

Numerik gewöhnlicher Differentialgleichungen

Oliver Ernst

Professur Numerische Mathematik

Wintersemester 2016/17



Mathematik!
TU Chemnitz

- Vorlesungswebseite:
`www.tu-chemnitz.de/mathematik/numa/lehre/ode-2014`
- Vorlesung: Prof. **Oliver Ernst**,
`oliver.ernst@mathematik.tu-chemnitz.de`
Mo 9:15 (ungerade Wochen) & Do 9:15 (wöchentlich)
- Übung: Dipl.-Math. oec. **Jan Blechschmidt**,
`jan.blechschmidt@mathematik.tu-chemnitz.de`
Mo 9:15 (gerade Wochen)
- Prüfung: 30 Minuten mündlich, Termin nach Vereinbarung.
- Modul M13: 6 LP, 180 AS.

- Eine **Differentialgleichung** beschreibt eine Beziehung zwischen einer gesuchten Funktion (abhängige Variable) und ihren Ableitungen bezüglich einer oder mehreren (unabhängigen) Variablen.
- Werden nur Ableitungen bezüglich einer Variablen betrachtet, so spricht man von einer **gewöhnlichen Differentialgleichung (GDG)** (ordinary differential equation, ODE), andernfalls von einer **partiellen Differentialgleichung (PDG)** (partial differential equation, PDE).
- Da explizite Lösungen nur in wenigen Ausnahmefällen zur Verfügung stehen, ist der Einsatz numerischer Methoden zur Lösung von GDGen unvermeidbar.
- Wir befassen uns hier ausschließlich mit numerischen Verfahren zur Lösung von GDGen: Zum Einen wegen der überragenden Rolle, die sie in vielen Anwendungen spielen, was (hoffentlich) im Laufe der Vorlesung klar wird; zum Zweiten aber auch wegen ihrer Bedeutung für die Entwicklung und Analyse numerischer Methoden zur Lösung von PDGen, welche sehr viel schwieriger zu lösen sind und im Mittelpunkt einer eigenen Vorlesung stehen.

- Begriff, Beispiele, theoretische und andere Hilfsmittel
- Numerische Methoden für GDGen
- Lineare Mehrschrittverfahren
- Runge-Kutta-Verfahren
- Schrittweitensteuerung
- Steife Systeme
- Exponentielle Integratoren
- Zufällige und stochastische Differentialgleichungen (wenn Zeit erlaubt)

- Harro Heuser. *Gewöhnliche Differentialgleichungen*. Vieweg-Teubner, 2009 (6. Auflage).
- Martin Braun. *Differentialgleichungen und ihre Anwendungen*. Springer-Verlag, 1979.
- Ernst Hairer, Syvert P. Nørsett, Gerhard Wanner. *Solving Ordinary Differential Equations I: Nonstiff Problems*. Springer, 1987.
Ernst Hairer and Gerhard Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*. Springer, 1991.
- John D. Lambert. *Numerical Methods for Ordinary Differential Systems: The Initial Value Problem* Wiley, 1991
- Arieh Iserles. *A First Course in the Numerical Analysis of Differential Equations*. Cambridge University Press, 1996.
- Uri M. Ascher, Linda R. Petzold. *Computer Methods for Ordinary Differential Equations and Differential-Algebraic Equations*. SIAM, 1998.
- David F. Griffiths and Desmond J. Higham. *Numerical Methods for Ordinary Differential Equations: Initial Value Problems*. Springer, 2010.

Siehe auch die (laufend ergänzte) [Literaturliste](#) auf der Webseite der Vorlesung.

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

- 2.1 Das Euler-Verfahren
- 2.2 Eine Sammlung von Beispielverfahren
- 2.3 Konvergenz, Konsistenz und Stabilität
- 2.4 Der Hauptsatz
- 2.5 Einschrittverfahren
- 2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

- 3.1 Begriffe
- 3.2 Konsistenzordnung linearer Mehrschrittverfahren

- 3.3 Die erste Dahlquist-Barriere
- 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
- 3.5 Prädiktor-Korrektor-Verfahren
- 3.6 Absolute Stabilität
- 3.7 BDF-Verfahren
- ① Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme

- 5.5 RKV für steife Probleme
- 5.6 Nichtlineare Stabilitätstheorie

6 Ausblick

- ① **Einleitung**
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

① Einleitung

1.1 Volterras Prinzip

1.2 Begriffe und theoretische Resultate

1.3 Lineare Differenzengleichungen

1.4 Matrixfunktionen

1.5 Systeme linearer Differentialgleichungen erster Ordnung

1.6 Die Fälschungen des Han van Meegeren

1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

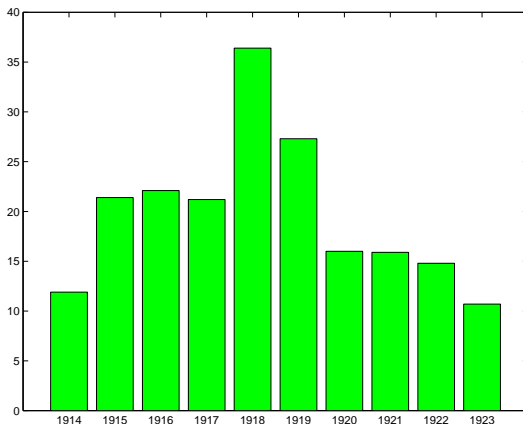
④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Volterras Prinzip

Umberto d'Anconas Beobachtung

Der Biologe Umberto d'Ancona (1896–1964) stellte 1925 den prozentualen Anteil der Haie am Gesamtfang (Speisefische und Haie) im Hafen von Triest fest:



Benachteiligt eingeschränkter Fischfang (1. Weltkrieg) die Speisefische?

Volterras Prinzip

Volterras Räuber-Beute-Modell

D'Ancona konsultierte den Mathematiker Volterra¹, der die Populationsdynamik wie folgt modellierte: Seien

$x(t)$: Beutepopulation zur Zeit t (Speisefische)

$y(t)$: Räuberpopulation zur Zeit t (Haie).

Ohne Räuber würde sich die Beute nach dem **Malthusianischen**² Gesetz

$$x'(t) = a x(t) \quad (\text{mit einer Konstanten } a > 0),$$

vermehrten, d.h. der Zuwachs wäre proportional zum Bestand bzw. das Wachstum wäre exponentiell

$$x(t) = x(0) \exp(at) \quad \text{für } t \geq 0$$

(eingeschränkt realistisch, falls Population nicht sehr dicht und ausreichend Nahrung vorhanden ist).

¹Vito Volterra (1860–1940)

²Thomas Malthus (1766–1834)

Volterras Prinzip

Interaktion von Räuber und Beute

Anzahl Räuber-Beute-Kontakte (pro Zeiteinheit):

$$b x(t) y(t) \quad (\text{mit einer Konstanten } b > 0).$$

$$\text{Insgesamt:} \quad x'(t) = a x(t) - b x(t) y(t).$$

$$\text{Analog:} \quad y'(t) = -c y(t) + d x(t) y(t), \quad \text{mit weiteren Konstanten } c, d > 0.$$

Wir erhalten ein **System zweier GDGen**.

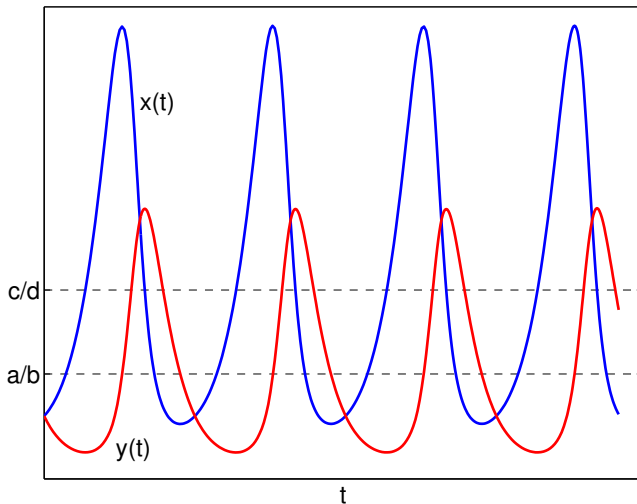
Man kann zeigen: dessen Lösungen sind **periodisch**: d.h. $\exists T > 0$ sodass

$$x(t + T) = x(t), \quad y(t + T) = y(t) \quad \text{für alle } t.$$

$$\text{Mittelwerte:} \quad \bar{x} := \frac{1}{T} \int_0^T x(t) dt = \frac{c}{d}, \quad \bar{y} := \frac{1}{T} \int_0^T y(t) dt = \frac{a}{b}.$$

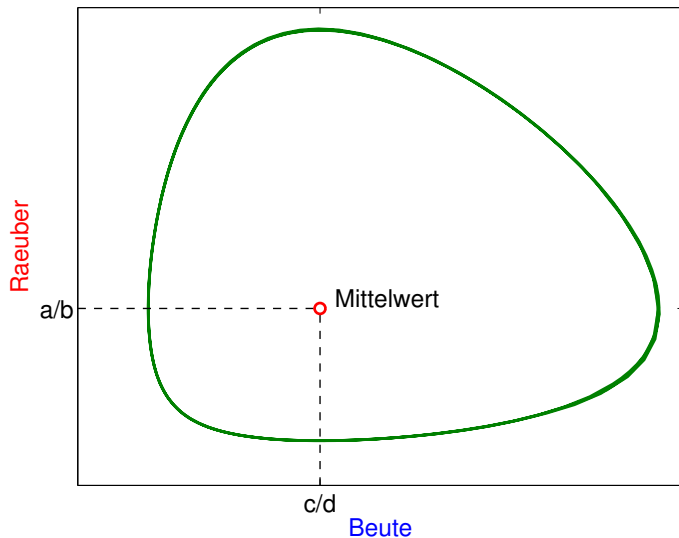
Volterra's Prinzip

Zeitlicher Verlauf der Populationen



Volterras Prinzip

Darstellung in der Phasebene



Volterras Prinzip

Auswirkung von Fischfang

Berücksichtige Fischfang:

$$\begin{aligned}x'(t) &= a x(t) - b x(t)y(t) - e x(t) &= (a - e) x(t) - b x(t)y(t), \\y'(t) &= -c y(t) + d x(t)y(t) - e y(t) &= -(c + e) y(t) + d x(t)y(t),\end{aligned}\quad (e > 0).$$

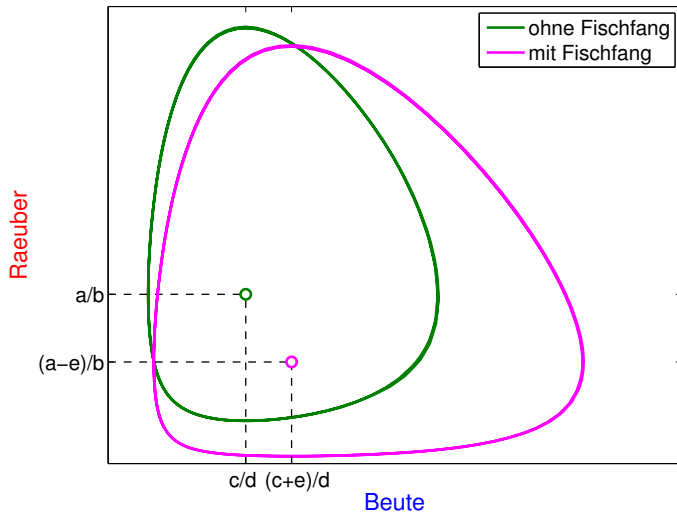
Gleiches System mit neuen Koeffizienten: $a \rightarrow a - e$ und $c \rightarrow c + e$.

$$\text{Mittelwerte:} \quad \frac{c + e}{d} > \frac{c}{d} \quad (\text{Beute}), \quad \frac{a - e}{b} < \frac{a}{b} \quad (\text{Räuber}).$$

Volterras Prinzip: Moderater Fischfang ($e < a$) steigert die durchschnittliche Zahl der Speisefische und reduziert die durchschnittliche Zahl der Haie.

Volterra's Prinzip

Darstellung in der Phasebene



① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Einen Ausdruck der Form

$$F(t, y, y', y'', \dots, y^{(n)}) = 0 \quad (\text{GDG})$$

mit einer Funktion $F : \mathbb{R}^{n+2} \supset M \rightarrow \mathbb{R}$ nennen wir eine **gewöhnliche Differentialgleichung** (GDG) **n -ter Ordnung**. Eine Funktion $y : \mathbb{R} \supset I \rightarrow \mathbb{R}$ heißt **Lösung** von (GDG) **über dem Intervall I** , wenn $y \in C^n(I)$ ist und für alle $t \in I$ gilt

$$F(t, y(t), y'(t), y''(t), \dots, y^{(n)}(t)) = 0.$$

- (GDG) besitzt die **Ordnung n** , weil n die Ordnung der höchsten auftretenden Ableitung ist.
- Sie heißt **gewöhnlich**, weil nur Ableitungen der gesuchten Funktion y nach einer Variablen auftreten.
- (GDG) heißt **implizit** — im Gegensatz zu einer **expliziten** GDG n -ter Ordnung, die nach der höchsten Ableitung von y aufgelöst ist:

$$y^{(n)} = f(t, y, y', \dots, y^{(n-1)}).$$

Wir werden fast ausschließlich **Systeme** von expliziten GDGen erster Ordnung betrachten (warum wir uns auf Systeme erster Ordnung beschränken können, wird später erklärt):

$$\begin{aligned}y_1' &= f_1(t, y_1, y_2, \dots, y_n) \\y_2' &= f_2(t, y_1, y_2, \dots, y_n) \\&\vdots = \vdots \\y_n' &= f_n(t, y_1, y_2, \dots, y_n)\end{aligned}\tag{DG}$$

mit den n unbekannten Funktionen y_1, y_2, \dots, y_n . Jedes System von n Funktionen

$$y_1 = y_1(t), \dots, y_n = y_n(t) \in C^1(I),$$

das (DG) für alle $t \in I$ erfüllt, heißt **Lösung** von (DG) über I .

Das System

$$\begin{aligned}y_1' &= 1 \\ y_2' &= 2y_1\end{aligned}$$

besitzt die Lösungen

$$y_1(t) = t + \alpha, \quad y_2(t) = t^2 + 2\alpha t + \beta \quad (\alpha, \beta \in \mathbb{R})$$

über $(-\infty, \infty)$.

Für eine eindeutige Lösung: **Anfangsbedingungen**, z.B.

$$y_1(0) = 1, \quad y_2(0) = 2.$$

Dann ist

$$y_1(t) = t + 1, \quad y_2(t) = t^2 + 2t + 2$$

die einzige Lösung.

Allgemein: Die Aufgabenstellung, eine Lösung von (DG) zu finden, die die Anfangsbedingung

$$y_1(t_0) = y_{0,1}, \dots, y_n(t_0) = y_{0,n} \quad (\text{AB})$$

erfüllt, heißt **Anfangswertproblem** (AWP) oder Anfangswertaufgabe für die gewöhnliche Differentialgleichung (DG).

Mit der Vektornotation

$$\mathbf{y} := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{f} := \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}, \quad \mathbf{y}_0 := \begin{bmatrix} y_{0,1} \\ \vdots \\ y_{0,n} \end{bmatrix},$$

ergibt sich die Kurzschreibweise

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}), \quad (\text{DG}')$$

$$\mathbf{y}(t_0) = \mathbf{y}_0. \quad (\text{AB}')$$

Bemerkung. GDGen höherer Ordnung lassen sich in (äquivalente) Systeme von GDGen erster Ordnung umschreiben:

Aus

$$y''' + 3y'' + y' = \sin(t)$$

wird etwa

$$\begin{bmatrix} y_1' \\ y_2' \\ y_3' \end{bmatrix} = \begin{bmatrix} y_2 \\ y_3 \\ -3y_3 - y_2 + \sin(t) \end{bmatrix}$$

mit den neuen Variablen

$$y_1 = y, \quad y_2 = y_1' = y', \quad y_3 = y_2' = y''.$$

Die explizite Abhängigkeit der rechten Seite von der unabhängigen Variable (hier t) kann durch Hinzunahme einer zusätzlichen Gleichung bzw. Komponente des Lösungsfunktionsvektors \mathbf{y} beseitigt werden:

$$y_4(t) = t \quad (\text{d.h. } y_4'(t) \equiv 1), \quad y_4(t_0) = t_0.$$

Im obigen Beispiel resultiert dies in der **autonomen** Differentialgleichung $\mathbf{y}' = \mathbf{f}(\mathbf{y})$, oder genauer:

$$\mathbf{y}'(t) = \mathbf{f}(\mathbf{y}(t)), \quad \mathbf{f}(\mathbf{y}) = \begin{bmatrix} y_2 \\ y_3 \\ -3y_3 - y_2 + \sin(y_4) \\ 1 \end{bmatrix}.$$

Satz 1.1 (Picard-Lindelöf)

Gegeben ist die Anfangswertaufgabe

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}), \quad \mathbf{y}(t_0) = \mathbf{y}_0. \quad (\text{AWP})$$

Die rechte Seite \mathbf{f} sei stetig im 'Quader'

$$Q := \{(t, \mathbf{y}) : |t - t_0| \leq a, \|\mathbf{y} - \mathbf{y}_0\| \leq b\}, \quad Q \subset \mathbb{R}^{n+1},$$

und es sei $M := \max\{\|\mathbf{f}(t, \mathbf{y})\| : (t, \mathbf{y}) \in Q\}$.

Außerdem erfülle \mathbf{f} in Q die **Lipschitz-Bedingung**

$$\|\mathbf{f}(t, \mathbf{y}) - \mathbf{f}(t, \tilde{\mathbf{y}})\| \leq L\|\mathbf{y} - \tilde{\mathbf{y}}\| \quad \forall (t, \mathbf{y}), (t, \tilde{\mathbf{y}}) \in Q. \quad (\text{Lip})$$

Dann besitzt das Problem (**AWP**) genau eine Lösung über $I := [t_0 - \alpha, t_0 + \alpha]$, wobei $\alpha = \min\{a, b/M\}$.

Emile Picard (1856–1941), Ernst Lindelöf (1870–1946), Rudolf Lipschitz (1832–1903).

Bemerkungen.

- (1) In der gesamten Vorlesung wird vorausgesetzt, dass die fundamentale Bedingung (Lip) erfüllt ist.
- (2) (AWP) besitzt in $[t_0 - a, t_0 + a]$ eine eindeutige Lösung, wenn \mathbf{f} die Bedingung (Lip) in $\tilde{Q} = \{(t, \mathbf{y}) : |t - t_0| \leq a, \|\mathbf{y}\| < \infty\}$ erfüllt.
- (3) Ist \mathbf{f} auf Q bez. \mathbf{y} stetig differenzierbar und bezeichnet $\mathbf{f}_y = [\partial f_i / \partial y_j]_{1 \leq i, j \leq n}$ die zugehörige Jacobi-Matrix, dann folgt aus dem Mittelwertsatz, dass die Voraussetzungen von Satz 1.1 mit

$$L = \sup_{(t, \mathbf{y}) \in Q} \|\mathbf{f}_y(t, \mathbf{y})\| < \infty$$

erfüllt sind.

- (4) (AWP) besitzt auch dann noch Lösungen, wenn \mathbf{f} nur als stetig auf Q vorausgesetzt wird (Existenzsatz von Peano³). Deren Eindeutigkeit ist aber nicht mehr gesichert.

³Giuseppe Peano (1858–1932)

Theoretische Grundlagen

Existenz- und Eindeutigkeit der Lösung: Beispiel

Beispiel:

$$y' = f(t, y) = \sqrt{y}, \quad y(0) = 0, \quad Q = \mathbb{R} \times [0, \infty),$$

mit den Lösungen

$$y_\lambda(t) = \begin{cases} 0, & 0 \leq t \leq \lambda, \\ (t - \lambda)^2/4, & t \geq \lambda, \end{cases} \quad (\lambda \geq 0).$$

Satz 1.2

Die Anfangswertaufgabe

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}), \quad \mathbf{y}(t_0) = \mathbf{y}_0,$$

erfülle die Voraussetzungen von Satz 1.1. Über eine weitere Anfangswertaufgabe

$$\mathbf{y}' = \tilde{\mathbf{f}}(t, \mathbf{y}), \quad \mathbf{y}(t_0) = \tilde{\mathbf{y}}_0,$$

setzen wir nur voraus, dass $\tilde{\mathbf{f}}$ stetig in Q ist. Sind dann \mathbf{y} und $\tilde{\mathbf{y}}$ Lösungen dieser Anfangswertaufgaben über dem Intervall I und gilt

$$\|\mathbf{y}_0 - \tilde{\mathbf{y}}_0\| \leq \gamma \quad \text{sowie} \quad \|\mathbf{f}(t, \mathbf{y}) - \tilde{\mathbf{f}}(t, \mathbf{y})\| \leq \delta \quad \forall (t, \mathbf{y}) \in Q,$$

so folgt für $t \in I$

$$\|\mathbf{y}(t) - \tilde{\mathbf{y}}(t)\| \leq \gamma e^{L(t-t_0)} + \frac{\delta}{L} \left(e^{L(t-t_0)} - 1 \right).$$

(vgl. [Heuser, Satz 13.1])

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Eine wichtige Rolle werden **lineare Differenzengleichungen**

$$y_{n+k} + \alpha_{k-1}y_{n+k-1} + \cdots + \alpha_0y_n = \beta_{n+k} \quad (n = 0, 1, 2, \dots) \quad (\text{DzG})$$

spielen.

- Genauer spricht man hier von einer linearen Differenzengleichung der **Ordnung** k mit **konstanten Koeffizienten** (die α 's hängen nicht von n ab). (O.B.d.A. sei $\alpha_0 \neq 0$).
- Die Gleichung heißt **homogen**, wenn $\beta_{n+k} = 0$ für alle n , andernfalls **inhomogen**.
- Jede Folge $\{y_n\}_n$, die (DzG) erfüllt, heißt eine **Lösung** von (DzG).
- Gibt man sich k Startwerte y_0, y_1, \dots, y_{k-1} (beliebig) vor, kann man sich mit (DzG) **rekursiv** eine solche Lösung berechnen.

Lineare Differenzengleichungen

Lösungsstruktur, homogene Gleichung

Lemma 1.3

Die Lösungsmenge einer homogenen linearen Differenzengleichung der Ordnung k ist ein Vektorraum der Dimension k .

Besitzt die Differenzengleichung darüber hinaus konstante Koeffizienten, so kann man eine Basis dieses Lösungsraums mit Hilfe der Nullstellen des zugehörigen **charakteristischen Polynoms**

$$p_k(\zeta) = \zeta^k + \alpha_{k-1}\zeta^{k-1} + \cdots + \alpha_1\zeta + \alpha_0$$

angeben: Bezeichnen λ_j ($1 \leq j \leq \ell$) die Nullstellen von p (mit Vielfachheiten m_j , $\sum_{j=1}^{\ell} m_j = k$), so bilden die k Folgen

$$(\lambda_j^n)_n, (n\lambda_j^{n-1})_n, \dots, (n(n-1)\dots(n-m_j+2)\lambda_j^{n-m_j+1})_n \quad (j = 1, 2, \dots, \ell)$$

eine solche Basis.

Lemma 1.4

Für die homogene Differenzengleichung

$$y_{n+k} + \alpha_{k-1}y_{n+k-1} + \cdots + \alpha_0y_n = 0 \quad (n = 0, 1, 2, \dots) \quad (*)$$

sind die folgenden drei Aussagen äquivalent :

- (1) Jede Lösung $\{y_n\}_n$ von $(*)$ ist beschränkt.
- (2) Für jede Lösung $\{y_n\}_n$ von $(*)$ ist $\{y_n/n\}_n$ eine Nullfolge.
- (3) Das zugehörige charakteristische Polynom p erfüllt die sogenannte **Stabilitätsbedingung**:

$$\begin{aligned} p(\lambda) = 0 &\Rightarrow |\lambda| \leq 1, \\ p(\lambda) = 0 \text{ und } |\lambda| = 1 &\Rightarrow \lambda \text{ ist einfach.} \end{aligned} \quad (\text{Stab})$$

Lineare Differenzengleichungen

Eine rekursive Abschätzung

Lemma 1.5

Es gebe Konstanten $M, K \geq 0$, so dass die ersten Glieder der Vektorfolge $(\mathbf{y}_n)_n$ die Ungleichung

$$\|\mathbf{y}_{n+1}\| \leq K\|\mathbf{y}_n\| + M \quad (n = 0, 1, \dots, n_0)$$

erfüllen. Dann gilt die Abschätzung

$$\|\mathbf{y}_{n+1}\| \leq K^{n+1}\|\mathbf{y}_0\| + \begin{cases} M \frac{K^{n+1} - 1}{K - 1}, & \text{für } K \neq 1, \\ (n+1)M, & \text{für } K = 1, \end{cases} \quad (n = 0, 1, \dots, n_0).$$

($\|\cdot\|$ bezeichnet eine beliebige Norm.)

Lineare Differenzengleichungen

Lösung der inhomogenen Gleichung

Gesucht ist eine explizite Darstellung der Lösung $(z_n)_n$ der inhomogenen Differenzengleichung

$$y_{n+k} + \alpha_{k-1}y_{n+k-1} + \cdots + \alpha_1y_{n+1} + \alpha_0y_n = \beta_{n+k}, \quad n \in \mathbb{N}_0,$$

die die k Anfangsbedingungen $z_n = y_n$ ($n = 0, 1, \dots, k-1$) erfüllt.

Antwort: Bezeichnen $(y_n^{(j)})_n$, $j = 0, 1, \dots, k-1$, die Lösungen der homogenen Gleichung

$$y_{n+k} + \alpha_{k-1}y_{n+k-1} + \cdots + \alpha_1y_{n+1} + \alpha_0y_n = 0 \quad (n = 0, 1, \dots),$$

die die Anfangsbedingungen $y_n^{(j)} = \delta_{n,j}$ (Kronecker-Symbol) ($n, j = 0, 1, \dots, k-1$) erfüllen, so ist

$$z_n = \sum_{j=0}^{k-1} y_j y_n^{(j)} + \sum_{j=0}^{n-k} \beta_{j+k} y_{n-j-1}^{(k-1)} \quad (n = 0, 1, \dots),$$

wobei $\beta_{n+k} = 0$ und $y_n^{(k-1)} = 0$ für $n < 0$ gesetzt wird.

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Matrixfunktionen

In diesem Abschnitt sei $A \in \mathbb{C}^{n \times n}$ stets eine quadratische Matrix. Außerdem sei eine Funktion

$$f : D \rightarrow \mathbb{C}, \quad D \subset \mathbb{C},$$

gegeben. Wir klären hier, wann und wie die Matrix

$$f(A) \in \mathbb{C}^{n \times n}$$

definiert ist, und wiederholen einige ihrer Eigenschaften. Im Zusammenhang mit GDGen von Interesse ist besonders $\exp(A)$, die **Exponentialfunktion** angewandt auf A .

Für einige elementare Funktionen f ist $f(A)$ kanonisch gegeben. Ist z.B. $f \in \mathcal{P}_m$ ein Polynom vom Grad m ,

$$f(\lambda) = \alpha_0 + \alpha_1 \lambda + \alpha_2 \lambda^2 + \cdots + \alpha_m \lambda^m \quad (\alpha_j \in \mathbb{C}, j = 0, 1, \dots, m),$$

so ist

$$f(A) = \alpha_0 I + \alpha_1 A + \alpha_2 A^2 + \cdots + \alpha_m A^m \in \mathbb{C}^{n \times n}.$$

Matrixfunktionen

Eigenschaften von $f(A)$ für Polynome f

Lemma 1.6

Sei $f \in \mathcal{P}_m$.

