*\| \leq c \|x^k - x^*\|^p, \quad \forall k = 0, 1, 2, \dots$$

ergibt sich für den Fehler $e_k := \|x^k - x^*\|$ die Beziehung

$$e_{k+1} \leq ce_k^p,$$

d.h. es ist

$$e_{k+1} = \mathcal{O}(e_k^p).$$

Dies verdeutlicht nochmals den Begriff der Ordnung eines Verfahrens.

Bemerkung 4.5. Im Sonderfall $p = 1$ (lineare Konvergenz) erhält man aus (4.4) die Abschätzung

$$(4.7) \quad e_{k+m} \leq c^m e_k, \quad m > 0.$$

$\{x_k\}$ mit einem unteren Index k bezeichnet.

Beispiel 4.6. Für $0 < q < 1$ sei x_k der Abschnitt der geometrischen Reihe

$$\begin{aligned} x_k &= \sum_{i=0}^k q^i, \\ \text{mit } x^* &= \lim_{k \rightarrow \infty} x_k = \frac{1}{1-q}, \\ &\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k} = q \end{aligned}$$

$\Rightarrow p = 1$: lineare Konvergenz.

4.3 Nichtlineare Gleichungen

In diesem Abschnitt betrachten wir den Sonderfall $D = \mathbb{R}$.

Hinweis: Zur Vermeidung von Mißverständnissen werden reelle Folgen $\{x_k\}$ mit einem unteren Index k bezeichnet.

4.3.1 Bisektionsverfahren

Wir suchen eine Nullstelle $x^* \in \mathbb{R}$ einer reellen stetigen Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x^*) = 0$. Das Bisektionsverfahren wird auch *Intervallhalbierungsverfahren* genannt. Der Name wird sofort aus dem Algorithmus ersichtlich.

Algorithmus 4.7. (Bisektionsverfahren) Gegeben sei eine stetige Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$ und Werte $a < b$ mit $f(a) \cdot f(b) < 0$ (d.h. $f(a)$ und $f(b)$ haben unterschiedliches Vorzeichen und damit existiert eine Nullstelle von $f(x)$ im Intervall (a, b)). Die gewünschte Genauigkeit sei durch ein $\varepsilon > 0$ gegeben.

1. Setze $k = 0$ (Zählindex) und $a_0 = a$, $b_0 = b$.

2. Setze $x_k = a_k + (b_k - a_k)/2$ (Intervallhalbierung).
3. Ist $f(x_k) = 0$ oder $(b_k - a_k)/2 < \varepsilon$ beende den Algorithmus.
4. Ist $f(x_k)f(a_k) < 0$, dann setze $a_{k+1} = a_k$, $b_{k+1} = x_k$.
Ist $f(x_k)f(a_k) > 0$, dann setze $a_{k+1} = x_k$, $b_{k+1} = b_k$.
Setze $k = k + 1$ und gehe zu Schritt 2.

Man kann die Punkte a_k und b_k als Intervallgrenzen der Intervalle $[a_k, b_k]$ verstehen, mit denen die Nullstelle durch immer weitere Halbierung eingeschachtelt wird. Daher stammt der Name *Bisektion* (=Zweiteilung).

Die Auswahlbedingung der neuen Werte a_{k+1} und b_{k+1} stellt sicher, daß $f(a_{k+1})$ und $f(b_{k+1})$ unterschiedliches Vorzeichen haben, daher muß sich immer eine Nullstelle zwischen den Werten befinden (daher auch die Voraussetzung der Stetigkeit).

Vorteilig beim Bisektionsverfahren ist:

- Es funktioniert für allgemeine stetige Funktionen.
- Es liefert immer ein Ergebnis (globale Konvergenz), wenn man geeignete Startwerte finden kann.
- Die Anzahl der Schritte bis zur gewünschten Genauigkeit hängt nur von a und b ab, aber nicht von f .

Leider konvergiert das Verfahren (vgl. auch Beispiel) nur sehr langsam und wird daher in der Praxis so gut wie nie eingesetzt.

Bemerkung 4.8. Wegen

$$|b_{k+1} - a_{k+1}| \leq \frac{1}{2}|b_k - a_k|$$

folgt aus (4.7) lineare Konvergenz.

Beispiel 4.9. Für $f(x) = x - \tan(x)$, berechnet man mit dem Bisektionsverfahren eine Nullstelle von f :

$$\begin{aligned} a &= 2, & b &= 4.6 \\ f(a) &\approx 4.18, & f(b) &\approx -4.26 \end{aligned}$$

Damit

$$\begin{aligned} x_5 &= 4.47812 & , & & f(x_5) &= 2.87 \cdot 10^{-1} \\ x_{20} &= 4.493410 & , & & f(x_{20}) &= -1.51 \cdot 10^{-5} \\ x^* &= x_{100} = 4.49340946 & , & & f(x_{100}) &= -1.72294616 \cdot 10^{-10}. \end{aligned}$$

4.3.2 Newton–Verfahren

Wir betrachten ein weiteres Verfahren zur Nullstellenbestimmung einer gegebenen Funktion f . Im Gegensatz zum Bisektionsverfahren benötigen wir nicht nur die Funktion f sondern auch ihre erste Ableitung.

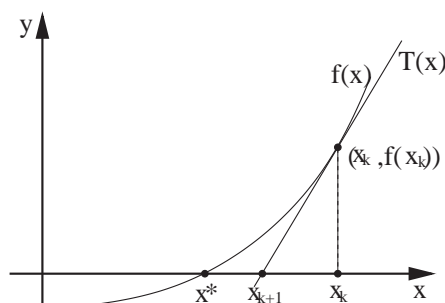


Abbildung 4.2: Zur Motivation des Newton–Verfahrens.

Sei x_k eine Näherung für x^* . Im Punkt $(x_k, f(x_k))$ wird eine Tangente

$$T(x) = f(x_k) + f'(x_k)(x - x_k)$$

an die Kurve $y = f(x)$ konstruiert und x_{k+1} als Nullstelle von T gewählt. Also

$$f(x_k) + f'(x_k)(x_{k+1} - x_k) = 0 \quad \Longleftrightarrow \quad x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

Eine Lösung existiert nur für $f'(x_k) \neq 0$. (Was bedeutet das anschaulich?)

Algorithmus 4.10. (Newton–Verfahren) Gegeben sei eine Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$, ihre Ableitung f' , ein Anfangswert x_0 und eine gewünscht Genauigkeit $\varepsilon > 0$.

1. Berechne $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$.
2. Ist $|x_{k+1} - x_k| < \varepsilon$, beende den Algorithmus, sonst setze $k = k + 1$ gehe zu 1.

Bemerkung 4.11. Das Newton–Verfahren ist eine Fixpunktiteration.

Satz 4.12 (Konvergenz des Newton–Verfahrens). Ist f zwei mal stetig differenzierbar, $f'(x) \neq 0$ und x_0 hinreichend nahe bei x^* , dann konvergiert das Newton–Verfahren quadratisch.

Beweis: Später.

Beispiel 4.13. Es sei $f(x) = x^2 - 2$ (d.h. Berechnung von $x^* = \sqrt{2} \approx 1.414213562373$) mit $f'(x) = 2x$. Die Iterationsvorschrift des Newton-Verfahrens ergibt hier

$$x_{k+1} = x_k - \frac{x_k^2 - 2}{2x_k} = \frac{1}{2}x_k + \frac{1}{x_k}$$

Wir starten mit $x_0 = 1$ und erhalten

k	x_k	Anzahl korrekter Dezimalstellen
0	<u>1.0</u>	1
1	<u>1.5</u>	1
2	<u>1.417</u>	3
3	<u>1.414216</u>	6
4	<u>1.414213562</u>	10

Die schnelle Konvergenz belegt das theoretische Resultat einer quadratischen Konvergenz.

Anmerkung: Newton hatte bereits 1669 ein Verfahren zur Berechnung einer Wurzel einer kubischen Gleichung entwickelt, das auf einen iterativen Linearisierungsprozeß hinausläuft. Er veröffentlichte sein Verfahren als Mittel zur Lösung der Keplerschen Gleichung:

$$E = e \cdot \sin(E) + \frac{2\pi}{U}t$$

zur Bahnbestimmung von Planeten. Gesucht ist die 'exzentrische Anomalie' E bei einer Umlaufzeit U , einer Zeit (in Tagen) t seit dem Periheldurchgang und einer numerischen Exzentrizität e der Bahnellipse. Joseph Raphson brachte um 1690 die Newtonschen Überlegungen für Polynome auf eine Form, die der heutigen Darstellung näher kommt. Man spricht deshalb häufig vom Newton-Raphson-Verfahren.

Es ist unmittelbar einzusehen, daß das Newton-Verfahren nicht immer konvergiert (vgl. Abbildung 4.3).

Bemerkung 4.14. Wesentlicher Nachteil des Newton-Verfahrens ist die Abhängigkeit von der Ableitung. Die Ableitung kann zwar auch numerisch berechnet werden, ist dann jedoch häufig anfällig gegenüber Rundungsfehlern.

4.3.3 Sekanten-Verfahren

Zur Vermeidung von Ableitungen betrachten wir das Sekanten-Verfahren, hier ergibt sich die neue Näherung x_{k+1} nicht nur aus x_k sondern auch aus x_{k-1} .

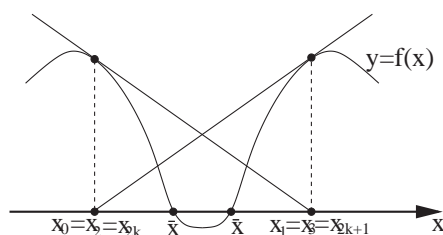


Abbildung 4.3: Keine Konvergenz beim Newton-Verfahren.

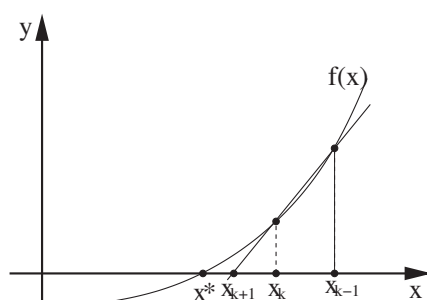


Abbildung 4.4: Zur Motivation des Sekanten-Verfahrens.

Die Sekantenmethode ist eine Vereinfachung des Newton-Verfahrens, wobei die Tangente durch die Sekante der letzten beiden Punkte ersetzt wird. Die Steigung der Sekante ergibt sich zu

$$(4.8) \quad \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}} \approx f'(x_k).$$

Für das Sekanten-Verfahren wird $f'(x_k)$ einfach durch (4.8) ersetzt.

Algorithmus 4.15. (Sekanten-Verfahren) Gegeben sei eine Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$, Anfangswerte x_0, x_1 und eine gewünscht Genauigkeit $\varepsilon > 0$. Setze $k = 1$.

1. Berechne $x_{k+1} = x_k - \frac{f(x_k)(x_k - x_{k-1})}{f(x_k) - f(x_{k-1})}$.
2. Ist $|x_{k+1} - x_k| < \varepsilon$, beende den Algorithmus, sonst setze $k = k + 1$ gehe zu 1.

Bemerkung 4.16. Das Sekanten-Verfahren ist eine Fixpunktiteration.

Satz 4.17 (Konvergenz des Sekanten-Verfahrens). Ist f zwei mal stetig differenzierbar, $f'(x^*) \neq 0$ und x_0, x_1 hinreichend nahe bei x^* , dann konvergiert das Sekanten-Verfahren mit der Ordnung $p = \frac{1}{2}(1 + \sqrt{5}) = 1.618 \dots$

Beweis: Später.

Beispiel 4.18. Wir betrachten erneut $f(x) = x^2 - 2$. Die Iterationsvorschrift des Sekanten-Verfahrens ergibt hier

$$x_{k+1} = x_k - \frac{x_k^2 - 2}{x_k + x_{k-1}}$$

Wir starten mit $x_0 = 1$, $x_1 = 2$ und erhalten

k	x_k	Anzahl korrekter Dezimalstellen
0	<u>1.0</u>	1
1	2.0	0
2	<u>1.3</u>	1
3	<u>1.43</u>	2
4	<u>1.414</u>	4
5	<u>1.414211</u>	6
6	<u>1.4142135627</u>	10

Das Sekanten-Verfahren startet zwar in diesem Beispiel recht langsam aufgrund der schlechten Startschätzung für x_1 , konvergiert aber später entsprechend schnell.

4.4 Konvergenz von Iterationsverfahren

4.4.1 Kontraktion

Sei $D \subset \mathbb{R}^n$ und $g : D \rightarrow \mathbb{R}^n$. Wir untersuchen die Frage, wann die Fixpunktiteration

$$(4.9) \quad x^{k+1} = g(x^k), \quad k \geq 0, \quad x^0 \in D \text{ gegeben,}$$

wohl definiert ist und gegen einen Fixpunkt $\bar{x} \in D$, konvergiert.

Definition 4.19. Die Abbildung $g : D \rightarrow \mathbb{R}^n$ heißt kontrahierend in D , falls es eine Zahl $0 \leq q < 1$ gibt mit

$$\|g(x) - g(y)\| \leq q\|x - y\| \quad \forall x, y \in D.$$

Für differenzierbare Abbildungen g kann ein einfaches Kriterium für die Kontraktion mit der Ableitung

$$g'(x) = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \cdots & \frac{\partial g_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial g_n}{\partial x_1} & \cdots & \frac{\partial g_n}{\partial x_n} \end{pmatrix}$$

gegeben werden. D heißt konvex, wenn für $x, y \in D$ gilt

$$\alpha x + (1 - \alpha)y \in D \quad \forall \alpha \in [0, 1].$$

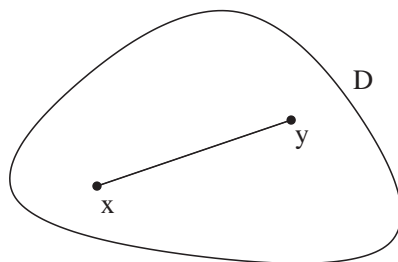


Abbildung 4.5: Konvexes Gebiet.

Satz 4.20. Sei D konvex, $g : D \rightarrow \mathbb{R}^n$ differenzierbar und sei

$$\sup_{x \in D} \|g'(x)\|_{\infty} \leq q < 1.$$

Dann ist g kontrahierend in D .

Beweis: Für zwei beliebige Punkte $x, y \in D$ betrachten wir $\varphi : [0, 1] \rightarrow \mathbb{R}^n$:

$$\varphi(\lambda) := g(\lambda x + (1 - \lambda)y), \quad \lambda \in [0, 1],$$

$$\varphi(1) = g(x), \quad \varphi(0) = g(y),$$

$$\varphi'(\lambda) = g'(\lambda x + (1 - \lambda)y)(x - y).$$

Aus dem Mittelwertsatz folgt:

$$|\varphi_i(1) - \varphi_i(0)| \leq \max_{0 \leq \lambda \leq 1} |\varphi'_i(\lambda)|, \quad i = 1, \dots, n.$$

$$\begin{aligned} \Rightarrow \|g(x) - g(y)\|_{\infty} &= \|\varphi(1) - \varphi(0)\|_{\infty} \\ &\leq \max_{0 \leq \lambda \leq 1} \|\varphi'(\lambda)\|_{\infty} \\ &= \max_{0 \leq \lambda \leq 1} \|g'(\lambda x + (1 - \lambda)y)(x - y)\|_{\infty} \\ &\leq \sup_{z \in D} \|g'(z)\|_{\infty} \|x - y\|_{\infty}. \end{aligned}$$

◇

Für $n = 1$ ist $D = [a, b]$ konvex und $g \in C^1[a, b]$ kontrahierend, falls

$$\max_{a \leq x \leq b} |g'(x)| = q < 1.$$

(Vgl. Graphik zum Fixpunkt).

4.4.2 Fixpunktsatz von Banach

Satz 4.21 (Fixpunktsatz von Banach). *Sei D abgeschlossen und $g : D \rightarrow \mathbb{R}^n$ kontrahierend in D mit $g(D) \subseteq D$. Dann konvergiert die Folge*

$$x^{k+1} = g(x^k), \quad k = 0, 1, 2, \dots, \quad x^0 \in D \text{ beliebig,}$$

gegen den eindeutig bestimmten Fixpunkt \bar{x} von g in D und es gilt:

$$(i) \quad \|\bar{x} - x^k\| \leq \frac{q}{1-q} \|x^k - x^{k-1}\| \leq \frac{q^k}{1-q} \|x^1 - x^0\|,$$

$$(ii) \quad \|\bar{x} - x^k\| \leq q \|\bar{x} - x^{k-1}\|.$$

Beweis: Wir zeigen zunächst, dass $\{x^k\}$ eine Cauchy-Folge ist. Für $k \geq 1$ gilt

$$\begin{aligned} \|x^{k+1} - x^k\| &= \|g(x^k) - g(x^{k-1})\| \leq q \|x^k - x^{k-1}\| \\ &\leq q^2 \|x^{k-1} - x^{k-2}\| \leq q^k \|x^1 - x^0\| \end{aligned}$$

und damit für $j > l$

$$\begin{aligned} \|x^j - x^l\| &= \left\| \sum_{k=l}^{j-1} (x^{k+1} - x^k) \right\| \leq \sum_{k=l}^{j-1} \|x^{k+1} - x^k\| \\ (4.10) \quad &\leq \sum_{k=l}^{j-1} q^k \|x^1 - x^0\| \\ &\leq q^l \frac{1}{1-q} \|x^1 - x^0\| \xrightarrow{l \rightarrow \infty} 0. \end{aligned}$$

Die x^k bilden damit eine Cauchy-Folge. Sei $\bar{x} \in D$ der Grenzwert der Folge x^k . Dann gilt

$$\begin{aligned} g(\bar{x}) &\stackrel{k \rightarrow \infty}{\longleftarrow} g(x^k) = x^{k+1} \stackrel{k \rightarrow \infty}{\longrightarrow} \bar{x}, \\ &\Rightarrow g(\bar{x}) = \bar{x}. \end{aligned}$$

In (4.10) ergibt der Grenzwert für $j \rightarrow \infty$:

$$\|\bar{x} - x^l\| \leq \frac{q^l}{1-q} \|x^1 - x^0\|.$$

Mit $l = 1$ folgt hieraus nach Ersetzen $x^0 \rightarrow x^{k-1}$

$$\|\bar{x} - x^k\| \leq \frac{q}{1-q} \|x^k - x^{k-1}\|.$$

Damit ist (i) gezeigt; (ii) folgt aus der Kontraktionseigenschaft:

$$\|\bar{x} - x^k\| = \|g(\bar{x}) - g(x^{k-1})\| \leq q \|\bar{x} - x^{k-1}\|.$$

Zu zeigen bleibt noch die Eindeutigkeit von \bar{x} .