- (a) Hat $A = \text{diag}(A_{1,1}, A_{2,2}, \dots, A_{k,k})$ Blockdiagonalstruktur mit quadratischen Diagonalblöcken

$$A_{j,j} \in \mathbb{C}^{n_j \times n_j}, \quad (j = 1, 2, \dots, k), \quad n_1 + n_2 + \dots + n_k = n,$$

dann gilt

$$f(A) = \text{diag}(f(A_{1,1}), f(A_{2,2}), \dots, f(A_{k,k})).$$

- (b) Ist $T \in \mathbb{C}^{n \times n}$ invertierbar und $B := TAT^{-1}$, dann gilt

$$f(B) = Tf(A)T^{-1}.$$

- (c) Ist λ ein Eigenwert von A mit zugehörigem Eigenvektor v , so ist $f(\lambda)$ ein Eigenwert von $f(A)$ mit zugehörigem Eigenvektor v :

$$Av = \lambda v \quad \implies \quad f(A)v = f(\lambda)v.$$

Matrixfunktionen

Beispiel: Potenzen eines Jordan-Blocks

Wir bestimmen $m_k(J)$ für das k -te Monom $m_k(\lambda) = \lambda^k$ und einen Jordan-Block⁴

$$J = J(\lambda) = \begin{bmatrix} \lambda & 1 & & & \\ & \lambda & 1 & & \\ & & \ddots & \ddots & \\ & & & \lambda & 1 \\ & & & & \lambda \end{bmatrix} \in \mathbb{C}^{n \times n}.$$

Eine elementare Rechnung zeigt, dass $m_k(J) = J^k$ eine obere Dreiecksmatrix mit Toeplitz-Struktur⁵ ist. Der Eintrag in der j -ten Diagonale ist

$$\binom{k}{j} \lambda^{k-j} = \frac{k(k-1) \cdots (k-j+1) \lambda^{k-j}}{j!} = \frac{m_k^{(j)}(\lambda)}{j!} \quad (j = 0, 1, \dots, n-1).$$

⁴Camille Jordan (1838–1922)

⁵Otto Toeplitz (1881–1940)

Matrixfunktionen

Beispiel: Potenzen eines Jordan-Blocks

Mit anderen Worten:

$$m_k(J) = J^k = \begin{bmatrix} m_k(\lambda) & m'_k(\lambda) & \cdots & \frac{m_k^{(n-2)}(\lambda)}{(n-2)!} & \frac{m_k^{(n-1)}(\lambda)}{(n-1)!} \\ & m_k(\lambda) & \cdots & \frac{m_k^{(n-3)}(\lambda)}{(n-3)!} & \frac{m_k^{(n-2)}(\lambda)}{(n-2)!} \\ & & \ddots & \vdots & \vdots \\ & & & m_k(\lambda) & m'_k(\lambda) \\ & & & & m_k(\lambda) \end{bmatrix}.$$

Jetzt sind wir in der Lage, $f(A)$ für beliebiges f zu definieren: Sei dazu $J_A = \text{diag}(J_1, J_2, \dots, J_k)$ die **Jordansche Normalform** von A , $A = TJ_AT^{-1}$. Die einzelnen Jordan-Blöcke $J_j = J_j(\lambda_j)$ seien $(n_j \times n_j)$ -Matrizen. Das charakteristische Polynom c_A von A hat dann die Form

$$c_A(\lambda) = \prod_{j=1}^k (\lambda - \lambda_j)^{n_j}.$$

Matrixfunktionen

Allgemeiner Fall

Wir sagen f ist auf A definiert, wenn f auf einer offenen Menge D definiert ist, die das Spektrum $\Lambda(A) = \{\lambda_1, \dots, \lambda_k\}$ von A enthält, und außerdem f in λ_j $(n_j - 1)$ -mal differenzierbar ist.

In diesem Fall setzen wir für $j = 1, 2, \dots, k$

$$f(J_j(\lambda_j)) := \begin{bmatrix} f(\lambda_j) & f'(\lambda_j) & \cdots & \frac{f^{(n_j-2)}(\lambda_j)}{(n_j-2)!} & \frac{f^{(n_j-1)}(\lambda_j)}{(n_j-1)!} \\ & f(\lambda_j) & \cdots & \frac{f^{(n_j-3)}(\lambda_j)}{(n_j-3)!} & \frac{f^{(n_j-2)}(\lambda_j)}{(n_j-2)!} \\ & & \ddots & \vdots & \vdots \\ & & & f(\lambda_j) & f'(\lambda_j) \\ & & & & f(\lambda_j) \end{bmatrix} \in \mathbb{C}^{n_j \times n_j}$$

und

$$f(A) := T \operatorname{diag}(f(J_1), f(J_2), \dots, f(J_k)) T^{-1}.$$

Bemerkungen.

- (1) Ist $f(\lambda) = \alpha_0 + \alpha_1\lambda + \cdots + \alpha_m\lambda^m$ ein Polynom, so gilt für die so definierte Matrix $f(A)$: $f(A) = \alpha_0 I + \alpha_1 A + \cdots + \alpha_m A^m$ (natürlich ist ein Polynom auf jeder Matrix definiert).
- (2) Ist $f = p/q$ eine rationale Funktion mit $p \in \mathcal{P}_m$ und $q \in \mathcal{P}_k$, so ist f genau dann auf A definiert, wenn kein Eigenwert von A eine Polstelle von f ist. In diesem Fall gilt $f(A) = p(A)[q(A)]^{-1} = [q(A)]^{-1}p(A)$.
- (3) Alle Aussagen von Lemma 1.6 gelten für jede beliebige Funktion f , wenn nur f auf A definiert ist.
- (4) Ist f analytisch in einer Umgebung von 0 und besitzt dort die Taylor⁶-Reihe $f(\lambda) = \sum_{j=0}^{\infty} \alpha_j \lambda^j$ mit Konvergenzradius $\tau > 0$ ($\tau = \infty$ ist erlaubt), so ist f auf jeder Matrix A mit Spektralradius $\rho(A) < \tau$ definiert und es gilt

$$f(A) = \sum_{j=0}^{\infty} \alpha_j A^j = \lim_{m \rightarrow \infty} \sum_{j=0}^m \alpha_j A^j.$$

⁶Brook Taylor (1685–1731)

Matrixfunktionen

Beispiele: Neumannsche Reihe, Exponentialfunktion

Ein bekanntes Beispiel für die letzte Bemerkung ist die **Neumannsche Reihe**⁷

$$(I - A)^{-1} = \sum_{j=0}^{\infty} A^j, \quad \text{falls } \rho(A) < 1.$$

Die **Exponentialfunktion einer Matrix** A kann z.B. auch durch

$$\exp(A) = \sum_{j=0}^{\infty} \frac{1}{j!} A^j$$

definiert werden. (Die Reihe konvergiert für jede Matrix A , weil die zugehörige skalare Reihe einen unendlichen Konvergenzradius besitzt.)

Bei GDGen ist es oft wichtig zu wissen, wie sich $\exp(tA)$ für $t \rightarrow \infty$ verhält. Die entscheidende Größe ist die **Spektralabszisse** $\alpha(A)$ von A :

$$\alpha(A) := \max\{\operatorname{Re}(\lambda) : \lambda \in \Lambda(A)\}.$$

⁷Carl Neumann (1832–1925)

Satz 1.7 (Asymptotisches Verhalten von $\exp(tA)$)

Sei $A \in \mathbb{C}^{n \times n}$.

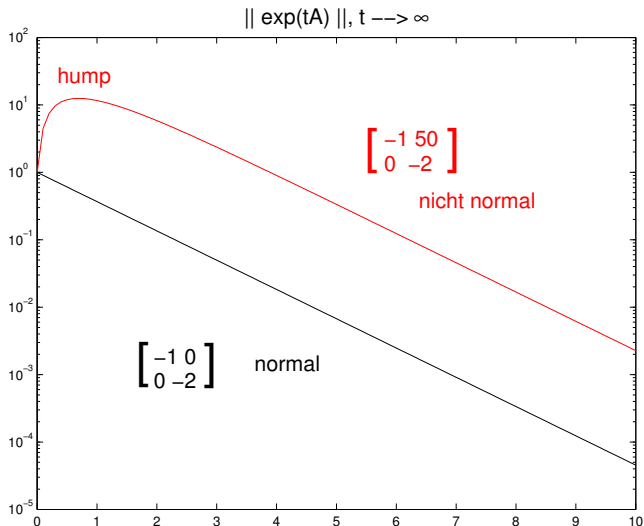
- (a) Es ist $\lim_{t \rightarrow \infty} \exp(tA) = 0$ genau dann, wenn $\alpha(A) < 0$ gilt.
- (b) Wenn $\alpha(A) > 0$ ist, so ist $\exp(tA)$ für $t \rightarrow \infty$ unbeschränkt.
Ist $\alpha(A) = 0$ und jeder Eigenwert λ von A mit $\operatorname{Re} \lambda = \alpha(A)$ halbeinfach, so ist $\exp(tA)$ für $t \rightarrow \infty$ beschränkt (aber i.A. nicht konvergent).
- (c) Es gilt $\|\exp(tA)\| \geq \exp(t\alpha(A))$ für jede Matrixnorm.
Ist A normal, dann gilt $\|\exp(tA)\|_2 = \exp(t\alpha(A))$.

Für normale Matrizen ist $\|\exp(tA)\|_2$ also eine streng monoton fallende Funktion von $t \geq 0$, wenn $\alpha(A) < 0$.

Ist A nicht normal, so beobachtet man die üblichen Nichtnormalitätseffekte (vollkommen analog zum Verhalten von $\|A^m\|_2$, $m \rightarrow \infty$, falls $\rho(A) < 1$).

Matrixfunktionen

Asymptotisches Verhalten der Matrix-Exponentialfunktion



Matrixfunktionen

Auswertung von Matrixfunktionen

Zum Abschluss soll noch ein Algorithmus zur Berechnung von $\exp(A)$ beschrieben werden.

Beachte: $\exp(A) = \lim_{m \rightarrow \infty} \sum_{j=0}^m A^j/j!$ ist nur geeignet, wenn $\rho(A)$ sehr klein ist; die Bestimmung von $\exp(A)$ über die Jordansche Normalform von A ist numerisch instabil oder zu aufwendig.

Der Algorithmus verwendet rationale Approximationen $(k/\ell)_{\exp}(\zeta) = p_{k,\ell}(\zeta)/q_{k,\ell}(\zeta)$ vom Typ (k, ℓ) (d.h. $p_{k,\ell} \in \mathcal{P}_k$, $q_{k,\ell} \in \mathcal{P}_\ell$) an die Exponentialfunktion, sog. **Padé-Approximationen**⁸. Diese sind eindeutig bestimmt durch die Vorschrift

$$\exp(\zeta) - p_{k,\ell}(\zeta)/q_{k,\ell}(\zeta) = O(\zeta^{k+\ell+1}) \quad \text{für } \zeta \rightarrow 0.$$

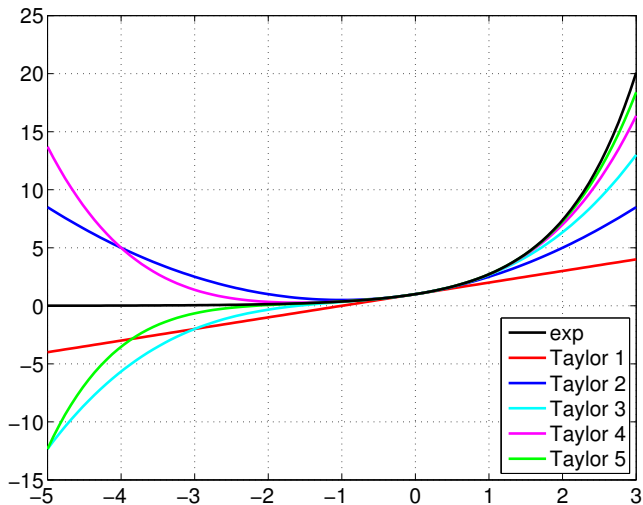
Beachte, dass hier die Taylor-Polynome für $\ell = 0$ als Spezialfall enthalten sind. Man kann die Padé-Approximationen im Fall der Exponentialfunktion explizit angeben:

$$p_{k,\ell}(\zeta) = \sum_{j=0}^k \frac{(k+\ell-j)!k!}{(k+\ell)!j!(k-j)!} \zeta^j, \quad q_{k,\ell}(\zeta) = \sum_{j=0}^{\ell} \frac{(k+\ell-j)! \ell!}{(k+\ell)!j!(\ell-j)!} \zeta^j.$$

⁸Henri Padé (1863–1953)

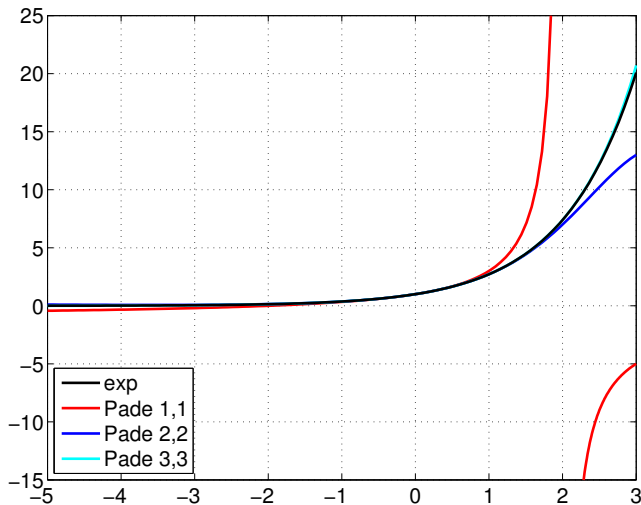
Matrixfunktionen

Taylorpolynome der Exponentialfunktion



Matrixfunktionen

Padé-Approximationen an die Exponentialfunktion



Satz 1.8 (Fehlerformel für Padé-Approximationen)

Für $k, \ell \in \mathbb{N}_0$ und $A \in \mathbb{C}^{n \times n}$ gilt

$$\begin{aligned} \exp(A) - (k, \ell)_{\exp}(A) \\ = \frac{(-1)^\ell}{(k + \ell)!} A^{k+\ell+1} [q_{k, \ell}(A)]^{-1} \int_0^1 u^k (1 - u)^\ell \exp((1 - u)A) \, du. \end{aligned}$$

Das bedeutet, dass $(k, \ell)_{\exp}(A)$

- sowohl für $k \rightarrow \infty$ bei festem ℓ ,
- als auch für $\ell \rightarrow \infty$ bei festem k ,
- als auch für $k \rightarrow \infty$ bei festem $k - \ell$

gegen $\exp(A)$ strebt.

Integrale über Matrizen $A = [a_{i,j}(u)]$ sind komponentenweise definiert:

$$\int A(u) \, du = \left[\int a_{i,j}(u) \, du \right].$$

Satz 1.8 zeigt, dass auch $(k, \ell)_{\exp}(A)$ nur dann eine akzeptable Näherung für $\exp(A)$ ist, wenn $\rho(A)$ nicht zu groß ist.

Daher verwendet man einen Trick: Kommutieren die Matrizen $A, B \in \mathbb{C}^{n \times n}$, so gilt

$$\exp(A + B) = \exp(A) \exp(B).$$

Insbesondere ist also

$$\exp(A) = \exp(A/m)^m \quad \text{für } m = 0, 1, 2, \dots$$

Das bedeutet, dass

$$E_{k,\ell} := [(k, \ell)_{\exp}(A/2^j)]^{2^j}$$

eine Approximation an $\exp(A)$ darstellt, bei der die Padé-Approximation an der Matrix $A/2^j$ ausgewertet wird, deren Spektralradius $\rho(A)/2^j$ man durch die Wahl von j steuern kann.

Die Berechnung von $E_{k,\ell}$ erfordert $j + \max\{k, \ell\}$ Multiplikationen mit A .

Lemma 1.9

Sei $\|A\|_\infty/2^j \leq 1/2$. Dann ist

$$\frac{\|\exp(A) - E_{k,\ell}\|_\infty}{\|\exp(A)\|_\infty} \leq \varepsilon(k, \ell) \|A\|_\infty \exp(\varepsilon(k, \ell) \|A\|_\infty)$$

mit

$$\varepsilon(k, \ell) = 2^{3-(k+\ell)} \frac{k!\ell!}{(k+\ell)!(k+\ell+1)!}.$$

Bei festem $d = \max\{k, \ell\}$ (Arbeitsaufwand zur Berechnung von $E_{k,\ell}$), wird $\varepsilon(k, \ell)$ durch die Wahl $k = \ell = d$ minimiert.

Algorithmus 1 : Berechnung von $\exp(A)$.

Gegeben : $A, \delta > 0$.

- 1 $j \leftarrow \max\{0, 1 + \text{floor}(\log_2 \|A\|_\infty)\}.$
 - 2 $A \leftarrow A/2^j.$
 - 3 Wähle ℓ minimal mit $\epsilon(\ell, \ell) \leq \delta.$
 - 4 $N \leftarrow I, Z \leftarrow I, X \leftarrow I, c \leftarrow 1.$
 - 5 **for** $m = 1$ **to** ℓ **do**
 - 6 $c \leftarrow c(\ell - m + 1)/((2\ell - m + 1)m).$
 - 7 $X \leftarrow AX; Z \leftarrow Z + cX; N \leftarrow N + (-1)^m cX.$
 - 8 Bestimme LU-Zerlegung von N und löse damit $NE = Z$ nach E auf.
 - 9 **for** $m = 1$ **to** j **do**
 - 10 $E \leftarrow EE$
-

- Dieser Algorithmus liefert eine Approximation $E \approx \exp(A)$ derart, dass

$$E = \exp(A + \Delta A), \quad \text{wobei} \quad \|\Delta A\|_{\infty} \leq \delta \|A\|_{\infty}.$$

- Seine Komplexität beträgt etwa $2(\ell + j + 1/3)n^3$ flops, vgl. [Moler & Van Loan, 2003]⁹ und [Ward, 1977]¹⁰
- Für große dünnbesetzte Matrizen ist unser Algorithmus jedoch ungeeignet (er verwendet die LU-Zerlegung einer Matrix der Dimension von A).
- Ähnlich wie bei linearen Gleichungssystemen, wo man selten an A^{-1} sondern vielmehr an $A^{-1}\mathbf{b}$, $\mathbf{b} \in \mathbb{C}^n$, interessiert ist, steht auch hier die Berechnung von $\exp(A)\mathbf{b}$ im Vordergrund. Bei großen dünn besetzten Problemen muss man auch dazu iterative Verfahren verwenden.

⁹C.B. Moler und C.F. Van Loan. *Nineteen dubious ways to compute the exponential of a matrix, twenty-five years later*. *SIAM Rev.* **45**, 3–49 (2003)

¹⁰R.C. Ward. *Numerical computation of the matrix exponential with accuracy estimate*. *SIAM J. Numer. Anal.* **14** (4) 600–610 (1977).

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Systeme linearer Differentialgleichungen erster Ordnung

Nur wenige Systeme der Form (DG) bzw. AWP der Form (AWP) kann man explizit lösen. Selbst **lineare Systeme erster Ordnung**,

$$y_1' = a_{1,1}(t)y_1 + a_{1,2}(t)y_2 + \cdots + a_{1,n}(t)y_n + b_1(t),$$

$$y_2' = a_{2,1}(t)y_1 + a_{2,2}(t)y_2 + \cdots + a_{2,n}(t)y_n + b_2(t),$$

$$\vdots = \vdots$$

$$y_n' = a_{n,1}(t)y_1 + a_{n,2}(t)y_2 + \cdots + a_{n,n}(t)y_n + b_n(t)$$

oder kürzer

$$\mathbf{y}' = A(t)\mathbf{y} + \mathbf{b}(t) \quad \text{mit } A(t) = [a_{i,j}(t)] \text{ und } \mathbf{b}(t) = [b_j(t)], \quad (\text{Lin})$$

gehören nur unter weiteren Einschränkungen zu diesen Ausnahmefällen.

Systeme linearer Differentialgleichungen erster Ordnung

Hinreichende Bedingung für Lösbarkeit

Sind die Funktionen $a_{i,j}(t)$, $b_j(t)$ stetig über einem Intervall I und ist $\|A(t)\| \leq L$ für alle $t \in I$ (was wir ab jetzt stets voraussetzen), so besitzt (Lin) nach Satz 1.1 für jede Wahl der Anfangsbedingungen

$$\mathbf{y}(t_0) = \mathbf{y}_0 \quad (t_0 \in I)$$

eine eindeutige Lösung.

Satz 1.10 (Lösungen linearer Systeme erster Ordnung)

Die Lösungen des **homogenen** Systems

$$\mathbf{y}' = A(t)\mathbf{y}$$

bilden einen n -dimensionalen Unterraum des $C^1(I)$. Die Differenz zweier Lösungen des **inhomogenen** Systems

$$\mathbf{y}' = A(t)\mathbf{y} + \mathbf{b}(t)$$

löst das zugehörige homogene System.

Systeme linearer Differentialgleichungen erster Ordnung

Konstante Koeffizienten, homogener Fall

Im Spezialfall **konstanter Koeffizienten**

$$a_{i,j}(t) = a_{i,j} \quad \text{für alle } t$$

lassen sich diese Lösungen angeben. Dazu betrachten wir zunächst den homogenen Fall, $\mathbf{b}(t) \equiv \mathbf{0}$: Es seien $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ die Einheitsvektoren im \mathbb{R}^n . Für $j = 1, 2, \dots, n$ löst

$$\mathbf{x}_j(t) := \exp(tA)\mathbf{u}_j$$

das AWP

$$\mathbf{y}' = A\mathbf{y}, \quad \mathbf{y}(0) = \mathbf{u}_j.$$

Darüber hinaus sind die Funktionen $\mathbf{x}_1(t), \mathbf{x}_2(t), \dots, \mathbf{x}_n(t)$ linear unabhängig und bilden deshalb eine Basis des Lösungsraums von $\mathbf{y}' = A\mathbf{y}$.

Schließlich ist die matrixwertige Funktion

$$X : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}, \quad X(t) := [\mathbf{x}_1(t) | \mathbf{x}_2(t) | \dots | \mathbf{x}_n(t)]$$

für alle $t \in \mathbb{R}$ invertierbar und löst das Anfangswertproblem

$$X'(t) = AX(t), \quad X(0) = I.$$

Systeme linearer Differentialgleichungen erster Ordnung

Konstante Koeffizienten, inhomogener Fall

Um das inhomogene Anfangswertproblem $\mathbf{y}' = A\mathbf{y} + \mathbf{b}(t)$ (der Einfachheit halber nehmen wir an, dass die Komponenten $b_j(t)$ auf ganz \mathbb{R} stetig sind), $\mathbf{y}(0) = \mathbf{y}_0$, zu lösen, bedient man sich einer Technik, die unter dem Namen **Variation der Konstanten** bekannt ist:

Die (eindeutige) Lösung ist gegeben durch

$$\mathbf{y}(t) = \sum_{j=1}^n \left[\int_{t_0}^t \frac{W_j(s)}{W(s)} ds + y_{0,j} \right] \mathbf{x}_j(t),$$

mit den **Wronski-Determinanten**¹¹

$$\begin{aligned} W(t) &= \det \begin{bmatrix} \mathbf{x}_1(t) & \mathbf{x}_2(t) & \cdots & \mathbf{x}_n(t) \end{bmatrix}, \\ W_j(t) &= \det \begin{bmatrix} \mathbf{x}_1(t) & \cdots & \mathbf{x}_{j-1}(t) & \mathbf{b}(t) & \mathbf{x}_{j+1}(t) & \cdots & \mathbf{x}_n(t) \end{bmatrix}. \end{aligned}$$

Bemerkung: Lautet die Anfangsbedingung $\mathbf{y}(t_0) = \mathbf{y}_0$, so müssen anstelle der Funktionen \mathbf{x}_j die Funktionen $\tilde{\mathbf{x}}_j(t) := \exp((t - t_0)A)\mathbf{u}_j$ verwendet werden.

¹¹ Josef Wronski (1778–1853)

Systeme linearer Differentialgleichungen erster Ordnung

Konstante Koeffizienten, inhomogener Fall

Eine weitere Lösungsdarstellung für inhomogene lineare Systeme mit konstanten Koeffizienten, ebenfalls unter der Bezeichnung **Variation der Konstanten** bekannt, lautet (mit Anfangsbedingung bei $t = t_0$)

$$\mathbf{y}(t) = \exp((t - t_0)A)\mathbf{y}_0 + \int_{t_0}^t \exp((t - \tau)A) \mathbf{b}(\tau) d\tau. \quad (1.1)$$

Beachte: Diese Formel gilt auch wenn \mathbf{b} neben t auch von \mathbf{y} abhängt, d.h.

$$\mathbf{b} = \mathbf{b}(t, \mathbf{y}(t)).$$

Systeme linearer Differentialgleichungen erster Ordnung

Linearisierung

Formel (1.1) gestattet auch die Lösung des **linearisierten Problems**: linearisiert man die Differentialgleichung $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ im Punkt (t_0, \mathbf{y}_0) , ergibt sich (multivariate Taylor-Entwicklung)

$$\mathbf{f}(t, \mathbf{y}) \approx \underbrace{\mathbf{f}(t_0, \mathbf{y}_0)}_{=: \mathbf{b}} + \underbrace{\mathbf{f}_t(t_0, \mathbf{y}_0)}_{=: \mathbf{a}}(t - t_0) + \underbrace{\mathbf{f}_y(t_0, \mathbf{y}_0)}_{=: \mathbf{A}}(\mathbf{y} - \mathbf{y}_0)$$

und, als Approximation in der Nähe des Linearisierungspunktes, das **linearisierte AWP**

$$\mathbf{y}'(t) = \mathbf{A}(\mathbf{y} - \mathbf{y}_0) + (t - t_0)\mathbf{a} + \mathbf{b}, \quad \mathbf{y}(t_0) = \mathbf{y}_0. \quad (1.2)$$

Formel (1.1) liefert als Lösung von (1.2)

$$\mathbf{y}(t) = \mathbf{y}_0 + (t - t_0)\mathbf{A}^{-1}\mathbf{a} + \left(e^{(t-t_0)\mathbf{A}} - \mathbf{I}\right)(\mathbf{A}^{-1}\mathbf{b} + \mathbf{A}^{-2}\mathbf{a}).$$

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Die Fälschungen des Han van Meegeren

Hintergrund

(vgl. [\[Braun, 1994\]](#))

- Im Mai 1945 entdeckten die Alliierten in der Kunstsammlung **Hermann Görings** ein bis dahin unbekanntes Gemälde von **Jan Vermeer van Delft** (1632–1675), nämlich „Christus und die Ehebrecherin“.
- Es dauerte nicht lange, bis der Maler **Han van Meegeren** als derjenige ermittelt wurde, der (über Mittelsmänner) dieses Bild an Göring verkauft hatte. Van Meegeren wurde wegen Kollaboration mit dem Feind verhaftet. Er behauptete daraufhin, dieses Bild sowie vier weitere mutmaßliche Vermeers, darunter „Christus und die Jünger in Emmaus“, selbst gemalt zu haben.
- Um diese Aussage zu bekräftigen, begann er im Gefängnis „Christus unter den Schriftgelehrten“ im Stil Vermeers zu malen. Er ging dabei sehr geschickt vor: Er kratzte von alten, wertlosen Gemälden die Farbe bis auf die Leinwand ab, vermischte die alte (und äußerst harte) Farbe mit Phenolformaldehyd, um mit ihr wieder malen zu können. Das fertige Bild wurde in einem Ofen erhitzt, wobei das Aldehyd zu Bakelit erstarrte.

Die Fälschungen des Han van Meegeren

Hintergrund

- Noch vor Vollendung seiner Arbeit erfuhr van Meegeren, dass die Anklage auf Kollaboration fallen gelassen wurde und er stattdessen ein Verfahren wegen Fälschung zu erwarten hatte. Er weigerte sich daraufhin, die Vermeer-Kopie zu vollenden.
- Weil u.A. in einigen der angeblichen Vermeers-Bilder Phenolformaldehyd nachgewiesen werden konnte (eine Substanz, die bis zum Ende des 19. Jahrhunderts völlig unbekannt war), wurde van Meegeren trotzdem am 12.10.1947 zu einem Jahr Gefängnis wegen Fälschung verurteilt. Er starb kurz darauf in der Haft.
- Dessen ungeachtet waren viele Experten immer noch der Meinung, dass es sich bei „Christus und die Jünger in Emmaus“ um einen echten Vermeer handelt (aufgrund der Expertise eines bekannten Kunsthistorikers erwarb die Rembrandt-Gesellschaft dieses Werk für 174.000 US-\$).
- Der Streit um die Authentizität dieses Gemäldes sollte schließlich 1967 von einer Forschergruppe an der Carnegie Mellon Universität (Pittsburgh, PA) entschieden werden.

Die Fälschungen des Han van Meegeren

Bleiweiß und radioaktiver Zerfall

Deren Analyse basierte auf der Tatsache, dass Künstler seit mehr als 2000 Jahren sog. Bleiweiß (Bleioxyd) verwenden, das kleine Bestandteile an radioaktivem Blei-210 und Radium-226 enthält. Um die Pittsburgher Analyse zu verstehen, sind elementare Kenntnisse über **radioaktiven Zerfall** erforderlich.

- Unter Radioaktivität versteht man den (ohne äußere Beeinflussung erfolgenden) Zerfall instabiler Atomkerne gewisser radioaktiver Substanzen.
- Für jede radioaktive Substanz gibt es eine charakteristische Übergangswahrscheinlichkeit λ (**Zerfallskonstante**), mit der im Mittel ein Atom pro Zeiteinheit zerfällt. Sind zur Zeit t also $N(t)$ radioaktive Atome vorhanden, so zerfallen im Zeitintervall $[t, t + \Delta t]$ durchschnittlich $\lambda N(t) \Delta t$ Atome.
- Für $\Delta t \rightarrow 0$ erhalten wir das **Zerfallsgesetz**

$$N'(t) = -\lambda N(t).$$

Die Fälschungen des Han van Meegeren

Halbwertszeit

Die Zahl der nach einer gewissen Zeit Δt , die seit dem Zeitpunkt t_0 verstrichen ist, noch vorhanden radioaktiven Atome ist deshalb

$$N(t_0 + \Delta t) = N(t_0) \exp(-\lambda \Delta t).$$

Die **Halbwertszeit**, d.h. die Zeitspanne, innerhalb der die Hälfte einer gegebenen Menge radioaktiver Atome zerfällt, ergibt sich damit [setze $N(t_0 + \Delta t)/N(t_0) = 1/2$] zu

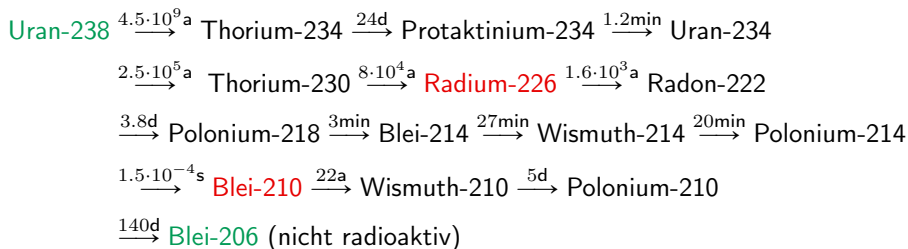
$$T_{1/2} = \Delta t = \log(2)/\lambda.$$

Die Fälschungen des Han van Meegeren

Zerfallsreihe Uran-238

Da sich die Zerfallsprodukte radioaktiver Stoffe weiter umwandeln, bis ein stabiles Endglied gebildet ist, entstehen sog. **Zerfallsreihen**.

Für uns ist die Zerfallsreihe von Uran-238 relevant¹²:

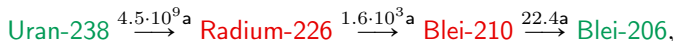


¹²Halbwertszeiten in Jahren [a], Tagen [d], Minuten [min] oder Sekunden [s])

Die Fälschungen des Han van Meegeren

Radioaktives Gleichgewicht

Innerhalb einer Zerfallsreihe stellt sich für die Zwischensubstanzen im Laufe der Zeit ein Gleichgewichtszustand ein, das **radioaktive Gleichgewicht**. Voraussetzung dafür ist, dass das Ausgangselement so langsam zerfällt, dass seine Menge als konstant (oder seine Zerfallskonstante als 0) betrachtet werden kann. Um dies zu verdeutlichen, werden wir die Zerfallsreihe von Uran-238 etwas ökonomisieren,



so dass sie nur noch aus vier Elementen besteht. Wir bezeichnen mit $N_1(t)$, $N_2(t)$, $N_3(t)$ die Anzahl der Atome von Uran-238, Radium-226 bzw. Blei-210 zur Zeit t und berechnen aus den angegebenen Halbwertszeiten die zugehörigen Zerfallskonstanten,

$$\lambda_1 = 1.54 \cdot 10^{-10}, \lambda_2 = 4.33 \cdot 10^{-4}, \lambda_3 = 3.09 \cdot 10^{-2} \text{ (gemessen in } a^{-1}\text{)}.$$

Die Fälschungen des Han van Meegeren

Vereinfachtes Zerfallssystem

Zu lösen ist damit das System

$$N_1'(t) = -\lambda_1 N_1(t),$$

$$N_2'(t) = -\lambda_2 N_2(t) + \lambda_1 N_1(t),$$

$$N_3'(t) = -\lambda_3 N_3(t) + \lambda_2 N_2(t),$$

oder kürzer $\mathbf{N}'(t) = A\mathbf{N}(t)$ (mit Anfangsbedingungen $\mathbf{N}(t_0) = \mathbf{N}_0$), wobei

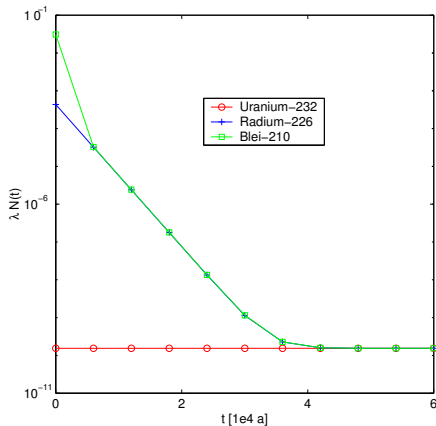
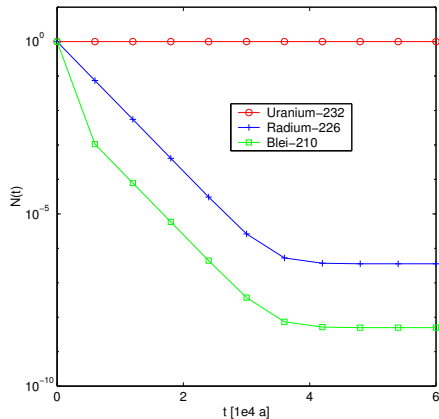
$$A = \begin{bmatrix} -\lambda_1 & 0 & 0 \\ \lambda_1 & -\lambda_2 & 0 \\ 0 & \lambda_2 & -\lambda_3 \end{bmatrix}.$$

Wir können die Lösung dieses AWP's explizit angeben:

$$\mathbf{N}(t) = \exp(tA)\mathbf{N}_0.$$

Die Fälschungen des Han van Meegeren

Radioaktives Gleichgewicht des vereinfachten Zerfallssystems



Die Fälschungen des Han van Meegeren

Radioaktives Gleichgewicht des vereinfachten Zerfallssystems

t [a]	$N_1(t)$	$N_2(t)$	$N_3(t)$
0	$1 \cdot 10^0$	$1 \cdot 10^0$	$1 \cdot 10^0$
10^1	$1 \cdot 10^0$	$1 \cdot 10^0$	$7 \cdot 10^{-1}$
10^3	$1 \cdot 10^0$	$6 \cdot 10^{-1}$	$9 \cdot 10^{-3}$
10^5	$1 \cdot 10^0$	$4 \cdot 10^{-7}$	$9 \cdot 10^{-9}$
10^6	$1 \cdot 10^0$	$4 \cdot 10^{-7}$	$5 \cdot 10^{-9}$
10^7	$1 \cdot 10^0$	$4 \cdot 10^{-7}$	$5 \cdot 10^{-9}$
10^8	$1 \cdot 10^0$	$4 \cdot 10^{-7}$	$5 \cdot 10^{-9}$
10^9	$9 \cdot 10^{-1}$	$3 \cdot 10^{-7}$	$4 \cdot 10^{-9}$
10^{11}	$2 \cdot 10^{-7}$	$7 \cdot 10^{-14}$	$1 \cdot 10^{-15}$
10^{12}	$1 \cdot 10^{-67}$	$5 \cdot 10^{-74}$	$6 \cdot 10^{-76}$

Die Fälschungen des Han van Meegeren

Gleichgewichtswerte und Halbwertszeiten

Natürlich gilt $\lim_{t \rightarrow \infty} N_j(t) = 0$ ($j = 1, 2, 3$) (warum?), aber für eine sehr lange Periode (etwa $10^5 \leq t \leq 10^8$) scheint sich ein Gleichgewicht einzustellen. Die „Gleichgewichtswerte“ sind (ziemlich genau) proportional zu den Halbwertszeiten bzw. umgekehrt proportional zu den Zerfallskonstanten. Für $t = 10^7$ gilt in unserem Beispiel

$$N_1(t)/N_2(t) = 2.812499 \cdots 10^6 = \lambda_2/\lambda_1,$$

$$N_1(t)/N_3(t) = 2.008927847142 \cdots 10^8 = \lambda_3/\lambda_1,$$

$$N_2(t)/N_3(t) = 7.1428571 \cdots 10^1 = \lambda_3/\lambda_2.$$

Um diese Phänomene zu untersuchen, diagonalisieren wir A , $AT = TD$ mit $D = \text{diag}(-\lambda_1, -\lambda_2, -\lambda_3)$ und

$$T = \begin{bmatrix} 1 & 0 & 0 \\ \frac{\lambda_1}{\lambda_2 - \lambda_1} & 1 & 0 \\ \frac{\lambda_1 \lambda_2}{(\lambda_2 - \lambda_1)(\lambda_3 - \lambda_1)} & \frac{\lambda_2}{\lambda_3 - \lambda_2} & 1 \end{bmatrix}.$$

Die Fälschungen des Han van Meegeren

Analyse in der Eigenbasis

Dann ergibt sich

$$\mathbf{N}(t) = \exp(tA)\mathbf{N}_0 = T \exp(tD)T^{-1}\mathbf{N}_0 = T \exp(tD)\tilde{\mathbf{N}}_0,$$

wobei wir $\tilde{\mathbf{N}}_0 = [\tilde{N}_1, \tilde{N}_2, \tilde{N}_3]^T := T^{-1}\mathbf{N}_0$ gesetzt haben. Entscheidend ist das Verhalten der einzigen Größe, die von t abhängt, nämlich von

$$\exp(tD) = \begin{bmatrix} \exp(-\lambda_1 t) & 0 & 0 \\ 0 & \exp(-\lambda_2 t) & 0 \\ 0 & 0 & \exp(-\lambda_3 t) \end{bmatrix}.$$

Für $t \in [10^5, 10^8]$ gelten $\exp(-\lambda_1 t) \in [0.999984 \dots, 1]$ und $\exp(-\lambda_2 t), \exp(-\lambda_3 t) \in [0, 1.5 \dots 10^{-19}]$. In diesem Zeitintervall gilt also

$$\mathbf{N}(t) = T \exp(tD)\tilde{\mathbf{N}}_0 \approx T \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \tilde{\mathbf{N}}_0 = \begin{bmatrix} \tilde{N}_1 \\ \frac{\lambda_1}{\lambda_2 - \lambda_1} \tilde{N}_1 \\ \frac{\lambda_1 \lambda_2}{(\lambda_2 - \lambda_1)(\lambda_3 - \lambda_1)} \tilde{N}_1 \end{bmatrix}.$$

Die Fälschungen des Han van Meegeren

Radioaktives Gleichgewicht in Bleiweiß

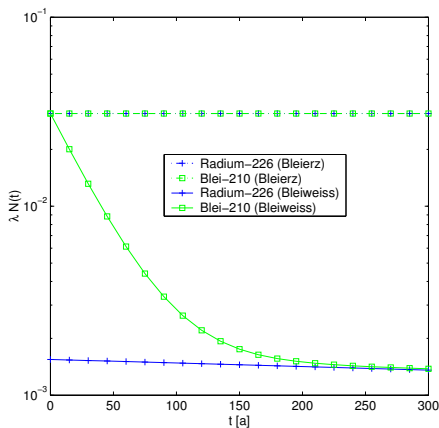
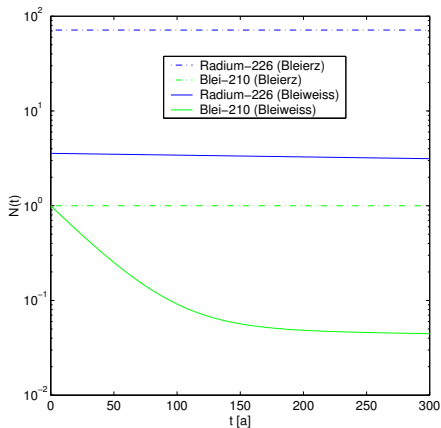
Beachtet man noch $\lambda_1 \ll \lambda_2 < \lambda_3$, so erhält man schließlich für $t \in [10^5, 10^8]$

$$[N_1(t), N_2(t), N_3(t)]^T \approx \tilde{N}_1[1, \lambda_1/\lambda_2, \lambda_1/\lambda_3]^T.$$

Zurück zur Altersbestimmung von Gemälden: Wie bereits erwähnt, enthalten fast alle Gemälde Bleiweiß und damit die radioaktiven Substanzen Radium-226 und Blei-210. Bleiweiß wird aus Blei gewonnen, welches wiederum durch Schmelzen von Bleierz entsteht. Bei diesem Schmelzvorgang werden 90–95% des Radiums und seiner Tochtersubstanzen mit der Schlacke entfernt, so dass Blei-210 von seinem Nachschub abgeschnitten ist und sich mit Radium-226 nicht mehr im radioaktiven Gleichgewicht befindet. Das Blei-210 zerfällt dann sehr schnell (da es eine kurze Halbwertszeit von $T_{1/2} = 22\text{a}$ besitzt), bis es mit den Resten von Radium-226 wieder im Gleichgewicht ist (nach ca. 200 Jahren).

Die Fälschungen des Han van Meegeren

Bleierz/Bleiweiß mit und ohne radioaktives Gleichgewicht



Die Fälschungen des Han van Meegeren

Hintergrund

Seien jetzt t_0 der Zeitpunkt, an dem das Bleiweiß hergestellt wurde, und (wie oben) $N_2(t)$, $N_3(t)$ die Mengen von Radium-226 bzw. Blei-210 (pro g Bleiweiß). Es gilt $N_3'(t) = -\lambda_3 N_3(t) + \lambda_2 N_2(t)$. Da wir uns nur für eine Zeitspanne von 300 Jahren interessieren und Radium-226 eine Halbwertszeit von $T_{1/2} = 1600$ Jahren besitzt, können wir annehmen, dass seine **Zerfallsrate (Aktivität)** $\rho_2 := \lambda_2 N_2(t)$ konstant ist. Die GDG vereinfacht sich zu $N_3'(t) = -\lambda_3 N_3(t) + \rho_2$ bzw. zu $N_3'(t) + \lambda_3 N_3(t) = \rho_2$. Multiplizieren wir mit dem **integrierenden Faktor** $\exp(\lambda_3 t)$, so ergibt sich

$$\frac{d}{dt} [\exp(\lambda_3 t) N_3(t)] = \rho_2 \exp(\lambda_3 t),$$

was zu

$$N_3(t) = \frac{\rho_2}{\lambda_3} [1 - \exp(-\lambda_3(t - t_0))] + N_3(t_0) \exp(-\lambda_3(t - t_0)) \quad (*)$$

führt. $N_3(t)$, λ_3 und ρ_2 sind bekannt (bzw. leicht zu messen). Wüssten wir die Größe von $N_3(t_0)$, könnten wir $t - t_0$ (und damit das Alter des Gemäldes) bestimmen.

Die Fälschungen des Han van Meegeren

Hintergrund

Natürlich ist es unmöglich, $N_3(t_0)$ ohne Kenntnis von t_0 zu ermitteln. Wir machen von der Tatsache Gebrauch, dass $N_3(t_0)$, also die Menge an Blei-210 zum Zeitpunkt der Herstellung des Bleiweißes, ein radioaktives Gleichgewicht mit dem Radium-226 im Bleierz bildete. Lösen wir also (*) nach der Zerfallsrate $\lambda_3 N_3(t_0)$ von Blei-210 zur Zeit t_0 auf,

$$\lambda_3 N_3(t_0) = \lambda_3 N_3(t) \exp(\lambda_3(t - t_0)) - \rho_2 [\exp(\lambda_3(t - t_0)) - 1],$$

und nehmen $t - t_0 = 300$ a an,

$$\begin{aligned}\lambda_3 N_3(t_0) &= \lambda_3 N(t) \exp(300\lambda_3) - \rho_2 [\exp(300\lambda_3) - 1] \\ &= \lambda_3 N(t) 2^{150/11} - \rho_2 [2^{150/11} - 1]\end{aligned}$$

$[\exp(300\lambda_3) = \exp(300 \log(2)/T_{1,2}) = \exp(300 \log(2)/22) = 2^{150/11}]$. Um $\lambda_3 N_3(t_0)$ zu berechnen, müssen wir die gegenwärtigen Zerfallsraten $\lambda_3 N_3(t)$ von Blei-210 bzw. ρ_2 von Radium-226 bestimmen, was für einige mutmaßliche Bilder Vermeers geschehen ist.

Die Fälschungen des Han van Meegeren

Hintergrund

	Zerfallsraten* von	
	Pb-210	Ra-226
„Christus und die Jünger in Emmaus“	8.5	0.80
„Fußwaschung“	12.6	0.26
„Die Notenleserin“	10.3	0.30
„Die Mandolinspielerin“	8.2	0.17
„Die Spitzenklöpplerin“	1.5	1.40
„Der Soldat und das lachende Mädchen“	5.2	6.00

(* pro Minute und pro g Bleiweiß)

Die Fälschungen des Han van Meegeren

Hintergrund



Die Fälschungen des Han van Meegeren

Hintergrund

Legende:

- (1,1) Han van Meegeren, „Christus und die Ehebrecherin“ (1941), ???
- (1,2) Han van Meegeren, „Christus und die Jünger in Emmaus“ (1936/37), Museum Boymans Van Beuningen, Rotterdam
- (1,3) Han van Meegeren, „Fußwaschung Christi“ (1941), Rijksmuseum, Amsterdam
- (1,4) Han van Meegeren, „Die Notenleserin“ (1935/36), Rijksmuseum, Amsterdam
- (2,1) Han van Meegeren, „Die Mandolinenspielerin“ (1935/36), Rijksmuseum, Amsterdam
- (2,2) Jan Vermeer, „Die Spitzenklöpplerin“ (ca. 1669/70), Louvre, Paris
- (2,3) Jan Vermeer, „Der Soldat und das lachende Mädchen“ (ca. 1658), Frick Collection, New York
- (2,4) Jan Vermeer, „Brieflesendes Mädchen am offenen Fenster“ (ca. 1659), Gemäldegalerie „Alte Meister“, Dresden

Quellen: <http://www.cacr.caltech.edu/roy/vermeer>, <http://www.mystudios.com/gallery/han>

Die Fälschungen des Han van Meegeren

Hintergrund

Für „Christus und die Jünger in Emmaus“ ergibt sich

$$\lambda_3 N_3(t_0) = 8.5 \cdot 2^{150/11} - 0.8 \left[2^{150/11} - 1 \right] \approx 98050.$$

Es bleibt die Frage, ob dies ein akzeptabler Wert für die Zerfallsrate von Blei-210 im radioaktiven Gleichgewicht ist. Man kann nachrechnen, dass, wenn das Blei zur Zeit der Gewinnung mit einer Zerfallsrate von 100 pro Minute und g Bleiweiß zerfiel, das Erz, aus dem es stammt, einen Urananteil von 0.014 % hatte. Dies ist eine sehr hohe Urankonzentration. Andererseits gibt es (seltene) Erze, deren Urangehalt bei 2–3 % liegt. Um sicher zu gehen, nennen wir $\lambda_3 N_3(t_0)$ deshalb unakzeptabel hoch, wenn

$$\lambda_3 N_3(t_0) > 100 \cdot 3 / 0.014 \approx 22000$$

gilt, was bei „Christus und die Jünger in Emmaus“ offenbar der Fall ist (zum Vergleich beträgt der entsprechende Wert bei der „Spitzenklöpplerin“ $\lambda_3 N_3(t_0) \approx 1275$).

① Einleitung

- 1.1 Volterras Prinzip
- 1.2 Begriffe und theoretische Resultate
- 1.3 Lineare Differenzengleichungen
- 1.4 Matrixfunktionen
- 1.5 Systeme linearer Differentialgleichungen erster Ordnung
- 1.6 Die Fälschungen des Han van Meegeren
- 1.7 Weitere Beispiele

② Numerische Methoden für Anfangswertprobleme

③ Lineare Mehrschrittverfahren

④ Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

Weitere Beispiele

Satellit im Kraftfeld von Erde und Mond

Wir betrachten die Bewegung eines Satelliten im Schwerfeld zweier großer Himmelskörper (z.B. Erde und Mond).

Annahmen:

- Die Bewegung aller drei Körper findet in einer Ebene statt; die beiden großen Körper rotieren in konstanter Entfernung und mit konstanter Winkelgeschwindigkeit um ihren gemeinsamen Schwerpunkt.
- Der Satellit hat somit keinen Einfluss auf die Bahnen von Erde und Mond.

Bezüglich eines mitrotierenden Koordinatensystems (in welchem Erde und Mond ruhen) mit Ursprung im gemeinsamen Schwerpunkt wird die Satellitenbahn $(x, y) = (x(t), y(t))$ beschrieben durch ein System zweier GDGen:

$$\begin{aligned}x'' &= x + 2y' - \mu' \frac{x + \mu}{[(x + \mu)^2 + y^2]^{3/2}} - \mu \frac{x - \mu'}{[(x - \mu')^2 + y^2]^{3/2}}, \\y'' &= y - 2x' - \mu' \frac{y}{[(x + \mu)^2 + y^2]^{3/2}} - \mu \frac{y}{[(x - \mu')^2 + y^2]^{3/2}}.\end{aligned}$$

Weitere Beispiele

Satellit im Kraftfeld von Erde und Mond

- $\mu = 1/82.45$ bezeichnet den Anteil der Mondmasse an der Gesamtmasse von Erde und Mond, $\mu' = 1 - \mu$ die der Erde.
- Als Längeneinheit wählen wir die Erde-Mond-Entfernung, wobei der Mond auf der positiven und die Erde auf der negativen reellen Achse platziert werden.
- Die Zeiteinheit ist gegeben durch die Winkelgeschwindigkeit der Rotation, genauer rotiert der Mond ein Mal um die Erde in 2π Zeiteinheiten.
- Bekannt als **restringiertes Dreikörperproblem**, da der dritte Körper die ersten beiden nicht beeinflusst.
- Anfangsbedingung zur Zeit $t = 0$:

Satellit in Position $(x(0), y(0)) = (1.2, 0)$

mit Geschwindigkeit $(x'(0), y'(0)) = (0, -1.05)$.

Umschreiben in System erster Ordnung:

$$y_1 = x, \quad y_2 = y, \quad y_3 = x', \quad y_4 = y',$$

führt auf

$$y_1' = y_3,$$

$$y_2' = y_4,$$

$$y_3' = y_1 + 2y_4 - \mu' \frac{y_1 + \mu}{[(y_1 + \mu)^2 + y_2^2]^{3/2}} - \mu \frac{y_1 - \mu'}{[(y_1 - \mu')^2 + y_2^2]^{3/2}},$$

$$y_4' = y_2 - 2y_3 - \mu' \frac{y_2}{[(y_1 + \mu)^2 + y_2^2]^{3/2}} - \mu \frac{y_2}{[(y_1 - \mu')^2 + y_2^2]^{3/2}},$$

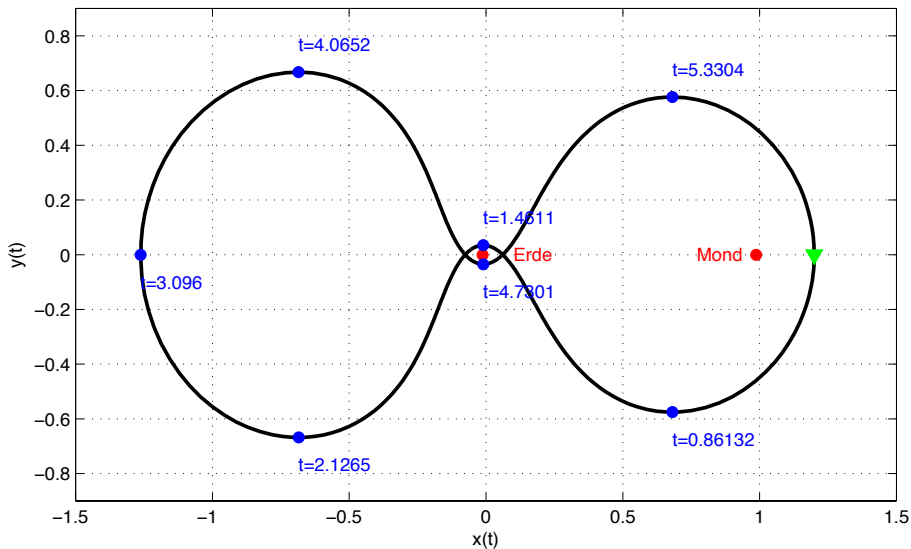
mit Anfangsbedingungen

$$y_1(0) = 1.2, \quad y_2(0) = y_3(0) = 0 \quad \text{and} \quad y_4(0) = -1.05.$$

Lösung $(x(t), y(t)) = (y_1(t), y_2(t))$: geschlossene Bahn mit Periode $T \approx 6.19$.

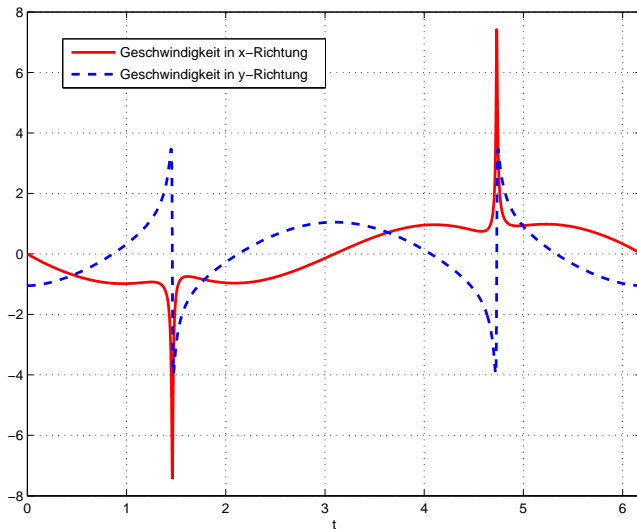
Weitere Beispiele

Satellit im Kraftfeld von Erde und Mond



Weitere Beispiele

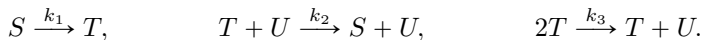
Satellit im Kraftfeld von Erde und Mond



Weitere Beispiele

Kinetik chemischer Reaktionsmechanismen

Drei Spezies S , T und U nehmen Teil an der autokatalytischen Reaktion



Zeitlicher Verlauf der Konzentrationen $y_1 = [S]$, $y_2 = [T]$, $y_3 = [U]$ beschrieben durch System GDGen (Massenwirkungsgesetz)

$$\begin{aligned} y_1' &= -k_1 y_1 + k_2 y_2 y_3, \\ y_2' &= k_1 y_1 - k_2 y_2 y_3 - 2k_3 y_2^2 + k_3 y_2^2 = k_1 y_1 - k_2 y_2 y_3 - k_3 y_2^2, \\ y_3' &= k_3 y_2^2. \end{aligned}$$

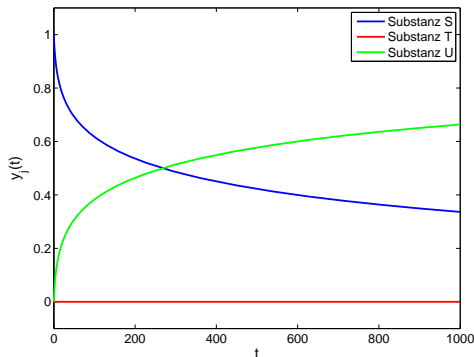
Reaktionsraten k_j sind ein Maß für die Geschwindigkeit mit denen die jeweilige Reaktion sich vollzieht. Sie differieren oft um mehrere Größenordnungen.

$$k_1 = 0.04, \quad k_2 = 10^4, \quad k_3 = 3 \cdot 10^7.$$

Weitere Beispiele

Kinetik chemischer Reaktionsmechanismen

Für die Anfangsbedingungen $y_1(0) = 1; y_2(0) = y_3(0) = 0$ erhalten wir:



Beachte: wegen $y_1'(t) + y_2'(t) + y_3'(t) = 0$ gilt für alle $t \geq t_0 = 0$, die **Erhaltungsgleichung**

$$y_1(t) + y_2(t) + y_3(t) = y_1(0) + y_2(0) + y_3(0) = 1.$$

Die Stabilitätsanalyse nichtlinearer dynamischer Systeme

$$\dot{\mathbf{u}} = \mathbf{f}(\mathbf{u}), \quad \mathbf{u}(0) = \mathbf{u}_0, \quad \mathbf{f} : \mathbb{C}^n \rightarrow \mathbb{C}^n, \quad (1.3)$$

geschieht meist durch Linearisierung um den (einen) stationären Zustand

$$\bar{\mathbf{u}} := \lim_{t \rightarrow \infty} \mathbf{u}(t).$$

Sofern ein solcher existiert ist er Lösung der Gleichung $\mathbf{f}(\mathbf{u}) = \mathbf{0}$.