Seien \bar{x}, x^* Fixpunkte von g :

$$\begin{aligned}\|\bar{x} - x^*\| &= \|g(\bar{x}) - g(x^*)\| \leq q\|\bar{x} - x^*\|, \quad q < 1, \\ \Rightarrow \quad \|\bar{x} - x^*\| &= 0.\end{aligned}$$

◇

Bemerkung 4.22. Wegen Teil (ii) konvergiert $\{x^k\}$ linear gegen \bar{x} .

Bemerkung 4.23. Die Schwierigkeiten bei der Anwendung des Kontraktionssatzes auf ein konkretes Problem bestehen darin:

- (a) man finde eine zugehörige kontrahierende Funktion $g : D \rightarrow \mathbb{R}^n$,
- (b) man prüfe $g(D) \subseteq D$.

Beispiel 4.24. Gesucht ist die Lösung \bar{x} der Gleichung

$$x = e^{-x} =: g(x), \quad x \in \mathbb{R}.$$

Auf das Intervall $D = [0.5, 0.69]$ trifft die Voraussetzung $g(D) \subset D$ zu. Als Kontraktionszahl q dient nach Satz 4.20 die Zahl

$$\max_{x \in D} |g'(x)| = e^{-0.5} = 0.606531 < 1.$$

Zum Startwert $x^{(0)} = 0.55 \in D$ berechnet man die Iterierten:

k	$x^{(k)}$	k	$x^{(k)}$	k	$x^{(k)}$
0	0.55000000	10	0.56708394	20	0.56714309
1	0.57694981	11	0.56717695	21	0.56714340
2	0.56160877	12	0.56712420	22	0.56714323
3	0.57029086	13	0.56715412	23	0.56714332
4	0.56536097	14	0.56713715	24	0.56714327

Mit der a priori Fehlerabschätzung $\|\bar{x} - x^{(k)}\| \leq \frac{q^k}{1-q} \|x^{(1)} - x^{(0)}\|$ kann die Anzahl k der Iteration geschätzt werden, die nötig sind, damit z. B. $\|\bar{x} - x^{(k)}\| \leq \varepsilon = 10^{-6}$ gilt. Man erhält

$$k \geq \log \left(\frac{\varepsilon(1-q)}{\|x^{(1)} - x^{(0)}\|} \right) / \log q = 22.3,$$

eine gegenüber der Tabelle leichte Überschätzung. Für den Wert $x^{(12)}$ erhält man die a posteriori Fehlerschranke

$$\|\bar{x} - x^{(12)}\| \leq \frac{q}{1-q} \|x^{(12)} - x^{(11)}\| = 8.3 \cdot 10^{-5}.$$

Beispiel 4.25. Erneut greifen wir das Beispiel $f(x) = x - \tan(x)$ auf.

Die Nullstelle \bar{x} wird in $D = [\pi, \frac{3}{2}\pi]$ gesucht. Die Funktion $g(x) = \tan x$ ist nicht kontrahierend wegen

$$g'(x) = \frac{1}{\cos^2 x} \geq 1.$$

Umformulierung:

$$x = \tan x = \tan(x - \pi) \Leftrightarrow \arctan x = x - \pi.$$

Setze nun

$$g(x) = \pi + \arctan x, \quad D = [\pi, \frac{3}{2}\pi].$$

Offenbar gilt $g(D) \subseteq D$ und

$$q := \max_{x \in D} |g'(x)| = \frac{1}{1 + \pi^2} \approx 0.092 < 1.$$

g ist also kontrahierend in D nach (4.20).

Für $\bar{x} = 4.4934094$ ist $g'(\bar{x}) = 0.04719$.

k	x^k	$\frac{q}{1-q} x^k - x^{k-1} $	$ \bar{x} - x^k $	$\frac{ \bar{x} - x^k }{ \bar{x} - x^{k-1} }$
0	3.14159265	–	–	–
1	4.40421991	0.1351	0.0892	–
2	4.48911945	0.008918	0.0043	0.0672
3	4.49320683	0.0004280	0.0002	0.0481
4	4.4933999	–	–	0.0472

4.4.3 Konvergenzsätze

Satz 4.26 (Lokaler Konvergenzsatz). Sei $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ mit $g(\bar{x}) = \bar{x}$. Ist g in einer Umgebung von \bar{x} stetig differenzierbar und $\|g'(\bar{x})\|_\infty < 1$, dann gibt es eine Umgebung D von \bar{x} , so dass das Iterationsverfahren

$$x^{k+1} = g(x^k), \quad x^0 \in D$$

gegen \bar{x} konvergiert.

Beweis: Sei D eine Kugel mit Radius r um \bar{x} mit $\|g'(x)\|_\infty \leq q < 1$ für $x \in D$. Für $x \in D$ gilt

$$\begin{aligned} \|g(x) - \bar{x}\|_\infty &= \|g(x) - g(\bar{x})\|_\infty \leq q\|x - \bar{x}\|_\infty \leq r \\ &\Rightarrow g(x) \in D. \end{aligned}$$

Damit ist g kontrahierend in D und es gilt $g(D) \subseteq D$. Mit dem Fixpunktsatz 4.21 folgt die Behauptung. \diamond

Als Anwendung erhält man im Falle $n = 1$ einfache Kriterien dafür, dass die Fixpunkt-Iteration ein Verfahren p -ter Ordnung ist.

Satz 4.27. Sei $g : \mathbb{R} \rightarrow \mathbb{R}$ eine C^p -Funktion mit $p \in \mathbb{N}_+$. Sei \bar{x} ein Fixpunkt von g mit

- (a) $|g'(\bar{x})| < 1$ für $p = 1$,
 (b) $g^{(i)}(\bar{x}) = 0$ ($i = 1, \dots, p - 1$) für $p > 1$.

Dann gibt es ein Intervall

$$I = [\bar{x} - \delta, \bar{x} + \delta], \quad \delta > 0,$$

so dass für alle $x_0 \in I$ die Iteration $x_{k+1} = g(x_k)$, $k = 0, 1, 2, \dots$ konvergent vom Grade p ist mit

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - \bar{x}|}{|x_k - \bar{x}|^p} = \frac{1}{p!} g^{(p)}(\bar{x}).$$

Beweis: Aus den Vor. (a),(b) folgt insbesondere $|g'(\bar{x})| < 1$. Der lokale Konvergenzsatz 4.26 sichert dann die (mindestens) lineare Konvergenz der Folge $x_{k+1} = g(x_k)$ für alle $x_0 \in I = [\bar{x} - \delta, \bar{x} + \delta]$, $\delta > 0$ geeignet.

Die Taylor-Entwicklung ergibt mit Vor. (b) und $\bar{x} = g(\bar{x})$:

$$x_{k+1} = \bar{x} + \frac{1}{p!} g^{(p)}(\bar{x})(x_k - \bar{x})^p + o(|x_k - \bar{x}|^p).$$

Hieraus folgt die Behauptung

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - \bar{x}|}{|x_k - \bar{x}|^p} = \frac{1}{p!} g^{(p)}(\bar{x})$$

\diamond

4.4.4 Konvergenz des Newton-Verfahrens

Als Anwendung betrachten wir das Newton-Verfahren

$$x_{k+1} = g(x_k) = x_k - \frac{f(x_k)}{f'(x_k)}.$$

Ist f eine C^3 -Funktion, so ist g eine C^2 -Funktion.

1. Fall: \bar{x} ist einfache Nullstelle von f , d.h. $f'(\bar{x}) \neq 0$: man berechnet

$$\begin{aligned} g'(x) &= \frac{f(x)f''(x)}{f'(x)^2}, & g'(\bar{x}) &= 0, \\ g''(\bar{x}) &= \frac{f''(\bar{x})}{f'(\bar{x})}. \end{aligned}$$

Also ist das Newton-Verfahren (mindestens) quadratisch konvergent mit der asymptotischen Fehlerkonstanten

$$c = \frac{1}{2} \frac{f''(\bar{x})}{f'(\bar{x})}.$$

2. Fall: \bar{x} sei m -fache Nullstelle von f , d.h. $f^{(i)}(\bar{x}) = 0$ für $i = 0, \dots, m-1$:

$$f(x) = (x - \bar{x})^m f_0(x), \quad f_0(\bar{x}) \neq 0$$

$$\Rightarrow g'(\bar{x}) = 1 - \frac{1}{m}.$$

Für $m > 1$ ist daher $g'(\bar{x}) \neq 0$ und das Newton-Verfahren ist nur linear konvergent.

Für das modifizierte Newton-Verfahren

$$x_{k+1} = g(x_k) := x_k - m \frac{f(x_k)}{f'(x_k)}$$

gilt jedoch $g'(\bar{x}) = 0$, also hat man quadratische Konvergenz.

4.5 Das Newton-Verfahren im \mathbb{R}^n

4.5.1 Herleitung des Newton-Verfahrens

Gegeben sei eine C^1 -Funktion $f : D \rightarrow \mathbb{R}^n$. Gesucht ist eine Nullstelle $\bar{x} \in D$ von f . Das Newton-Verfahren zur Berechnung von \bar{x} ist die folgende Fixpunktiteration:

$$(4.11) \quad \boxed{x^{k+1} = x^k - (f'(x^k))^{-1} f(x^k), \quad k \geq 0, \quad x^0 \in D \text{ gegeben.}}$$

Das Newton-Verfahren im \mathbb{R}^n läßt sich auf verschiedene Weise erklären:

- (1) Verallgemeinerung des Newton-Verfahrens mittels Taylor-Entwicklung:

Es gilt

$$0 = f(\bar{x}) = f(x^k) + f'(x^k)(\bar{x} - x^k) + o(\|\bar{x} - x^k\|).$$

Vernachlässigt man $o(\|\bar{x} - x^k\|)$ und ersetzt den unbekanntem Punkt \bar{x} durch x^{k+1} , so erhält man

$$0 = f(x^k) + f'(x^k)(x^{k+1} - x^k)$$

und daraus (4.11).

- (2) Anwendung des lokalen Konvergenzsatzes 4.26:
 $f(\bar{x}) = 0$ gilt genau dann, wenn \bar{x} Fixpunkt von

$$g(x) := x + A(x)f(x)$$

ist mit einer geeignet zu wählenden regulären (n, n) C^1 -Matrix $A(x)$. Nach Satz 4.26 ist g kontrahierend, falls $\|g'(\bar{x})\|_\infty < 1$ ist. Wegen $f(\bar{x}) = 0$ gilt

$$g'(\bar{x}) = I + A(\bar{x})f'(\bar{x}).$$

Wählen wir nun

$$A(\bar{x}) = -(f'(\bar{x}))^{-1}$$

so ist $g'(\bar{x}) = 0$. Da \bar{x} unbekannt ist setzen wir

$$A(x) = -(f'(x))^{-1}$$

d.h.

$$g(x) = x - (f'(x))^{-1}f(x).$$

Die Fixpunktiteration $x^{k+1} = g(x^k)$ ergibt gerade (4.11). Satz 4.26 sichert wegen $g'(\bar{x}) = 0$ die lokale Konvergenz.

Beispiel 4.28. *Gesucht ist die Lösung des Systems*

$$f(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \end{pmatrix} = \begin{pmatrix} x_1^2 + x_2^2 - 1 \\ x_1 \end{pmatrix} = 0$$

Die Ableitung $f'(x)$ ist gegeben durch

$$f'(x) = \begin{pmatrix} 2x_1 & 2x_2 \\ 1 & 0 \end{pmatrix},$$

wodurch sich die Inverse zu

$$(f'(x))^{-1} = \begin{pmatrix} 0 & 1 \\ \frac{1}{2x_2} & -\frac{x_1}{x_2} \end{pmatrix}$$

ergibt. Die Iterationsvorschrift lautet somit

$$x^{k+1} = x^k - \begin{pmatrix} 0 & 1 \\ \frac{1}{2x_2^k} & -\frac{x_1^k}{x_2^k} \end{pmatrix} \begin{pmatrix} (x_1^k)^2 + (x_2^k)^2 - 1 \\ x_1^k \end{pmatrix}.$$

Mit $x^0 = (1, 1)$ ergibt sich

k	x_1^k	x_2^k	Max. Anzahl korrekter Dezimalstellen
0	1	<u>1</u>	0
1	<u>0</u>	<u>1,5</u>	1
2	<u>0</u>	<u>1.08</u>	2
3	<u>0</u>	<u>1.003</u>	3
4	<u>0</u>	<u>1.000005</u>	6

4.5.2 Praktische Realisierung

Bemerkung 4.29. *Praktisch benutzt man in höheren Dimensionen das Newton-Verfahren in der Form*

$$f'(x^k)(x^{k+1} - x^k) = -f(x^k).$$

So muß anstatt der Invertierung von $f'(x^k)$ (n^3 Operation) nur noch ein LGS gelöst werden ($\frac{1}{3}n^3$ Operationen).

Beispiel 4.30. $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$,

$$f(x) = \begin{pmatrix} 10^4 x_1 x_2 - 1 \\ e^{-x_1} + e^{-x_2} - 1.0001 \end{pmatrix}, \quad f'(x) = \begin{pmatrix} 10^4 x_2 & 10^4 x_1 \\ -e^{-x_1} & -e^{-x_2} \end{pmatrix}.$$

$$x^0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad f(x^0) = \begin{pmatrix} -1 \\ 0,36 \end{pmatrix},$$

$$\bar{x} = x^{13} = \begin{pmatrix} 1.0981595 \times 10^{-5} \\ 9.10614 \end{pmatrix}.$$

Die Matrix

$$f'(\bar{x}) = \begin{pmatrix} 9.1 \times 10^4 & 0.11 \\ -1 & -1.1 \times 10^{-4} \end{pmatrix}$$

hat die Kondition

$$\|f'(\bar{x})\|_\infty \cdot \|(f'(\bar{x}))^{-1}\|_\infty = \mathcal{O}(10^9).$$

Bei der Berechnung von $f_2(x)$ entsteht Auslösung; $f_2(x)$ läßt sich in folgender Gestalt besser berechnen:

$$\begin{aligned} e^{-x_1} + e^{-x_2} - 1.0001 &= (e^{-x_1} - 1) + (e^{-x_2} - 10^{-4}) \\ &\approx (-x_1 + (x_1)^2/2) + (e^{-x_2} - 10^{-4}). \end{aligned}$$

4.5.3 Newton–Kantorovich

Zum Nachweis der lokalen quadratischen Konvergenz des Newton-Verfahrens (4.11) benötigen wir den folgenden Hilfssatz:

Hilfssatz 4.31. *Sei $D_0 \subset D$ konvex. Es gebe $\gamma > 0$ mit*

$$\|f'(x) - f'(y)\| \leq \gamma \|x - y\| \quad \text{für } x, y \in D_0.$$

Dann gilt

$$\|f(x) - f(y) - f'(y)(x - y)\| \leq \frac{\gamma}{2} \|x - y\|^2 \quad \forall x, y \in D_0$$

Beweis: Definiere die differenzierbare Funktion $\varphi : [0, 1] \rightarrow \mathbb{R}^n$ durch

$$\begin{aligned} \varphi(t) &:= f(y + t(x - y)), & x, y \in D_0, \\ \varphi'(t) &= f'(y + t(x - y))(x - y) \in \mathbb{R}^n. \end{aligned}$$

Mit der Voraussetzung folgt

$$\begin{aligned} \|\varphi'(t) - \varphi'(0)\| &= \|(f'(y + t(x - y)) - f'(y))(x - y)\| \\ &\leq \gamma t \|x - y\| \|x - y\|. \end{aligned}$$

Es ist

$$\begin{aligned} \Delta &:= f(x) - f(y) - f'(y)(x - y) \\ &= \varphi(1) - \varphi(0) - \varphi'(0) = \int_0^1 (\varphi'(t) - \varphi'(0)) dt, \\ \Rightarrow \|\Delta\| &\leq \int_0^1 \|\varphi'(t) - \varphi'(0)\| dt \\ &\leq \gamma \|x - y\|^2 \int_0^1 t dt = \frac{\gamma}{2} \|x - y\|^2. \end{aligned}$$

◇

Satz 4.32 (Newton-Kantorovich). *Es sei eine offene Menge $D \subseteq \mathbb{R}^n$ gegeben, ferner eine konvexe Menge D_0 mit $\overline{D_0} \subseteq D$ und $f : D \rightarrow \mathbb{R}^n$ sei eine für alle $x \in D_0$ differenzierbare und für alle $x \in D$ stetige Funktion.*

Für ein $x^0 \in D_0$ gebe es positive Konstanten $r, \alpha, \beta, \gamma, h$ mit:

$$\begin{aligned} S_r(x^0) &:= \{x \mid \|x - x^0\| < r\} \subseteq D_0, \\ h &:= \alpha\beta\gamma/2 < 1, \\ r &:= \alpha/(1 - h). \end{aligned}$$

$f(x)$ habe die Eigenschaften

- (a) $\|f'(x) - f'(y)\| \leq \gamma \|x - y\|$ für alle $x, y \in D_0$
 (Lipschitz-Bedingung für f'),
- (b) $f'(x)^{-1}$ existiert und es gilt
 $\|(f'(x))^{-1}\| \leq \beta$ für alle $x \in D_0$,
- (c) $\|(f'(x^0))^{-1}f(x^0)\| \leq \alpha$.

Dann gilt

- (i) ausgehend von x^0 ist jedes

$$x^{k+1} = x^k - (f'(x^k))^{-1}f(x^k), \quad k \geq 0$$

wohldefiniert und es gilt $x^k \in S_r(x^0)$ für alle $k \geq 0$.

- (ii) $\bar{x} = \lim_{k \rightarrow \infty} x^k$ existiert und es gilt

$$\bar{x} \in \overline{S_r(x^0)} \text{ und } f(\bar{x}) = 0.$$

- (iii) $\|\bar{x} - x^k\| \leq \alpha \frac{h^{2^k-1}}{1-h^{2^k}}$ für alle $k \geq 0$.

Wegen $0 < h < 1$ ist also das Newton-Verfahren mindestens quadratisch konvergent.