Das dynamische System (1.3) heißt **lokal stabil** in $\bar{\mathbf{u}}$, falls $\epsilon > 0$ existiert mit

$$\lim_{t \rightarrow \infty} \mathbf{u}(t) = \bar{\mathbf{u}}, \quad \text{sofern } \|\mathbf{u}(0) - \bar{\mathbf{u}}\| < \epsilon.$$

In vielen Fällen lässt sich die Frage nach der lokalen Stabilität von (1.3) durch Analyse der Linearisierung

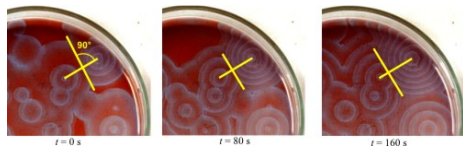
$$\mathbf{f}(\mathbf{u}) \approx \mathbf{f}(\bar{\mathbf{u}}) + \mathbf{A}(\mathbf{u} - \bar{\mathbf{u}}), \quad \mathbf{A} := \mathbf{f}'(\bar{\mathbf{u}})$$

klären, d.h. durch die Realteile der Eigenwerte von \mathbf{A} .

Weitere Beispiele

Chemische Reaktionskinetik: der Brusselator

Die Belousov-Zhabotinsky-Reaktion ist ein Beispiel für einen sog. **chemischen Oszillator**, bei dem sich zeitliche Oszillationen in einem chemischen Reaktionsmechanismus zeigen.



BZ-Reaktion in einer Petri-Schale,
Wellenfront in gelb markiert.

Ein mathematisches Modell der BZ-Reaktion ist der sog. **Brusselator**¹³, einer Evolutionsgleichung der örtlichen Variation (in einer Raumkoordinate $r \in (0, 1)$) der Konzentration zweier miteinander reagierender Spezies x und y :

$$\partial_t x = D_1 \partial_{rr} x + A - (B + 1)x + x^2 y, \quad x(0, t) = x(1, t) = A.$$

$$\partial_t y = D_2 \partial_{rr} y + Bx - x^2 y, \quad y(0, t) = y(1, t) = \frac{B}{A},$$

$$x(r, 0) = x_0(r), \quad y(r, 0) = y_0(r).$$

¹³Brussels + Oszillator, Ilya Prigogine FU Brüssel

Weitere Beispiele

Chemische Reaktionskinetik: der Brusselator

Ein stationärer Zustand ist gegeben durch

$$\bar{x} = A, \quad \bar{y} = \frac{B}{A}.$$

Die Jacobi-Matrix an dieser Stelle ist gegeben durch

$$\mathbf{J} = \begin{bmatrix} D_1 \partial_{rr} + B - 1 & A^2 \\ -B & D_2 \partial_{rr} - A^2 \end{bmatrix},$$

was nach Ortsdiskretisierung eine große, dünn besetzte Matrix ergibt.

Bifurkationsproblem: Ab welchem Wert von B setzt periodisches Verhalten ein?
Hierbei überqueren zwei konjugiert-komplexe Eigenwerte von \mathbf{J} die imaginäre Achse.

(vgl. [\[Hairer, Norsett & Wanner, 1987, Abschnitt I.16\]](#))

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren

Wir betrachten das AWP

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}), \quad \mathbf{y}(t_0) = \mathbf{y}_0. \quad (\text{AWP})$$

Unter den Voraussetzungen von Satz 1.1 besitzt es eine eindeutige Lösung, sagen wir über dem Intervall I .

Wir wollen diese Lösung $\mathbf{y}(t)$ für $t \in [t_0, t_{\text{end}}] \subseteq I$ durch das **Euler-Verfahren**¹⁴, auch **Euler-Cauchy-Verfahren**¹⁵ oder **Polygonzug-Verfahren**, den Prototyp eines numerischen Verfahrens zur Lösung von AWPen, approximieren. Wie alle numerischen Methoden, die hier besprochen werden, basiert das Euler-Verfahren auf der Idee der **Diskretisierung**: Statt (AWP) über $[t_0, t_{\text{end}}]$ zu lösen, geben wir uns damit zufrieden, Näherungswerte für die Lösung auf einer diskreten Teilmenge

$$\{t_n : n = 0, 1, \dots, N\} \subset [t_0, t_{\text{end}}]$$

zu bestimmen.

¹⁴Leonhard Euler (1707–1783)

¹⁵Augustin Louis Cauchy (1789–1857)

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren

Wähle $N \in \mathbb{N}$, $h = (t_{\text{end}} - t_0)/N$ und definiere

$$t_n := t_0 + nh, \quad n = 0, 1, \dots, N,$$

d.h. $t_0 < t_1 < \dots < t_{n-1} < t_N = t_{\text{end}}$. Die Zahl $h > 0$ heißt **Schrittweite**. (Nur der Bequemlichkeit wählen wir die Schrittweite h zunächst **konstant**.)

Bezeichnet nun \mathbf{y}_n einen Näherungswert für $\mathbf{y}(t_n)$, $n = 0, 1, \dots, N-1$, dann ist (falls $\mathbf{y} \in C^2[t_0, t_{\text{end}}]$)

$$\begin{aligned} \mathbf{y}(t_{n+1}) &= \mathbf{y}(t_n + h) = \mathbf{y}(t_n) + h\mathbf{y}'(t_n) + \frac{1}{2}h^2\mathbf{y}''(\xi) \\ &\approx \mathbf{y}(t_n) + h\mathbf{y}'(t_n) = \mathbf{y}(t_n) + h\mathbf{f}(t_n, \mathbf{y}(t_n)) \approx \mathbf{y}_n + h\mathbf{f}(t_n, \mathbf{y}_n). \end{aligned}$$

Euler-Verfahren

\mathbf{y}_0 gegeben,

$$\mathbf{y}_{n+1} := \mathbf{y}_n + h\mathbf{f}(t_n, \mathbf{y}_n), \quad n = 0, 1, \dots, N-1.$$

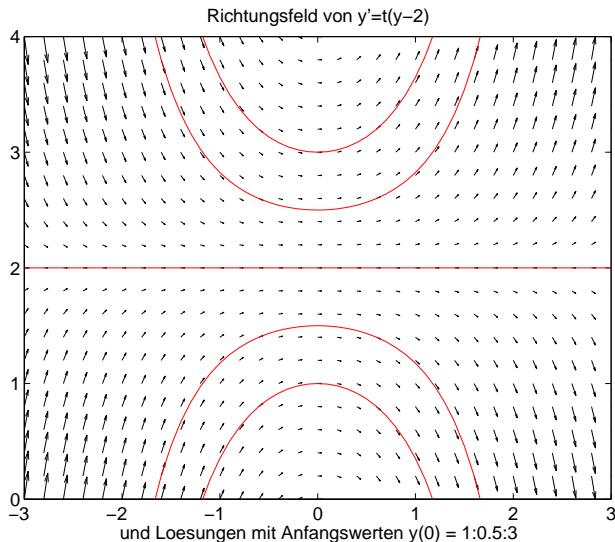
Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren

- Zur Veranschaulichung einer GDG 1. Ordnung $y' = f(t, y)$ wird oft das assoziierte **Richtungsfeld** herangezogen: In jedem Punkt (t, y) wird ein Pfeil gezeichnet, der von dort aus in Richtung der Steigung $y' = f(t, y)$ weist.
- Der Graph der Lösung des AWP's $y' = f(t, y), y(t_0) = y_0$, muss einerseits das Richtungsfeld respektieren (d.h. die Tangenten an den Graphen sind Elemente des Richtungsfelds), andererseits den Punkt (t_0, y_0) enthalten.

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren



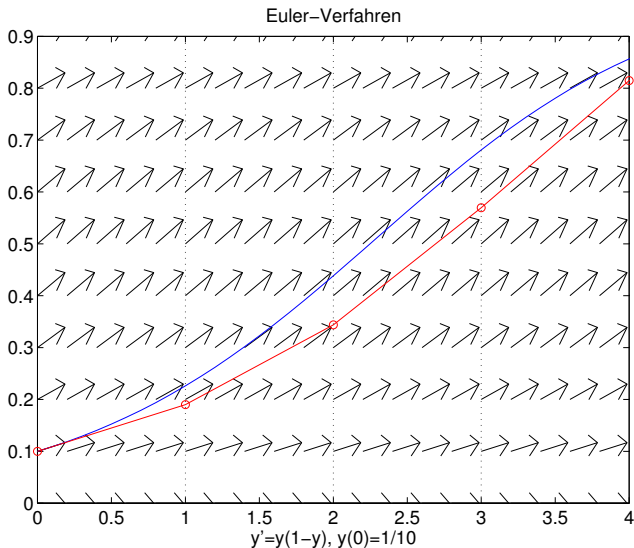
Die nächste Abbildung zeigt, wie das Euler-Verfahren die **logistische Gleichung**

$$y' = y(1 - y), \quad y(0) = \frac{1}{10},$$

löst. Statt der exakten Trajektorie zu folgen (was natürlich unmöglich ist), produziert das Euler-Verfahren eine „stückweise lineare Lösung“ (einen Polygonzug). An der Stelle $t_0 = 0$ arbeitet das Euler-Verfahren mit der richtigen Steigung $f(t_0, y_0) = 9/100$, bereits an der Stelle $t = 1$ ist die Steigung „falsch“. In späteren Schritten entfernt man sich (potentiell) immer weiter von der exakten Lösung.

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren



Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren

Konvergiert das Euler-Verfahren, d.h. strebt die Näherungslösung für $h \rightarrow 0$ gegen die exakte Lösung $\mathbf{y}(t)$?

Formal: Zu jedem Wert von $h > 0$ gehört eine Folge von Näherungen

$$\mathbf{y}_n = \mathbf{y}_n(h), \quad n = 0, 1, \dots, N(h) := \text{floor}((t_{\text{end}} - t_0)/h).$$

Das Verfahren heißt **konvergent** (in $[t_0, t_{\text{end}}]$), wenn gilt

$$\lim_{h \rightarrow 0+} \max_{0 \leq n \leq N(h)} \|\mathbf{y}_n(h) - \mathbf{y}(t_n)\| = 0. \quad (\text{Konv})$$

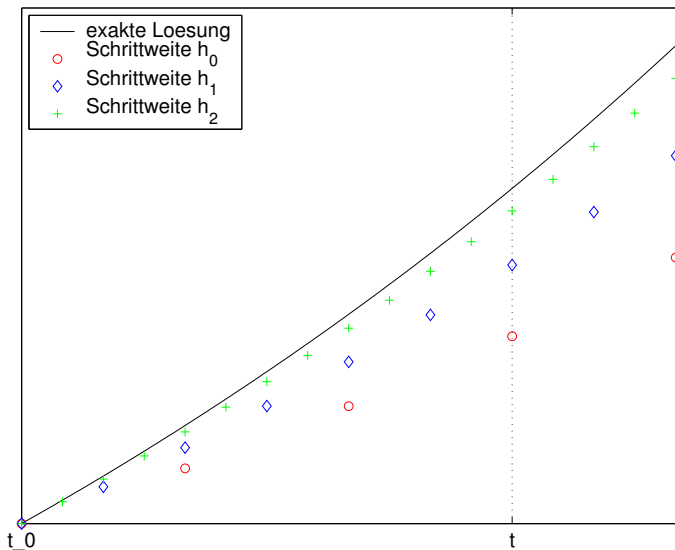
Satz 2.1

Unter den Voraussetzungen von Satz 1.1 konvergiert das Euler-Verfahren.
Genauer:

$$\max_{0 \leq n \leq N(h)} \|\mathbf{y}_n(h) - \mathbf{y}(t_n)\| = O(h) \quad \text{für } h \rightarrow 0.$$

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren



Numerische Methoden für Anfangswertprobleme

Modifiziertes und verbessertes Euler-Verfahren

Die Idee dieser beiden Verfahren macht man sich leicht im Richtungsfeld klar.

Modifiziertes Euler-Verfahren

y_0 gegeben,

$$y_{n+1} := y_n + h\bar{f}\left(t_n + \frac{1}{2}h, y_n + \frac{1}{2}hf(t_n, y_n)\right) \quad (n = 0, 1, \dots, N-1).$$

Verbessertes Euler-Verfahren (Verfahren von Heun)

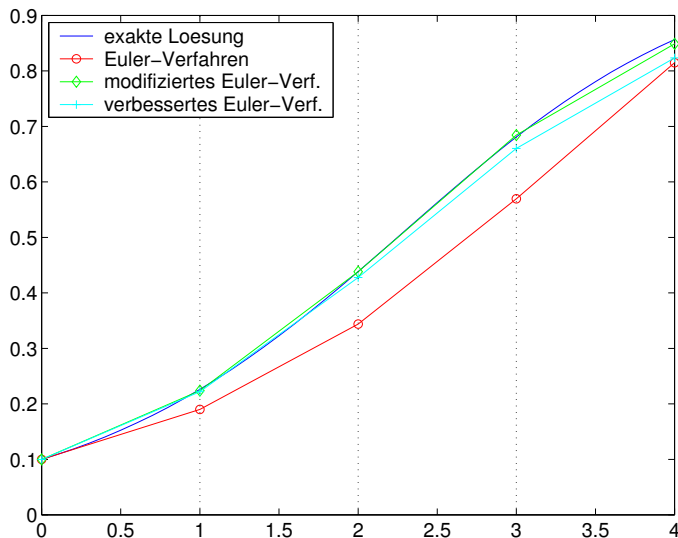
y_0 gegeben,

$$y_{n+1} := y_n + \frac{1}{2}h[f(t_n, y_n) + f(t_n + h, y_n + hf(t_n, y_n))] \\ (n = 0, 1, \dots, N-1).$$

Auch das modifizierte und das verbesserte Euler-Verfahren konvergieren: Eine triviale Modifikation des Beweises von Satz 2.1 zeigt, dass für diese Verfahren sogar $\max_{0 \leq n \leq N(h)} \|y_n(h) - y(t_n)\| = O(h^2)$ gilt.

Numerische Methoden für Anfangswertprobleme

Das Euler-Verfahren



① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Numerische Methoden für Anfangswertprobleme

Eine Sammlung von Beispielverfahren

Wir werden nun eine Reihe numerischer Verfahren zur Lösung von AWPen angeben (nicht alle davon sind für einen praktischen Einsatz geeignet), um später auf eine gewisse Anzahl von Beispielen zurückgreifen zu können.

Beispiel 1

$$\begin{aligned} \mathbf{y}_{n+1} - \mathbf{y}_n &= \frac{1}{4}h(\mathbf{k}_1 + 3\mathbf{k}_3) & \text{mit } \mathbf{k}_1 &= \mathbf{f}(t_n, \mathbf{y}_n), \\ & & \mathbf{k}_2 &= \mathbf{f}(t_n + \frac{1}{3}h, \mathbf{y}_n + \frac{1}{3}h\mathbf{k}_1), \\ & & \mathbf{k}_3 &= \mathbf{f}(t_n + \frac{2}{3}h, \mathbf{y}_n + \frac{2}{3}h\mathbf{k}_2). \end{aligned} \quad (2.1)$$

- Dieses Verfahren ist ein **Einschrittverfahren** (zur Berechnung von \mathbf{y}_{n+1} ist nur die Kenntnis von \mathbf{y}_n erforderlich).
- Es ist **explizit**, d.h. nach \mathbf{y}_{n+1} aufgelöst.
- Es gehört zur Klasse der **Runge-Kutta-Verfahren**¹⁶.

¹⁶Carl Runge (1856–1927), Martin Kutta (1867–1944)

Numerische Methoden für Anfangswertprobleme

Eine Sammlung von Beispielverfahren

Beispiel 2

$$\begin{aligned} \mathbf{y}_{n+1} - \mathbf{y}_n &= \frac{1}{2}h(\mathbf{k}_1 + \mathbf{k}_2) \quad \text{mit } \mathbf{k}_1 = \mathbf{f}(t_n, \mathbf{y}_n), \\ \mathbf{k}_2 &= \mathbf{f}(t_n + h, \mathbf{y}_n + \frac{1}{2}h\mathbf{k}_1 + \frac{1}{2}h\mathbf{k}_2) \end{aligned} \quad (2.2)$$

ist ebenfalls ein Einschrittverfahren der Runge-Kutta-Klasse. Im Gegensatz zu (2.1) ist es aber **implizit**, d.h. um \mathbf{y}_{n+1} zu bestimmen, muss ein (im Allgemeinen nicht-lineares) Gleichungssystem gelöst werden.

Beispiel 3

$$\mathbf{y}_{n+2} - \mathbf{y}_{n+1} = \frac{1}{3}h[3\mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}) - 2\mathbf{f}(t_n, \mathbf{y}_n)] \quad (2.3)$$

ist ein explizites **Zweischrittverfahren**: Man benötigt \mathbf{y}_n und \mathbf{y}_{n+1} , um \mathbf{y}_{n+2} zu bestimmen. Es gehört zur Familie der **linearen Mehrschrittverfahren**. Neben \mathbf{y}_0 ist hier zusätzlich \mathbf{y}_1 (ein sog. **Anlaufstück**) erforderlich, was man sich üblicherweise mit Hilfe eines Einschrittverfahrens beschafft.

Numerische Methoden für Anfangswertprobleme

Eine Sammlung von Beispielverfahren

Auch

Beispiel 4

$$\mathbf{y}_{n+2} + \mathbf{y}_{n+1} - 2\mathbf{y}_n = \frac{1}{4}h [\mathbf{f}(t_{n+2}, \mathbf{y}_{n+2}) + 8\mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}) + 3\mathbf{f}(t_n, \mathbf{y}_n)] \quad (2.4)$$

ist ein lineares Mehr- (genauer: Zwei-) Schrittverfahren. Es ist implizit.

Beispiel 5

$$\mathbf{y}_{n+3} + \frac{1}{4}\mathbf{y}_{n+2} - \frac{1}{2}\mathbf{y}_{n+1} - \frac{3}{4}\mathbf{y}_n = \frac{1}{8}h[19\mathbf{f}(t_{n+2}, \mathbf{y}_{n+2}) + 5\mathbf{f}(t_n, \mathbf{y}_n)] \quad (2.5)$$

stellt ein explizites lineares Dreischrittverfahren dar.

Beispiel 6

$$\begin{aligned} \mathbf{y}_{n+2} - \mathbf{y}_n &= h[\mathbf{f}(t_{n+2}, \mathbf{y}_{n+2}^*) + \mathbf{f}(t_n, \mathbf{y}_n)] \\ &\text{mit} \\ \mathbf{y}_{n+2}^* - 3\mathbf{y}_{n+1} + 2\mathbf{y}_n &= \frac{1}{2}h[\mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}) - 3\mathbf{f}(t_n, \mathbf{y}_n)] \end{aligned} \tag{2.6}$$

ist ein **Prädiktor-Korrektor-Verfahren**, bei dem ein implizites Zweischnittverfahren (der Korrektor) mit einem expliziten Zweischnittverfahren (dem Prädiktor) kombiniert wird. Das zusammengesetzte Verfahren ist explizit.

Numerische Methoden für Anfangswertprobleme

Eine Sammlung von Beispielverfahren

Alle Verfahrensbeispiele haben die Struktur

$$\sum_{j=0}^k \alpha_j \mathbf{y}_{n+j} = h \Phi_f(\mathbf{y}_{n+k}, \mathbf{y}_{n+k-1}, \dots, \mathbf{y}_n, t_n; h) \quad (\text{V})$$

(wir normieren im Folgenden durch $\alpha_k := 1$) mit

$$\Phi_{f \equiv 0}(\mathbf{y}_{n+k}, \mathbf{y}_{n+k-1}, \dots, \mathbf{y}_n, t_n; h) \equiv \mathbf{0} \quad (\text{V}_1)$$

und

$$\|\Phi_f(\mathbf{y}_{n+k}, \dots, \mathbf{y}_n, t_n; h) - \Phi_f(\mathbf{y}_{n+k}^*, \dots, \mathbf{y}_n^*, t_n; h)\| \leq M \sum_{j=0}^k \|\mathbf{y}_{n+j} - \mathbf{y}_{n+j}^*\|. \quad (\text{V}_2)$$

Die Eigenschaft (V₂) ist eine Folge der Lipschitz-Stetigkeit von f (vgl. Satz 1.1), die immer vorausgesetzt wird.

Wir werden ausschließlich numerische Verfahren untersuchen, die die Struktur (V) besitzen und dabei den Bedingungen (V₁) und (V₂) genügen.

① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Definition 2.2 (Konvergenz)

Das Verfahren (V) heißt **konvergent**, wenn

$$\lim_{\substack{h \rightarrow 0 \\ t=t_0+nh}} \mathbf{y}_n = \lim_{\substack{h \rightarrow 0 \\ t=t_0+nh}} \mathbf{y}_n(h) = \mathbf{y}(t)$$

gilt, und zwar

- für alle AWPes, die den Voraussetzungen von Satz 1.1 genügen ($\mathbf{y}(t)$ bezeichnet die Lösung eines solchen AWPes),
- gleichmäßig für alle $t \in [t_0, t_{\text{end}}]$,
- für alle Lösungen $\{\mathbf{y}_n(h)\} = \{\mathbf{y}_n\}$ von (V) mit Anfangswerten $\mathbf{y}_0(h), \dots, \mathbf{y}_{k-1}(h)$, die $\lim_{h \rightarrow 0} \mathbf{y}_j(h) = \mathbf{y}_0$, $j = 0, \dots, k-1$, erfüllen.

Äquivalent:

$$\lim_{h \rightarrow 0} \max_{0 \leq n \leq N} \|\mathbf{y}(t_n) - \mathbf{y}_n(h)\| = 0.$$

Numerische Methoden für Anfangswertprobleme

Lokaler Diskretisierungsfehler, Residuum

Setzt man die exakte Lösung in (V) ein, so werden linke und rechte Seite nicht übereinstimmen. Es ergibt sich ein **Residuum**

$$\mathbf{R}_{n+k} := \sum_{j=0}^k \alpha_j \mathbf{y}(t_{n+j}) - h\Phi_f(\mathbf{y}(t_{n+k}), \mathbf{y}(t_{n+k-1}), \dots, \mathbf{y}(t_n), t_n; h).$$

\mathbf{R}_{n+k} ist eng verknüpft mit dem **lokalen Diskretisierungsfehler** \mathbf{T}_{n+k} .

Unter der **Lokalisierungsannahme**

$$\mathbf{y}_{n+j} = \mathbf{y}(t_{n+j}) \quad \text{für } j = 0, 1, \dots, k-1$$

liefert (V) in Schritt $n+k$:

$$\hat{\mathbf{y}}_{n+k} + \sum_{j=0}^{k-1} \alpha_j \mathbf{y}(t_{n+j}) = h\Phi_f(\hat{\mathbf{y}}_{n+k}, \mathbf{y}(t_{n+k-1}), \dots, \mathbf{y}(t_n), t_n; h).$$

Für die exakte Lösung gilt hingegen

$$\mathbf{y}(t_{n+k}) + \sum_{j=0}^{k-1} \alpha_j \mathbf{y}(t_{n+j}) = h\Phi_f(\mathbf{y}(t_{n+k}), \dots, \mathbf{y}(t_n), t_n; h) + \mathbf{R}_{n+k}.$$

Definiert man nun

$$\mathbf{T}_{n+k} := \mathbf{y}(t_{n+k}) - \hat{\mathbf{y}}_{n+k}$$

(**Vorsicht:** die Definition ist nicht einheitlich in der Literatur), so folgt

$$\begin{aligned} \mathbf{T}_{n+k} = & h [\Phi_f(\mathbf{y}(t_{n+k}), \dots, \mathbf{y}(t_n), t_n; h) \\ & - \Phi_f(\hat{\mathbf{y}}_{n+k}, \mathbf{y}(t_{n+k-1}), \dots, \mathbf{y}(t_n), t_n; h)] + \mathbf{R}_{n+k}. \end{aligned}$$

In speziellen Fällen (lineare Mehrschrittverfahren) kann man die rechte Seite mit dem Mittelwertsatz weiter bearbeiten. Immer folgt aus (V₂)

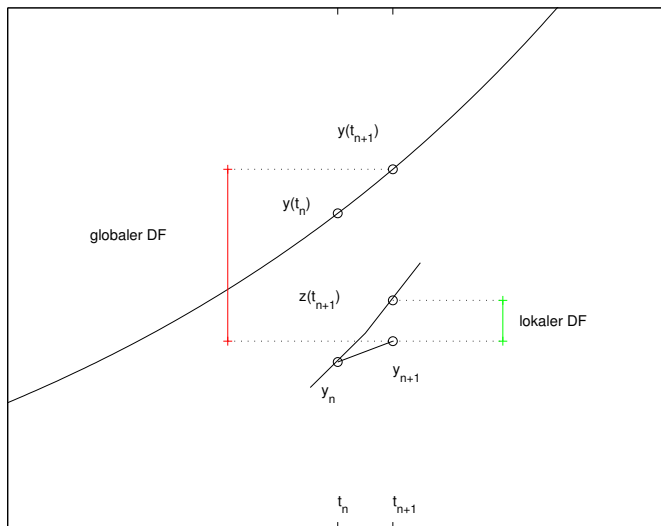
$$\|\mathbf{T}_{n+k}\| \leq hM \|\mathbf{T}_{n+k}\| + \|\mathbf{R}_{n+k}\|$$

und damit

$$(1 - hM) \|\mathbf{T}_{n+k}\| \leq \|\mathbf{R}_{n+k}\|.$$

Numerische Methoden für Anfangswertprobleme

Lokaler Diskretisierungsfehler, Residuum



Definition 2.3 (Konsistenz)

Das Verfahren (V) heißt **konsistent** [mit (AWP)], wenn

$$\lim_{\substack{h \rightarrow 0 \\ t=t_0+nh}} \frac{1}{h} \mathbf{R}_{n+k} = \mathbf{0}$$

für alle AWP-e gilt, die den Voraussetzungen von Satz 1.1 genügen.

Satz 2.4 (Konsistenzbedingung)

Das Verfahren (V) ist genau dann konsistent, wenn gilt

$$\sum_{j=0}^k \alpha_j = 0 \quad (K_1)$$

$$\text{und} \quad \left[\sum_{j=0}^k j \alpha_j \right] \mathbf{f}(t_n, \mathbf{y}(t_n)) = \Phi_{\mathbf{f}}(\mathbf{y}(t_n), \dots, \mathbf{y}(t_n), t_n; 0). \quad (K_2)$$

Das zum Verfahren (V) gehörende **erste charakteristische Polynom** ist durch

$$\rho(\zeta) = \alpha_k \zeta^k + \alpha_{k-1} \zeta^{k-1} + \cdots + \alpha_1 \zeta + \alpha_0$$

definiert. Hiermit läßt sich Satz 2.4 kompakter formulieren:

Satz 2.5 (Konsistenzbedingung')

Das Verfahren (V) ist genau dann konsistent, wenn gilt

$$\rho(1) = 0 \quad (K'_1)$$

$$\text{und} \quad f(t_n, \mathbf{y}(t_n)) = \frac{\Phi_f(\mathbf{y}(t_n), \dots, \mathbf{y}(t_n), t_n; 0)}{\rho'(1)}. \quad (K'_2)$$

Numerische Methoden für Anfangswertprobleme

Beispiel zu Konsistenz

Wir betrachten als Beispiel das AWP

$$y' = -y, \quad y(0) = 1 \quad \text{mit Lösung } y(t) = \exp(-t).$$

Das implizite Zweischrittverfahren

$$\begin{aligned} y_{n+2} - (1 + \alpha)y_{n+1} + \alpha y_n \\ = \frac{1}{2}h[f(t_{n+2}, y_{n+2}) + (1 - \alpha)f(t_{n+1}, y_{n+1}) - \alpha f(t_n, y_n)] \end{aligned}$$

ist konsistent (für jedes $\alpha \in \mathbb{R}$):

$$\begin{aligned} \rho(1) &= 1 - (1 + \alpha) + \alpha = 0, \\ \Phi_f(y(t_n), y(t_n), y(t_n), t_n; 0)/\rho'(1) &= \frac{1}{2}(2 - 2\alpha)f(t_n, y(t_n))/(2 - (1 + \alpha)) \\ &= f(t_n, y(t_n)). \end{aligned}$$

Lösung mit Anfangswerten $y_0 = y_1 = 1$:

$$(1 + \frac{1}{2}h)y_{n+2} - [1 + \alpha - \frac{1}{2}(1 - \alpha)h]y_{n+1} + \alpha(1 - \frac{1}{2}h)y_n = 0 \quad (*)$$

(Differenzengleichung zweiter Ordnung mit konstanten Koeffizienten).

Allgemeine Lösung: $y_n = c_1 \xi_1^n + c_2 \xi_2^n$, $\xi_1 = \alpha$ und $\xi_2 = (1 - \frac{1}{2}h)/(1 + \frac{1}{2}h)$.