Beweis: zu (i):

Zunächst wird $x^k \in S_r(x^0)$, $k \geq 0$ induktiv gezeigt. Für $k = 1$ ist

$$x^1 = x^0 - f'(x^0)^{-1}f(x^0) \Rightarrow \|x^1 - x^0\| \leq \alpha < \frac{\alpha}{1-h} = r$$

Seien $x^0, \dots, x^k \in S_r(x^0)$:

$$\begin{aligned} \|x^{k+1} - x^k\| &= \| -f'(x^k)^{-1}f(x^k) \| \leq \beta \|f(x^k)\| \\ &= \beta \| \underbrace{f(x^k) - f(x^{k-1}) - f'(x^{k-1})(x^k - x^{k-1})}_{=0 \text{ nach Def.}} \| \\ &\leq \frac{1}{2} \beta \gamma \|x^k - x^{k-1}\|^2 \text{ nach Hilfssatz 4.31} \end{aligned}$$

Hiermit zeigen wir nun induktiv

$$(4.12) \quad \|x^{k+1} - x^k\| \leq \alpha h^{2^k-1}.$$

Für $k = 0$ ist dies bereits gezeigt (s.o.). Ist die Abschätzung für $k \geq 0$ richtig, so (wegen $h = \frac{1}{2}\alpha\beta\gamma$) auch für $k + 1$, denn

$$\|x^{k+1} - x^k\| \leq \frac{\beta\gamma}{2} \|x^k - x^{k-1}\|^2 \leq \frac{\beta\gamma}{2} \alpha^2 h^{2^k-2} = \alpha h^{2^k-1}.$$

Nun folgt mit (4.12)

$$\begin{aligned}\|x^{k+1} - x^0\| &\leq \|x^{k+1} - x^k\| + \|x^k - x^{k-1}\| + \dots + \|x^1 - x^0\| \\ &\leq \alpha(1 + h + h^3 + h^7 + \dots + h^{2^k-1}) < \alpha/(1-h) = r\end{aligned}$$

und daher $x^{k+1} \in S_r(x^0)$.

Zu (ii) und (iii):

$\{x^k\}$ ist eine Cauchy-Folge, denn für $m \geq n$ hat man nach (4.12)

$$\begin{aligned}\|x^{m+1} - x^n\| &\leq \|x^{m+1} - x^m\| + \|x^m - x^{m-1}\| + \dots + \|x^{n+1} - x^n\| \\ &\leq \alpha h^{2^n-1}(1 + h^{2^n} + (h^{2^n})^2 + \dots) \\ &< \frac{\alpha h^{2^n-1}}{1-h^{2^n}} < \varepsilon\end{aligned}$$

für genügend großes $n \geq N(\varepsilon)$, da $0 < h < 1$. Also existiert

$$\lim_{k \rightarrow \infty} x^k =: \bar{x} \in \overline{S_r(x^0)},$$

und für $m \rightarrow \infty$ ergibt sich die Abschätzung (iii). Zu zeigen ist noch $f(\bar{x}) = 0$: Mit der Dreiecksungleichung und $\|x - y\| \leq r$ folgt aus (a) für $K := \gamma r + \|f'(x^0)\|$

$$\begin{aligned}\|f(x^k)\| &= \|f'(x^k)(x^{k+1} - x^k)\| \leq K\|x^{k+1} - x^k\| \rightarrow 0 \text{ für } k \rightarrow \infty \\ &\Rightarrow f(\bar{x}) = 0, \quad \text{da } f \text{ stetig in } \bar{x} \in D \text{ ist. } \diamond\end{aligned}$$

Bemerkung 4.33. Sei \bar{x} eine isolierte Nullstelle von f , so dass die Vor. (a), (b) in einer Umgebung D_0 von \bar{x} erfüllt sind. Für x^0 in einer genügend kleinen Umgebung E (Einzugsbereich) von \bar{x} ist dann α in Vor. (c) hinreichend klein, so dass die Vor. von Satz 4.32 erfüllt sind. Dann konvergiert $\{x^k\}$ quadratisch gegen \bar{x} für $x^0 \in E$.

Bemerkung 4.34. Für eine C^2 -Funktion f ist die Lipschitz-Bedingung in Vor. (a) erfüllt.

4.5.4 Erweiterungen

4.5.4.1 Approximation von $f'(x)$ durch Differenzen

Die Berechnung von $f'(x)$ kann sehr zeitraubend bzw. explizit nicht möglich sein. Man kann $f'(x^k)$ z.B. numerisch approximieren, indem man $f(x)$ in der Nähe von x^k geeignet auswertet.

$$\left. \frac{\partial f_i(x)}{\partial x_j} \right|_{x=x^k} \approx \frac{f_i(x^k + h e_j) - f_i(x^k)}{h}, \quad h > 0 \text{ klein}$$

Bemerkung 4.35. Diese Methode erfordert n Berechnungen von f für jeden Schritt des Newton-Verfahrens und dies ist für n groß ungünstig.

4.5.4.2 λ -Strategie, Modifiziertes Newton-Verfahren

Durch Einführung eines konvergenzerzeugenden Faktors $0 < \lambda \leq 1$ kann der Einzugsbereich des Newton-Verfahrens vergrößert werden. Sei

$$d^k = f'(x^k)^{-1} f(x^k) \in \mathbb{R}^n.$$

Gewöhnliches Newton-Verfahren:

$$x^{k+1} = x^k - d^k.$$

Modifiziertes Newton-Verfahren:

$$(4.13) \quad \boxed{x^{k+1} = x^k - \lambda_k d^k \quad 0 < \lambda_k \leq 1}$$

Zur Bestimmung der konvergenzerzeugenden Faktoren λ_k vergleiche man STOER, §5.4.

Kapitel 5

Interpolation

5.1 Einführung und Aufgabenstellung

Häufig tritt die Situation auf, daß statt einer Funktion nur einige diskrete Daten (x_j, f_j) , $j = 0, \dots, n$ gegeben sind (z.B. Messwerte eines Experimentes). Historisch trat das Problem bei der Berechnung von zusätzlichen Funktionswerten zwischen tabellierten Werten auf (z.B. für \sin , \cos , \log). Heute ist es ein häufig auftretendes Problem sowohl in der Mathematik als auch in vielen Anwendungen.

Gesucht ist eine Funktion $\Phi : \mathbb{R} \rightarrow \mathbb{R}$, für die die Gleichungen

$$(5.1) \quad \Phi(x_j) = f_j, \quad j = 0, \dots, n$$

gelten. Sind z.B. $f_j = f(x_j)$ Messdaten einer unbekanntes Funktion f , so soll Φ möglichst nahe an f liegen. Zudem soll Φ leicht auswertbar sein. Ohne Einschränkung verlangen wir $x_0 < x_1 < \dots < x_n$.

Beispiel 5.1. (*Lineare Interpolationsprobleme*)

- *Interpolation durch Polynome*

$$\Phi(x) = \sum_{j=0}^n a_j x^j,$$

- *Trigonometrische Interpolation* (Es gilt $-1 = i^2$, $e^{ix} = \cos x + i \sin x$)

$$\Phi(x) = \sum_{j=0}^n a_j e^{ijx} = a_0 + \sum_{j=1}^n a_j \cos(jx) + i \sum_{j=1}^n a_j \sin(jx)$$

- *Kubische Spline-Interpolation*

$$\Phi(x) \in C^2[x_0, x_n] \quad \text{und} \quad \Phi(x) \in C^3[x_j, x_{j+1}]$$

Beispiel 5.2. (*Nichtlineare Interpolationsprobleme*)

- *Rationale Interpolation*

$$\Phi(x) = \frac{\sum_{j=0}^n a_j x^j}{\sum_{j=0}^m b_j x^j}$$

- *Interpolation durch Exponentialsummen*

$$\Phi(x) = \sum_{j=0}^n a_j e^{\lambda_j x}$$

In diesem Kapitel werden wir uns mit linearen Interpolationsproblemen beschäftigen.

5.2 Polynominterpolation

5.2.1 Existenz und Eindeutigkeit der Polynominterpolation

Polynominterpolation ist eine einfache und effiziente Möglichkeit der Interpolation. Es bezeichne Π_n die Menge aller reellen oder komplexen Polynome vom Grade $\leq n$. Es gilt der

Satz 5.3 (Eindeutigkeit). *Zu beliebigen $n+1$ Stützstellen (x_j, f_j) , $j = 0, \dots, n$, $x_j \neq x_k$ für $j \neq k$ gibt es genau ein Polynom $P \in \Pi_n$*

$$P(x) = a_0 + a_1 x + \dots + a_n x^n$$

mit

$$(5.2) \quad P(x_j) = f_j, \quad j = 0, \dots, n.$$

Beweis: Die Zerlegung von (5.2) ergibt das LGS

$$\sum_{k=0}^n a_k x_j^k = f_j, \quad j = 0, \dots, n,$$

d.h.

$$(5.3) \quad \begin{pmatrix} 1 & x_0 & \cdots & x_0^n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^n \end{pmatrix} \begin{pmatrix} a_0 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} f_0 \\ \vdots \\ f_n \end{pmatrix}$$

Die Determinante der sogenannten *Vandermonde-Matrix* in (5.3) ist

$$\prod_{i=1}^n \left(\prod_{j=i+1}^n (x_i - x_j) \right)$$

und damit ungleich Null, wenn die x_i paarweise verschieden sind. Also ist die Matrix invertierbar und das Gleichungssystem eindeutig lösbar. \diamond

Bemerkung 5.4. *Prinzipiell ist es somit möglich aus dem LGS die Koeffizienten zu bestimmen, jedoch mit dem Wissen eines Aufwandes von $O(n^3)$ Operationen recht teuer.*

Aus diesem Grund betrachten wir andere Techniken.

5.2.2 Interpolationsformel von Lagrange

Für die $n + 1$ *Lagrange-Polynome* vom Grad n

$$L_i(x) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}$$

gilt offensichtlich

$$L_i(x_j) = \begin{cases} 1, & \text{für } i = j \\ 0, & \text{für } i \neq j \end{cases}$$

Damit ist unser gesuchtes Polynom gegeben durch

$$\boxed{P(x) = \sum_{i=0}^n f_i L_i(x)}, \quad \text{da } P(x_j) = f_j.$$

Beispiel 5.5. *Wir betrachten die Daten $(3, 68)$, $(2, 16)$, $(5, 352)$. Die Lagrange-Polynome sind gegeben durch*

$$L_0(x) = \frac{(x-2)(x-5)}{(3-2)(3-5)} = -\frac{1}{2}(x-2)(x-5)$$

$$L_1(x) = \frac{(x-3)(x-5)}{(2-3)(2-5)} = \frac{1}{3}(x-3)(x-5)$$

$$L_2(x) = \frac{(x-2)(x-3)}{(5-2)(5-3)} = \frac{1}{6}(x-2)(x-3)$$

Damit erhalten wir

$$\begin{aligned} P(x) &= -68 \cdot \frac{1}{2}(x-2)(x-5) + 16 \cdot \frac{1}{3}(x-3)(x-5) + 352 \cdot \frac{1}{6}(x-2)(x-3) \\ &= 30x^2 - 98x + 92 \end{aligned}$$

Bemerkung 5.6. Ein Abzählen der notwendigen Operationen ergibt den Aufwand $O(n^2)$.

Die Lagrange-Darstellung ist praktisch, wenn man Messwerte f_i nachträglich verändern will, da diese explizit in der Polynomdarstellung auftauchen. Sie ist aber i.A. unpraktisch, wenn man Datenpunkte hinzufügen will, da man alle L_i neu berechnen muss.

5.2.3 Der Algorithmus von Aitken und Neville

Gesucht ist ein numerisch sparsamer Algorithmus zur Berechnung von $P(x)$ an einigen wenigen Stellen (z.B. nur einmalig) x .

Für $i_0, \dots, i_k \in \{0, 1, \dots, n\}$ sei $P_{i_0, \dots, i_k} \in \Pi_k$ das Interpolationspolynom zu x_{i_0}, \dots, x_{i_k} ; insbesondere

$$P_i(x) \equiv f_i, \quad P_{0, \dots, n}(x) = P(x).$$

5.2.3.1 Rekursionsformel von Aitken

Es gilt die Rekursionsformel von Aitken

$$(5.4) \quad P_{i_0, \dots, i_k}(x) = \frac{(x-x_{i_0})P_{i_1, \dots, i_k}(x) - (x-x_{i_k})P_{i_0, \dots, i_{k-1}}(x)}{x_{i_k} - x_{i_0}}$$

Beweis: Das Polynom $Q(x) \in \Pi_k$ auf der rechten Seite von (5.4) erfüllt

$$\begin{aligned} Q(x_{i_0}) &= P_{i_0, \dots, i_{k-1}}(x_{i_0}) &&= f_{i_0} \\ Q(x_{i_k}) &= P_{i_1, \dots, i_k}(x_{i_k}) &&= f_{i_k} \\ Q(x_{i_j}) &= f_{i_j}, \quad j = 1, \dots, k-1. \end{aligned}$$

Wegen der Eindeutigkeit der Polynom-Interpolation folgt dann (5.4). \diamond

5.2.3.2 Variante von Neville

Es sei x fest. Erzeuge die Werte $P_{i-k, \dots, i}(x)$ ($k \leq i$) nach dem Schema

	$k = 0$	1	2	3
x_0	$f_0 = P_0(x)$			
		\searrow $P_{0,1}(x)$		
x_1	$f_1 = P_1(x)$		\searrow $P_{0,1,2}(x)$	
		\searrow $P_{1,2}(x)$		\searrow $P_{0,1,2,3}(x)$
x_2	$f_2 = P_2(x)$		\searrow $P_{1,2,3}(x)$	
		\searrow $P_{2,3}(x)$		
x_3	$f_3 = P_3(x)$			

Die Auswertung dieses Ausdrucks für festes x geschieht mit dem Horner-Schema:

$$P(x) = [\dots (a_n(x - x_{n-1}) + a_{n-1})(x - x_{n-2}) + \dots + a_1](x - x_0) + a_0$$

d.h. rekursiv

$$\begin{aligned} b_n &= a_n, \\ b_i &= b_{i+1}(x - x_i) + a_i, \quad i = n - 1, \dots, 0, \\ P(x) &= b_0. \end{aligned}$$

Für die Abschnittspolynome

$$Q_k(x) = \sum_{i=0}^k a_i \prod_{j=0}^{i-1} (x - x_j)$$

folgt

$$(a) \quad Q_k(x) = P_{0,1,\dots,k}(x), \quad k = 0, \dots, n,$$

$$(b) \quad a_k \text{ ist Koeffizient von } x^k \text{ in } P_{0,1,\dots,k}(x).$$

Die Koeffizienten a_k werden nicht mittels der Bedingung $P(x_j) = f_j$ berechnet, sondern mit Hilfe des Differenzenschemas:

$$\begin{array}{l} x_0 \left| \begin{array}{l} [f_0] \\ \searrow \\ [f_0, f_1] \\ \searrow \\ [f_0, f_1, f_2] \end{array} \\ x_1 \left| \begin{array}{l} [f_1] \\ \searrow \\ [f_1, f_2] \end{array} \\ x_2 \left| \begin{array}{l} [f_2] \end{array} \end{array} \right. \right.$$

Die hier auftretenden „Dividierten Differenzen“ $[f_i, \dots, f_k]$ sind rekursiv definiert durch

$$(a) \quad [f_i] = f_i,$$

$$(b) \quad [f_i, \dots, f_k] = \frac{[f_{i+1}, \dots, f_k] - [f_i, \dots, f_{k-1}]}{x_k - x_i},$$

z.B.

$$\begin{aligned} [f_0, f_1] &= \frac{f_1 - f_0}{x_1 - x_0} \\ [f_0, f_1, f_2] &= \frac{[f_1, f_2] - [f_0, f_1]}{x_2 - x_0}. \end{aligned}$$

Satz 5.7. *Es gilt:*

$$P_{0,\dots,k}(x) = \sum_{i=0}^k [f_0, \dots, f_i] \prod_{j=0}^{i-1} (x - x_j).$$

Beweis: Durch Induktion bzgl. k : Für $k = 0$ ist die Behauptung offensichtlich richtig: Die Aussage gelte für $k - 1 \geq 0$.

$$P_{0,\dots,k}(x) = P_{0,\dots,k-1}(x) + a(x - x_0) \cdots (x - x_{k-1}).$$

Zu zeigen ist $a = [f_0, \dots, f_k]$.

$$\begin{aligned} \text{Koeffizient von } x^k \text{ in } P_{0,\dots,k}(x) &: a \\ \text{Koeffizient von } x^{k-1} \text{ in } P_{0,\dots,k-1}(x) &: [f_0, \dots, f_{k-1}] \\ \text{Koeffizient von } x^{k-1} \text{ in } P_{1,\dots,k}(x) &: [f_1, \dots, f_k] \end{aligned}$$

(nach Induktionsvoraussetzung)

Nach der Formel von Aitken (5.4) ist

$$P_{0,\dots,k}(x) = \frac{(x - x_0)P_{1,\dots,k}(x) - (x - x_k)P_{0,\dots,k-1}(x)}{x_k - x_0}.$$

Der Koeffizient von x^k auf der rechten Seite ist

$$a = \frac{[f_1, \dots, f_k] - [f_0, \dots, f_{k-1}]}{x_k - x_0} = [f_0, \dots, f_k].$$

↑
Definition

◇

Beispiel 5.8. (vgl. Lagrange-Form)

Das Differenzenschema lautet

$$\begin{array}{rcl} [f_0] = \boxed{68} & & \\ & \searrow & \\ & [f_0, f_1] = \frac{16-68}{2-3} = \boxed{52} & \\ [f_1] = 16 & & \searrow [f_0, f_1, f_2] = \frac{52-112}{3-5} = \boxed{30} \\ & \searrow & \\ & [f_1, f_2] = \frac{16-352}{2-5} = 112 & \\ [f_2] = 352 & & \end{array}$$

Damit erhalten wir das Interpolationspolynom (vgl. (5.5))

$$P(x) = P_{0,2}(x) = 68 + 52(x - 3) + 30(x - 3)(x - 2) = 30x^2 - 98x + 92$$

Satz 5.9. Für eine beliebige Permutation i_0, \dots, i_n von $0, \dots, n$ gilt

$$[f_{i_0}, \dots, f_{i_n}] = [f_0, \dots, f_n].$$

Bemerkung 5.10. Die Berechnung der dividierten Differenzen kann offensichtlich mit $O(n^2)$ Rechenoperationen erfolgen. Die Bestimmung des Polynoms P hingegen benötigt nur $O(n)$ Operationen, wenn man zur Auswertung das Horner-Schema verwendet.

Mit dem Newton-Schema können leicht zusätzliche Datenpunkte mit in die Interpolation aufgenommen werden. Außerdem erlaubt das Horner-Schema eine effiziente Berechnung des Polynoms, wenn selbiges sehr oft für verschiedene x auszuwerten ist.

5.2.5 Interpolationsfehler

Wir betrachten den Fehler $f(x) - P(x)$ zwischen einer Funktion f und dem Interpolationspolynom P zu den Stützstellen (x_j, f_j) , $f_j = f(x_j)$, $j = 0, \dots, n$. Es gilt der

Satz 5.11. *Sei $f \in C^{n+1}[a, b]$ und $\bar{x}, x_0, \dots, x_n \in [a, b]$. Dann gibt es ein $\xi \in [a, b]$ mit*

$$(5.7) \quad f(\bar{x}) - P(\bar{x}) = \frac{f^{(n+1)}(\xi)}{(n+1)!} L(\bar{x}), \quad L(x) := (x-x_0)(x-x_1) \dots (x-x_n)$$

Beweis: Betrachte für $\bar{x} \neq x_j$ die Funktion

$$F(x) := f(x) - P(x) - \frac{f(\bar{x}) - P(\bar{x})}{L(\bar{x})} L(x) \in C^{n+1}[a, b]$$

F hat die $n+2$ Nullstellen \bar{x}, x_0, \dots, x_n in $[a, b]$. Nach dem Satz von Rolle (Zwischen zwei Nullstellen einer Funktion $F(x)$ liegt (mindestens) eine Nullstelle ihrer Ableitung $F'(x)$) hat $F^{(n+1)}(x)$ mindestens eine Nullstelle $\xi \in [a, b]$:

$$0 = F^{(n+1)}(\xi) = f^{(n+1)}(\xi) - \frac{f(\bar{x}) - P(\bar{x})}{L(\bar{x})} (n+1)!$$

\Rightarrow Behauptung. ◇

Bemerkung 5.12. *Für den Interpolationsfehler erhalten wir somit die Abschätzung*

$$(5.8) \quad |f(x) - P(x)| \leq \max_{y \in [a, b]} \left| \frac{f^{(n+1)}(y)}{(n+1)!} L(x) \right|.$$

Bemerkung 5.13. *Wegen $|x - x_i| \leq b - a$ folgt sofort*

$$(5.9) \quad |f(x) - P(x)| \leq \max_{y \in [a, b]} |f^{(n+1)}(y)| \frac{(b-a)^{n+1}}{(n+1)!}.$$

Beispiel 5.14. *Es sei $f(x) = \sin(x)$ auf $[0, 2\pi]$. Aus*

$$f^{(1)}(x) = \cos(x), \quad f^{(2)}(x) = -\sin(x), \quad f^{(3)}(x) = -\cos(x), \dots$$

folgt $|f^{(k)}(y)| \leq 1, \forall y \in \mathbb{R}$. Mit den äquidistanten Stützstellen $x_i = 2\pi i/n$ folgt aus (5.9)

$$|f(x) - P(x)| \leq \frac{(2\pi)^{n+1}}{(n+1)!}.$$

Man sieht leicht, daß bereits für kleine n eine sehr gute Übereinstimmung der Funktionen erwartet werden kann.

5.2.6 Konvergenz

Sei

$$\Delta_m := \{a = x_0^{(m)} < x_1^{(m)} < \dots < x_{n_m}^{(m)} = b\}, \quad m = 0, 1, \dots,$$

eine Folge von Intervallteilungen von $[a, b]$.

$$\|\Delta_m\| = \max_i |x_{i+1}^{(m)} - x_i^{(m)}|.$$

$P_{\Delta_m}(x)$ bezeichne das interpolierende Polynom von f bzgl. Δ_m .

Problem: Gilt $\lim_{m \rightarrow \infty} P_{\Delta_m}(x) = f(x)$ für

$$\lim_{m \rightarrow \infty} \|\Delta_m\| = 0 ?$$

Antwort: i.A. nicht richtig.

Beispiel 5.15. (Beispiel von Runge)

$$f(x) = (1 + 25x^2)^{-1} \text{ in } [-1, 1]$$

$f(x)$ ist bzgl. $x \in \mathbb{C}$ keine ganze Funktion: vgl. Satz 5.17.

f wird an den Stellen $x_j = -1 + \frac{2j}{n}$, $j = 0, \dots, n$, interpoliert durch ein Polynom

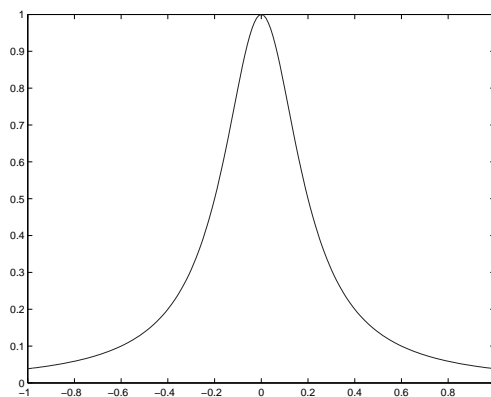


Abbildung 5.1: Beispiel von Runge.

P_n .

n	$\ f - P_n\ _\infty$
1	0.96
5	0.43
13	1.07
19	8.57

Wir wählen nun die Knoten x_0, \dots, x_n so, dass $\|L\|_\infty = \max_{-1 \leq x \leq 1} |L(x)|$ möglichst klein wird, vgl. (5.8). Man erhält

$$L(x) = 2^{-n} T_{n+1}(x), \quad x \in [-1, 1],$$

wobei die Tschebyscheff-Polynome T_n wie folgt definiert sind:

$$T_n(x) = \cos n\theta, \quad x = \cos \theta, \quad 0 \leq \theta \leq 2\pi.$$

Es gilt $\|T_n\|_\infty \leq 1$ und die Nullstellen von $L(x) = 2^{-n} T_{n+1}(x)$ sind

$$x_j = \cos \frac{(j + \frac{1}{2})\pi}{n+1}, \quad j = 0, \dots, n.$$

Bei dieser Wahl der Knoten ergibt sich die Fehlerabschätzung

n	$\ f - P_n\ _\infty$
1	0.93
5	0.56
13	0.12
19	0.04

Die Verbesserung ist erheblich, die Approximation aber immer noch unbefriedigend.

Ohne Beweis geben wir noch an:

Satz 5.16. Zu jeder Folge $\{\Delta_m\}$ gibt es $f \in C[a, b]$, so dass $\{P_{\Delta_m}\}$ nicht gleichmäßig gegen f konvergiert.

Satz 5.17. Sei f eine ganze Funktion. Dann gilt $P_{\Delta_m} \rightarrow f$ gleichmäßig für alle $\{\Delta_m\}$ mit $\|\Delta_m\| \rightarrow 0$.

Bemerkung 5.18. $P(x) = P_{0, \dots, n}(x)$ oszilliert i.A. stark für großes n ; die Interpolation ist dann unbrauchbar. Diese Schwierigkeiten werden bei der Spline-Interpolation vermieden.

5.3 Trigonometrische Interpolation

5.3.1 Diskrete Fouriertransformation

Gegeben seien Stützstellen (x_k, f_k) , $k = 0, \dots, n-1$ mit

$$x_k = \frac{2\pi}{n} k$$

Gesucht ist ein *Trigonometrisches Polynom*

$$P(\omega) = p(x) = \sum_{k=0}^{n-1} \beta_k \omega^k$$

mit

$$\omega := e^{ix} = \cos(x) + i \cdot \sin(x), \quad i = \sqrt{-1}$$

so dass

$$p(x_k) = f_k, \quad k = 0, 1, \dots, n - 1.$$

Bemerkung 5.19. *Ausdrücke der Form $e^{\frac{2\pi i}{n} \cdot k}$ heißen n te Einheitswurzeln $k = 0, \dots, n - 1$. Sie liegen alle auf dem Einheitskreis.*

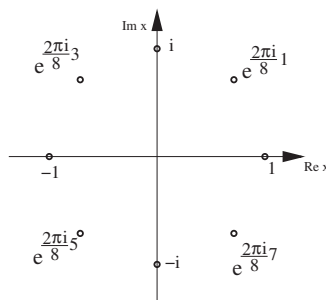


Abbildung 5.2: Einheitswurzeln, hier am Beispiel $n = 8$.

Nach Satz 5.3 ist die Interpolationsaufgabe

$$(5.10) \quad P(\omega_k) = \sum_{j=0}^{n-1} \beta_j \omega_k^j = f_k, \quad \omega_k = e^{i \cdot x_k}, \quad (\text{Fourier-Synthese})$$

eindeutig lösbar. Die Interpolationsaufgabe (5.10) ist schreibbar als

$$f = W\beta$$

mit

$$f = \underbrace{\begin{pmatrix} f_0 \\ \vdots \\ f_{n-1} \end{pmatrix}}_{=:f} = \underbrace{\begin{pmatrix} 1 & \omega_0 & \cdots & \omega_0^{n-1} \\ \vdots & \vdots & & \vdots \\ 1 & \omega_{n-1} & \cdots & \omega_{n-1}^{n-1} \end{pmatrix}}_{=:W} \underbrace{\begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{n-1} \end{pmatrix}}_{=: \beta}$$

Es gilt die Beziehung

$$(5.11) \quad W\bar{W}^T = nI, \quad \text{mit } \bar{W} \text{ konjugiert komplexe Matrix zu } W,$$

da

$$(W\bar{W}^T)_{k,l} = \sum_{j=0}^{n-1} \omega_k^j \cdot \bar{\omega}_l^j = \sum_{j=0}^{n-1} e^{\frac{2\pi i}{n} \cdot k \cdot j} \cdot e^{-\frac{2\pi i}{n} \cdot l \cdot j} = \sum_{j=0}^{n-1} e^{\frac{2\pi i}{n} \cdot j(k-l)} = \sum_{j=0}^{n-1} q^j$$

mit $q = e^{\frac{2\pi i}{n} \cdot (k-l)}$. Dies ist eine geometrische Reihe und wir erhalten weiter

$$\sum_{j=0}^{n-1} q^j = \begin{cases} \frac{1-q^n}{1-q}, & \text{falls } q \neq 1 \\ n, & \text{falls } q = 1 \end{cases} \stackrel{q^n=1}{=} \begin{cases} 0, & \text{falls } k \neq l \\ n, & \text{falls } k = l. \end{cases}$$

Damit ist die Inverse von W aus (5.11) einfach zu berechnen:

$$(5.12) \quad W^{-1} = \frac{1}{n} \bar{W}^T \quad \Rightarrow \quad \boxed{\beta = \frac{1}{n} \bar{W}^T f}$$

In Komponentenschreibweise:

$$(5.13) \quad \boxed{\beta_j = \frac{1}{n} \sum_{k=0}^{n-1} f_k \cdot e^{-\frac{2\pi i}{n} \cdot j \cdot k}}$$

Formel (5.13) wird *Diskrete Fouriertransformation der Länge n* genannt.

Bemerkung 5.20. Der Aufwand zur Berechnung aller Komponenten beträgt offensichtlich $O(n^2)$.

5.3.2 Trigonometrische Interpolation

Der Zusammenhang mit der Überschrift dieses Abschnittes ergibt sich durch den nachfolgenden Satz.

Satz 5.21 (Reelle Fouriertransformation). *Definiert man*

$$a_j := \frac{2}{n} \sum_{k=0}^{n-1} f_k \cdot \cos\left(\frac{2\pi k j}{n}\right), \quad b_j := \frac{2}{n} \sum_{k=0}^{n-1} f_k \cdot \sin\left(\frac{2\pi k j}{n}\right), \quad j = 0, \dots, n-1.$$

und setzt für ungerades $n = 2m + 1$

$$\Psi(x) := \frac{a_0}{2} + \sum_{j=1}^m (a_j \cos(jx) + b_j \sin(jx))$$

bzw. für gerades $n = 2m$

$$\Psi(x) := \frac{a_0}{2} + \sum_{j=1}^{m-1} (a_j \cos(jx) + b_j \sin(jx)) + \frac{a_m}{2} \cos(mx)$$

so gilt

$$\Psi(x_k) = f_k, \quad k = 0, \dots, n-1$$

mit x_k aus (5.3.1).

Beweis: Zunächst notieren wir einige interessante Beziehungen.

Wegen $\omega_k^n = 1$ folgt für die komplexen Koeffizienten:

$$(5.14) \quad \beta_{n-j} = \frac{1}{n} \sum_{k=0}^{n-1} f_k \omega_k^{j-n} = \frac{1}{n} \sum_{k=0}^{n-1} f_k \omega_k^j.$$

Der Zusammenhang der reellen Fourierkoeffizienten mit den komplexen Koeffizienten ergibt sich wegen (5.13) zu

$$(5.15) \quad \beta_j = \frac{1}{2}(a_j - ib_j), \quad \text{und} \quad \beta_{n-j} = \frac{1}{2}(a_j + ib_j).$$

Hierraus folgt unmittelbar die Beziehung

$$(5.16) \quad \beta_j \omega_k^j + \beta_{n-j} \omega_k^{n-j} = a_j \cos jx_k + b_j \sin jx_k.$$

Nach diesen Vorüberlegungen wenden wir uns nun dem eigentlichen Beweis zu. Für gerades $n = 2m$ ist

$$b_0 = 0, \quad b_m = 0, \quad \omega_k^m = \cos(mx_k),$$

und daher gilt (wegen $f = W\beta$) mit (5.14)–(5.16)

$$\begin{aligned} f_k = \sum_{j=0}^{n-1} \beta_j \omega_k^j &= \underbrace{\beta_0}_{\frac{a_0}{2}} + \sum_{j=1}^{m-1} \underbrace{(\beta_j \omega_k^j + \beta_{n-j} \omega_k^{n-j})}_{a_j \cos(jx_k) + b_j \sin(jx_k)} + \underbrace{\beta_m \omega_k^m}_{\frac{a_m}{2} \cos(mx)} \\ &= \psi(x_k). \end{aligned}$$

Ebenso folgt die Behauptung für ungerades $n = 2m + 1$. ◇

Bemerkung 5.22. Der Aufwand zur Berechnung aller Komponenten beträgt $O(n^2)$.

5.3.3 Schnelle Fourier–Transformation (FFT)

Cooley–Tukey haben eine Möglichkeit gefunden die diskrete Fouriertransformation effizienter durchzuführen. Sei hierzu $q = e^{\frac{-2\pi i}{n}}$ und n gerade, $n = 2m$, dann ist (5.13) darstellbar als

$$(5.17) \quad \beta_k = \frac{1}{n} \sum_{j=0}^{n-1} f_j \cdot q^{j \cdot k}.$$

Es gilt $q^n = 1$ und $q^{n/2} = -1$.

Idee: Spalte Summe in (5.17) nach geraden und ungeraden Indizes auf.

$$\begin{aligned}\beta_k &= \frac{1}{n} \left(\sum_{l=0}^{m-1} f_{2l} \cdot q^{(2l)k} + \sum_{l=0}^{m-1} f_{2l+1} \cdot q^{(2l+1)k} \right) \\ &= \frac{1}{2} \left(\frac{1}{n/2} \sum_{l=0}^{m-1} f_{2l} \cdot (q^2)^{lk} + q^k \frac{1}{n/2} \sum_{l=0}^{m-1} f_{2l+1} \cdot (q^2)^{lk} \right) \\ &=: \frac{1}{2} (g_k + q^k u_k)\end{aligned}$$

Wegen $q^2 = e^{\frac{2\pi i}{n/2}}$ lassen sich g_k und u_k als Fouriertransformation der Länge $\frac{n}{2}$ berechnen. Weiterhin folgt durch einfaches Nachrechnen

$$(5.18) \quad \begin{aligned}g_{k+m} &= g_k \\ u_{k+m} &= u_k\end{aligned}$$

und damit für $k = 0, 1, \dots, m-1$

$$(5.19) \quad \beta_k = \frac{1}{2} (g_k + q^k u_k)$$

und

$$(5.20) \quad \beta_{k+m} = \frac{1}{2} (g_k + q^{k+m} u_k) = \frac{1}{2} (g_k - q^k u_k).$$

Bemerkung 5.23. Die g_k aus (5.19) können somit bei (5.20) direkt ohne zusätzlichen Rechenaufwand wieder verwendet werden.

Anstelle von 1 mal n^2 hat man jetzt nur noch 2 mal $(n/2)^2$, also $\frac{1}{2}n^2$ als Aufwand.

Sei nun $n = 2^p$. Die Idee ist es nun die zuvor dargestellte Berechnungsvorschrift mehrfach anzuwenden, also auch g_k und u_k gemäß (5.20) zu bestimmen. Sei M_p die Anzahl an Multiplikationen/Divisionen, die wir für die FFT benötigen. Wir vernachlässigen die Berechnung von q^k . Es gilt

$$\begin{aligned}M_0 &= 0 \\ M_{p+1} &= 2M_p + \frac{n}{2} + \frac{n}{2} = 2M_p + 2^p\end{aligned}$$

also

$$M_p = 2^{p-1} + 2 \cdot 2^{p-2} + \dots + \underbrace{2 \cdot 2 \cdot \dots \cdot 2}_{p-1 \text{ mal}} = p \cdot 2^{p-1} = \frac{1}{2} \cdot p \cdot 2^p = \boxed{\frac{1}{2} \log_2 n \cdot n}$$

Damit gilt der

Satz 5.24 (Aufwand FFT). Eine Fouriertransformation der Länge $n = 2^p$ kann mit $O(n \cdot \log_2 n)$ komplexen Multiplikationen berechnet werden.

Bemerkung 5.25. Die Anzahl an Additionen/Subtraktionen ergibt sich ebenfalls zu $O(n \cdot \log_2 n)$.

Bemerkung 5.26. Es gilt:

$$O(n) < O(n \cdot \log_2 n) < O(n^2).$$

5.3.4 Anwendungen

Die FT findet Ihre Anwendung z.B. in der *Signalverarbeitung*. Eine wichtige Anwendung ist die *Frequenzanalyse*, hier haben die Fourierkoeffizienten der trigonometrischen Polynome (im Gegensatz zu herkömmlichen Koeffizienten) eine klare physikalische Bedeutung: Jeder der Fourierkoeffizienten a_k bzw. b_k gehört zu einer bestimmten Frequenz im Sinus bzw. Kosinus. Die Anwendbarkeit der FT ergibt sich z.B. wenn Signale durch Störungen (Rauschen) beeinflusst werden. Die FT glättet und das Originalresultat kann häufig reproduziert werden, bzw. die Hauptfrequenzen des Signals können besser identifiziert werden. Dies nutzt man dann z.B. beim *Filtern* aus. Hat man einmal ein Signal zerlegt und stehen die Fourierkoeffizienten (die eine bestimmte Frequenz beschreiben) zur Verfügung kann man durch Vergrößern oder Verkleinern der einzelnen Koeffizienten spezielle Frequenzen unterschiedlich stark betonen oder sogar ganz herausnehmen (Filtern). Weitere Anwendungen sind die Zerlegung eines Signals in verschiedenen Frequenzbereiche oder Datenkompression (Null-Setzen der Koeffizienten mit geringem Betrag). Da die Anwendungen alle in Echtzeit durchgeführt werden müssen ist eine effiziente Implementierung im Sinne der FFT offensichtlich.