Numerische Methoden für Anfangswertprobleme

Beispiel zu Konsistenz

Störe nun (*) (Rundungsfehler, $\delta > 0$):

$$\begin{aligned}(1 + \tfrac{1}{2}h)z_{n+2} - [1 + \alpha - \tfrac{1}{2}(1 - \alpha)h]z_{n+1} + \alpha(1 - \tfrac{1}{2}h)z_n &= h\delta \\ z_0 &= 1 + \delta, \quad z_1 = 1 + \delta.\end{aligned}\tag{*}'$$

Allgemeine Lösung:

$$z_n = c_1 \xi_1^n + c_2 \xi_2^n + \begin{cases} \delta/(1 - \alpha), & \text{falls } \alpha \neq 1, \\ n\delta, & \text{falls } \alpha = 1. \end{cases}$$

Fall 1: $\alpha \neq 1$

$$z_n = \frac{1}{\gamma} [\mu(\delta) \xi_1^n + \nu(\delta) \xi_2^n] + \frac{\delta}{1 - \alpha}$$

mit

$$\begin{aligned}\gamma &= 1 - \alpha - h(1 + \alpha)/2, & \mu(\delta) &= h[\alpha\delta/(1 - \alpha) - 1], \\ \nu(\delta) &= [1 - \alpha(1 + \delta)](1 + h/2).\end{aligned}$$

Numerische Methoden für Anfangswertprobleme

Beispiel zu Konsistenz

Ersetze δ durch δ^* , ergibt Lösung $\{z_n^*\}$ statt $\{z_n\}$.

Fall 1a: $-1 \leq \alpha < 1$,

$$|z_n - z_n^*| \leq \left[\frac{h/(1-\alpha) + 1 + h/2}{|1-\alpha - h(1+\alpha)/2|} + \frac{1}{1-\alpha} \right] |\delta - \delta^*|$$

für alle $h \leq h_0 = 2(1-\alpha)/(1+\alpha)$.

Fall 1b: $|\alpha| > 1$ $\{z_n - z_n^*\}$ ist unbeschränkt für $h \rightarrow 0$.

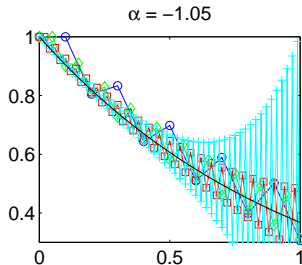
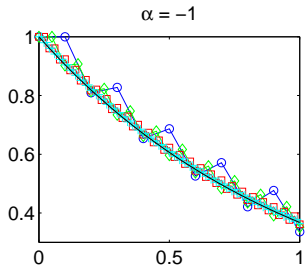
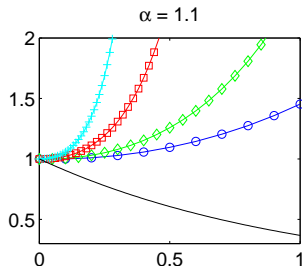
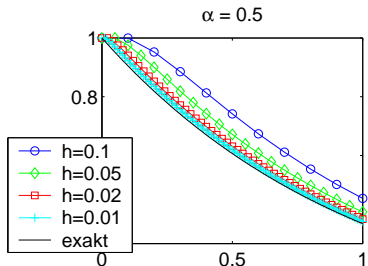
Fall 2: $\alpha = 1$

$$z_n = 1 + \frac{h-2}{2h}\delta + \frac{h+2}{2h}\delta \left(\frac{1-h/2}{1+h/2} \right)^n + n\delta$$

und $\{z_n - z_n^*\}$ ist unbeschränkt für $h \rightarrow 0$.

Numerische Methoden für Anfangswertprobleme

Beispiel zu Konsistenz



Definition 2.6 (Stabilität)

Seien $\{\delta_n^{(\ell)}\}_{n=0,1,\dots,N}$, $\ell = 1, 2$, zwei beliebige Störungen der Differenzengleichung (V) und seien $\{z_n^{(\ell)}\}$, $\ell = 1, 2$, die Lösungen von

$$\sum_{j=0}^k \alpha_j \mathbf{y}_{n+j} = h[\Phi_f(\mathbf{y}_{n+k}, \mathbf{y}_{n+k-1}, \dots, \mathbf{y}_n, t_n; h) + \delta_{n+k}^{(\ell)}], \quad (\text{V}_\delta)$$
$$\mathbf{y}_j^{(\ell)} = \mathbf{y}_j(h) + \delta_j^{(\ell)} \quad (j = 0, 1, \dots, k-1).$$

Die Methode (V) heißt **stabil**, wenn es positive Konstanten S und h_0 gibt mit

$$\max_{0 \leq n \leq N} \|z_n^{(1)} - z_n^{(2)}\| \leq S \max_{0 \leq n \leq N} \|\delta_n^{(1)} - \delta_n^{(2)}\|$$

für alle $h \in (0, h_0]$ und alle zulässigen f .

Bemerkungen.

- (1) Stabilität bedeutet, dass die Differenzengleichung **sachgemäß gestellt** ist (vgl. Satz 1.1 und 1.2, wo gezeigt wird, dass die Bedingung (**Lip**) garantiert, dass das (**AWP**) sachgemäß gestellt ist).
- (2) Stabilität ist ein Muss für Rechnungen in Gleitpunktarithmetik.

Das Verfahren (**V**) erfüllt die **Wurzelbedingung**, wenn sämtliche Nullstellen ξ seines ersten charakteristischen Polynoms betragsmäßig kleiner oder gleich 1 sind, und zusätzlich aus $|\xi| = 1$ stets folgt, dass ξ eine einfache Nullstelle ist.

Satz 2.7 (Stabilität \Leftrightarrow Wurzelbedingung)

Das Verfahren (**V**) ist genau dann stabil, wenn es die Wurzelbedingung erfüllt.

① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Numerische Methoden für Anfangswertprobleme

Der Hauptsatz

Satz 2.8 (Dahlquist)

Das Verfahren (V) ist genau dann konvergent, wenn es konsistent und stabil ist.
Oder kürzer:

Konsistenz & Stabilität \Leftrightarrow Konvergenz.

Bemerkungen.

- (1) Der Beweis zeigt, dass der globale Diskretisierungsfehler $\max_{0 \leq n \leq N} \|\mathbf{y}(t_n) - \mathbf{y}_n(h)\|$ bei stabilen Verfahren von der Ordnung p ist, wenn

$$\max_{0 \leq n \leq N-k} R_{n+k}/h = O(h^p) \quad (\text{Konsistenzordnung } p) \text{ und}$$

$$\max_{0 \leq n \leq k-1} \|\mathbf{y}(t_n) - \mathbf{y}_n(h)\| = O(h^p) \quad (\text{Anfangswerte})$$

gelten.

- (2) Ob ein Verfahren konvergent ist, kann jetzt rein algebraisch nachgeprüft werden (i.W. über die Nullstellen des ersten charakteristischen Polynoms).

Germund Dahlquist (1925–2005)

① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Numerische Methoden für Anfangswertprobleme

Einschrittverfahren

Ein **Einschrittverfahren** hat die Form

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h\Phi_f(\mathbf{y}_{n+1}, \mathbf{y}_n, t_n; h). \quad (\text{ESV})$$

Es ist immer stabil, also genau dann konvergent, wenn es konsistent ist. Ist

$$\frac{1}{h}\mathbf{R}_{n+1} = \frac{\mathbf{y}(t_{n+1}) - \mathbf{y}(t_n)}{h} - \Phi_f(\mathbf{y}(t_{n+1}), \mathbf{y}(t_n), t_n; h) = O(h^p) \quad (h \rightarrow 0),$$

so besitzt es (mindestens) die **Konsistenzordnung** p . Das Euler-Verfahren besitzt die Konsistenzordnung 1, denn:

$$\frac{\mathbf{y}(t_{n+1}) - \mathbf{y}(t_n)}{h} - \mathbf{f}(t_n, \mathbf{y}(t_n)) = \frac{h\mathbf{y}'(t_n) + \frac{1}{2}h^2\mathbf{y}''(\xi_n)}{h} - \mathbf{y}'(t_n) = \frac{1}{2}h\mathbf{y}''(\xi_n).$$

Ist \mathbf{f} p -mal stetig differenzierbar, so kann man Einschrittverfahren der Konsistenzordnung p mit der Methode des **Taylor-Abgleichs** konstruieren:

Numerische Methoden für Anfangswertprobleme

Einschrittverfahren

Man entwickelt die (unbekannte) Lösung $\mathbf{y}(t_n + h)$ formal nach Potenzen von h ,

$$\mathbf{y}(t_n + h) = \mathbf{y}(t_n) + h\mathbf{y}'(t_n) + \frac{1}{2}h^2\mathbf{y}''(t_n) + \cdots,$$

und nutzt aus, dass $\mathbf{y}'(t) = \mathbf{f}(t, \mathbf{y}(t))$ gilt, z.B.

$$\mathbf{y}'(t_n) = \mathbf{f}(t_n, \mathbf{y}(t_n)),$$

$$\begin{aligned}\mathbf{y}''(t_n) &= \mathbf{f}_t(t, \mathbf{y}(t))|_{t=t_n} + \mathbf{f}_y(t, \mathbf{y}(t))\mathbf{y}'(t)|_{t=t_n} \\ &= \mathbf{f}_t(t_n, \mathbf{y}(t_n)) + \mathbf{f}_y(t_n, \mathbf{y}(t_n))\mathbf{f}(t_n, \mathbf{y}(t_n)).\end{aligned}$$

Bricht man etwa nach dem zweiten Term ab, so ergibt sich mit

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h \left[\mathbf{f}(t_n, \mathbf{y}_n) + \frac{1}{2}h (\mathbf{f}_t(t_n, \mathbf{y}_n) + \mathbf{f}_y(t_n, \mathbf{y}_n)\mathbf{f}(t_n, \mathbf{y}_n)) \right]$$

ein Verfahren zweiter Ordnung.

Numerische Methoden für Anfangswertprobleme

Einschrittverfahren

Es ist klar, dass man dieses Verfahren (oder entsprechende Verfahren für $p > 2$) nur dann verwenden kann, wenn f einfach zu differenzieren ist (oder mittels automatische Differentiation).

Für die skalare GDG $y'(t) = ty(t)$ ergibt sich z.B.

$$y_{n+1} = y_n \left[1 + \frac{1}{2}h^2 + ht_n + \frac{1}{2}h^2t_n^2 \right]$$

im Fall von $p = 2$ und

$$y_{n+1} = y_n \left[1 + \frac{1}{2}h^2 + ht_n + \frac{1}{2}h^2t_n^2 + \frac{1}{2}h^3t_n + \frac{1}{6}h^3t_n^3 \right]$$

im Fall von $p = 3$.

Für $p = 1$ erhält man stets das Euler-Verfahren.

Für dieses Beispiel ergibt sich

mit Anfangsbedingung $y(0) = 1$, d.h. mit exakter Lösung $y(t) = \exp(t^2/2)$

als (normalisierter) globaler Diskretisierungsfehler

$$h^{-p} \max_{0 \leq n \leq N} |y(t_n) - y_n|$$

für die Taylor-Verfahren der Ordnung $p \in \{1, 2, 3\}$:

N	h	$p = 1$	$p = 2$	$p = 3$
10	1.e-1	1.016e+0	4.301e-1	3.267e-1
100	1.e-2	1.090e+0	4.757e-1	3.541e-1
1000	1.e-3	1.098e+0	4.804e-1	3.569e-1
10000	1.e-4	1.099e+0	4.808e-1	3.548e-1

① Einleitung

② Numerische Methoden für Anfangswertprobleme

2.1 Das Euler-Verfahren

2.2 Eine Sammlung von Beispielerfahren

2.3 Konvergenz, Konsistenz und Stabilität

2.4 Der Hauptsatz

2.5 Einschrittverfahren

2.6 Numerische Experimente

③ Lineare Mehrschrittverfahren

① Runge-Kutta-Verfahren

⑤ Steife Differentialgleichungen

⑥ Ausblick

Wir werden jetzt die Verfahren aus Abschnitt 2 (ab Seite 104) an folgendem AWP testen:

$$y_1' = y_2$$

$$y_1(0) = \frac{1}{2},$$

$$y_2' = \frac{y_2(y_2 - 1)}{y_1}$$

mit den Anfangsbedingungen

$$y_2(0) = -3.$$

Die Lösung wird für $t \in [0, 1]$ gesucht.

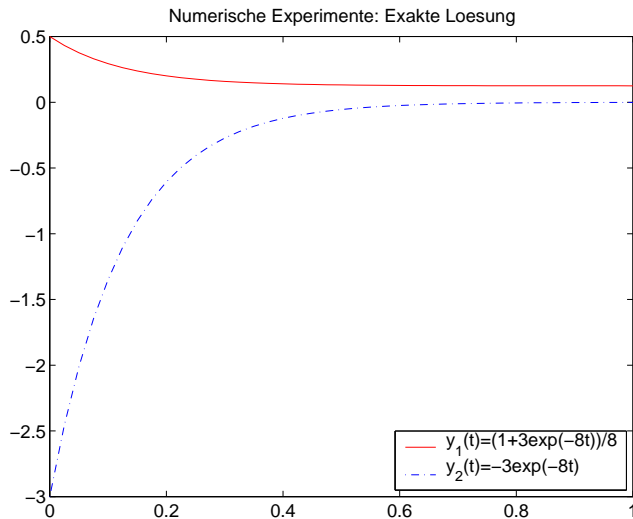
Die exakte Lösung lautet

$$y_1(t) = [1 + 3 \exp(-8t)]/8,$$

$$y_2(t) = -3 \exp(-8t).$$

Numerische Methoden für Anfangswertprobleme

Numerische Experimente



zu Beispiel 4:

Für $\rho(\zeta) = \zeta^2 + \zeta - 2 = (\zeta - 1)(\zeta + 2)$ gelten

$$\rho(1) = 0 \text{ und } \frac{\Phi_f(\mathbf{y}(t_n), \mathbf{y}(t_n), \mathbf{y}(t_n), t_n; 0))}{\rho'(1)} = \frac{3\mathbf{f}(t_n, \mathbf{y}(t_n))}{3} = \mathbf{f}(t_n, \mathbf{y}(t_n)).$$

Die Methode ist also konsistent, aber instabil.

t	$h = 0.1$	$h = 0.05$	$h = 0.025$	$h = 0.0125$
0.2	2.6e-02	8.2e-03	9.0e-03	1.5e-01
0.4	1.4e-01	2.1e-01	4.1e+00	1.8e+04
0.6	9.0e-01	5.9e+00	1.9e+03	2.2e+11
0.8	6.2e+00	1.7e+02	8.6e+05	Inf
1.0	4.2e+01	4.7e+03	4.0e+10	

zu Beispiel 3:

Für $\rho(\zeta) = \zeta^2 - \zeta = (\zeta - 1)\zeta$ gelten

$$\rho(1) = 0 \text{ und } \frac{\Phi_f(\mathbf{y}(t_n), \mathbf{y}(t_n), \mathbf{y}(t_n), t_n; 0))}{\rho'(1)} = \frac{\mathbf{f}(t_n, \mathbf{y}(t_n))/3}{1} = \frac{\mathbf{f}(t_n, \mathbf{y}(t_n))}{3}.$$

Die Methode ist also stabil, aber inkonsistent.

t	$h = 10^{-1}$	$h = 10^{-2}$	$h = 10^{-3}$	$h = 10^{-4}$
0.2	1.3e+00	1.1e+00	1.2e+00	1.2e+00
0.4	1.1e+00	9.1e+00	9.2e+00	9.2e+00
0.6	8.0e-01	5.9e-01	5.9e-01	5.9e-01
0.8	5.5e-01	3.7e-01	3.5e-01	3.5e-01
1.0	3.8e-01	2.2e-01	2.1e-01	2.1e-01

Die Methode ist aber konsistent mit dem AWP

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y})/3, \quad \mathbf{y}(0) = \mathbf{y}_0.$$

Für den Abstand zu dessen exakter Lösung

$$\tilde{\mathbf{y}}(t) = [(1 + 3 \exp(-8t/3))/8, \quad -3 \exp(-8t/3)]^\top$$

ergibt sich:

t	$h = 10^{-1}$	$h = 10^{-2}$	$h = 10^{-3}$	$h = 10^{-4}$
0.2	3.0e-01	3.6e-02	3.8e-03	3.8e-04
0.4	4.0e-01	4.4e-02	4.4e-03	4.4e-04
0.6	3.9e-01	3.9e-02	3.9e-03	3.9e-04
0.8	3.4e-01	3.1e-02	3.1e-03	3.1e-04
1.0	2.8e-01	2.3e-02	2.3e-03	2.2e-04

zu Beispiel 5: Für

$$\rho(\zeta) = \zeta^3 + \zeta^2/4 - \zeta/2 - 3/4 = (\zeta - 1)(\zeta + 5/8 - \sqrt{23}i/8)(\zeta + 5/8 + \sqrt{23}i/8)$$

gelten

$$\rho(1) = 0 \text{ und } \frac{\Phi_f(\mathbf{y}(t_n), \mathbf{y}(t_n), \mathbf{y}(t_n), t_n; 0))}{\rho'(1)} = \frac{3\mathbf{f}(t_n, \mathbf{y}(t_n))}{3} = \mathbf{f}(t_n, \mathbf{y}(t_n)).$$

Die Methode ist also konsistent und stabil, damit konvergent.

t	$h = 0.1 \uparrow$	$h = 0.05 \uparrow$	$h = 0.025 \downarrow$	$h = 0.0125 \downarrow$
0.2		8.4e-03	9.3e-04	1.1e-04
0.4	2.6e-01	4.1e-02	2.4e-04	4.5e-05
0.6	1.5e+00	1.2e-01	1.6e-04	1.4e-05
0.8	8.1e+00	3.3e-01	2.1e-05	3.8e-06
1.0	4.4e+01	9.1e-01	6.8e-06	9.6e-07

zu Beispiel 6: Die Methode ist konsistent und stabil, also konvergent.

t	$h = 10^{-1}$	$h = 10^{-2}$	$h = 10^{-3}$	10^{-4}
0.2	9.0e-01	6.7e-03	2.7e-05	1.7e-07
0.4	4.0e+00	8.9e-02	1.5e-04	4.7e-06
0.6	2.3e+01	1.7e+00	3.4e-03	1.2e-04
0.8	1.4e+02	3.3e+01	8.2e-02	2.9e-03
1.0	7.9e+03	6.6e+02	2.0e+00	7.1e-02
1.2				1.7e+00
1.4				4.3e+01

Numerische Methoden für Anfangswertprobleme

Numerische Experimente

zu Beispiel 1: Die Methode ist konsistent und stabil, also konvergent.

t	$h = 0.4 \uparrow$	$h = 0.2 \downarrow$	$h = 0.1 \downarrow$	$h = 0.05 \downarrow$
0.2		6.2e-01	3.9e-02	3.6e-03
0.4	7.8e+00	1.2e-01	1.5e-02	1.4e-03
0.6		2.5e-02	4.5e-03	4.4e-04
0.8	2.0e+01	5.0e-03	1.2e-03	1.2e-04
1.0		1.0e-03	2.9e-04	3.0e-05
1.2	5.0e+01			

zu Beispiel 2: Die Methode ist konsistent und stabil, also konvergent.

t	$h = 0.8$ ↓	$h = 0.4$ ↓	$h = 0.2$ ↓	$h = 0.1$ ↓
0.2			2.7e-01	5.5e-02
0.4		8.2e-01	8.6e-02	2.1e-02
0.6			2.1e-02	6.2e-03
0.8	1.6e+00	1.6e-01	4.6e-03	1.6e-03
1.0			9.6e-04	3.8e-04
1.6	8.3e-01	8.6e-03		
2.4	4.3e-01	4.6e-04		

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Verfahren der Bauart

$$\sum_{j=0}^k \alpha_j \mathbf{y}_{n+j} = h \sum_{j=0}^k \beta_j \mathbf{f}_{n+j}, \quad \text{wobei } \mathbf{f}_{n+j} := \mathbf{f}(t_{n+j}, \mathbf{y}_{n+j}), \quad (\text{LMV})$$

heißen **lineare Mehrschrittverfahren**, genauer **lineare k -Schritt-Verfahren**.

- O.B.d.A. sei $\alpha_k = 1$ und $(\alpha_0, \beta_0) \neq (0, 0)$.
- Falls $\beta_k = 0$, ist (LMV) explizit, sonst implizit.
- Bei impliziten Verfahren muss in jedem Zeitschritt ein (i.Allg. nichtlineares) Gleichungssystem der Form

$$\mathbf{y}_{n+k} = h\beta_k \mathbf{f}(t_{n+k}, \mathbf{y}_{n+k}) + \sum_{j=0}^{k-1} (h\beta_j \mathbf{f}_{n+j} - \alpha_j \mathbf{y}_{n+j}) = \mathbf{g}(\mathbf{y}_{n+k}) + \mathbf{c}$$

gelöst werden. Wegen

$$\|\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w})\| = h|\beta_k| \|\mathbf{f}(t_{n+k}, \mathbf{v}) - \mathbf{f}(t_{n+k}, \mathbf{w})\| \leq h|\beta_k| L \|\mathbf{v} - \mathbf{w}\|$$

besitzt dies eine eindeutige Lösung, wenn $h|\beta_k| L < 1$, die mit **Fixpunktiteration** bestimmt werden kann.

Das Polynom

$$\sigma(\zeta) := \beta_0 + \beta_1\zeta + \cdots + \beta_k\zeta^k \in \mathcal{P}_k$$

heißt **zweites charakteristisches Polynom** von (LMV) und

$$\mathcal{L}(\mathbf{z}(t); h) := \sum_{j=0}^k [\alpha_j \mathbf{z}(t + jh) - h\beta_j \mathbf{z}'(t + jh)], \quad \mathbf{z} \in C^1(I)$$

der mit (LMV) assoziierte **Differenzenoperator**.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lemma 3.1

Ist z genügend oft differenzierbar, so gilt

$$\mathcal{L}(z(t); h) = C_0 z(t) + C_1 z'(t)h + \dots + C_q z^{(q)}(t)h^q + \dots$$

$$\text{mit } C_0 = \sum_{j=0}^k \alpha_j = \rho(1), \quad C_1 = \sum_{j=0}^k [j\alpha_j - \beta_j] = \rho'(1) - \sigma(1)$$

$$\text{und } C_q = \sum_{j=0}^k \left[\frac{j^q}{q!} \alpha_j - \frac{j^{q-1}}{(q-1)!} \beta_j \right] \quad (q = 2, 3, \dots).$$

Für die Koeffizienten der analogen Entwicklung

$$\mathcal{L}(z(t + \tau h); h) = \sum_{j=0}^{\infty} D_j z^{(j)}(t + \tau h) h^j$$

$$\text{gelten } C_q = \sum_{j=0}^q \frac{\tau^j}{j!} D_{q-j} \quad (q = 0, 1, \dots).$$

Lineare Mehrschrittverfahren

Konsistenzordnung

- Der lineare Differenzenoperator \mathcal{L} entspricht im Wesentlichen dem bekannten Residuum: $\mathbf{R}_{n+k} = \mathcal{L}(\mathbf{y}(t_n); h)$.
- Das lineare Mehrschrittverfahren (LMV) besitzt die genaue **Konsistenzordnung** p , wenn

$$C_0 = C_1 = \dots = C_p = 0 \quad \text{und} \quad C_{p+1} \neq 0$$

gelten. Mit den Bezeichnungen von Lemma 3.1 ist das äquivalent zu

$$D_0 = D_1 = \dots = D_p = 0 \quad \text{und} \quad D_{p+1} \neq 0.$$

C_{p+1} ($= D_{p+1}$) heißt dann die **Fehlerkonstante** des Verfahrens.

- Beachte, dass (LMV) genau dann konsistent ist (mit anderen Worten: seine Konsistenzordnung beträgt mindestens $p = 1$), wenn $\rho(1) = 0$ und $\rho'(1) = \sigma(1)$ erfüllt sind.
- (LMV) ist damit genau dann konvergent, wenn ρ die Wurzelbedingung erfüllt und $\rho(1) = 0$ sowie $\rho'(1) = \sigma(1)$ gelten (was insbesondere $\rho'(1) = \sigma(1) \neq 0$ impliziert).

Satz 3.2

Für jedes lineare k -Schritt-Verfahren sind die folgenden fünf Aussagen äquivalent:

- (a) Das k -Schritt-Verfahren besitzt (mindestens) die Konsistenzordnung p .
- (b) $q! C_q = \sum_{j=0}^k [j^q \alpha_j - q j^{q-1} \beta_j] = 0 \quad (q = 0, 1, \dots, p)$.
- (c) Das k -Schritt-Verfahren ist konsistent mit $y' = y$, $y(0) = 1$, von (mindestens) der Ordnung p .
- (d) Die Funktion

$$\frac{\rho(\zeta)}{\log \zeta} - \sigma(\zeta)$$

hat in $\zeta = 1$ eine (mindestens) p -fache Nullstelle.

- (e) $\mathcal{L}(z(t); h) = 0$ für alle Polynome $z \in \mathcal{P}_p$.

Lineare Mehrschrittverfahren

Peano-Kern

Nach der Definition der Fehlerkonstanten C_{p+1} eines LMV wissen wir lediglich dass

$$\mathcal{L}(y(t_n); h) = h^{p+1} C_{p+1} y^{(p+1)}(t_n) + O(h^{p+2}) \quad (h \rightarrow 0).$$

Die Frage, wann auch eine Darstellung der Form $\mathcal{L}(y(t_n); h) = h^{p+1} C_{p+1} y^{(p+1)}(\xi)$ mit $t_n \leq \xi \leq t_n + h$ möglich ist, führt auf die Darstellung mittels **Peano-Kern**.

Lemma 3.3

Das lineare k -Schritt-Verfahren (**LMV**) zur Lösung von (**AWP**) besitze die Konsistenzordnung p . Für Funktionen $\mathbf{y} \in C^{(p+1)}(I)$ gilt

$$\mathcal{L}(\mathbf{y}(t); h) = h^{p+1} \int_0^k G(\tau) \mathbf{y}^{(p+1)}(t + \tau h) d\tau \quad (3.1)$$

schreiben mit der **Peano-Kernfunktion**

$$G(\tau) = \sum_{j=0}^k \left[\alpha_j \frac{(j - \tau)_+^p}{p!} - \beta_j \frac{(j - \tau)_+^{p-1}}{(p-1)!} \right], \quad u_+^k := \max\{0, u\}^k.$$

Bemerkungen:

- (1) Durch früheres Abbrechen der Taylor-Reihen im Beweis von Lemma 3.3 erhält man entsprechende Darstellungen

$$\mathcal{L}(\mathbf{y}(t); h) = h^{q+1} \int_0^k G_q(\tau) \mathbf{y}^{(q+1)}(t + \tau h) d\tau, \quad 1 \leq q \leq p$$

mit

$$G_q(\tau) = \sum_{j=0}^k \left[\alpha_j \frac{(j - \tau)_+^q}{q!} - \beta_j \frac{(j - \tau)_+^{q-1}}{(q-1)!} \right].$$

- (2) $G_q(\tau) = 0$ für $\tau \in \mathbb{R} \setminus (0, k)$.
- (3) G_q ist $(q-2)$ -mal stetig differenzierbar und $G'_q(\tau) = -G_{q-1}(\tau)$ (für $q = 2$ stückweise zu verstehen).
- (4) $G_1(s)$ ist stückweise linear mit Sprüngen der Höhe β_j an den Stellen $\tau_j = j$, $j = 0, \dots, k$ und Steigung $-(\alpha_j + \alpha_{j+1} + \dots + \alpha_k)$ im Intervall $(j, j+1)$.

Satz 3.4

Für den Differenzenoperator eines LMV der Konsistenzordnung p und $y \in C^{(p+1)}(I)$ gilt

$$|\mathcal{L}(y(t_n); h)| \leq h^{p+1} G Y$$

mit $Y := \max_{t \in I} |y^{(p+1)}(t)|$ sowie, falls der Peano-Kern $G(\tau)$ das Vorzeichen in $[0, k]$ nicht wechselt,

$$G = |C_{p+1}| = \left| \int_0^k G(\tau) \, d\tau \right|$$

und andernfalls

$$G = \int_0^k |G(\tau)| \, d\tau.$$

LMV ohne Vorzeichenwechsel des Peano-Kerns sind beispielsweise die Familie der Adams-Bashforth bzw. Adams-Moulton Verfahren.

Satz 3.4 erlaubt Abschätzungen des lokalen Diskretisierungsfehlers:

Beispiel: Für das stabile Zweischrittverfahren

$$\mathbf{y}_{n+2} - \mathbf{y}_{n+1} = \frac{h}{12}(5\mathbf{f}_{n+2} + 8\mathbf{f}_{n+1} - \mathbf{f}_n) \quad (3.2)$$

der Konsistenzordnung 3 ($C_0 = C_1 = C_2 = C_3 = 0$, $C_4 = -1/24$) erhalten wir

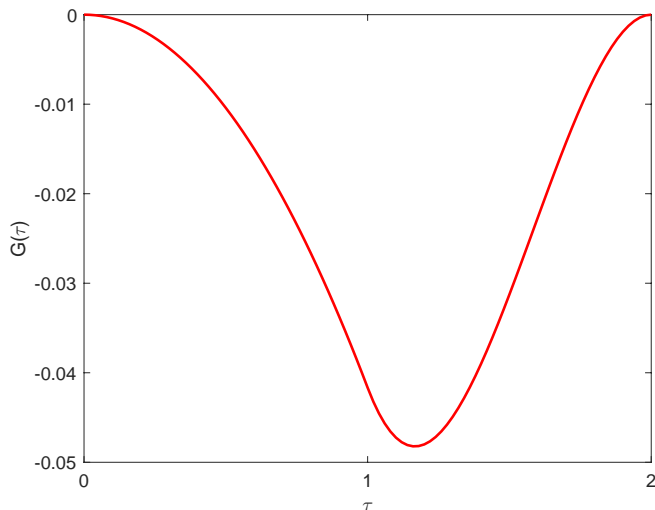
$$\begin{aligned} G(\tau) &= \frac{1}{24}(0 - \tau)_+^2 - \frac{1}{6}(1 - \tau)_+^3 - \frac{3}{4}(1 - \tau)_+^2 + \frac{1}{6}(2 - \tau)_+^3 - \frac{5}{24}(2 - \tau)_+^2 \\ &= \begin{cases} -\frac{\tau^2}{24} & 0 \leq \tau \leq 1, \\ -\frac{1}{6}\tau^3 + \frac{19}{24}\tau^2 - \frac{7}{6}\tau + \frac{1}{2} & 1 \leq \tau \leq 2. \end{cases} \end{aligned}$$

Wie man leicht nachprüft gilt $G(\tau) \leq 0$ für $\tau \in [0, 2]$.

Lineare Mehrschrittverfahren

Peano-Kern, Beispiel

Peano-Kern $G(\tau)$ des Verfahrens (3.2):



Lineare Mehrschrittverfahren

Peano-Kern, Beispiel

Ist $y \in C^4(I)$, so gilt nach Satz 3.4

$$\begin{aligned} \left| \frac{1}{h} R_{n+2} \right| &= \left| \frac{1}{h} \mathcal{L}(y(t_n); h) \right| = h^3 \left| \int_0^2 G(\tau) y^{(4)}(t_n + \tau h) d\tau \right| \\ &\leq h^3 |C_4| \max_{t \in I} |y^{(4)}(t)| = \frac{h^3}{24} \max_{t \in I} |y^{(4)}(t)|. \end{aligned}$$

Alternativ folgt mit der **Hölder-Ungleichung**

$$\left| \frac{1}{h} R_{n+2} \right| \leq h^3 \left[\int_0^2 |G(\tau)|^\mu d\tau \right]^{1/\mu} \left[\int_0^2 |y^{(4)}(t_n + \tau h)|^\nu d\tau \right]^{1/\nu},$$

falls $1/\mu + 1/\nu = 1$.

Lineare Mehrschrittverfahren

Peano-Kern, Beispiel

Für $\mu = 1$, $\nu = \infty$ ergibt sich (entspricht Satz 3.4)

$$\begin{aligned} \left| \frac{1}{h} R_{n+2} \right| &\leq -h^3 \left[\int_0^2 G(\tau) d\tau \right] \max_{t_n \leq t \leq t_n + 2h} |y^{(4)}(t)| \\ &= \frac{h^3}{24} \max_{t_n \leq t \leq t_n + 2h} |y^{(4)}(t)| \end{aligned}$$

sowie für $\mu = \infty$, $\nu = 1$:

$$\begin{aligned} \left| \frac{1}{h} R_{n+2} \right| &\leq h^3 \max_{0 \leq \tau \leq 2} |G(\tau)| \left[\int_0^2 |y^{(4)}(t_n + \tau h)| d\tau \right] \\ &= \frac{125}{2592} h^3 \int_0^2 |y^{(4)}(t_n + \tau h)| d\tau. \end{aligned}$$

Lineare Mehrschrittverfahren

Peano-Kern, globaler Fehler

Mit Hilfe von Satz 3.4 kann man auch Abschätzungen für den globalen Diskretisierungsfehler angeben. Sei

$$e_n := y(t_n) - \tilde{y}_n, \quad \tilde{y}_n = \tilde{y}_n(h),$$

wobei die Folge $(\tilde{y}_n)_{n \in \mathbb{N}_0}$ die durch Rundungsfehler gestörte Differenzengleichung

$$\sum_{j=0}^k \alpha_j \tilde{y}_{n+j} = h \sum_{j=0}^k \beta_j f(t_{n+j}, \tilde{y}_{n+j}) + \underbrace{\vartheta_n K h^{q+1}}_{\text{lokaler Rundungsfehler}}, \quad \|\vartheta_n\|_\infty \leq 1$$

erfüllt. Wir setzen

$$A := \sum_{j=0}^k |\alpha_j|, \quad B := \sum_{j=0}^k |\beta_j|.$$

sowie für den Fehler der Startwerte

$$E := \max_{0 \leq j \leq k-1} \|e_j\|.$$

Lineare Mehrschrittverfahren

Peano-Kern, globaler Fehler

Bezeichnet $(\gamma_j)_{j \in \mathbb{N}_0}$ die Folge der Koeffizienten der Potenzreihe

$$\frac{1}{\zeta^k \rho(\zeta^{-1})} = \frac{1}{\alpha_k + \alpha_{k-1}\zeta + \cdots + \alpha_0\zeta^k} = \sum_{j=0}^{\infty} \gamma_j \zeta^j,$$

so gilt, sofern ρ die Wurzelbedingung erfüllt ist,

$$\Gamma := \sup_{j \geq 0} |\gamma_j| < \infty.$$

Wir setzen ferner

$$\Gamma^* := \frac{\Gamma}{1 - h|\beta_k|L}$$

mit der Lipschitz-Konstanten L der rechten Seite $f(t, \mathbf{y})$.

Lineare Mehrschrittverfahren

Peano-Kern, globaler Fehler

Unter der Annahme, dass

$$h|\beta_k|L < 1 \quad (3.3)$$

lässt sich dann zeigen¹⁷, dass

$$\|e_n\| \leq \Gamma^* [kAE + (t_n - t_0)(\frac{1}{h}\mathcal{L}(\mathbf{y}(t_n); h) + h^q K)] \exp(\Gamma^* BL(t_n - t_0)).$$

Mit Hilfe der Peano-Kern-Darstellung für den lokalen Fehler $\mathcal{L}(\mathbf{y}(t_n); h)$ aus Satz 3.4 ergibt sich

$$\|e_n\| \leq \Gamma^* [kAE + (t_n - t_0)(h^p GY + h^q K)] \exp(\Gamma^* BL(t_n - t_0)). \quad (3.4)$$

Details in [Lambert, 1991; §3.6], [Henrici, 1962; §5.3–4] .

¹⁷Hier wird wie schon in Kapitel 2 auf die Lösungsdarstellung für inhomogene lineare Differenzengleichungen zurückgegriffen.

Lineare Mehrschrittverfahren

Bemerkungen

- (1) (3.3) ist für explizite LMV stets erfüllt, für implizite ist es eine hinreichende Bedingung für die eindeutige Lösbarkeit der Verfahrensgleichung.
- (2) (3.4) zeigt den Einfluss von lokalem Fehler, Startfehler und Rundfehler; Insbesondere implizieren $E = O(h^p)$, Konsistenzordnung p sowie $q \geq p$ für den Gesamtfehler $\|e_n\| = O(h^p)$.
- (3) Für den Startfehler findet keine Fehlerakkumulation statt, d.h. für die Anlaufrechnung genügt ein Verfahren der Konsistenzordnung $p - 1$.
- (4) Das Rundungsfehlermodell Kh^{q+1} ist unrealistisch; wir beobachten jedoch, dass man eine h -Potenz beim Übergang zum globalen Fehler verliert. Realistischer wäre eine absolute Schranke ϵ für den (lokalen) Rundungsfehler, was im globalen auf einen Anteil $\sim \epsilon/h$ führt. (Konsequenz?)
- (5) Ist die Wurzelbedingung nicht erfüllt, so ist Γ nicht endlich.
- (6) (3.4) fördert das Verständnis der Fehlerentwicklung, ist in der Praxis aber kaum von Nutzen. Obwohl die einzelnen Terme – zumindest a posteriori – geschätzt werden können, wird der globale Fehler durch (3.4) oft stark überschätzt.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

- Was ist die maximale Ordnung eines **konvergenten** linearen k -Schritt-Verfahrens?
- $2k + 2$ freie Parameter $\{\alpha_j, \beta_j\}_{j=0}^k$,
 $2k + 1$ nach Normierung, $2k$ für ein explizites Verfahren.
- Konsistenz der Ordnung p führt auf $p + 1$ homogene lineare Gleichungen für die Koeffizienten. Bis zu welcher Ordnung p liegt auch Stabilität vor? Erste Vermutung: $p = 2k$ [$p = 2k - 1$] im impliziten [expliziten] Fall?
- 1956 beantwortet in

CONVERGENCE AND STABILITY IN THE NUMERICAL INTEGRATION OF ORDINARY DIFFERENTIAL EQUATIONS

GERMUND DAHLQUIST

1. Introduction and summary

1.1. Statement of the problem. Consider a class of difference equations

$$(1.1) \quad \alpha_k y_{n+k} + \alpha_{k-1} y_{n+k-1} + \dots + \alpha_0 y_n = h(\beta_k f_{n+k} + \dots + \beta_0 f_n),$$

Lineare Mehrschrittverfahren

Maximale Konsistenzordnung

Nach Satz 3.2 ist die Konsistenzordnung eines LMV von mindestens p äquivalent damit, dass die Funktion

$$\varphi(\zeta) = \frac{\rho(\zeta)}{\log \zeta} - \sigma(\zeta), \quad (|\arg \zeta| < \pi, \log 1 = 0)$$

an der Stelle $\zeta = 1$ eine Nullstelle der Vielfachheit mindestens p besitzt.

Als einfache Folgerung hieraus lässt sich ein zweites charakteristisches Polynom σ bei gegebenem ersten charakteristischem Polynom ρ optimal wählen:

Satz 3.5

Sei $\rho \in \mathcal{P}_k$ mit $\rho(1) = 0$ sowie $\ell \in \{0, 1, \dots, k\}$. Dann gibt es genau ein Polynom $\sigma \in \mathcal{P}_\ell$ sodass das zugehörige LMV die Konsistenzordnung mindestens ℓ besitzt.

Von praktischem Interesse lediglich

$\ell = k - 1$	führt auf bestmögliches explizites Verfahren,
$\ell = k$	" " implizites " .

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

Ausgangspunkt der Analyse der maximalen Konsistenzordnung eines stabilen LMV sind wieder dessen charakteristischen Polynome $\rho(\zeta)$ und $\sigma(\zeta)$. Folgende Variablensubstitution ist hierbei hilfreich:

$$\zeta = \frac{z+1}{z-1}, \quad z = \frac{\zeta+1}{\zeta-1}, \quad (3.5)$$

sowie die hierdurch bestimmten Polynome

$$R(z) = \left(\frac{z-1}{2} \right)^k \rho(\zeta) = \sum_{j=0}^k a_j z^j, \quad S(z) = \left(\frac{z-1}{2} \right)^k \sigma(\zeta) = \sum_{j=0}^k b_j z^j.$$

Lemma 3.6

Für ein stabiles LMV mit mindestens Konsistenzordnung $p = 0$ gilt

- (a) $a_k = 0$ sowie $a_{k-1} = 2^{1-k} \rho'(1) \neq 0$.
- (b) Alle von Null verschiedenen Koeffizienten von $R(z)$ besitzen das gleiche Vorzeichen.

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

Lemma 3.7

Ein LMV besitzt genau dann die Konsistenzordnung p , wenn

$$R(z) \left(\log \frac{z+1}{z-1} \right)^{-1} - S(z) = C_{p+1} \left(\frac{2}{z} \right)^{p-k} + O \left(\left(\frac{2}{z} \right)^{p-k+1} \right), \quad z \rightarrow \infty. \quad (3.6)$$

Lemma 3.8

Für die Koeffizienten der Laurent-Reihe

$$\left(\log \frac{z+1}{z-1} \right)^{-1} = \frac{z}{2} - \frac{\mu_1}{z} - \frac{\mu_3}{z^3} - \frac{\mu_5}{z^5} - \dots \quad (3.7)$$

gilt $\mu_{2j+1} > 0$ für alle $j \geq 0$.

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

Satz 3.9 (Dahlquist-Barriere)

Für die Konsistenzordnung p eines stabilen linearen k -Schritt-Verfahrens gilt

$$p \leq \begin{cases} k + 2, & \text{falls } k \text{ gerade,} \\ k + 1, & \text{falls } k \text{ ungerade,} \\ k, & \text{falls } \beta_k / \alpha_k \leq 0. \end{cases}$$

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

Optimalität bei der Ordnung LMV hat strukturelle Eigenschaften zur Folge:

Satz 3.10

Stabile LMV mit (maximaler) Konsistenzordnung $k + 2$ sind symmetrisch, d.h.

$$\alpha_j = -\alpha_{k-j} \quad \text{und} \quad \beta_j = \beta_{k-j}, \quad j = 0, \dots, k. \quad (3.8)$$

Beachte: Bei (stabilen) symmetrischen LMV gilt $\rho(\zeta) = -\zeta^k \rho(1/\zeta)$. Mit ζ ist somit auch $1/\zeta$ Nullstelle von ρ , d.h. alle Nullstellen von ρ liegen auf dem Einheitskreis und sind somit einfach.

Korollar 3.11

Ist k gerade, dann ist ein stabiles lineares k -Schritt-Verfahren mit optimaler Konsistenzordnung $k + 2$ nur **schwach stabil**, d.h. alle Nullstellen des ersten charakteristischen Polynoms haben Betrag 1.

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

(Jedes) stabile k -Schritt-Verfahren ($k = 2\ell$) mit Konsistenzordnung $k+2$ kann man wie folgt konstruieren:

(1) Setze

$$\rho(\zeta) := (\zeta - 1)(\zeta + 1) \prod_{j=1}^{(k-2)/2} (\zeta - \zeta_j)(\zeta - \bar{\zeta}_j)$$

mit paarweise verschiedenen ζ_j , $|\zeta_j| = 1$, $\operatorname{Im} \zeta_j > 0$.

(2) Bestimme die ersten Koeffizienten der Taylor-Entwicklung

$$\frac{\left(\frac{z-1}{2}\right)^k \rho\left(\frac{z+1}{z-1}\right)}{\log \frac{z+1}{z-1}} = \sum_{j=0}^{\infty} b_j z^j \quad \text{und setze} \quad S(z) := \sum_{j=0}^k b_j z^j.$$

(3) Setze

$$\sigma(\zeta) := (\zeta - 1)^k S\left(\frac{\zeta + 1}{\zeta - 1}\right).$$

Lineare Mehrschrittverfahren

Die erste Dahlquist-Barriere

- Das einzige stabile Zweischrittverfahren der Ordnung 4 ist die **Simpson-Regel**¹⁸

$$y_{n+2} - y_n = \frac{h}{3}(f_{n+2} + 4f_{n+1} + f_n).$$

- Für $k = 4$ ist z.B.

$$y_{n+4} - y_n = \frac{h}{90}(56f_{n+4} - 31f_{n+3} + 96f_{n+2} - 31f_{n+1} + 56f_n)$$

ein stabiles Verfahren der Ordnung 6.

- In der Praxis spielen diese Beispiele (wie alle linearen 2ℓ -Schritt-Verfahren der Ordnung $2\ell + 2$) keine Rolle (vgl. dazu Abschnitt 6).

¹⁸Thomas Simpson (1710–1761)

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lineare Mehrschrittverfahren

Die Verfahren von Adams-Bashforth und Adams-Moulton

Die Idee der **Adams-Bashforth-Verfahren**¹⁹:

$$\mathbf{y}(t_{n+1}) - \mathbf{y}(t_n) = \int_{t_n}^{t_{n+1}} \mathbf{y}'(t) dt = \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{y}(t)) dt.$$

Ersetze $\mathbf{f}(t, \mathbf{y}(t))$ durch ein Polynom $\mathbf{q}_{k-1}(t) \in \mathcal{P}_{k-1}$, das die k Datenpaare

$$(t_n, \mathbf{f}_n), (t_{n-1}, \mathbf{f}_{n-1}), \dots, (t_{n-k+1}, \mathbf{f}_{n-k+1})$$

interpoliert. In der Lagrange²⁰-Darstellung ist dieses durch

$$\mathbf{q}_{k-1}(t) = \sum_{j=0}^{k-1} \mathbf{f}_{n-j} \prod_{\substack{\ell=0 \\ \ell \neq j}}^{k-1} \frac{t - t_{n-\ell}}{t_{n-j} - t_{n-\ell}}$$

gegeben.

¹⁹ John Couch Adams (1819–1892), Francis Bashforth, 1819–1912)

²⁰ Joseph-Louis Lagrange (1736–1813)

Lineare Mehrschrittverfahren

Die Verfahren von Adams-Bashforth und Adams-Moulton

k -te Adams-Bashforth-Formel

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h \sum_{j=0}^{k-1} \beta_{k,j} \mathbf{f}_{n-j} \quad (\text{A-B})$$

$$\text{mit } \beta_{k,j} = \frac{1}{h} \int_{t_n}^{t_{n+1}} \prod_{\substack{\ell=0 \\ \ell \neq j}}^{k-1} \frac{t - t_{n-\ell}}{t_{n-j} - t_{n-\ell}} dt = \int_0^1 \prod_{\substack{\ell=0 \\ \ell \neq j}}^{k-1} \frac{s + \ell}{\ell - j} ds.$$

Die **Adams-Moulton-Verfahren**²¹ konstruiert man fast genauso. Der einzige Unterschied besteht darin, dass ein Interpolationspolynom q_k vom Grad k zu den $(k+1)$ Daten

$$(t_{n+1}, \mathbf{f}_{n+1}), (t_n, \mathbf{f}_n), (t_{n-1}, \mathbf{f}_{n-1}), \dots, (t_{n-k+1}, \mathbf{f}_{n-k+1})$$

als Approximation an $\mathbf{f}(t, \mathbf{y}(t))$ verwendet wird.

²¹Forest Ray Moulton (1872–1952),

Lineare Mehrschrittverfahren

Die Verfahren von Adams-Bashforth und Adams-Moulton

k -te Adams-Moulton-Formel:

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h \sum_{j=0}^k \beta_{k,j}^* \mathbf{f}_{n+1-j} \quad \text{mit} \quad \beta_{k,j}^* = \int_{-1}^0 \prod_{\substack{\ell=0 \\ \ell \neq j}}^k \frac{s + \ell}{\ell - j} ds. \quad (\text{A-M})$$

Satz 3.12

Das Adams-Bashforth-Verfahren (A-B) ist ein explizites lineares k -Schritt-Verfahren. Es ist stabil und besitzt die Konsistenzordnung k .

Das Adams-Moulton-Verfahren (A-M) ist ein implizites lineares k -Schritt-Verfahren. Es ist stabil und besitzt die Konsistenzordnung $k + 1$.

Lineare Mehrschrittverfahren

Die Verfahren von Adams-Bashforth und Adams-Moulton

Koeffizienten für Adams-Bashforth-Verfahren:

k						
1	1					
2	$\frac{3}{2}$	$-\frac{1}{2}$				
3	$\frac{23}{12}$	$-\frac{16}{12}$	$\frac{5}{12}$			
4	$\frac{55}{24}$	$-\frac{59}{24}$	$\frac{37}{24}$	$-\frac{9}{24}$		
5	$\frac{1901}{720}$	$-\frac{2774}{720}$	$\frac{2616}{720}$	$-\frac{1274}{720}$	$\frac{251}{720}$	
6	$\frac{4277}{1440}$	$-\frac{7923}{1440}$	$\frac{9982}{1440}$	$-\frac{7298}{1440}$	$\frac{2877}{1440}$	$-\frac{475}{1440}$

Beispielsweise ist

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{h}{2}(3\mathbf{f}_n - \mathbf{f}_{n-1})$$

das Adams-Bashforth-Verfahren für $k = 2$.

Lineare Mehrschrittverfahren

Die Verfahren von Adams-Bashforth und Adams-Moulton

Koeffizienten für Adams-Moulton-Verfahren:

k							
1	$\frac{1}{2}$	$\frac{1}{2}$					
2	$\frac{5}{12}$	$\frac{8}{12}$	$-\frac{1}{12}$				
3	$\frac{9}{24}$	$\frac{19}{24}$	$-\frac{5}{24}$	$\frac{1}{24}$			
4	$\frac{251}{720}$	$\frac{646}{720}$	$-\frac{264}{720}$	$\frac{106}{720}$	$-\frac{19}{720}$		
5	$\frac{475}{1440}$	$\frac{1427}{1440}$	$-\frac{798}{1440}$	$\frac{482}{1440}$	$-\frac{173}{1440}$	$\frac{27}{1440}$	
6	$\frac{19087}{60480}$	$\frac{65112}{60480}$	$-\frac{46461}{60480}$	$\frac{37504}{60480}$	$-\frac{20211}{60480}$	$\frac{6312}{60480}$	$-\frac{863}{60480}$

Beispielsweise ist

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{h}{2}(\mathbf{f}_{n+1} + \mathbf{f}_n) \quad (\text{Trapezregel})$$

das Adams-Moulton-Verfahren für $k = 1$ und

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{h}{12}(5\mathbf{f}_{n+1} + 8\mathbf{f}_n - \mathbf{f}_{n-1})$$

das für $k = 2$.

Lineare Mehrschrittverfahren

Die Verfahren von Nyström und Milne-Simpson

Natürlich kann man auch in

$$\mathbf{y}(t_{n+k}) - \mathbf{y}(t_{n+k-\ell}) = \int_{t_{n+k-\ell}}^{t_{n+k}} \mathbf{y}'(t) dt = \int_{t_{n+k-\ell}}^{t_{n+k}} \mathbf{f}(t, \mathbf{y}(t)) dt$$

($\ell = 1, 2, \dots$) den Integrand durch ein Interpolationspolynom ersetzen, um lineare Mehrschrittverfahren zu konstruieren. (Für $\ell = 1$ ergeben sich die Adams-Formeln.) Für $\ell = 2$ erhält man so die expliziten **Nyström-Verfahren**, z.B. die **Mittelpunktsregel**

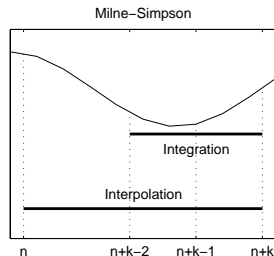
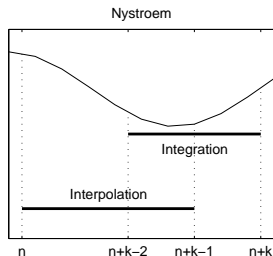
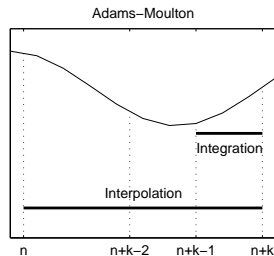
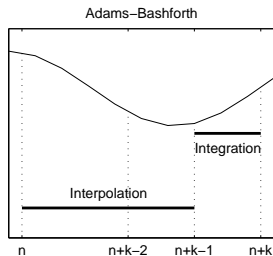
$$\mathbf{y}_{n+2} - \mathbf{y}_n = 2h\mathbf{f}_{n+1},$$

bzw. die impliziten **Milne-Simpson-Verfahren**, wie etwa die **Simpson-Regel**

$$\mathbf{y}_{n+2} - \mathbf{y}_n = \frac{h}{3} (\mathbf{f}_{n+2} + 4\mathbf{f}_{n+1} + \mathbf{f}_n).$$

Lineare Mehrschrittverfahren

Interpolation und Integration bei den Adams-artigen LMV



- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lineare Mehrschrittverfahren

Prädiktor-Korrektor-Verfahren

Löst man ein AWP durch ein implizites Mehrschrittverfahren, so muss in jedem Zeitschritt ein Gleichungssystem der Form

$$\mathbf{y}_{n+k} = h\beta_k \mathbf{f}(t_{n+k}, \mathbf{y}_{n+k}) + \mathbf{c}$$

(vgl. Abschnitt 1) gelöst werden. Ist $h|\beta_k|L < 1$, so konvergiert die Fixpunktiteration

$$\mathbf{y}_{n+k}^{(\nu)} = h\beta_k \mathbf{f}(t_{n+k}, \mathbf{y}_{n+k}^{(\nu-1)}) + \mathbf{c} \quad \nu = 1, 2, \dots$$

gegen die (eindeutige) Lösung. Einen Startwert $\mathbf{y}_{n+k}^{(0)}$ kann man mit einem expliziten Verfahren berechnen. Man verwendet dazu ein Verfahren gleicher Ordnung.

Lineare Mehrschrittverfahren

Prädiktor-Korrektor-Verfahren

Kombiniert man etwa die $(k + 1)$ -te Adams-Bashforth-Formel (Prädiktor) mit der k -ten Adams-Moulton-Formel (Korrektor), so liest sich ein Zeitschritt des resultierenden **Prädiktor-Korrektor-Verfahrens** wie folgt:

$$(P): \quad \mathbf{y}_{n+1}^{(0)} = \mathbf{y}_n + h \sum_{j=0}^k \beta_{k+1,j} \mathbf{f}_{n-j},$$

For $\nu = 0, 1, 2, \dots$:

$$\mathbf{f}_{n+1}^{(\nu)} = \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}^{(\nu)}),$$

$$(K): \quad \mathbf{y}_{n+1}^{(\nu+1)} = \mathbf{y}_n + h \beta_{k,0}^* \mathbf{f}_{n+1}^{(\nu)} + h \sum_{j=1}^k \beta_{k,j}^* \mathbf{f}_{n+1-j}.$$

Man bricht die Iteration (K) ab, wenn $\|\mathbf{y}_{n+1}^{(\nu+1)} - \mathbf{y}_{n+1}^{(\nu)}\|$ „genügend klein“ ist. Dann setzt man $\mathbf{y}_{n+1} = \mathbf{y}_{n+1}^{(\nu+1)}$ und (für weitere Zeitschritte) $\mathbf{f}_{n+1} = \mathbf{f}_{n+1}^{(\nu)}$ (oder alternativ: $\mathbf{f}_{n+1} = \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}^{(\nu+1)})$).

Lineare Mehrschrittverfahren

Prädiktor-Korrektor-Verfahren

- In der Praxis kann man so nur selten vorgehen, da völlig unklar ist, wie viele Schritte von (K) durchgeführt werden müssen, um das Abbruchkriterium zu erfüllen. Stattdessen wird man nur eine feste (kleine) Zahl μ von diesen Iterationsschritten durchführen.
- Bezeichnen p bzw. p^* die Konsistenzordnungen von Prädiktor und Korrektor (in unserem Beispiel $p = p^* = k$), dann ist die Konsistenzordnung des zusammengesetzten Verfahrens gleich der des Korrektors, wenn $p \geq p^*$ oder $\mu > p^* - p$ gilt.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lineare Mehrschrittverfahren

Absolute Stabilität, einführendes Beispiel

Der Begriff **absolute Stabilität** befasst sich anstelle der Stabilität eines Verfahrens im Grenzwert $h \rightarrow 0$, mit dessen Verhalten für lange Integrationsintervalle bei fester Schrittweite $h > 0$.

Qualitativ: wie klein muss man bei gegebener AWA und gegebenem Verfahren die Schrittweite wählen, damit die numerische Approximation der Lösung sich zumindest qualitativ richtig verhält.