5.4 Spline-Interpolation

Spline-Funktionen, die sich stückweise aus Polynomen zusammensetzen, verbinden den Vorteil einer glatten oszillationsfreien Interpolation mit denjenigen, die der Umgang mit Polynomen niedrigen Grades mit sich bringt.

5.4.1 Polynom-Splines

Sei $\Delta = \{a = x_0 < x_1 < \dots < x_n = b\}$ eine Zerlegung des Intervalls $[a, b] \subset \mathbb{R}$ mit inneren Knoten x_1, \dots, x_{n-1} und Randknoten x_0, x_n .

Definition 5.27. Eine Funktion $s : [a, b] \rightarrow \mathbb{R}$ heißt Polynom-Spline vom Grad l , $l = (0, 1, 2, \dots)$ zur Zerlegung Δ , wenn sie folgende Eigenschaften besitzt:

$$(a) s \in C^{l-1}[a, b]$$

$$(b) s \in \Pi_l \text{ für } x_j \leq x < x_{j+1}, \quad j = 0, 1, \dots, n-1.$$

Hierbei ist $C^{-1}[a, b]$ der Raum der auf $[a, b]$ stückweise stetigen Funktionen. Die Menge aller Polynom-Splines vom Grade l zur Zerlegung Δ bezeichnen wir mit $S_l(\Delta)$. Fortan wird schlechthin von Splines gesprochen.

Beispiel 5.28. (Lineare Splines): Sind $(n+1)$ Punkte

$$(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$$

gegeben, so stellt der Polygonzug durch diese Punkte einen Spline $s \in S_1(\Delta)$ dar.

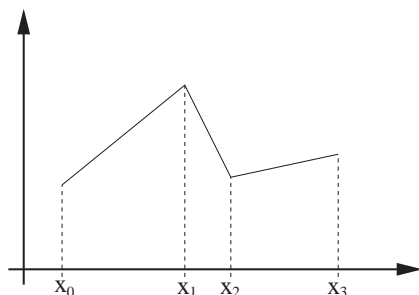


Abbildung 5.3: Polynomspline vom Grade 1.

Beispiel 5.29. (Quadratische Splines): Bei äquidistanten Knoten $x_j = a + jh$, $j = 0, \dots, n$, ist

$$s(x) = \frac{1}{2h^2} \begin{cases} (x - x_j)^2 & , x_j \leq x < x_{j+1} \\ h^2 + 2h(x - x_{j+1}) - 2(x - x_{j+1})^2, & x_{j+1} \leq x < x_{j+2} \\ (x_{j+3} - x)^2 & , x_{j+2} \leq x < x_{j+3} \end{cases}$$

ein quadratischer B-Spline $s \in S_2(\Delta)$.

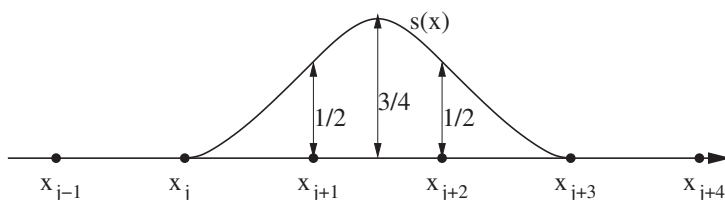


Abbildung 5.4: Polynomspline vom Grade 2.

Beispiel 5.30. Die Funktionen $q_j : [a, b] \rightarrow \mathbb{R}$, $j = 0, 1, \dots, n - 1$,

$$q_j(x) = (x - x_j)_+^l := \begin{cases} (x - x_j)^l & , x \geq x_j \\ 0 & , x < x_j \end{cases}$$

sind Splines vom Grade l zu Δ . Man beachte, dass q_j keine Polynome in $[a, b]$ sind.

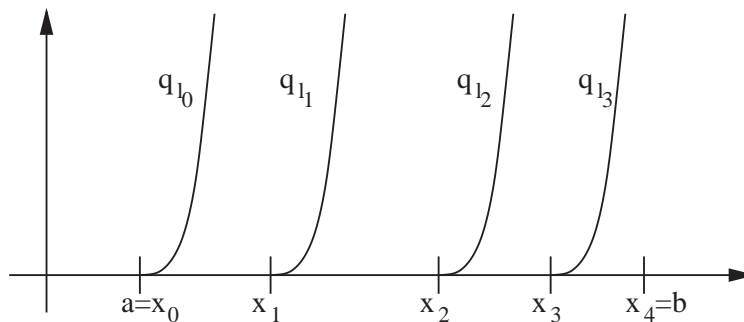


Abbildung 5.5: Splines vom Grade l .

Der Spline-Raum $S_l(\Delta)$ ist ein linearer Teilraum von $C^{l-1}(\Delta)$. Der folgende Satz gibt Auskunft über eine Basis von $S_l(\Delta)$.

Satz 5.31. Die Menge $S_l(\Delta)$ ist ein linearer Raum der Dimension $n + l$. Die Elemente $p_i(x) = x^i$, ($i = 0, \dots, l$), $q_{lj}(x) = (x - x_j)_+^l$ ($j = 1, \dots, n - 1$) bilden eine Basis von $S_l(\Delta)$.

Beweis: Wir haben zu zeigen, dass es zu $s \in S_l(\Delta)$ eine eindeutige Darstellung

$$s(x) = \sum_{i=0}^l a_i x^i + \sum_{j=1}^{n-1} b_j (x - x_j)_+^l, \quad x \in [a, b],$$

gibt. Dies erkennt man durch Induktion bzgl. des Index j der Zerlegung Δ . Im Intervall $I_1 = [x_0, x_1]$ ist s ein Polynom $s(x) = a_0 + a_1 x + \dots + a_l x^l$, also gilt die Darstellung

$$s(x) = \sum_{i=0}^l a_i x^i + \sum_{j=1}^{k-1} b_j (x - x_j)_+^l$$

für $k = 1$ auf $I_k = [x_0, x_k]$. Wir betrachten nun

$$d(x) := s(x) - \sum_{i=0}^l a_i x^i - \sum_{j=1}^{k-1} b_j (x - x_j)_+^l.$$

Dann ist $d \in C^{l-1}(I_{k+1})$ und $d(x) = 0$ für $x \in I_k$. Außerdem ist $d \in \Pi_l$ in $[x_k, x_{k+1}]$. Also genügt d auf $[x_k, x_{k+1}]$ der Differentialgleichung

$$\begin{aligned} d^{(l+1)}(x) &= 0, & x_k \leq x \leq x_{k+1} \\ d^{(i)}(x_k) &= 0, & i = 0, \dots, l-1 \end{aligned}$$

Die Lösung dieser Anfangswertaufgabe ist

$$d(x) = b_k(x - x_k)_+^l, \quad x_k \leq x, \quad b_k \in \mathbb{R}.$$

Damit ist die Behauptung für den Index $k+1$ gezeigt. Für $k=n$ bilden daher die $n+l$ linear unabhängigen Elemente

$$p_i(x) = x^i, \quad (i = 0, \dots, l), \quad q_j(x) = (x - x_j)_+^l, \quad (j = 1, \dots, n-1)$$

eine Basis von $S_l(\Delta)$. ◇

Wir untersuchen nun die folgende Interpolationsaufgabe:

Bestimme zu Stützwerten $y_j, j = 0, \dots, n$ einen interpolierenden Spline $s \in S_l(\Delta)$ mit

$$s(x_j) = y_j, \quad j = 0, \dots, n.$$

Die Stützwerte y_j sind dabei als Funktionswerte $y_j = f(x_j)$ einer hinreichend glatten Funktion f aufzufassen. Wegen $\dim S_l(\Delta) = n+l$ können über die $n+1$ Interpolationsbedingungen hinaus noch

$$n+l - (n+1) = l-1$$

freie Parameter bestimmt werden. Für ungerades $l = 2m-1$ ist dies eine gerade Anzahl $2m-2$ von Parametern, welche sich symmetrisch in den Randknoten x_0, x_n anordnen lassen.

Wir beschränken uns nachfolgend auf den in der Praxis wichtigen Fall $m=2$ der kubischen Splines, d.h. Polynom-Splines mit $l=3$, der sich besonders häufig in den Anwendungen wiederfindet.

5.4.2 Kubische Splines

5.4.2.1 Einführung und Aufgabenstellung

Nach Definition ist ein Spline $s \in C^2[x_0, x_n]$ auf jedem der Intervalle $[x_{k-1}, x_k]$, $k = 1, \dots, n$ durch

$$s(x) = s_k(x) = a_k + b_k(x - x_{k-1}) + c_k(x - x_{k-1})^2 + d_k(x - x_{k-1})^3$$

gegeben. Neben den Interpolationsbedingungen

$$(5.21) \quad s(x_j) = f(x_j), \quad j = 0, \dots, n$$

erfordert die Glattheit $s \in C^2[x_0, x_n]$ (an den Nahtstellen sollen keine Knicke auftreten und die Krümmungen sollen stetig ineinander übergehen) die zusätzlichen Bedingungen

$$(5.22) \quad s'_k(x_k) = s'_{k+1}(x_k) \quad \text{und} \quad s''_k(x_k) = s''_{k+1}(x_k), \quad k = 1, \dots, n-1.$$

Die $2n - 2$ linearen Gleichungen in (5.22) ergeben dann zusammen mit den $2n$ linearen Gleichungen aus (5.21) zusammen $4n - 2$ lineare Gleichungen für die $4n$ Unbekannten $a_k, b_k, c_k, d_k, k = 1, \dots, n$. Das Problem ist somit zwar lösbar, jedoch existieren ∞ viele Lösungen. Der Grund liegt in den Randpunkten x_0 und x_n , in denen wir keine Glattheitseigenschaften fordern mußten.

Um dennoch eine eindeutige Lösung erhalten zu können, müssen wir zwei weitere Bedingungen fordern, die üblicherweise an den Randpunkten verlangt werden. Hierzu gibt es verschiedenen Möglichkeiten:

(a) Natürliche Endbedingungen:

$$s''(a) = 0, \quad s''(b) = 0$$

(b) Hermite Endbedingungen:

$$(5.23) \quad s'(a) = f'(a), \quad s'(b) = f'(b)$$

(c) Periodische Endbedingungen:

$$s^{(i)}(a) = s^{(i)}(b), \quad i = 0, 1, 2,$$

falls f periodisch mit

$$f^{(i)}(a) = f^{(i)}(b), \quad i = 0, 1, 2.$$

5.4.2.2 Existenz und Eindeutigkeit

Man kann nun zeigen, dass die Interpolationsaufgabe (5.21) zusammen mit einer der Bedingungen aus (5.23) eindeutig lösbar ist. Zusätzlich erfüllt der interpolierende Spline eine Minimum-Norm-Eigenschaft bzgl. der Norm in $C^2[a, b]$.

$$(5.24) \quad \|f\|_2 := \left(\int_a^b (f''(x))^2 dx \right)^{\frac{1}{2}}, \quad f \in C^2[a, b].$$

Für diese Norm gilt

$$\begin{aligned}\|f\|_2 = 0 &\Leftrightarrow f''(x) = 0 \quad \text{für } x \in [a, b] \\ &\Leftrightarrow f \text{ linear in } [a, b].\end{aligned}$$

Satz 5.32 (Existenz, Eindeutigkeit und Extremaleigenschaft von Splines). Sei $f \in C^2[a, b]$. Dann gibt es genau einen Spline $s \in S_3(\Delta)$, der (5.21) und eine der Interpolationsbedingungen (5.23) erfüllt. Dieser interpolierende Spline genügt der Minimum-Norm-Bedingung

$$0 \leq \|f - s\|_2^2 = \|f\|_2^2 - \|s\|_2^2.$$

Beweis: Mit einem Spline $s \in S_3(\Delta)$ berechnet man für die Abweichung $d(x) = f(x) - s(x)$:

$$\begin{aligned}\|f - s\|_2^2 &= \int_a^b (f''(x) - s''(x))^2 dx \\ &= \|f\|_2^2 - \|s\|_2^2 - 2 \int_a^b d''(x) s''(x) dx.\end{aligned}$$

Da nur $s \in C^2[a, b]$ gilt, müssen wir für die partielle Integration aufspalten:

$$\begin{aligned}\int_a^b d''(x) s''(x) dx &= \sum_{j=1}^n \int_{x_{j-1}}^{x_j} d''(x) s''(x) dx \\ &= \sum_{j=1}^n \left\{ [d'(x) s''(x) - d(x) s^{(3)}(x)]_{x_{j-1}}^{x_j} + \int_{x_{j-1}}^{x_j} d(x) s^{(4)}(x) dx \right\}.\end{aligned}$$

Nun ist $s^{(4)}(x) \equiv 0$ auf $[x_{j-1}, x_j]$. Die Interpolationsforderungen (5.21), (5.23) bewirken gerade, dass

$$\sum_{j=1}^n [d'(x) s''(x) - d(x) s^{(3)}(x)]_{x_{j-1}}^{x_j} = [d'(x) s''(x) - d(x) s^{(3)}(x)]_a^b = 0$$

in den Fällen (5.23)(a)-(c). Damit ist die Minimum-Norm-Bedingung

$$0 \leq \|f - s\|_2^2 = \|f\|_2^2 - \|s\|_2^2$$

gezeigt, mit der sich die Eindeutigkeit von s folgendermaßen ergibt:

Ist \tilde{s} ein weiterer interpolierender Spline, so kann man $f = \tilde{s}$ in der letzten Ungleichung nehmen:

$$0 \leq \|\tilde{s} - s\|_2^2 = \|\tilde{s}\|_2^2 - \|s\|_2^2.$$

Durch Vertauschung von s und \tilde{s} sieht man $\|s - \tilde{s}\|_2 = 0$. Also ist $s - \tilde{s}$ linear. Wegen $s(x) - \tilde{s}(x) = 0$ für $x = a$ und $x = b$ muss dann $s = \tilde{s}$ gelten.

Zum Nachweis der Existenz greifen wir auf die Basis-Darstellung in Satz 5.31 zurück. Die Interpolationsforderungen stellen ein LGS in $(n + 3)$ Unbekannten $a_0, \dots, a_3, b_1, \dots, b_{n-1}$ dar. In den Fällen (5.23)(a)-(c) wird das System homogen, wenn $f \equiv 0$ zu interpolieren ist. Dann ist aber $s = 0$ interpolierender Spline und nach den obigen Überlegungen auch der einzige. \diamond

5.4.2.3 Geometrische und mechanische Interpretation

Die Extremaleigenschaft des kubischen Splines

$$\int_a^b (s''(x))^2 dx \leq \int_a^b (d''(x))^2 dx$$

erlaubt die folgende geometrische und mechanische Interpretation:

Die Krümmung $k(x)$ einer Kurve $y = f(x)$ in der (x,y) -Ebene ist gegeben durch

$$k(x) = \frac{f''(x)}{(1 + (f'(x))^2)^{3/2}} .$$

Unter der Annahme $|f'(x)| \ll 1$ wird die mittlere Gesamtkrümmung

$$\|k\|_2^2 \approx \int_a^b (f''(x))^2 dx .$$

Der kubische Spline minimiert also die Norm $\|k\|_2$ unter allen interpolierenden Funktionen.

Das Biegemoment $M(x)$ eines homogenen, isotropen Stabes, dessen Biegelinie durch $y = f(x)$ beschrieben wird, ist $M(x) = c_1 k(x)$, $c_1 > 0$. Die Biege-Energie ist dann näherungsweise

$$E(f) = c_2 \int_a^b M(x)^2 dx \approx c_3 \int_a^b (f''(x))^2 dx .$$

Wird ein gebogener Stab durch Lager in "Interpolationspunkten" fixiert, so wird die minimale Biege-Energie durch einen kubischen Spline realisiert. Außerhalb von $[a, b]$, wo der Stab nicht fixiert ist, nimmt er die spannungsfreie "natürliche" Lage $s''(x) = 0$ an. In diesem Sinne sind die Endbedingungen $s''(a) = 0$, $s''(b) = 0$ in (5.23)(a) als "natürlich" zu verstehen.