Beispiel 1: Für die AWA

$$y'(t) = -\sin t, \quad y(0) = 1 \quad (3.9)$$

mit Lösung $y(t) = \cos t$ (und $L = 0$) beträgt das Residuum R_{n+1} des expliziten Euler-Verfahrens an der Stelle t_n

$$R_{n+1} = \mathcal{L}(y(t_n); h) = -\frac{h^2}{2} y''(t_n) + O(h^3) = \frac{h^2}{2} \cos t_n + O(h^3),$$

sodass der globale Fehler für $t \in [t_0, t_{\text{end}}]$ beschränkt ist durch (vgl. Beweis Satz 2.1)

$$\max_{t_0 \leq n h \leq t_{\text{end}}} |e_n| \leq h \max_{t_0 \leq t \leq t_{\text{end}}} |\cos t| = h.$$

Lineare Mehrschrittverfahren

Absolute Stabilität, einführendes Beispiel

Für einen Fehler $|e| \leq 10^{-3}$ bei der Integration bis $t_{\text{end}} = 2$ müßte also eine Schrittweite von $h = 10^{-3}$ ausreichen. Man erhält für $N = 2000$

$$y_N = -4.166014 \cdot 10^{-1}, \quad |y(2) - y_N| = 4.547667 \cdot 10^{-4}.$$

Beispiel 2: Wir modifizieren obige AWA zu

$$y'(t) = \lambda(y - \cos t) - \sin t, \quad y(0) = 1 \quad (3.10)$$

mit derselben Lösung $y(t) = \cos t$. Man rufe sich in Erinnerung, dass beim expliziten Euler-Verfahren der globale Fehler der Rekursion

$$|e_{n+1}| \leq (1 + hL)|e_n| + O(h^2)$$

genügt, mit L der Lipschitz-Konstanten der rechten Seite. Für $\lambda = -10$ erhalten wir nun

$$y_N = -4.170721 \cdot 10^{-1}, \quad |y(2) - y_N| = 1.611611 \cdot 10^{-5}.$$

Lineare Mehrschrittverfahren

Absolute Stabilität, einführendes Beispiel

Beispiel 3: Wir betrachten obiges Beispiel für $\lambda = -2100$. In diesem Fall ergibt sich

$$y_N = 1.597768 \cdot 10^{76}, \quad |y(2) - y_N| = 1.452516 \cdot 10^{76}.$$

Verschiedene Schrittweiten für dieses AWP liefern

h	Fehler bei $t_{\text{end}} = 2$
0.001	$1.7e + 76$
0.000976	$3.1e + 36$
0.00095	$8.6e - 04$
0.0008	$7.3e - 04$
0.0004	$3.6e - 04$

Lineare Mehrschrittverfahren

Absolute Stabilität

- Bekanntlich streben die Lösungen $\mathbf{y}(t)$ von $\mathbf{y}' = A\mathbf{y}$, $A \in \mathbb{R}^{m \times m}$ (konstant) gegen $\mathbf{0}$ für $t \rightarrow \infty$, wenn $\operatorname{Re} \lambda < 0$ für alle Eigenwerte λ von A gilt.
- Wir suchen Bedingungen an ein numerisches Verfahren (zunächst LMV), so dass die Näherungslösungen dasselbe asymptotische Verhalten besitzen.
- Dazu eine Bezeichnung: Seien ρ und σ die charakteristischen Polynome eines LMV; dann heißt

$$\pi(\zeta; \hat{h}) := \rho(\zeta) - \hat{h}\sigma(\zeta), \quad \hat{h} = h\lambda$$

Stabilitätspolynom des Verfahrens.

Lemma 3.13

Es seien $\{y_n\}$ die Näherungen eines linearen k -Schritt-Verfahrens

$$\sum_{j=0}^k \alpha_j \mathbf{y}_{n+j} = h \sum_{j=0}^k \beta_j \mathbf{f}_{n+j} \quad (n = 0, 1, 2, \dots)$$

zur Lösung von $\mathbf{y}' = A\mathbf{y}$ (inkl. Anfangsbedingungen). Bei festem h gilt

$$\lim_{n \rightarrow \infty} \|\mathbf{y}_n\| = 0$$

genau dann, wenn alle Nullstellen von $\pi(\zeta; \hat{h}) = \pi(\zeta; h\lambda)$ (als Polynom in ζ betrachtet) betragsmäßig echt kleiner als 1 sind und zwar für jedes $\lambda \in \Lambda(A)$.

Lineare Mehrschrittverfahren

Absolute Stabilität

- Das Verfahren heißt **absolut stabil** für \hat{h} , wenn alle Nullstellen ζ von $\pi(\cdot; \hat{h})$ die Beziehung $|\zeta| < 1$ erfüllen.
- Die Menge

$$\mathcal{R}_A := \{\hat{h} \in \mathbb{C} : \pi(\cdot; \hat{h}) \text{ hat nur Nullstellen in } |\zeta| < 1\}$$

heißt **Stabilitätsgebiet** des Verfahrens.

- Das Verfahren heißt **A-stabil** (absolut stabil), wenn \mathcal{R}_A die linke Halbebene $\{\operatorname{Re} \zeta < 0\}$ enthält.

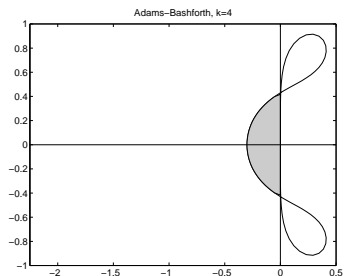
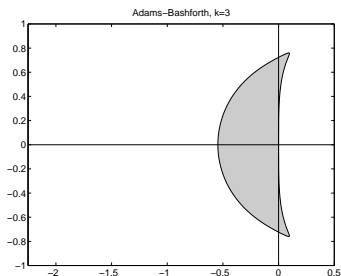
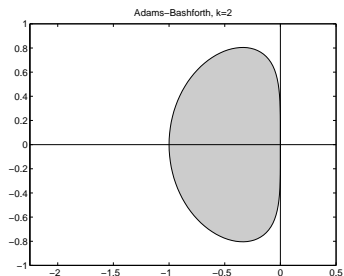
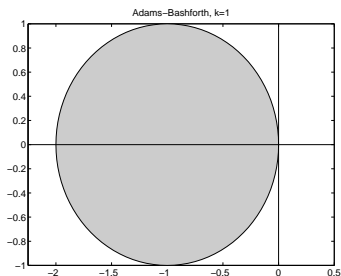
Bemerkungen.

- (1) Für jedes konvergente LMV gibt es ein $\hat{h}_0 > 0$, so dass $\mathcal{R}_A \cap [0, \hat{h}_0] = \emptyset$.
- (2)

$$\partial \mathcal{R}_A \subseteq \left\{ \hat{h} \in \mathbb{C} : \hat{h} = \frac{\rho(e^{i\phi})}{\sigma(e^{i\phi})}, 0 \leq \phi \leq 2\pi \right\}.$$

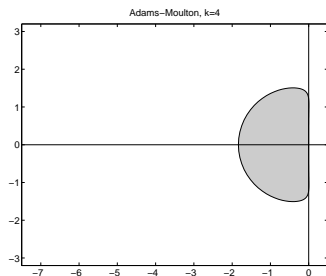
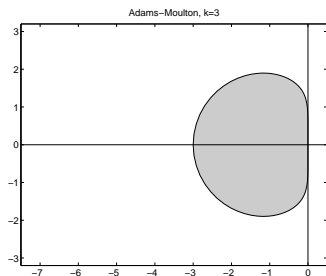
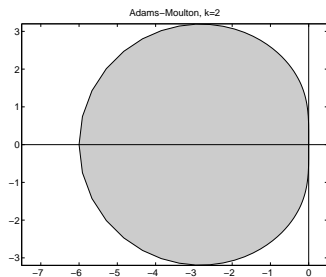
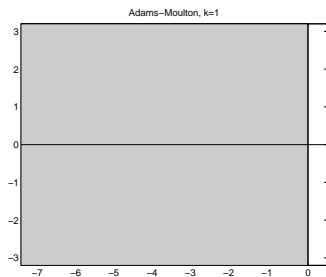
Lineare Mehrschrittverfahren

Absolute Stabilität



Lineare Mehrschrittverfahren

Absolute Stabilität



Lineare Mehrschrittverfahren

Zweite Dahlquist-Barriere

Ein LMV mit charakteristischen Polynomen $\rho(\zeta)$ und $\sigma(\zeta)$ heißt **irreduzibel**, falls $\rho(\zeta)$ und $\sigma(\zeta)$ keine gemeinsamen Nullstellen besitzen.

Lemma 3.14

Ist ein lineares Mehrschrittverfahren A-stabil, so gilt

$$\operatorname{Re} \frac{\rho(\zeta)}{\sigma(\zeta)} \geq 0 \quad \text{für } |\zeta| \geq 1. \quad (3.11)$$

Für irreduzible LMV gilt auch die Umkehrung, d.h. (3.11) impliziert A-Stabilität.

Satz 3.15 (Dahlquist, 1963)

Für die Konsistenzordnung p eines A-stabilen LMV gilt $p \leq 2$. Gilt $p = 2$, so gilt für die Fehlerkonstante des Verfahrens $C \leq -\frac{1}{12}$. Die Trapezregel ist das einzige A-stabile LMV mit Fehlerkonstante $C = -\frac{1}{12}$.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
 - 3.1 Begriffe
 - 3.2 Konsistenzordnung linearer Mehrschrittverfahren
 - 3.3 Die erste Dahlquist-Barriere
 - 3.4 Die Verfahren von Adams-Bashforth und Adams-Moulton
 - 3.5 Prädiktor-Korrektor-Verfahren
 - 3.6 Absolute Stabilität
 - 3.7 BDF-Verfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen

Lineare Mehrschrittverfahren

BDF-Verfahren

Idee der **BDF-Verfahren** (Backward Differentiation Formulas), oder auch **Gear-Verfahren**²²:

Um eine Näherung für $\mathbf{y}(t_{n+1})$ zu gewinnen, approximieren wir $t \mapsto \mathbf{f}(t, \mathbf{y}(t))$ durch das Interpolationspolynom $P_k \in \mathcal{P}_k$ mit den $k+1$ Wertepaaren

$$(t_j, \mathbf{f}(t_j, \mathbf{y}_j)) = (t_j, \mathbf{f}_j), \quad j = n - k + 1, \dots, n + 1.$$

Es ergibt sich

$$\mathbf{y}'(t) = \mathbf{f}(t, \mathbf{y}(t)) \approx P_k(t) = \sum_{j=0}^k \ell_j \left(\frac{t_{n+1} - t}{h} \right) \mathbf{y}_{n+1-j} \quad \text{mit} \quad \ell_j(t) = \prod_{\substack{s=0 \\ s \neq j}}^k \frac{t - s}{j - s}.$$

Jetzt approximieren wir $\mathbf{y}'(t_{n+1}) \approx P'_k(t_{n+1}) = \sum_{j=0}^k \left(-\frac{1}{h}\right) \ell'_j(0) \mathbf{y}_{n+1-j}$ und setzen diese Näherung ein in

$$P'_k(t_{n+1}) \approx \mathbf{y}'(t_{n+1}) = \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}) \approx \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}) = \mathbf{f}_{n+1}$$

²²Charles William Gear, (1935–)

Lineare Mehrschrittverfahren

BDF-Verfahren

und erhalten die Verfahrensgleichung eines linearen k -Schritt-Verfahrens

$$\sum_{j=0}^k (-\ell'_j(0)) \mathbf{y}_{n+1-j} = h \mathbf{f}_{n+1},$$

das **BDF(k)-Verfahren**. In Standardform:

$$\mathbf{y}_{n+1} + \sum_{j=1}^k \frac{\ell'_j(0)}{\ell'_0(0)} \mathbf{y}_{n+j-1} = -\frac{h}{\ell'_0(0)} \mathbf{f}_{n+1}.$$

Nach Konstruktion besitzt es die Konsistenzordnung k .

Satz 3.16 (Stabilität von BDF-Verfahren)

Das BDF(k)-Verfahren ist genau dann stabil (und damit konvergent), wenn $k \leq 6$.

Lineare Mehrschrittverfahren

BDF-Verfahren

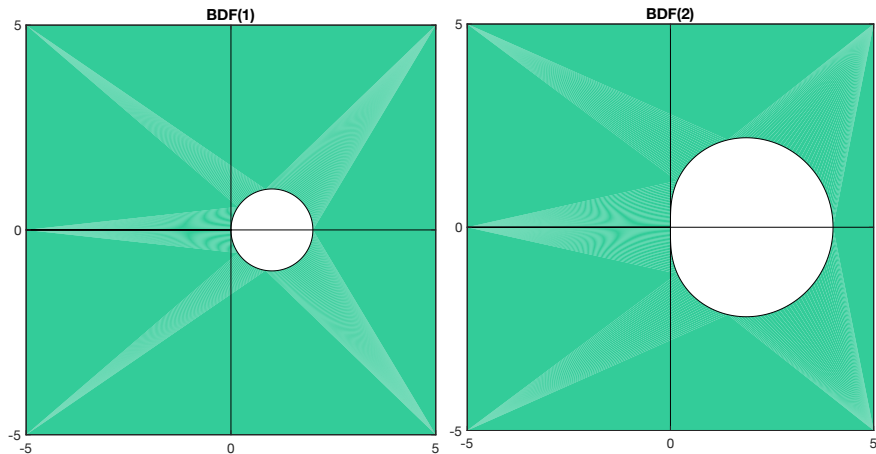
Koeffizienten und Fehlerkonstanten der BDF-Verfahren:

k	α_6	α_5	α_4	α_3	α_2	α_1	α_0	β_k	C_{k+1}
1						1	-1	1	$-\frac{1}{2}$
2					1	$-\frac{4}{3}$	$\frac{1}{3}$	$\frac{2}{3}$	$-\frac{2}{9}$
3				1	$-\frac{18}{11}$	$\frac{9}{11}$	$-\frac{2}{11}$	$\frac{6}{11}$	$-\frac{3}{22}$
4			1	$-\frac{48}{25}$	$\frac{36}{25}$	$-\frac{16}{25}$	$\frac{3}{25}$	$\frac{12}{25}$	$-\frac{12}{125}$
5		1	$-\frac{300}{137}$	$\frac{300}{137}$	$-\frac{200}{137}$	$\frac{75}{137}$	$-\frac{12}{137}$	$\frac{60}{137}$	$-\frac{10}{137}$
6	1	$-\frac{360}{147}$	$\frac{450}{147}$	$-\frac{400}{147}$	$\frac{225}{147}$	$-\frac{72}{147}$	$\frac{10}{147}$	$\frac{60}{147}$	$-\frac{20}{343}$

Für $k = 1$ erhält man das **implizite Euler-Verfahren** $\mathbf{y}_{n+1} = \mathbf{y}_n + h\mathbf{f}(t_{n+1}, \mathbf{y}_{n+1})$.
BDF-Verfahren zeichnen sich durch „große“ Stabilitätsbereiche aus.

Lineare Mehrschrittverfahren

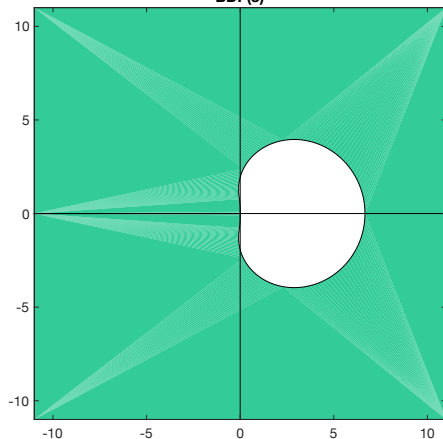
BDF-Verfahren



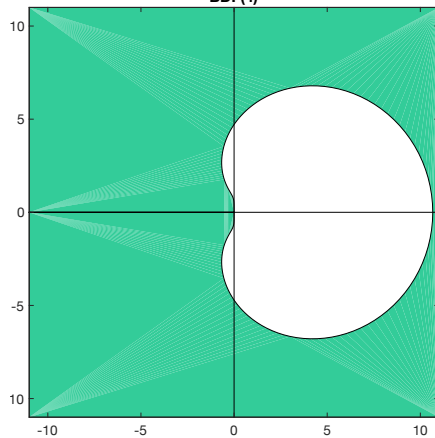
Lineare Mehrschrittverfahren

BDF-Verfahren

BDF(3)

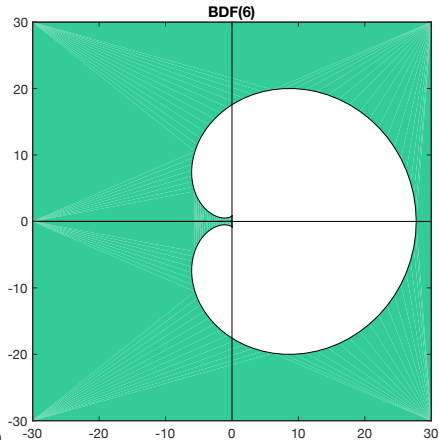
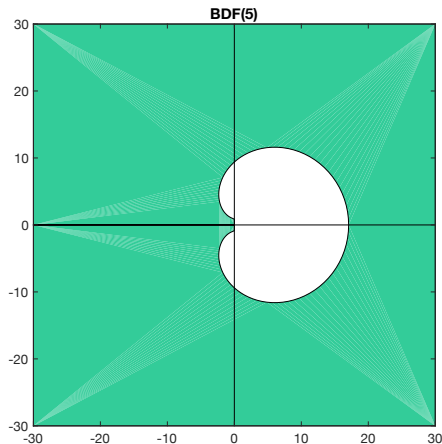


BDF(4)



Lineare Mehrschrittverfahren

BDF-Verfahren



- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ① Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ① Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Ausgangspunkt wie immer (Substitution: $s = t + \tau h$, $0 \leq \tau \leq 1$)

$$\begin{aligned} \mathbf{y}(t+h) &= \mathbf{y}(t) + [\mathbf{y}(t+h) - \mathbf{y}(t)] = \mathbf{y}(t) + \int_t^{t+h} \mathbf{y}'(s) \, ds \\ &= \mathbf{y}(t) + h \int_0^1 \mathbf{y}'(t + \tau h) \, d\tau. \end{aligned}$$

Approximiere durch Quadraturformel

$$\int_0^1 g(\tau) \, d\tau \approx \sum_{j=1}^m \beta_j g(\gamma_j). \quad (*)$$

Damit zumindest $g \equiv 1$ exakt integriert wird, fordern wir $\sum_{j=1}^m \beta_j = 1$.

Daraus folgt

$$\begin{aligned}\mathbf{y}(t+h) &\approx \mathbf{y}(t) + h \sum_{j=1}^m \beta_j \mathbf{y}'(t + \gamma_j h) \\ &= \mathbf{y}(t) + h \sum_{j=1}^m \beta_j \mathbf{f}(t + \gamma_j h, \mathbf{y}(t + \gamma_j h)).\end{aligned}\tag{RK-1}$$

Problem: $\mathbf{y}(t + \gamma_j h) = \mathbf{y}(t) + h \int_0^{\gamma_j} \mathbf{y}'(t + \tau h) d\tau$ sind unbekannt. Näherungen wieder durch Quadraturformeln, aber mit den alten Knoten γ_j ($j = 1, \dots, m$) aus (*) (sonst würden sich neue „Unbekannte“ $\mathbf{y}(t + \text{Knoten} \cdot h)$ ergeben).

$$\int_0^{\gamma_j} g(\tau) d\tau \approx \sum_{\ell=1}^m \alpha_{j,\ell} g(\gamma_\ell) \quad (j = 1, \dots, m).\tag{**}$$

Damit zumindest $g \equiv 1$ exakt integriert wird, fordern wir

$$\sum_{\ell=1}^m \alpha_{j,\ell} = \gamma_j \quad (j = 1, \dots, m).$$

Damit ergibt sich

$$\begin{aligned}\mathbf{y}(t + \gamma_j h) &\approx \mathbf{y}(t) + h \sum_{\ell=1}^m \alpha_{j,\ell} \mathbf{y}'(t + \gamma_\ell h) \\ &= \mathbf{y}(t) + h \sum_{\ell=1}^m \alpha_{j,\ell} \mathbf{f}(t + \gamma_\ell h, \mathbf{y}(t + \gamma_\ell h)).\end{aligned}\tag{RK-2}$$

Abkürzung: $\tilde{\mathbf{k}}_j := \mathbf{f}(t + \gamma_j h, \mathbf{y}(t + \gamma_j h))$ ($j = 1, \dots, m$).

(RK-2): $\tilde{\mathbf{k}}_j \approx \mathbf{f}\left(t + \gamma_j h, \mathbf{y}(t) + h \sum_{\ell=1}^m \alpha_{j,\ell} \tilde{\mathbf{k}}_\ell\right)$ ($j = 1, \dots, m$).

(RK-1): $\mathbf{y}(t + h) \approx \mathbf{y}(t) + h \sum_{j=1}^m \beta_j \tilde{\mathbf{k}}_j$.

m -stufiges Runge-Kutta-Verfahren (RKV)

$$\begin{aligned}\mathbf{y}_{n+1} &= \mathbf{y}_n + h \sum_{j=1}^m \beta_j \mathbf{k}_j \quad \text{mit} \\ \mathbf{k}_j &= \mathbf{f}\left(t_n + \gamma_j h, \mathbf{y}_n + h \sum_{\ell=1}^m \alpha_{j,\ell} \mathbf{k}_\ell\right) \quad (j = 1, \dots, m).\end{aligned}\tag{RKV}$$

Runge-Kutta-Verfahren

Konstruktion

Butcher-Tableau

(John Charles Butcher, *1933)

γ_1	$\alpha_{1,1}$	\cdots	$\alpha_{1,m}$
\vdots	\vdots		\vdots
γ_m	$\alpha_{m,1}$	\cdots	$\alpha_{m,m}$
	β_1	\cdots	β_m

Beispiele.

0	0	0
1	1	0
	$1/2$	$1/2$

symbolisiert ein zweistufiges explizites RKV, nämlich das verbesserte Euler-Verfahren.

(Ein RKV ist **explizit**, wenn $\alpha_{j,\ell} = 0 \ \forall j \leq \ell$ gilt.)

symbolisiert ein zweistufiges implizites RKV:

$$\begin{array}{c|cc} 0 & 1/4 & -1/4 \\ 2/3 & 1/4 & 5/12 \\ \hline & 1/4 & 3/4 \end{array}$$

$$\mathbf{k}_1 = \mathbf{f} \left(t_n, \mathbf{y}_n + \frac{1}{4}h\mathbf{k}_1 - \frac{1}{4}h\mathbf{k}_2 \right),$$

$$\mathbf{k}_2 = \mathbf{f} \left(t_n + \frac{2}{3}h, \mathbf{y}_n + \frac{1}{4}h\mathbf{k}_1 + \frac{5}{12}h\mathbf{k}_2 \right),$$

(„zwei“ i.A. nichtlineare Gleichungen für \mathbf{k}_1 und \mathbf{k}_2)

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{4}h(\mathbf{k}_1 + 3\mathbf{k}_2).$$

(Das Verfahren (2.2) aus Abschnitt 2 ist ein weiteres implizites zweistufiges RKV.)

Runge-Kutta-Verfahren

Konstruktion

0	0	0	0
1/2	1/2	0	0
1	-1	2	0
<hr/>			
	1/6	4/6	1/6

symbolisiert ein dreistufiges explizites RKV.

Verfahren dritter Ordnung von Kutta

$$\mathbf{k}_1 = \mathbf{f}(t_n, \mathbf{y}_n),$$

$$\mathbf{k}_2 = \mathbf{f}\left(t_n + \frac{1}{2}h, \mathbf{y}_n + \frac{1}{2}h\mathbf{k}_1\right),$$

$$\mathbf{k}_3 = \mathbf{f}(t_n + h, \mathbf{y}_n - h\mathbf{k}_1 + 2h\mathbf{k}_2),$$

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{6}h(\mathbf{k}_1 + 4\mathbf{k}_2 + \mathbf{k}_3).$$

Runge-Kutta-Verfahren

Konstruktion

0	0	0	0
1/3	1/3	0	0
2/3	0	2/3	0
	1/4	0	3/4

symbolisiert ein dreistufiges explizites RKV.

Verfahren dritter Ordnung von Heun

$$\mathbf{k}_1 = \mathbf{f}(t_n, \mathbf{y}_n),$$

$$\mathbf{k}_2 = \mathbf{f}\left(t_n + \frac{1}{3}h, \mathbf{y}_n + \frac{1}{3}h\mathbf{k}_1\right),$$

$$\mathbf{k}_3 = \mathbf{f}\left(t_n + \frac{2}{3}h, \mathbf{y}_n + \frac{2}{3}h\mathbf{k}_2\right),$$

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{4}h(\mathbf{k}_1 + 3\mathbf{k}_3).$$

(Vgl. (2.1) aus Abschnitt 2.)

Runge-Kutta-Verfahren

Konstruktion

0	0	0	0	0
1/2	1/2	0	0	0
1/2	0	1/2	0	0
1	0	0	1	0
<hr/>				
	1/6	2/6	2/6	1/6

symbolisiert ein vierstufiges explizites RKV.

Klassisches Runge-Kutta-Verfahren

$$\mathbf{k}_1 = \mathbf{f}(t_n, \mathbf{y}_n),$$

$$\mathbf{k}_2 = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{y}_n + \frac{1}{2}h\mathbf{k}_1),$$

$$\mathbf{k}_3 = \mathbf{f}(t_n + \frac{1}{2}h, \mathbf{y}_n + \frac{1}{2}h\mathbf{k}_2),$$

$$\mathbf{k}_4 = \mathbf{f}(t_n + h, \mathbf{y}_n + h\mathbf{k}_3),$$

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{6}h(\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4).$$

Alternative Form von RKV

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h \sum_{j=1}^m \beta_j \mathbf{f}(t_n + \gamma_j h, \tilde{\mathbf{y}}_j) \quad (\text{RKV}^*)$$

$$\text{mit } \tilde{\mathbf{y}}_j = \mathbf{y}_n + h \sum_{\ell=1}^m \alpha_{j,\ell} \mathbf{f}(t_n + \gamma_\ell h, \tilde{\mathbf{y}}_\ell) \quad (j = 1, \dots, m).$$

Setze $\mathbf{k}_j = \mathbf{f}(t_n + \gamma_j h, \tilde{\mathbf{y}}_j)$.

- $\mathbf{k}_j \approx \mathbf{y}'(t_n + \gamma_j h)$
- $\tilde{\mathbf{y}}_j \approx \mathbf{y}(t_n + \gamma_j h)$.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Jedes RKV hat die Form

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h\Phi_f(\mathbf{y}_n, t_n; h) \quad \text{mit} \quad \Phi_f(\mathbf{y}_n, t_n; h) = \sum_{j=1}^m \beta_j \mathbf{k}_j.$$

Es ist ein Einschrittverfahren ($\rho(\zeta) = \zeta - 1$), also stabil und (vgl. Abschnitt 3) genau dann konsistent, wenn

$$\Phi_f(\mathbf{y}(t_n), t_n; 0) = \mathbf{f}(t_n, \mathbf{y}(t_n))\rho'(1)$$

erfüllt ist, was hier zu $\sum_{j=1}^m \beta_j = 1$ äquivalent ist.

Ein RKV ist deshalb genau dann konvergent, wenn

$$\sum_{j=1}^m \beta_j = 1$$

gilt.

Runge-Kutta-Verfahren

Konsistenzordnung

Um die Konsistenzordnung eines RKVs zu bestimmen (oder um m -stufige RKV mit möglichst hoher Konsistenzordnung zu konstruieren), sind wie im Fall der Taylor-Verfahren (siehe Abschnitt 5) komplizierte Rechnungen erforderlich. Wir untersuchen als Beispiel explizite dreistufige RKV,

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ \gamma_2 & \gamma_2 & 0 & 0 \\ \gamma_3 & \gamma_3 - \alpha_{3,2} & \alpha_{3,2} & 0 \\ \hline & \beta_1 & \beta_2 & \beta_3 \end{array},$$

und entwickeln

$$\frac{1}{h} \mathbf{R}_{n+1} = \frac{\mathbf{y}(t_{n+1}) - \mathbf{y}(t_n)}{h} - \Phi_f(\mathbf{y}(t_n), t_n; h) = \frac{\mathbf{y}(t_{n+1}) - \mathbf{y}(t_n)}{h} - \sum_{j=1}^3 \beta_j \mathbf{k}_j$$

nach Potenzen von h (unter der Voraussetzung, dass \mathbf{y} bzw. f genügend oft differenzierbar sind).