5.4.2.4 Die Berechnung von Spline-Funktionen

Zu berechnen sei die Spline-Funktion $s(x)$ mit $s(x_j) = y_j$, $j = 0, \dots, n$, welche zusätzlich eine der Eigenschaften (a),(b),(c) hat. Wir setzen

$$\begin{aligned} h_j &:= x_j - x_{j-1}, \quad j = 1, \dots, n, \\ M_j &:= s''(x_j) \quad j = 0, \dots, n, \quad (\text{Momente}). \end{aligned}$$

Da s'' linear in $[x_{j-1}, x_j]$ ist, gilt

$$s''(x) = \frac{1}{h_j}(M_j(x - x_{j-1}) + M_{j-1}(x_j - x)), \quad x_{j-1} \leq x \leq x_j.$$

Durch Integration erhält man für $x \in [x_{j-1}, x_j]$

$$\begin{aligned} s'(x) &= \frac{1}{2h_j}(M_j(x - x_{j-1})^2 - M_{j-1}(x_j - x)^2) + a_j, \\ s(x) &= \frac{1}{6h_j}(M_j(x - x_{j-1})^3 + M_{j-1}(x_j - x)^3) + a_j(x - x_{j-1}) + b_j \end{aligned}$$

mit $a_j, b_j \in \mathbb{R}$. Für die Koeffizienten a_j, b_j erhält man aus $s(x_{j-1}) = y_{j-1}$, $s(x_j) = y_j$ die Gleichungen

$$\begin{aligned} M_{j-1} \frac{h_j^2}{6} + b_j &= y_{j-1}, \\ M_j \frac{h_j^2}{6} + a_j h_j + b_j &= y_j, \end{aligned}$$

und daraus

$$\begin{aligned} b_j &= y_{j-1} - M_{j-1} \frac{h_j^2}{6}, \\ a_j &= \frac{y_j - y_{j-1}}{h_j} - \frac{h_j}{6}(M_j - M_{j-1}). \end{aligned}$$

Damit ergibt sich

$$\begin{aligned} s'(x_j^-) &= \frac{1}{2h_j} M_j h_j^2 + a_j = \frac{y_j - y_{j-1}}{h_j} + \frac{h_j}{3} M_j + \frac{h_j}{6} M_{j-1}, \\ s'(x_j^+) &= \frac{y_{j+1} - y_j}{h_{j+1}} - \frac{h_{j+1}}{3} M_j - \frac{h_{j+1}}{6} M_{j+1}. \end{aligned}$$

Wegen $s'(x_j^-) = s'(x_j^+)$ folgt dann

$$(5.25) \quad \boxed{\mu_j M_{j-1} + M_j + \lambda_j M_{j+1} = d_j, \quad j = 1, \dots, n-1}$$

mit

$$\begin{aligned} \mu_j &:= \frac{h_j}{2(h_j + h_{j+1})}, \\ \lambda_j &:= \frac{h_{j+1}}{2(h_j + h_{j+1})}, \quad \mu_j + \lambda_j = \frac{1}{2}, \\ d_j &:= \frac{3}{h_j + h_{j+1}} \left(\frac{y_{j+1} - y_j}{h_{j+1}} - \frac{y_j - y_{j-1}}{h_j} \right). \end{aligned}$$

Fall (5.23)(a): (Natürliche Endbedingungen)

Vorgegeben: Werte M_0, M_n

Die "natürliche" Bedingung $M_0 = M_n = 0$ ist darin enthalten. Dann stellt (5.25) ein LGS für M_1, \dots, M_{n-1} dar:

$$\begin{pmatrix} 1 & \lambda_1 & & & \\ \mu_2 & 1 & \lambda_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \mu_{n-2} & 1 & \lambda_{n-2} \\ & & & \mu_{n-1} & 1 \end{pmatrix} \begin{pmatrix} M_1 \\ M_2 \\ \vdots \\ \vdots \\ M_{n-1} \end{pmatrix} = \begin{pmatrix} d_1 - \mu_1 M_0 \\ d_2 \\ \vdots \\ d_{n-2} \\ d_{n-1} - \lambda_{n-1} M_n \end{pmatrix}$$

Die Matrix ist tridiagonal und wegen $\mu_j + \lambda_j = \frac{1}{2}$ diagonal-dominant, also LR-zerlegbar nach Satz 3.11. Das LGS kann daher mit Algorithmus (3.26) gelöst werden. Bei äquidistanten Knoten x_j gilt $\mu_j = \lambda_j = \frac{1}{4}$, also ist die Matrix symmetrisch und damit positiv definit.

Fall (5.23)(b): (Hermite Endbedingungen)

Vorgegeben: $s'(a) = y'_0, s'(b) = y'_n$.

Aus der Darstellung von $s'(x)$ folgt

$$\begin{aligned} \frac{h_1}{3} M_0 + \frac{h_1}{6} M_1 &= \frac{y_1 - y_0}{h_1} - y'_0 \\ \frac{h_n}{6} M_{n-1} + \frac{h_n}{3} M_n &= y'_n - \frac{y_n - y_{n-1}}{h_n}. \end{aligned}$$

Mit

$$\begin{aligned} \lambda_0 = \frac{1}{2}, \quad d_0 &= \frac{3}{h_1} \left(\frac{y_1 - y_0}{h_1} - y'_0 \right), \\ \mu_n = \frac{1}{2}, \quad d_n &= \frac{3}{h_n} \left(y'_n - \frac{y_n - y_{n-1}}{h_n} \right) \end{aligned}$$

erhalten wir ein LGS und den $(n + 1)$ Unbekannten M_0, \dots, M_n :

$$\begin{pmatrix} 1 & \lambda_0 & & & \\ \mu_1 & 1 & \lambda_1 & & \\ & \ddots & \ddots & \ddots & \\ & & \mu_{n-1} & 1 & \lambda_{n-1} \\ & & & \mu_n & 1 \end{pmatrix} \begin{pmatrix} M_0 \\ \vdots \\ \vdots \\ M_n \end{pmatrix} = \begin{pmatrix} d_0 \\ \vdots \\ \vdots \\ d_n \end{pmatrix}.$$

Die Matrix ist ebenfalls tridiagonal und diagonal-dominant, also LR-zerlegbar.

Beispiel 5.33. Für die nicht äquidistanten Knoten

j	x_j	y_j	M_j
0	0	0	0.022181
1	8.2	0.5	-0.000665
2	14.7	1.0	-0.010253
3	17.0	1.1	-0.006909
4	21.1	1.2	-0.000613
5	35.0	1.4	-0.000691
6	54.1	1.5	-0.000040
7	104	1.6	-0.000014
8	357	1.7	0.000004

$$y'_0 = 0.0012566,$$

$$y'_8 = 0.0001$$

liefert Algorithmus (3.26) die angegebenen Werte M_j .

Fall (5.23)(c): (Periodische Endbedingungen)

Hier ist

$$M_0 = M_n, \quad \text{da } s''(a) = s''(b)$$

$$y_0 = y_n, \quad \text{da } s(a) = s(b).$$

Die weitere Gleichung $s'(a) = s'(b)$ ergibt eine Beziehung

$$\mu_n M_{n-1} + M_n + \lambda_n M_{n+1} = d_n,$$

wenn man setzt $h_{n+1} = h_1$, $M_{n+1} = M_1$, $y_{n+1} = y_1$. Für M_1, \dots, M_n ist das LGS zu lösen

$$\begin{pmatrix} 1 & \lambda_1 & & & \mu_1 \\ \mu_2 & 1 & \lambda_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \mu_{n-1} & 1 & \lambda_{n-1} \\ \lambda_n & & & \mu_n & 1 \end{pmatrix} \begin{pmatrix} M_1 \\ \vdots \\ \vdots \\ M_n \end{pmatrix} = \begin{pmatrix} d_1 \\ \vdots \\ \vdots \\ d_n \end{pmatrix}.$$

Die Matrix ist nicht mehr tridiagonal, aber immer noch diagonal-dominant, also LR-zerlegbar. Bei symmetrischen Matrizen kann jedoch das CHOLESKY-Verfahren (5.3) durch eine kleine zusätzliche Betrachtung modifiziert werden, auf die wir hier jedoch nicht weiter eingehen.

5.4.2.5 Konvergenzeigenschaften

Entgegen dem Grenzverhalten von Interpolations-Polynomen konvergieren Spline-Funktionen gegen die Funktion, die sie interpolieren, bei Verfeinerung der Unterteilungen Δ . Sei

$$\Delta_m = \{a = x_0^{(m)} < x_1^{(m)} < \dots < x_{n_m}^{(m)} = b\}$$

eine Folge von Unterteilungen des Intervalls $[a, b]$. Mit

$$\|\Delta_m\| := \max_j (x_{j+1}^{(m)} - x_j^{(m)})$$

gilt:

Satz 5.34. Sei $f \in C^4[a, b]$ mit $L = \|f^{(4)}\|_\infty$ und sei Δ_m eine Zerlegungsfolge von $[a, b]$ mit

$$\sup_{m,j} \frac{\|\Delta_m\|}{x_{j+1}^{(m)} - x_j^{(m)}} \leq K < +\infty.$$

Seien s_m die zu f gehörigen Spline-Funktionen mit

$$\begin{aligned} s_m(\xi) &= f(\xi) & \text{für } \xi \in \Delta_m, \\ s'_m(x) &= f'(x) & \text{für } x = a, b. \end{aligned}$$

Dann gibt es von Δ_m unabhängige Konstanten C_i (≤ 2), so dass für $x \in [a, b]$ gilt

$$|f^{(i)}(x) - s_m^{(i)}(x)| \leq C_i L K \|\Delta_m\|^{4-i}, \quad i = 0, 1, 2, 3.$$

Beweis: Sei

$$\Delta := \Delta_m = \{a = x_0 < \dots < x_n = b\}$$

eine feste Zerlegung. Für die Momente $M_j = s''(x_j)$ einer Spline-Funktion $s(x)$ mit $s'(x) = f'(x)$ für $x = x_0, x_n$ gilt nach (5.25) die Gleichung

$$AM = d$$

mit

$$\begin{aligned} \lambda_0 &= \mu_n = \frac{1}{2}, \\ \lambda_j &= \frac{1}{2} \frac{h_{j+1}}{h_j + h_{j+1}}, & \mu_j &= \frac{1}{2} - \lambda_j, \quad j = 1, \dots, n-1 \\ d_0 &= \frac{3}{h_1} \left(\frac{y_1 - y_0}{h_1} - f'(x_0) \right), \\ d_n &= \frac{3}{h_n} \left(f'(x_n) - \frac{y_n - y_{n-1}}{h_n} \right) \\ d_j &= \frac{3}{h_j + h_{j+1}} \left(\frac{y_{j+1} - y_j}{h_{j+1}} - \frac{y_j - y_{j-1}}{h_j} \right) \quad j = 1, \dots, n-1. \end{aligned}$$

Man zeigt leicht (Übung), dass die Matrix A die Eigenschaft hat

$$(5.26) \quad \|x\|_\infty \leq 2\|Ax\|_\infty \quad \text{für } x \in \mathbb{R}^n.$$

Für die Vektoren

$$\begin{aligned} F &:= (f''(x_0), f''(x_1), \dots, f''(x_n))^T \\ r &:= d - AF = A(M - F) \end{aligned}$$

wird zunächst gezeigt

$$(5.27) \quad \|r\|_\infty \leq \frac{3}{8}L \|\Delta\|^2.$$

Beweis von (5.27):

Die Taylor-Entwicklung von f und f'' um x_0 ergibt

$$\begin{aligned} r_0 &= d_0 - f''(x_0) - \frac{1}{2}f''(x_1) \\ &= \frac{3}{h_1} \left(\frac{f(x_1) - f(x_0)}{h_1} - f'(x_0) \right) - f''(x_0) - \frac{1}{2}f''(x_1) \\ &= \frac{h_1^2}{8} f^{(4)}(\tau_1) - \frac{h_1^2}{4} f^{(4)}(\tau_2) \quad \text{mit } \tau_1, \tau_2 \in [x_0, x_1]. \\ \Rightarrow \quad |r_0| &\leq \frac{3}{8}L \|\Delta\|^2. \end{aligned}$$

Analog erhält man für

$$r_n = d_n - \frac{1}{2}f''(x_{n-1}) - f''(x_n)$$

die Abschätzung

$$|r_n| \leq \frac{3}{8}L \|\Delta\|^2.$$

Entsprechend ergibt sich durch die Taylor-Entwicklung um x_j für $j = 1, \dots, n-1$:

$$\begin{aligned} r_j &= d_j - \mu_j f''(x_{j-1}) - f''(x_j) - \lambda_j f''(x_{j+1}) \\ &= \frac{1}{2(h_j + h_{j+1})} \left[\frac{h_{j+1}^3}{4} f^{(4)}(\tau_1) + \frac{h_j^3}{4} f^{(4)}(\tau_2) - \frac{h_j^3}{4} f^{(4)}(\tau_3) - \frac{h_{j+1}^3}{2} f^{(4)}(\tau_4) \right] \\ &\quad \text{mit } \tau_1, \dots, \tau_4 \in [x_{j-1}, x_{j+1}]. \end{aligned}$$

Also

$$|r_j| \leq \frac{3}{8}L \frac{h_j^3 + h_{j+1}^3}{h_j + h_{j+1}} \leq \frac{3}{8}L \|\Delta\|^2.$$

Insgesamt gilt

$$\|r\|_\infty \leq \frac{3}{8}L \|\Delta\|^2$$

und damit wegen $r = A(M - F)$ und (5.26)

$$(5.28) \quad \|M - F\|_\infty \leq 2\|r\|_\infty \leq \frac{3}{4}L \|\Delta\|^2.$$

Wir zeigen nun die Behauptung des Satzes für $\underline{i = 3}$:

Für $x \in [x_{j-1}, x_j]$ ist

$$\begin{aligned} s^{(3)}(x) - f^{(3)}(x) &= \frac{M_j - M_{j-1}}{h_j} - f^{(3)}(x) \\ &= \frac{M_j - f''(x_j)}{h_j} - \frac{M_{j-1} - f''(x_{j-1})}{h_j} \\ &\quad + \frac{f''(x_j) - f''(x) - (f''(x_{j-1}) - f''(x))}{h_j} - f^{(3)}(x) \end{aligned}$$

Taylor-Entwicklung um x ergibt mit der Abschätzung (5.28)

$$(5.29) \quad |s^{(3)}(x) - f^{(3)}(x)| \leq \frac{3}{2}L \frac{\|\Delta\|^2}{h_j} + \frac{L}{2} \frac{\|\Delta\|^2}{h_j} \leq 2LK\|\Delta\|,$$

da $\|\Delta\|/h_j \leq K$ nach Voraussetzung.

Die Behauptung für $\underline{i = 2}$ folgt so:

Für $x \in [a, b]$ gibt es $x_j = x_j(x)$ mit $|x_j(x) - x| \leq \frac{1}{2}\|\Delta\|$.

Mit

$$f''(x) - s''(x) = f''(x_j(x)) - s''(x_j(x)) + \int_{x_j(x)}^x (f^{(3)}(t) - s^{(3)}(t))dt$$

erhält man wegen (5.29) und $K \geq 1$:

$$\begin{aligned} |f''(x) - s''(x)| &\leq \frac{3}{4}L \|\Delta\|^2 + LK \|\Delta\|^2 \\ &\leq \frac{7}{4}LK \|\Delta\|^2 \end{aligned}$$

Nun zeigen wir die Behauptung für $\underline{i=1}$:

Es gilt

$$f(x_j) = s(x_j), \quad j = 0, \dots, n.$$

Außer $\xi_0 := a$, $\xi_{n+1} := b$ gibt es daher nach dem Satz von Rolle n Punkte $\xi_j \in (x_{j-1}, x_j)$ mit

$$f'(\xi_j) = s'(\xi_j), \quad j = 0, \dots, n + 1.$$

Zu jedem $x \in [a, b]$ kann man also $\xi_j(x)$ wählen mit

$$|\xi_j(x) - x| \leq \|\Delta\|.$$

Damit erhält man

$$|f'(x) - s'(x)| = \left| \int_{\xi_j(x)}^x (f''(t) - s''(t)) dt \right| \leq \frac{7}{4} LK \|\Delta\|^3$$

Analog zum Fall $i = 1$ ergibt sich schließlich die Behauptung für $i = 0$:

$$|f(x) - s(x)| = \left| \int_{x_j(x)}^x (f'(t) - s'(t)) dt \right| \leq \frac{7}{4} LK \|\Delta\|^4$$

◇

5.5 Numerische Differentiation

Häufig werden in Anwendungen Ableitungen benötigt, die nicht mehr oder nur sehr aufwendig 'von Hand' zu berechnen sind. Die Interpolation bietet eine einfache Möglichkeit auf numerischem Wege die Ableitungen zu approximieren. Anstelle die Funktion selber an einer Stelle x abzuleiten, approximieren wir die Funktion um x durch geeignete Polynome und leiten dann einfach das so erhaltene Polynom an der Stelle x ab. Hierbei nutzen wir die leichte Ableitbarkeit von Polynomen aus.

Beispiel 5.35. Gesucht sei $f'(x_1)$ an einer Stelle x_1 für eine gegebene Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$. Sei $h > 0$ eine kleine Zahl und $x_0 = x_1 - h$, $x_2 = x_1 + h$. Mittels der Interpolationsformel nach Lagrange erhalten wir ein Interpolationspolynom vom Grad 2 mit $f_i = f(x_i)$, $i = 0, 1, 2$, das nach leichter Umformung übergeht in

$$(5.30) \quad f(x) \approx P(x) = \frac{f_0(x-x_1)(x-x_2)}{2h^2} - 2 \frac{f_1(x-x_0)(x-x_2)}{2h^2} + \frac{f_2(x-x_0)(x-x_1)}{2h^2}$$

Einmaliges Ableiten ergibt

$$f'(x) \approx P'(x) = \frac{(f_2 - 2f_1)(x-x_0) + (f_0 + f_2)(x-x_1) + (f_0 - 2f_1)(x-x_2)}{2h^2}$$

Auswerten von $P'(x)$ an der Stelle x_1 ergibt dann:

$$(5.31) \quad \boxed{f'(x_1) \approx P'(x_1) = \frac{f_2 - f_0}{2h}}$$

Beispiel 5.36. Soll auch noch die zweite Ableitung $f''(x_1)$ berechnet werden, so leiten wir (5.30) erneut ab und erhalten:

$$(5.32) \quad \boxed{f''(x_1) \approx P''(x_1) = \frac{f_2 - 2f_1 + f_0}{h^2}}$$

Natürlich lassen sich auf diesem Wege eine Reihe weiterer Formeln herleiten. Besteht der Bedarf nach höheren Ableitungen ist das Interpolationspolynom von entsprechend höherem Grad zu wählen. Es ist jedoch zu beachten, dass wir die Interpolation durch Polynome von höherem Grade bereits als numerisch instabil festgestellt hatten, so dass damit zu rechnen ist, dass auch die numerische Differentiation dann instabil werden kann.

Kapitel 6

Integration

6.1 Einführung und Aufgabenstellung

Die Integration von Funktionen ist eine elementare mathematische Operation, die in vielen Formeln benötigt wird. Im Gegensatz zur Ableitung, die für praktisch alle mathematischen Funktionen explizit analytisch berechnet werden kann, gibt es viele Funktionen, deren Integrale man nicht explizit angeben kann. Verfahren zur numerischen Integration (man spricht auch von *Quadratur*) spielen daher eine wichtige Rolle, sowohl als eigenständige Algorithmen als auch als Basis für andere Anwendungen.

Das Problem läßt sich hierbei sehr leicht beschreiben. Für eine Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$ soll das Integral

$$(6.1) \quad I := \int_a^b f(x) dx$$

auf einem Intervall $[a, b]$ berechnet werden.

6.2 Newton–Cotes–Formeln

Die *Stützstellen* x_i seien äquidistant und enthalten die Randpunkte des Intervalls, d.h.

$$(6.2) \quad x_i = a + i \cdot h, \quad h = \frac{b - a}{n}, \quad i = 0, \dots, n.$$

Sei P_n das interpolierende Polynom mit

- a) Grad $P_n \leq n$,

$$\text{b) } P_n(x_i) = f_i := f(x_i), \quad i = 0, \dots, n.$$

Als *Approximation für I* nehmen wir den Ausdruck

$$(6.3) \quad I \approx I_n := \int_a^b P_n(x) dx = \sum_{i=0}^n a_i f(x_i).$$

In der Form von Lagrange lautet P_n :

$$P_n(x) = \sum_{i=0}^n f_i L_i(x), \quad L_i(x) = \prod_{\substack{k=0 \\ k \neq i}}^n \frac{x - x_k}{x_i - x_k}$$

Damit folgt

$$I_n = \sum_{i=0}^n f_i \underbrace{\int_a^b L_i(x) dx}_{=: a_i}$$

Die Koeffizienten in (6.3) ergeben sich daher zu

$$(6.4) \quad a_i = \int_a^b L_i(x) dx = \int_a^b \prod_{\substack{k=0 \\ k \neq i}}^n \frac{x - x_k}{x_i - x_k} dx$$

Mit der Substitution $x = a + s \cdot h$, $s \in [0, n]$, $dx = h \cdot ds$ erhalten wir die *Formeln von Newton-Cotes*:

$$(6.5) \quad \boxed{\begin{aligned} I_n &= \sum_{i=0}^n a_i f(x_i), \quad f_i = f(a + i \cdot h) \\ a_i &= h \int_0^n \prod_{\substack{k=0 \\ k \neq i}}^n \frac{s - k}{i - k} ds =: h A_i, \quad A_i \in \mathbb{Q} \end{aligned}}$$

Beispiel 6.1. $n = 1$ (Trapezregel): Es ist $h = b - a$, $A_0 = A_1 = \frac{1}{2}$ und damit

$$I_1 = \frac{h}{2} (f(a) + f(b))$$

Beispiel 6.2. $n = 2$ (Simpsonregel): Es ist $h = \frac{b-a}{2}$, $A_0 = \frac{1}{3}$, $A_1 = \frac{4}{3}$, $A_2 = \frac{1}{3}$ und damit

$$I_2 = \frac{h}{3} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

Bemerkung 6.3. Allgemein gilt:

$$\sum_{i=0}^n A_i = n$$

(d.h. $\sum_{i=0}^n a_i = 1$) und $A_i = A_{n-i}$ (Symmetrie).

Die folgende Tabelle enthält die Koeffizienten A_i für $n \leq 4$:

n	A_0	A_1	A_2	A_3	A_4	Bezeichnung
1	$\frac{1}{2}$	$\frac{1}{2}$				Trapezregel
2	$\frac{1}{3}$	$\frac{4}{3}$	$\frac{1}{3}$			Simpson–Regel
3	$\frac{3}{8}$	$\frac{9}{8}$	$\frac{9}{8}$	$\frac{3}{8}$		Newton’sche $\frac{3}{8}$ –Regel
4	$\frac{14}{45}$	$\frac{64}{45}$	$\frac{24}{45}$	$\frac{64}{45}$	$\frac{14}{45}$	Milne–Regel

Bemerkung 6.4. Für $n \geq 8$ können negative Gewichte A_i auftreten und die Formeln werden dann aufgrund von Rundungsfehlereinfluß unbrauchbar.

Es gilt der Fehler

$$(6.6) \quad R_n(f) = I - I_n = \int_a^b f(x) dx - \int_a^b P_n(x) dx.$$

Offensichtlich ist $R_n(f) = 0$, falls f Polynom vom Grade $\leq n$ ist. Der folgende Satz gibt eine Abschätzung für den Fehler

Satz 6.5 (Fehlerabschätzung).

1. Für $f \in C^{n+1}[a, b]$ gilt

$$|R_n(f)| \leq h^{n+2} c_n \max_{[a,b]} |f^{(n+1)}(x)|$$

mit

$$c_n = \frac{1}{(n+1)!} \int_0^n \prod_{i=0}^n |s-i| ds.$$

2. Für n gerade und $f \in C^{n+2}[a, b]$ gilt

$$|R_n(f)| \leq h^{n+3} c_n^* \max_{[a,b]} |f^{(n+2)}(x)|, \quad c_n^* = \frac{n}{2} c_n$$

Beweis: Wir beweisen nur den ersten Teil der Aussage: Nach der Restgliedformel der Polynominterpolation (5.7) existiert ein $\xi \in [a, b]$ mit

$$f(x) - P_n(x) = \frac{L(x)}{(n+1)!} f^{(n+1)}(\xi), \quad L(x) = \prod_{i=0}^n (x - x_i).$$

Hieraus folgt

$$(6.7) \quad |R_n(f)| \leq \frac{1}{(n+1)!} \int_a^b |L(x)| dx \cdot \max_{[a,b]} |f^{(n+1)}(x)|.$$

Mit der Substitution $x = a + s \cdot h$, $s \in [0, n]$ ergibt sich

$$\int_a^b |L(x)| dx = \int_a^b \prod_{i=0}^n |x - x_i| dx = h^{n+2} \int_0^n \prod_{i=0}^n |s - i| ds.$$

Zusammen mit (6.7) erhält man die Behauptung. \diamond

Da das Maximum hoher Ableitungen von f sehr schwer zu bestimmen ist, sind die Formeln zur praktischen Fehlerabschätzung i.a. unbrauchbar. Ihr Nutzen liegt in der Information, mit welcher Potenz von h der Fehler abfällt.

Beispiel 6.6. $n = 1$ (Trapezregel):

$$|R_1(f)| \leq \frac{h^3}{12} \max_{[a,b]} |f^{(2)}(x)|.$$

Beispiel 6.7. $n = 2$ (Simpson-Regel):

$$|R_2(f)| \leq \frac{h^5}{90} \max_{[a,b]} |f^{(4)}(x)|.$$

Bemerkung 6.8. Bei geradem n gewinnt man durch den Übergang zu $n+1$ keine Potenz von h .

Beispiel 6.9. Es sei $I = \int_0^1 e^x dx = e - 1 = 1.7182\dots$

$$I_1 = \frac{1}{2}(1 + e) = 1.8591\dots$$

$$I_2 = \frac{1}{6}(1 + 4e^{1/2} + e) = 1.7189\dots$$

$$I_3 = \frac{1}{8}(1 + 3e^{1/3} + 3e^{2/3} + e) = 1.7185\dots$$

6.3 Zusammengesetzte Newton–Cotes–Formeln

Wie beschrieben kann man höhere Genauigkeiten nicht durch immer höhere Polynomgrade erwirken. Es tauchen also die gleichen prinzipiellen Probleme auf, wie bei der Interpolation durch Polynome. Dies ist nicht weiter verwunderlich, da Interpolationspolynome ja den hier beschriebenen Verfahren zu Grunde liegen. Formeln höherer Genauigkeit kann man konstruieren, indem man (ähnlich wie bei der Splineinterpolation) die oben angegebenen Regeln auf Teilintervalle anwendet. Bei der Integration fällt jedoch der mühsame Weg der zusätzlichen

Glattheitsbedingungen an den Nahtstellen weg, da wir ja nicht an einer schönen Approximation der Funktion, sondern 'nur' an einer guten Approximation des Integrals interessiert sind.

Wir beschränken uns in diesem Abschnitt auf die Herleitung der zusammengesetzten Trapezregel, eine Übertragung auf die anderen Newton–Cotes Formeln verläuft analog.

6.3.1 Zusammengesetzte Trapezregel

Wir betrachten wieder äquidistante Stützstellen (6.2). Die Anwendung der Trapezregel auf das Teilintervall $[x_i, x_{i+1}]$ ergibt die Approximation

$$\int_{x_i}^{x_{i+1}} f(x) dx \approx \frac{h}{2} (f(x_i) + f(x_{i+1})), \quad i = 0, 1, \dots, n-1.$$

Durch Summation erhalten wir die *zusammengesetzte Trapezregel*:

(6.8)

$$\begin{aligned} \int_a^b f(x) dx \approx T(h) &:= \sum_{i=0}^{n-1} \frac{h}{2} (f(x_i) + f(x_{i+1})) \\ &= h \left[\frac{f(a)}{2} + f(a+h) + \dots + f(b-h) + \frac{f(b)}{2} \right]. \end{aligned}$$

Die zusammengesetzte Trapezregel ist in Abbildung 6.1 dargestellt.

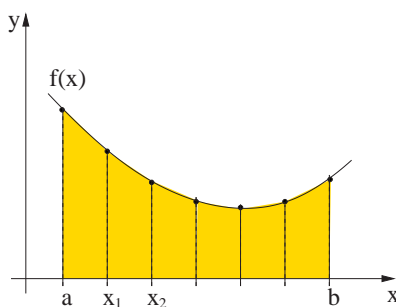


Abbildung 6.1: Zusammengesetzte Trapezregel.

Der Gesamtfehler ergibt sich aus der Summation der Fehler (bzw. aus Beispiel 6.6) zu

$$\left| T(h) - \int_a^b f(x) dx \right| \leq n \cdot \frac{h^3}{12} \max_{[a,b]} |f^{(2)}(x)| = (b-a) \cdot \frac{h^2}{12} \max_{[a,b]} |f^{(2)}(x)|,$$

also ein Fehler in der Größenordnung $O(h^2)$.

6.3.2 Verfeinerung der zusammengesetzten Trapezregel

Wir wollen in diesem Abschnitt einige weitere Besonderheiten aufweisen, wie aus unscheinbaren Zusammenhängen numerisches Kapital geschlagen werden kann.

Eine ebenso anschauliche Approximation des Integrals (6.1), welche der Riemannschen Summe entspricht, bildet die *Mittelpunktsumme*

$$(6.9) \quad \int_a^b f(x) dx \approx M(h) := h \sum_{i=0}^{n-1} f\left(x_{i+\frac{1}{2}}\right), \quad x_{i+\frac{1}{2}} := a + \left(i + \frac{1}{2}\right) h.$$

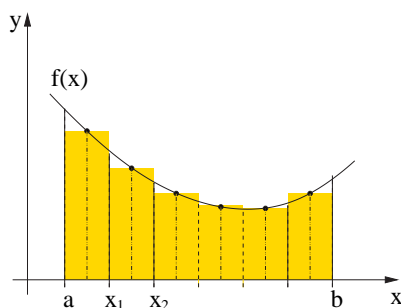


Abbildung 6.2: Zusammengesetzte Mittelpunktregel.

$M(h)$ stellt die Fläche unterhalb der Treppenkurve in Abbildung 6.2 dar und ist von gleicher Fehlerordnung $O(h^2)$ wie die zusammengesetzte Trapezregel. Aus (6.8) und (6.9) folgt unmittelbar die Relation

$$(6.10) \quad T\left(\frac{h}{2}\right) = \frac{1}{2} [T(h) + M(h)]$$

Bemerkung 6.10. Die Relation (6.10) erlaubt eine Verbesserung der zusammengesetzten Trapezregel durch sukzessive Halbierung der Schrittweite in der Weise, daß zu bereits berechneten Näherungen $T(h)$ auch noch $M(h)$ berechnet wird. Bei jeder Halbierung der Schrittweite h wird der Rechenaufwand (gemessen mit der Anzahl der Funktionsauswertungen) etwa verdoppelt, doch werden die schon berechneten Funktionswerte auf ökonomische Weise wieder verwendet. Die sukzessive Halbierung kann etwa dann abgebrochen werden, wenn für eine gegebene Fehlerschranke $\varepsilon > 0$ $|T(h) - M(h)| < \varepsilon$ gilt. Dann ist der Fehler $|T(\frac{h}{2}) - I|$ im allgemeinen höchstens gleich ε .

6.4 Die Gaußsche Integrationsmethode

6.4.1 Orthogonalpolynome

Der Raum $V = C[a, b]$ mit dem inneren Produkt

$$(6.11) \quad \langle f, g \rangle = \int_a^b f(x)g(x)w(x)dx, \\ w \in C(a, b), \quad w(x) > 0 \quad \text{für } a < x < b,$$

ist ein sogenannter Prä-Hilbertraum. Wendet man das SCHMIDT'sche Orthogonalisierungsverfahren auf die Monome $u_n(x) = x^n$ an, so erhält man Orthogonalpolynome mit Höchstkoeffizient $a_n = 1$:

$$\tilde{p}_n \in \tilde{\Pi}_n := \left\{ x^n + \sum_{k=0}^{n-1} a_k x^k \right\} \quad (n = 0, 1, 2, \dots), \\ \langle \tilde{p}_i, \tilde{p}_k \rangle = 0 \quad \text{für } i \neq k, \quad \langle \tilde{p}_i, \tilde{p}_i \rangle = 1 \quad \text{bei Orthonormalpolynomen}$$

Die Polynome $\tilde{p}_0 \dots, \tilde{p}_n$ bilden eine Basis von Π_n .

Darüberhinaus gilt der bemerkenswerte

Satz 6.11 (Nullstellensatz). *Das Orthogonalpolynom $\tilde{p}_n \in \tilde{\Pi}_n$ hat in (a, b) genau n einfache Nullstellen.*

Beweis: Seien $x_1 \dots, x_m$ ($m \geq 0$) die Vorzeichenwechsel von \tilde{p}_n in (a, b) . Wir zeigen, dass $m = n$ gilt. Mit dem Polynom

$$q(x) = \prod_{i=1}^m (x - x_i) \in \Pi_m$$

hat das Polynom $q \cdot \tilde{p}_n$ konstantes Vorzeichen. Für $m < n$ würde folgen

$$\langle q, \tilde{p}_n \rangle = \int_a^b q(x)\tilde{p}_n(x)w(x)dx = 0,$$

also $q \cdot \tilde{p}_n \equiv 0$ in (a, b) : Widerspruch. ◇

Beispiel 6.12. (Legendre–Polynome): *Es sei $[a, b] = [-1, 1]$, $w = 1$. Man prüft leicht nach, dass die Legendre–Polynome*

$$(6.12) \quad \tilde{L}_n(x) = \frac{n!}{(2n)!} \frac{d^n (x^2 - 1)^n}{dx^n} = x^n + \dots, \quad n \geq 0$$

ein Orthogonalsystem bilden mit

$$\langle \tilde{L}_n, \tilde{L}_n \rangle = \frac{2n+1}{2} \frac{1}{(2^n n!)^2}.$$

Zum Beispiel ist

$$\tilde{L}_1(x) = x, \quad \tilde{L}_2(x) = x^2 - \frac{1}{3}, \quad \tilde{L}_3(x) = x^3 - \frac{3}{5}x.$$

Beispiel 6.13. (Tschebychev–Polynome): Es sei $[a, b] = [-1, 1]$, $w(x) = 1/\sqrt{1-x^2}$:

Die Tschebychev-Polynome $T_n \in \Pi_n$ werden rekursiv definiert durch

$$(6.13) \quad \begin{aligned} T_{n+1}(x) &= 2xT_n(x) - T_{n-1}(x) \quad (n = 1, 2, \dots), \\ T_0(x) &= 1, \quad T_1(x) = x. \end{aligned}$$

Mit der Substitution

$$x = \cos \theta, \quad \theta = \arccos x$$

gelangt man zur Darstellung

$$T_n(x) = \cos(n\theta), \quad \theta = \arccos x,$$

denn die Rekursionsformel

$$\cos((n+1)\theta) = 2 \cos \theta \cos(n\theta) - \cos((n-1)\theta)$$

entspricht gerade der Rekursion (6.13). Damit bestätigt man die Orthogonalität bzgl. $w(x) = 1/\sqrt{1-x^2}$:

$$(6.14) \quad \int_{-1}^1 T_i(x)T_k(x) \frac{dx}{\sqrt{1-x^2}} = \int_0^\pi \cos(i\theta) \cos(k\theta) \frac{\sin \theta}{\sin \theta} d\theta = \begin{cases} 0 & \text{für } i \neq k \\ \pi & \text{für } i = k = 0 \\ \frac{\pi}{2} & \text{für } i = k \neq 0. \end{cases}$$

Die Darstellung $T_n(x) = \cos(n\theta)$ zeigt, dass $T_n(x)$ die Nullstellen (Tschebychev–Abszissen)

$$x_k = \cos\left(\frac{2k-1}{n} \frac{\pi}{2}\right) \in (-1, 1), \quad k = 1, \dots, n,$$

und die Extremalstellen

$$x_k^{(l)} = \cos\left(\frac{k\pi}{n}\right) \quad (k = 0, 1, \dots, n), n \geq 1,$$

mit

$$T_n(x_k^{(l)}) = (-1)^k,$$

besitzt. Z.B. liefert die Rekursion (6.13)

$$\begin{aligned} T_2(x) &= 2x^2 - 1, & T_3(x) &= 4x^3 - 3x, \dots, \\ T_n(x) &= 2^{n-1} x^n - \dots \end{aligned}$$

Zu normierten Polynomen $\tilde{T}_n \in \tilde{\Pi}_n$ gelangt man durch die Normierung

$$\tilde{T}_n(x) = \frac{1}{2^{n-1}} T_n(x).$$

Weitere Orthogonalpolynome bzgl. anderer Gewichte $w(x)$ finden sich in M. AB-
RAMOVITZ, F. STEGUN: Handbook of Mathematical Functions.

6.4.2 Gaußintegration

Sei $f \in C[a, b]$ und sei $w \in C[a, b]$ eine positive Gewichtsfunktion mit $w(x) > 0$ für $x \in (a, b)$. Wir suchen eine Integrationsformel für das Integral

$$(6.15) \quad I(f) = \int_a^b w(x) f(x) dx.$$

Im folgenden benutzen wir die Bezeichnungen:

$$\begin{aligned} \Pi_j &: \quad \text{Polynome vom Grade } \leq j \quad (j = 0, 1, 2, \dots) \\ \tilde{\Pi}_j &: \quad = \{p \in \Pi_j \mid p(x) = x^j + a_{j-1}x^{j-1} + \dots + a_0\}. \end{aligned}$$

Eine Integrationsformel für $I(f)$ der Form

$$(6.16) \quad G_n(f) = \sum_{i=1}^n A_i f(x_i)$$

hat die $2n$ freien Parameter A_i und x_i . Die Formeln von Newton–Cotes (6.5) mit äquidistanten Stützstellen sind exakt in Π_{n-1} . Wir wollen nun fordern, dass

$$G_n(f) = I(f) \quad \text{für alle } f \in \Pi_{2n-1},$$

d.h. $G_n(f)$ ist exakt in Π_{2n-1} . Dies ergibt gerade $2n$ Bedingungen für die $2n$ Parameter. Der folgende Satz zeigt, dass diese Forderung maximal ist.

Satz 6.14. *Es gibt keine Formel $G_n(f)$ des Typs (6.16), die in Π_{2n} exakt ist.*

Beweis: Annahme: $G_n(f) = \int_a^b w(x)f(x)dx$ für $f \in \Pi_{2n}$.

Mit

$$f := \prod_{i=1}^n (x - x_i)^2 \in \Pi_{2n}$$

erhält man einen Widerspruch wegen

$$G_n(f) = 0 \neq \int_a^b w(x)f(x)dx > 0.$$

◇

Zur Konstruktion einer in Π_{2n-1} exakten Formel $G_n(f)$ benutzt man die zur Gewichtsfunktion w gehörenden Orthogonalpolynome \tilde{p}_n bzgl. des Skalarproduktes (vgl. (6.11))

$$\langle f, g \rangle = \int_a^b f(x)g(x)w(x)dx, \quad f, g \in C[a, b].$$

Das Polynom $\tilde{p}_n \in \tilde{\Pi}_n$ hat nach Satz 6.11 n Nullstellen $x_1, \dots, x_n \in (a, b)$. Damit gelangen wir zum Hauptresultat dieses Abschnittes.

Satz 6.15. *Seien x_1, \dots, x_n die Nullstellen des Orthogonalpolynoms \tilde{p}_n und sei*

$$L_i(x) = \prod_{\substack{k=1 \\ k \neq i}}^n \frac{x - x_k}{x_i - x_k}, \quad i = 1, \dots, n.$$

Dann ist die Integrationsformel

$$G_n(f) = \sum_{i=1}^n A_i f(x_i), \quad A_i := \int_a^b w(x)L_i(x)^2 dx$$

exakt in Π_{2n-1} und es gilt

$$A_i > 0.$$

Beweis: Nach der Interpolationsformel von Lagrange gilt

$$f(x) = \sum_{i=1}^n L_i(x)f(x_i) \quad \text{für alle } f \in \Pi_{n-1}.$$

Daher ist

$$G_n(f) = \sum_{i=1}^n \int_a^b w(x)L_i(x)dx f(x_i)$$

exakt in Π_{n-1} . Ein Polynom $f \in \Pi_{2n-1}$ faktorisieren wir in

$$f = q \cdot \tilde{p}_n + r \quad \text{mit } q, r \in \Pi_{n-1}$$

$$\begin{aligned} \Rightarrow I(f) &= \int_a^b w(x)f(x)dx = \underbrace{\int_a^b w(x)q(x)\tilde{p}_n(x)dx}_{=0, \text{ (da } q \in \Pi_{n-1})} + \int_a^b w(x)r(x)dx \\ &= G_n(r) \quad (\text{da } r \in \Pi_{n-1}) \\ &= G_n(r) + \underbrace{G_n(q \cdot \tilde{p}_n)}_{=0, \text{ (da } \tilde{p}_n(x_i)=0, i=1, \dots, n)} \\ &= G_n(r + q \cdot \tilde{p}_n) \\ &= G_n(f). \end{aligned}$$

Also ist $G_n(f)$ exakt in Π_{2n-1} und die Formeln für die Gewichte A_i folgen so:
Es ist $L_i^2 \in \Pi_{2n-2}$, also

$$\int_a^b w(x)L_i(x)^2 dx = G_n(L_i^2) = \sum_{k=1}^n A_k L_i(x_k)^2 = \sum_{k=1}^n A_k \delta_{ik} = A_i.$$

◇

Beispiel 6.16. Sei $[a, b] = [-1, 1]$, $w(x) \equiv 1$,

$$\tilde{p}_n(x) = \frac{n!}{(2n)!} \frac{d^n}{dx^n} (x^2 - 1)^n, \quad n = 0, 1, 2, \dots,$$

Legendre-Polynome bis auf Normierungsfaktoren:

$$\tilde{p}_1(x) = x, \quad \tilde{p}_2(x) = x^2 - \frac{1}{3}, \quad \tilde{p}_3(x) = x^3 - \frac{3}{5}x$$

n	x_1	x_2	x_3	A_1	A_2	A_3
1	0			2		
2	$-\sqrt{\frac{1}{3}}$	$+\sqrt{\frac{1}{3}}$		1	1	
3	$-\sqrt{\frac{3}{5}}$	0	$\sqrt{\frac{3}{5}}$	$\frac{5}{9}$	$\frac{8}{9}$	$\frac{5}{9}$

$$I = \int_{-1}^1 e^x dx = 2.350402.$$

Die Simpson-Regel liefert

$$I_2 = 2.362054.$$

Dagegen ist mit gleich vielen Funktionsauswertungen

$$G_3 = 2.350337.$$

Für den Fehler der Gauß'schen Integrationsmethode gilt die folgende Abschätzung: (vgl. Stoer I, S. 126).

Satz 6.17. Sei $f \in C^{2n}[a, b]$ und sei $G_n(f)$ die in Satz 6.15 definierte Integrationsformel, dann gilt

$$I(f) - G_n(f) = \frac{f^{(2n)}(\xi)}{(2n)!} \langle p_n, p_n \rangle$$

mit einem $\xi \in (a, b)$.

Die Gauß-Formeln liefern im Vergleich zu den Newton-Cotes-Formeln bzw. den Extrapolationsverfahren die genauesten Resultate (gemessen an der Zahl der Funktionsauswertungen). Im Gegensatz zu Extrapolationsverfahren können jedoch beim Übergang von einem Index n zu $n + 1$ die bis dahin berechneten Funktionswerte $f(x_i)$ nicht weiter verwendet werden. Daher sind in der Praxis Extrapolationsverfahren vorzuziehen.

6.5 Integration und Extrapolation

6.5.1 Euler-Maclaurin'sche Summenformel

Durch Anwendung der Extrapolation auf die zusammengesetzte Trapezregel wollten wir nun Formeln konstruieren, deren Fehler mit einer hohen Potenz von h abfällt. Grundlage dieser Extrapolationsverfahren ist die folgende asymptotische Entwicklung von $T(h)$ nach Potenzen von h^2 .

Satz 6.18 (Euler-Maclaurin'sche Summenformel). Für $f \in C^{2m+2}[a, b]$ gilt die Entwicklung

$$T(h) = \tau_0 + \tau_1 h^2 + \tau_2 h^4 + \dots + \tau_m h^{2m} + \alpha_{m+1}(h) h^{2m+2}$$

mit

$$1. \tau_0 = \int_a^b f(x) dx,$$

$$2. \tau_k = \frac{(-1)^{k+1} B_k}{(2k)!} (f^{(2k-1)}(b) - f^{(2k-1)}(a)), \quad k = 1, \dots, m$$

$$3. \alpha_{m+1}(h) = \frac{1}{(2m+2)!} \int_a^b f^{(2m+2)}(x) K_{2m+2} \left(\frac{x-a}{h} \right) dx$$

mit $K_{2m+2} \in C[0, n]$ und

$$\int_a^b K_{2m+2} \left(\frac{x-a}{h} \right) dx = (-1)^m B_{m+1} (b-a).$$

Hierbei sind B_k die Bernoulli-Zahlen

$$B_1 = \frac{1}{6}, \quad B_2 = \frac{1}{30}, \quad B_3 = \frac{1}{42}, \dots$$

Beweis: Zum Beweis vergleiche man etwa Stoer 1. ◇

6.5.2 Anwendung der Extrapolation auf die Integration

Nach dem vorigen Satz gilt

$$T(h) = \underbrace{\tau_0 + \tau_1 h^2 + \dots + \tau_m h^{2m}}_{\text{Polynom vom Grade } m \text{ in } h^2} + \alpha_{m+1}(h) h^{2m+2}.$$

Es interessiert die Größe

$$\tau_0 = \int_a^b f(x) dx = \lim_{h \rightarrow 0} T(h).$$

Idee der Extrapolation: Zu $m+1$ Schrittweiten

$$h_0 = b-a, \quad h_1 = \frac{h_0}{n_1}, \dots, \quad h_m = \frac{h_0}{n_m}; \quad n_i < n_{i+1} \quad (n_i \in \mathbb{N})$$

bestimme man Trapezsummen

$$T_{i0} := T(h_i), \quad i = 0, \dots, m,$$

und dann durch Interpolation dasjenige Polynom in h^2

$$(6.17) \quad \tilde{T}_{mm}(h) := a_0 + a_1 h^2 + \dots + a_m h^{2m}$$

mit

$$\tilde{T}_{mm}(h_i) = T(h_i), \quad i = 0, \dots, m.$$

Dann ist

$$\tilde{T}_{mm}(0) = a_0 \approx \tau_0 \quad (\text{Extrapolation}).$$

Beispiel 6.19. $h_0 = b - a$, $h_1 = \frac{b-a}{2}$. Dann ist

$$T_{11} := \tilde{T}_{11}(0) = L_0(0)T(h_0) + L_1(0)T(h_1)$$

mit

$$L_i(x) = \prod_{\substack{k=0 \\ k \neq i}}^1 \frac{x - x_k}{x_i - x_k}, \quad x_i = h_i^2, \quad i = 0, 1.$$

Für $x = 0$ erhält man

$$L_0(0) = \frac{-h_1^2}{h_0^2 - h_1^2} = -\frac{1}{3}, \quad L_1(0) = \frac{-h_0^2}{h_1^2 - h_0^2} = \frac{4}{3}$$

und daher

$$\begin{aligned} T_{11} &= -\frac{1}{3} \underbrace{\frac{b-a}{2} (f(a) + f(b))}_{=T(h_0)} + \frac{4}{3} \underbrace{\frac{b-a}{2} \left(\frac{f(a)}{2} + f\left(\frac{a+b}{2}\right) + \frac{f(b)}{2} \right)}_{=T(h_1)} \\ &= \frac{h_1}{3} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) \end{aligned}$$

und dies ist die Simpson-Regel!

Die Berechnung des Wertes $\tilde{T}_{mm}(0) = a_0$ in (6.17) erfolgt mit dem Algorithmus von Neville; vgl. (5.5). Dazu sei $1 \leq k \leq i \leq m$ und $\tilde{T}_{ik}(h)$ dasjenige Polynom in h^2 mit

$$\tilde{T}_{ik}(h_j) = T_{j0} := T(h_j) \quad \text{für } j = i - k, i - k + 1, \dots, i.$$

Die Rekursion für $T_{ik} := \tilde{T}_{ik}(0)$ ergibt sich aus dem Algorithmus von Neville mit $x_i = h_i^2$, $x = 0$, zu

$$T_{ik} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\left(\frac{h_{i-k}}{h_i}\right)^2 - 1}, \quad 1 \leq k \leq i \leq m.$$

Die Berechnung des Tableaus mit den Größen T_{ik} erfolgt spaltenweise, z. B.

$$\begin{array}{c|cccc} h_0 & T_{00} & & & \\ h_1 & T_{10} & & T_{11} & \\ h_2 & T_{20} & \searrow & T_{21} & T_{22} \\ h_3 & T_{30} & & T_{31} & T_{32} \searrow \\ & & & & T_{33} \end{array}$$

Beispiel 6.20. $I = \int_0^1 e^x dx = 1.718281828$, $m = 3$

h_i	T_{i0}	T_{i1}	T_{i2}	T_{i3}
1	1.859140914			
$\frac{1}{2}$	1.753931092	1.718861151		
$\frac{1}{4}$	1.727221904	1.718318841	1.718282687	
$\frac{1}{8}$	1.720518792	1.718284155	1.718281842	1.718281828

In der Praxis haben sich zwei Schrittweitenfolgen bewährt:

$$\begin{aligned} \text{Romberg-Folge: } & h_0 = b - a, \quad h_i = \frac{h_0}{2^i}, \quad i = 0, 1, \dots \\ \text{Bulirsch-Folge: } & h_0 = b - a, \quad h_1 = \frac{h_0}{2}, \quad h_2 = \frac{h_0}{3}, \quad h_3 = \frac{h_0}{4}, \\ & h_4 = \frac{h_0}{6}, \quad h_5 = \frac{h_0}{8}, \dots \end{aligned}$$

Bemerkung 6.21. Die Bulirsch-Folge hat den Vorteil, dass bei ihr die Rechenarbeit für die Berechnung neuer $T(h_i)$ nicht so rasch ansteigt wie bei der Romberg-Folge. Bei der praktischen Durchführung beachte man, dass bei der Berechnung von $T(h_{i+1})$ auf die schon bei $T(h_i)$ berechneten Funktionswerte zurückgegriffen wird, vgl. etwa die Verfeinerung der zusammengesetzten Trapezregel mittels der Mittelpunktschrittweite.

6.5.3 Integrationsfehler

Im folgenden soll ein Ausdruck für den Fehler

$$T_{mm} - \int_a^b f(x) dx$$

angegeben werden.

Hilfssatz 6.22. Seien x_i , $i = 0, \dots, m$, paarweise verschiedene Zahlen und sei

$$L_i(x) = \prod_{\substack{k=0 \\ k \neq i}}^m \frac{x - x_k}{x_i - x_k}, \quad i = 0, \dots, m.$$

Dann gilt

$$\sum_{i=0}^m x_i^j L_i(x) = \begin{cases} 1 & , \quad j = 0 \\ 0 & , \quad j = 1, \dots, m \\ (-1)^m x_0 \dots x_m & , \quad j = m + 1 \end{cases}$$

Beweis: Man setze $x = 0$ in den beiden folgenden Identitäten ein:

$$\begin{aligned} x^j &\equiv \sum_{i=0}^m x_i^j L_i(x), & j = 0, \dots, m, \\ x^{m+1} &\equiv \sum_{i=0}^m x_i^{m+1} L_i(x) + (x - x_0)(x - x_1) \dots (x - x_m). \end{aligned}$$

Die zweite Identität folgt daraus, dass die rechte Seite gleich der linken Seite in den Punkten $x_i, i = 0, \dots, m$, ist und die Koeffizienten von x^{m+1} auf beiden Seiten übereinstimmen; vgl. auch die Restgliedformel (5.7). \diamond

Die Substitution $x = h^2, x_i = h_i^2$, in Hilfssatz 6.22 ergibt die Beziehung

$$(6.18) \quad \sum_{i=0}^m h_i^{2j} L_i(0) = \begin{cases} 1 & , j = 0 \\ 0 & , j = 1, \dots, m \\ (-1)^m h_0^2 h_1^2 \dots h_m^2 & , j = m + 1 \end{cases}.$$

Das Polynom $\tilde{T}_{mm}(h)$ in (6.17) interpolierte die Werte $T(h_i), i = 0, \dots, m$. Also gilt nach der Formel von Lagrange

$$(6.19) \quad T_{mm} = \tilde{T}_{mm}(0) = \sum_{i=0}^m L_i(0) T(h_i).$$

Die asymptotische Entwicklung in Satz 6.18) von $T(h)$ ergab

$$T(h) = \int_a^b f(x) dx + \tau_1 h^2 + \dots + \tau_m h^{2m} + \alpha_{m+1}(h) h^{2m+2},$$

$$\alpha_{m+1}(h) = \frac{1}{(2m+2)!} \int_a^b f^{(2m+2)}(x) K_{2m+2}\left(\frac{x-a}{h}\right) dx,$$

$$(6.20) \quad \int_a^b K_{2m+2}\left(\frac{x-a}{h}\right) dx = (-1)^m B_{m+1}(b-a).$$

Mit (6.18), (6.19) folgt dann

$$(6.21) \quad T_{mm} = \int_a^b f(x) dx + \frac{1}{(2m+2)!} \int_a^b f^{(2m+2)}(x) K(x) dx,$$

$$(6.22) \quad K(x) := \sum_{i=0}^m L_i(0) h_i^{2m+2} K_{2m+2} \left(\frac{x-a}{h_i} \right).$$

Man kann zeigen: Die Funktion $K(x)$ hat gleiches Vorzeichen in $[a, b]$ für die Romberg-Folge und die Bulirsch-Folge h_i . Daher gilt

$$\int_a^b f^{(2m+2)}(x) K(x) dx = f^{(2m+2)}(\xi) \int_a^b K(x) dx, \quad \xi \in [a, b],$$

und mit (6.18), (6.20) und (6.22) haben wir

$$\begin{aligned} \int_a^b K(x) dx &= \sum_{i=0}^m L_i(0) h_i^{2m+2} \int_a^b K_{2m+2} \left(\frac{x-a}{h_i} \right) dx \\ &= (-1)^m h_0^2 h_1^2 \dots h_m^2 (-1)^m B_{m+1} (b-a) \\ &= (b-a) h_0^2 \dots h_m^2 B_{m+1} \end{aligned}$$

Insgesamt ergibt sich dann aus (6.21) und den vorigen Beziehungen der Fehler

$$(6.23) \quad \boxed{T_{mm} - \int_a^b f(x) dx = (b-a) h_0^2 \dots h_m^2 \frac{B_{m+1}}{(2m+2)!} f^{(2m+2)}(\xi)}.$$

Bei Interpolation mit den Schrittweiten h_{i-k}, \dots, h_i erhält man auf ähnliche Weise

$$(6.24) \quad T_{ik} - \int_a^b f(x) dx = (b-a) h_{i-k}^2 \dots h_i^2 \frac{B_{k+1}}{(2k+2)!} f^{(2k+2)}(\xi).$$

Für $k=0$ gewinnt man hieraus die Abschätzung für den Gesamtfehler der zusammengesetzten Trapezregel zurück wegen $B_1 = \frac{1}{6}$.

Wegen (6.24) verhält sich der Fehler von T_{im} in der $(i+1)$ -ten Spalte des Tableaus wie h_{i-m}^{2m+2} , also wie der Fehler eines Verfahrens $(2m+2)$ -ter Ordnung. Aus Gründen der Auslöschung geht man in der Praxis nicht über $m=6$ hinaus. Man beendet die Rechnung, falls das erste Mal

$$|T_{i,6} - T_{i+1,6}| \leq \epsilon \cdot s$$

erfüllt ist, wobei

ϵ : gewünschte relative Genauigkeit,

s : grober Näherungswert von $\int_a^b |f(x)| dx$.

Literaturverzeichnis

- [1] Büskens, C. *Numerische Mathematik für Naturwissenschaftler und Ingenieure*, Skript, Universität Bayreuth, Bayreuth, 2004.
- [2] Cryer, C.W. *Praktische Mathematik I*, Skript, Universität Münster, Münster, 1989.
- [3] Deuffhard, P., Hohmann, A. *Numerische Mathematik I*, de Gruyter, Berlin, New York, 1993.
- [4] Hämmerlin, G., Hoffmann, K.H. *Numerische Mathematik*, Springer-Verlag, Berlin Heidelberg, 1991.
- [5] Maurer, H. *Numerische Mathematik*, Skript, Universität Münster, Münster, 1990.
- [6] Schaback, R., Werner, H. *Numerische Mathematik*, Springer-Verlag, Berlin, Heidelberg, New York, 1992.
- [7] Schwarz, H.R. *Numerische Mathematik*, B.G. Teubner, Stuttgart, 1993.
- [8] Stoer, J. *Einführung in die Numerische Mathematik I*, Springer-Verlag, Berlin Heidelberg, 1983.