Für skalare AWPes ergibt sich mit den Abkürzungen

$$F := f_t + f_y f \quad \text{und} \quad G := f_{tt} + 2f_{ty}f + f_{yy}f^2$$

(alle Ableitungen von f werden an der Stelle $(t_n, y(t_n))$ ausgewertet) die Beziehung

$$\frac{y(t_{n+1}) - y(t_n)}{h} = f + \frac{1}{2}Fh + \frac{1}{6}(G + f_y F)h^2 + O(h^3).$$

Andererseits ist

$$k_1 = f(t_n, y(t_n)) = f,$$

$$k_2 = f(t_n + h\gamma_2, y(t_n) + h\gamma_2 k_1) = f + h\gamma_2 F + \frac{1}{2}h^2\gamma_2^2 G + O(h^3),$$

$$\begin{aligned} k_3 &= f(t_n + h\gamma_3, y(t_n) + h(\gamma_3 - \alpha_{3,2})k_1 + h\alpha_{3,2}k_2) \\ &= f + h\gamma_3 F + h^2(\gamma_2\alpha_{3,2}Ff_y + \frac{1}{2}\gamma_3^2 G) + O(h^3). \end{aligned}$$

Das bedeutet:

$$\begin{aligned}\frac{1}{h} \mathbf{R}_{n+1} = & \left[1 - \sum_{j=1}^3 \beta_j \right] f + \left[\frac{1}{2} - \beta_2 \gamma_2 - \beta_3 \gamma_3 \right] Fh \\ & + \left[\left(\frac{1}{3} - \beta_2 \gamma_2^2 - \beta_3 \gamma_3^2 \right) \frac{1}{2} G + \left(\frac{1}{6} - \beta_3 \gamma_2 \alpha_{3,2} \right) F f_y \right] h^2 + O(h^3).\end{aligned}$$

Folgerungen:

1. Das Euler-Verfahren ist das einzige einstufige explizite RKV der Ordnung 1 ($\beta_1 = 1$). Es gibt kein einstufiges explizites RKV höherer Ordnung.
2. Die zweistufigen expliziten RKV der Ordnung 2 sind durch

$$\beta_1 + \beta_2 = 1 \quad \text{und} \quad \beta_2 \gamma_2 = \frac{1}{2}$$

charakterisiert. Beispiele sind das modifizierte ($\beta_1 = 0, \beta_2 = 1, \gamma_2 = \frac{1}{2}$) und das verbesserte Euler-Verfahren ($\beta_1 = \beta_2 = \frac{1}{2}, \gamma_2 = 1$). Kein explizites zweistufiges RKV besitzt die Ordnung 3.

3. Explizite dreistufige RKV der Ordnung 3 sind durch die vier Gleichungen

$$\begin{aligned}\beta_1 + \beta_2 + \beta_3 &= 1, & \beta_2\gamma_2^2 + \beta_3\gamma_3^2 &= \frac{1}{3}, \\ \beta_2\gamma_2 + \beta_3\gamma_3 &= \frac{1}{2}, & \beta_3\gamma_2\alpha_{3,2} &= \frac{1}{6}\end{aligned}$$

charakterisiert. (Man kann zeigen, dass keine dieser Methoden die Ordnung 4 besitzt.) Beispiele sind das Verfahren von Heun ($\beta_1 = \frac{1}{4}$, $\beta_2 = 0$, $\beta_3 = \frac{3}{4}$, $\gamma_2 = \frac{1}{3}$, $\gamma_3 = \alpha_{3,2} = \frac{2}{3}$) und das Verfahren von Kutta ($\beta_1 = \frac{1}{6}$, $\beta_2 = \frac{2}{3}$, $\beta_3 = \frac{1}{6}$, $\gamma_2 = \frac{1}{2}$, $\gamma_3 = 1$, $\alpha_{3,2} = 2$).

4. Ähnliche (kompliziertere) Rechnungen zeigen, dass es eine zweiparametrische Familie expliziter vierstufiger RKV der Ordnung 4 gibt, von denen keines die Ordnung 5 besitzt. Ein Beispiel ist das klassische Runge-Kutta-Verfahren.

Runge-Kutta-Verfahren

Konsistenzordnung

Weitere Beispiele sind

0	0	0	0	0
1/3	1/3	0	0	0
2/3	-1/3	1	0	0
1	1	-1	1	0
<hr/>				
	1/8	3/8	3/8	1/8

(3/8-Regel)

0	0	0	0	0
2/5	2/5	0	0	0
3/5	-3/20	3/4	0	0
1	19/44	-15/44	40/44	0
<hr/>				
	55/360	125/360	125/360	55/360

(Formel von Kuntzmann).

- Die oben beschriebene Methode, die Ordnung eines RKVs zu bestimmen, wird für Verfahren höherer Ordnung schnell unübersichtlich: Die Koeffizienten eines expliziten Verfahrens der Ordnung 3 müssen 4 Gleichungen erfüllen (s.o.), während bei einem Verfahren der Ordnung 8 bereits 200 nicht-lineare Gleichungen überprüft werden müssen.
- Die sog. Butcher-Theorie²³ erleichtert mit Hilfe graphentheoretischer Bäume die Buchhaltung bei den partiellen Ableitungen von f und erlaubt eine elegante Berechnung der Ordnung eines gegebenen RKVs (sie liefert aber keine Methode, ein Verfahren mit gewünschter Ordnung zu konstruieren).
- Wir beschränken uns hier darauf, *notwendige Ordnungsbedingungen* abzuleiten, die sich aus den speziellen AWPen

$$y' = y + t^{\ell-1}, \quad y(0) = 0, \quad \ell \in \mathbb{N},$$

ergeben.

²³J. C. Butcher, *The Numerical Analysis of Ordinary Differential Equations. Runge-Kutta and General Linear Methods*. John Wiley & Sons, Chichester 1987

Satz 4.1 (Notwendige Ordnungsbedingungen für RKV)

Das durch das Butcher-Tableau

$$\begin{array}{c|c} \mathbf{c} & A \\ \hline & \mathbf{b}^\top \end{array}$$

definierte RKV besitze die Ordnung p . Dann gelten

$$\mathbf{b}^\top A^k C^{\ell-1} \mathbf{e} = \frac{(\ell-1)!}{(\ell+k)!} = \frac{1}{\ell(\ell+1)\dots(\ell+k)}$$

für $\ell = 1, 2, \dots, p$ und $k = 0, 1, \dots, p - \ell$.

Dabei sind $\mathbf{b} := [\beta_1, \beta_2, \dots, \beta_m]^\top$, $A := [\alpha_{j,\nu}]_{1 \leq j, \nu \leq m}$,
 $C := \text{diag}(\gamma_1, \gamma_2, \dots, \gamma_m)$, und $\mathbf{e} := [1, 1, \dots, 1]^\top \in \mathbb{R}^m$.

Spezialfälle der notwendigen Bedingungen aus Satz 4.1 sind (für $k = 0$)

$$\mathbf{b}^\top C^{\ell-1} \mathbf{e} = \sum_{j=1}^m \beta_j \gamma_j^{\ell-1} = \frac{1}{\ell} \quad \text{für } \ell = 1, 2, \dots, p$$

sowie (für $\ell = 1$ mit $k \leftarrow k + 1$)

$$\mathbf{b}^\top A^{k-1} \mathbf{e} = \frac{1}{k!} \quad \text{für } k = 1, 2, \dots, p.$$

Bemerkung. Ein explizites m -stufiges RKV besitzt höchstens die Konsistenzordnung m , denn hier ist $A^m = O$ (A ist echte untere Dreiecksmatrix). Für die optimale Ordnung $p(m)$ eines expliziten m -stufigen RKVs gilt sogar $p(m) \leq m - 1$ falls $m \geq 5$, genauer:

m	1	2	3	4	5	6	7	8	9	10	11	12
$p(m)$	1	2	3	4	4	5	6	6	7	7	8	9

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Wir wenden ein m -stufiges RKV auf die Testgleichung $y' = \lambda y$ an und erhalten

$$y_{n+1} = \left[1 + \hat{h} \mathbf{b}^\top (I_m - \hat{h} A)^{-1} \mathbf{e} \right] y_n =: R(\hat{h}) y_n, \quad (\hat{h} = \lambda h)$$

so dass (bei festem h)

$$\lim_{n \rightarrow \infty} y_n = 0 \text{ (für alle } y_0) \Leftrightarrow |R(\hat{h})| < 1.$$

In Analogie zu Abschnitt 6 definieren wir den **Stabilitätsbereich** eines RKVs durch

$$\mathcal{R}_A := \{ \hat{h} \in \mathbb{C} : |R(\hat{h})| < 1 \}.$$

Für ein beliebiges m -stufiges RKV gilt

$$R(\hat{h}) = 1 + \hat{h} \mathbf{b}^\top (I_m - \hat{h} A)^{-1} \mathbf{e} = 1 + \sum_{j=1}^{\infty} \hat{h}^j \mathbf{b}^\top A^{j-1} \mathbf{e}.$$

- Besitzt das Verfahren die Ordnung p , so folgt

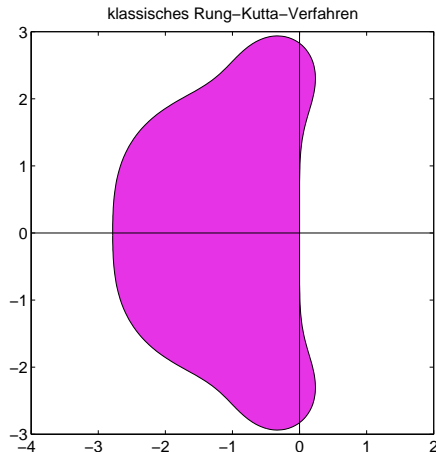
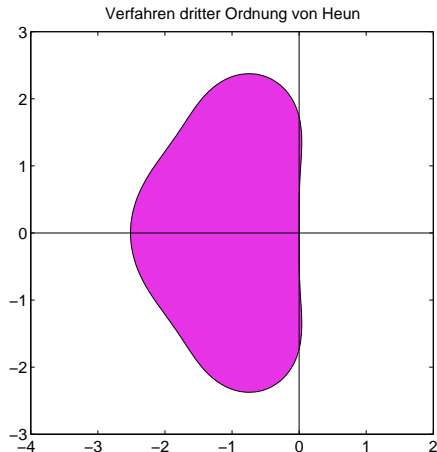
$$R(\hat{h}) = \sum_{j=0}^p \frac{1}{j!} \hat{h}^j + \sum_{j=p+1}^{\infty} \hat{h}^j \mathbf{b}^{\top} A^{j-1} \mathbf{e}.$$

- Ist das RKV explizit, so folgt

$$R(\hat{h}) = 1 + \sum_{j=1}^m \hat{h}^j \mathbf{b}^{\top} A^{j-1} \mathbf{e}.$$

- Insbesondere hängt der Stabilitätsbereich eines m -stufigen expliziten RKVs der Ordnung m ($1 \leq m \leq 4$) wegen $R(\hat{h}) = \sum_{j=0}^m \frac{1}{j!} \hat{h}^j$ nicht von den Koeffizienten des Verfahrens ab.
- Außerdem besitzt kein explizites RKV einen unbeschränkten Stabilitätsbereich (denn R ist in diesem Fall ein Polynom).

Beispiele für Stabilitätsgebiete zweier RKV:



- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Runge-Kutta-Verfahren

Eingebettete Runge-Kutta-Verfahren

- Kein Verfahren zur Lösung von AWPen arbeitet in der Praxis mit einer *konstanten* Schrittweite.
- Man wird vielmehr versuchen, die Schrittweite an das Verhalten der Lösung y anzupassen (ändert sich y in einem Bereich schnell, so ist dort eine kleine Schrittweite angebracht; in Bereichen, in denen y kaum variiert, ist eine größere Schrittweite ausreichend).
- Wir werden hier eine **Schrittweitensteuerung** vorstellen, die zum Ziel hat, den Konsistenzfehler $K_{n+1} := \frac{1}{h} R_{n+1}$ (wird in der Literatur oft lokaler Diskretisierungsfehler genannt, vgl. Abschnitt 3) zu kontrollieren:

$$\|K_n\| \sim \text{tol}, \quad n = 1, 2, \dots,$$

mit einer vorgebenen Toleranz tol .

- Bei Systemen von DGen (insbesondere dann, wenn die Lösungskomponenten von unterschiedlicher Größenordnung sind) wird man für jede Komponente eine eigene absolute Fehlertoleranz und global eine relative Fehlertoleranz festsetzen.

Runge-Kutta-Verfahren

Eingebettete Runge-Kutta-Verfahren

Das folgende Lemma besagt, dass mit dem Konsistenzfehler auch der (eigentlich interessierende) globale Diskretisierungsfehler kontrolliert wird.

Lemma 4.2

Für den globalen Diskretisierungsfehler $e_n = \mathbf{y}(t_n) - \mathbf{y}_n$ eines Einschrittverfahrens gilt

$$\|e_n\| \leq (t_n - t_0) \kappa_n \exp(M(t_n - t_0)).$$

Dabei ist $\kappa_n := \max\{\|\mathbf{K}_j\| : j = 0, 1, \dots, n\}$ und M die Lipschitzkonstante der Verfahrensfunktion (vgl. (V₂) aus Abschnitt 2).

Um den Konsistenzfehler zu schätzen, verwendet man **zwei Methoden unterschiedlicher Konsistenzordnungen** (sagen wir p und q mit $p < q$), um \mathbf{y}_n aus \mathbf{y}_{n-1} zu berechnen:

$$\mathbf{y}_n = \mathbf{y}_{n-1} + h\Phi_f(\mathbf{y}_{n-1}, t_{n-1}; h) \quad \text{bzw.}$$

$$\hat{\mathbf{y}}_n = \mathbf{y}_{n-1} + h\hat{\Phi}_f(\mathbf{y}_{n-1}, t_{n-1}; h).$$

Für die zugehörigen Konsistenzfehler gelten (Lokalisierungsannahme):

$$\mathbf{K}_n = \frac{\mathbf{y}(t_n) - \mathbf{y}(t_{n-1})}{h} - \Phi_f(\mathbf{y}(t_{n-1}), t_{n-1}; h) = O(h^p),$$

$$\hat{\mathbf{K}}_n = \frac{\mathbf{y}(t_n) - \mathbf{y}(t_{n-1})}{h} - \hat{\Phi}_f(\mathbf{y}(t_{n-1}), t_{n-1}; h) = O(h^q).$$

Daraus folgt

$$\mathbf{K}_n - \hat{\mathbf{K}}_n = \hat{\Phi}_f(\mathbf{y}(t_{n-1}), t_{n-1}; h) - \Phi_f(\mathbf{y}(t_{n-1}), t_{n-1}; h) = \frac{1}{h}(\hat{\mathbf{y}}_n - \mathbf{y}_n) + O(h^p).$$

Wegen $\mathbf{K}_n - \hat{\mathbf{K}}_n = \mathbf{K}_n(1 + O(h^{q-p})) \sim \mathbf{K}_n$ erhalten wir aus

$$\frac{1}{h} \|\mathbf{y}_n - \hat{\mathbf{y}}_n\| \sim \|\mathbf{K}_n\|$$

eine (grobe) Schätzung für $\|\mathbf{K}_n\|$.

- Ist $\frac{1}{h} \|\mathbf{y}_n - \hat{\mathbf{y}}_n\| > \text{tol}$, so wird die Schrittweite h verworfen und mit

$$\left(\frac{\tilde{h}}{h}\right)^p = \alpha \frac{h \text{ tol}}{\|\mathbf{y}_n - \hat{\mathbf{y}}_n\|} \quad (*)$$

eine neue Schrittweite \tilde{h} bestimmt (α ist hier ein Sicherheitsfaktor, etwa $\alpha = 0.9$).

- Ausgehend von \mathbf{y}_{n-1} werden jetzt neue Näherungen \mathbf{y}_n und $\hat{\mathbf{y}}_n$ (an der Stelle $t_{n-1} + \tilde{h}$) berechnet. Diesen Prozess wiederholt man so lange, bis $\frac{1}{\tilde{h}} \|\mathbf{y}_n - \hat{\mathbf{y}}_n\| \leq \text{tol}$ erfüllt ist. Dann wird $(*)$ verwendet, um eine neue (größere) Schrittweite für den nächsten Schritt ($n \rightarrow n+1$) vorzuschlagen.
- Die Wahl von \tilde{h} nach $(*)$ motiviert sich folgendermaßen:

$$\begin{array}{ll} \text{benutzte Schrittweite } h: & \frac{1}{h} \|\mathbf{y}_n - \hat{\mathbf{y}}_n\| \sim \|\mathbf{K}_n\| = ch^p + O(h^{p+1}) \sim ch^p, \\ \text{erwünschte Schrittweite } \tilde{h}: & \text{tol} = \|\mathbf{K}_n\| = c\tilde{h}^p + O(\tilde{h}^{p+1}) \sim c\tilde{h}^p. \end{array}$$

Runge-Kutta-Verfahren

Eingebettete Runge-Kutta-Verfahren

- Um den Aufwand in Grenzen zu halten, verwendet man zur Berechnung von y_n und \hat{y}_n zwei RKV (verschiedener Ordnungen), deren Butcher-Tableaus sich nur im Vektor b unterscheiden (d.h. A und c sind für beide Verfahren gleich, so dass die Größen k_j nur einmal berechnet werden müssen).
- Man spricht von **eingebetteten RKV** und schreibt

$$\begin{array}{c|c} c & A \\ \hline & b^\top \\ \hline & \hat{b}^\top \end{array}, \quad \text{z.B.} \quad \begin{array}{c|cc} & 0 & 0 \\ & 1 & 0 \\ \hline & 1 & 0 \\ \hline & 1/2 & 1/2 \end{array}.$$

Im Beispiel wird ein RKV der Ordnung 1 (das Euler-Verfahren) in ein RKV der Ordnung 2 (das verbesserte Euler-Verfahren) eingebettet.

Runge-Kutta-Verfahren

Eingebettete Runge-Kutta-Verfahren

Ein populäres Beispiel ist die **Fehlberg 4(5)-Formel**:

0	0	0	0	0	0	0
$\frac{1}{4}$	$\frac{1}{4}$	0	0	0	0	0
$\frac{3}{8}$	$\frac{3}{32}$	$\frac{9}{32}$	0	0	0	0
$\frac{12}{13}$	$\frac{1932}{2197}$	$-\frac{7200}{2197}$	$\frac{7296}{2197}$	0	0	0
1	$\frac{439}{216}$	-8	$\frac{3680}{513}$	$-\frac{845}{4104}$	0	0
$\frac{1}{2}$	$-\frac{8}{27}$	2	$-\frac{3544}{2565}$	$\frac{1859}{4104}$	$-\frac{11}{40}$	0
	$\frac{25}{216}$	0	$\frac{1408}{2565}$	$\frac{2197}{4104}$	$-\frac{1}{5}$	0
	$\frac{16}{135}$	0	$\frac{6656}{12825}$	$\frac{28561}{56430}$	$-\frac{9}{50}$	$\frac{2}{55}$

Hier werden zwei sechsstufige RKV der Ordnungen 4 bzw. 5 kombiniert.

Runge-Kutta-Verfahren

Eingebettete Runge-Kutta-Verfahren

Ein weiteres Beispiel ist das Verfahren von **Dormand-Prince 4(5)**:

0	0	0	0	0	0	0	0
$\frac{1}{5}$	$\frac{1}{5}$	0	0	0	0	0	0
$\frac{3}{10}$	$\frac{3}{40}$	$\frac{9}{40}$	0	0	0	0	0
$\frac{4}{5}$	$\frac{44}{45}$	$-\frac{56}{15}$	$\frac{32}{9}$	0	0	0	0
$\frac{8}{9}$	$\frac{19372}{6561}$	$-\frac{25360}{2187}$	$\frac{64448}{6561}$	$-\frac{212}{729}$	0	0	0
1	$\frac{9017}{3168}$	$-\frac{355}{33}$	$\frac{46732}{5247}$	$\frac{49}{176}$	$-\frac{5103}{18656}$	0	0
1	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	0
	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	0
	$\frac{5179}{57600}$	0	$\frac{7571}{16695}$	$\frac{393}{640}$	$-\frac{92097}{339200}$	$\frac{187}{2100}$	$\frac{1}{40}$

Hier wird ein sechstufiges RKV der Ordnung 4 in ein siebenstufiges RKV der Ordnung 5 eingebettet.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Runge-Kutta-Verfahren

Implizite und halb-implizite Verfahren

- Ist die Matrix A eines m -stufigen RKVs **keine** echte untere Δ -Matrix (ist das RKV also **implizit**), so muss in jedem Zeitschritt ein nicht-lineares Gleichungssystem der Form

$$\begin{array}{rcl} \mathbf{k}_1 & = & \mathbf{f}(t_n + \gamma_1 h, h \sum_{\ell=1}^m \alpha_{1,\ell} \mathbf{k}_\ell) \\ \vdots & & \vdots \\ \mathbf{k}_m & = & \mathbf{f}(t_n + \gamma_m h, h \sum_{\ell=1}^m \alpha_{m,\ell} \mathbf{k}_\ell) \end{array} \quad (\square)$$

gelöst werden.

- Dieses System hat also mn Unbekannte, wenn $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ aus n Gleichungen besteht.
- Mit Hilfe des Banachschen Fixpunktsatzes erkennt man, dass (\square) für genügend kleine h eindeutig lösbar ist.
- In der Praxis wird man dieses System aber nicht mit der Fixpunktiteration, sondern mit einem Newton- bzw. Quasi-Newton-Verfahren lösen.

Runge-Kutta-Verfahren

Implizite und halb-implizite Verfahren

- Ist A eine untere (aber keine echte untere) Δ -Matrix, so nennt man das zugehörige RKV **halb-implizit**. Das System (\square) zerfällt dann in m Systeme mit jeweils n Unbekannten.
- Implizite RKV werden oft mit Hilfe von **Gauß-Quadraturformeln** konstruiert. Dies sind Formeln der Bauart

$$\int_a^b g(\tau) \, d\tau = \sum_{j=1}^m \beta_j g(\gamma_j) + R_m(g).$$

Hier werden die Gewichte β_j und Knoten γ_j so gewählt, dass $R_m(p) = 0$ für Polynome p möglichst hohen Grades d erfüllt ist. Man kann zeigen (vgl. Numerik I), dass eine optimale Wahl auf $d = 2m - 1$ führt (man sagt die Quadraturformel hat **Exaktheitsgrad** $2m - 1$).

- Die zugehörigen RKV (auch sie werden **Gauß-Formeln** genannt) haben die Ordnung $2m$. Beachte, dass kein m -stufiges RKV eine höhere Ordnung besitzen kann. Warum?

Runge-Kutta-Verfahren

Implizite und halb-implizite Verfahren

Für $m = 1$ ergibt sich die **implizite Mittelpunktsregel** $\frac{1/2}{1} \mid \frac{1/2}{1}$,
welche die Ordnung 2 besitzt.

Für $m = 2$ und $m = 3$ ergeben sich die Gauß-Formeln

$$\begin{array}{c|cc} \frac{3-\sqrt{3}}{6} & \frac{1}{4} & \frac{3-2\sqrt{3}}{12} \\ \frac{3+\sqrt{3}}{6} & \frac{3+2\sqrt{3}}{12} & \frac{1}{4} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array} \quad \text{bzw.} \quad \begin{array}{c|ccc} \frac{5-\sqrt{15}}{10} & \frac{5}{36} & \frac{10-3\sqrt{15}}{45} & \frac{25-6\sqrt{15}}{180} \\ \frac{1}{2} & \frac{10+3\sqrt{15}}{72} & \frac{2}{9} & \frac{10-3\sqrt{15}}{72} \\ \frac{5+\sqrt{15}}{10} & \frac{25+6\sqrt{15}}{180} & \frac{10+3\sqrt{15}}{45} & \frac{5}{36} \\ \hline & \frac{5}{18} & \frac{4}{9} & \frac{5}{18} \end{array}$$

mit den Konsistenzordnungen 4 bzw. 6.

Bei **Gauß-Radau-Integrationsformeln** wählt man einen Knoten als (entweder linken oder rechten) Endpunkt des Integrationsintervalls.

Runge-Kutta-Verfahren

Implizite und halb-implizite Verfahren

- Die übrigen Knoten und alle Gewichte werden so bestimmt, dass sich ein möglichst hoher Exaktheitsgrad ergibt.
- Man kann zeigen, dass eine Gauß-Radau-Formel mit m Knoten den Exaktheitsgrad $2m - 2$ besitzt. Daher haben die zugehörigen impliziten RKV die Konsistenzordnung $2m - 1$.
- Zu dieser Klasse gehören die Verfahren

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array} \quad \text{und} \quad \begin{array}{c|cc} \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ \hline 1 & \frac{3}{4} & \frac{1}{4} \\ \hline & \frac{3}{4} & \frac{1}{4} \end{array}$$

- Schließlich kann man noch beide Enden des Integrationsintervalls als Knoten wählen und die übrigen Daten so bestimmen, dass die zugehörige Integrationsformel den Exaktheitsgrad $2m - 3$ (bzw. das zugehörige implizite RKV die Konsistenzordnung $2m - 2$) besitzt. Man spricht von **Gauß-Lobatto-Formeln**. Ein Beispiel ist die Trapezregel (für $m = 2$).

Runge-Kutta-Verfahren

Implizite und halb-implizite Verfahren

- Ein Vorteil von impliziten gegenüber expliziten RKV ist ihr wesentlich größerer Stabilitätsbereich (wird ausführlicher im nächsten Kapitel diskutiert).
- Wir betrachten die Trapezregel

$$\begin{array}{c|cc} 0 & 0 & \\ 1 & 1/2 & 1/2 \\ \hline & 1/2 & 1/2 \end{array} \quad \text{d.h.} \quad \begin{aligned} \mathbf{k}_1 &= \mathbf{f}(t_n, \mathbf{y}_n), \\ \mathbf{k}_2 &= \mathbf{f}(t_{n+1}, \mathbf{y}_n + h(\mathbf{k}_1 + \mathbf{k}_2)/2), \\ \mathbf{y}_{n+1} &= \mathbf{y}_n + h(\mathbf{k}_1 + \mathbf{k}_2)/2, \end{aligned}$$

oder kürzer:

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h(\mathbf{f}(t_n, \mathbf{y}_n) + \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}))/2.$$

Die zugehörige Stabilitätsfunktion ist $R(\hat{h}) = (1 + \hat{h}/2)/(1 - \hat{h}/2)$ und es gilt:

$$|R(\hat{h})| < 1 \Leftrightarrow |1 + \hat{h}/2| < |1 - \hat{h}/2| \Leftrightarrow \operatorname{Re} \hat{h} < 0.$$

Die Trapezregel ist daher **A-stabil**.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
 - 1.1 Konstruktion
 - 1.2 Konsistenzordnung
 - 1.3 Absolute Stabilität
 - 1.4 Eingebettete Runge-Kutta-Verfahren
 - 1.5 Implizite und halb-implizite Verfahren
 - 1.6 Kollokationsmethoden
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Kollokationsmethoden sind spezielle implizite RKV, die – auf Grund ihrer Konstruktion – sehr viel leichter zu analysieren sind als allgemeine RKV.

Mit gegebenen

$$0 \leq \gamma_1 < \gamma_2 < \dots < \gamma_m \leq 1 \text{ setzen wir } t^{(j)} := t_n + \gamma_j h \quad (j = 1, 2, \dots, m)$$

und suchen ein „Polynom“ \mathbf{p} vom Grad $\leq m$ (bei Systemen von k DGen ist das ein „Vektor“ $\mathbf{p} = [p_1, p_2, \dots, p_k]^\top$ aus k Polynomen vom Grad höchstens m), das die Interpolationsbedingungen

$$\mathbf{p}(t_n) = \mathbf{y}_n, \quad \mathbf{p}'(t^{(j)}) = \mathbf{f}(t^{(j)}, \mathbf{p}(t^{(j)})) \quad (j = 1, 2, \dots, m)$$

erfüllt. Die Näherung \mathbf{y}_{n+1} an der Stelle t_{n+1} wird dann definiert durch

$$\mathbf{y}_{n+1} := \mathbf{p}(t_{n+1}).$$

\mathbf{p}' ist ein Polynom vom Grad $m - 1$, das durch die letzten m dieser Interpolationsbedingungen eindeutig bestimmt ist. In Lagrange-Form besitzt es die Darstellung

$$\mathbf{p}'(t_n + \tau h) = \sum_{j=1}^m \ell_j(t_n + \tau h) \mathbf{k}_j$$

$$\text{mit } \ell_j(t_n + \tau h) = \prod_{j \neq i=1}^m \frac{\tau - \gamma_i}{\gamma_j - \gamma_i} \quad \text{und} \quad \mathbf{k}_j := \mathbf{p}'(t^{(j)}).$$

Jetzt folgt für jedes $i \in \{1, 2, \dots, m\}$

$$\begin{aligned} \mathbf{p}(t^{(i)}) - \mathbf{p}(t_n) &= \int_{t_n}^{t^{(i)}} \mathbf{p}'(t_n + sh) \, ds = \int_{t_n}^{t^{(i)}} \sum_{j=1}^m \ell_j(t_n + sh) \mathbf{k}_j \, ds \\ &= h \sum_{j=1}^m \left(\int_0^{\gamma_i} \ell_j(r) \, dr \right) \mathbf{k}_j =: h \sum_{j=1}^m \alpha_{i,j} \mathbf{k}_j. \end{aligned}$$

Analog:

$$\mathbf{p}(t_{n+1}) - \mathbf{p}(t_n) = h \sum_{j=1}^m \left(\int_0^1 \ell_j(r) \, dr \right) \mathbf{k}_j =: h \sum_{j=1}^m \beta_j \mathbf{k}_j.$$

Daraus folgt

$$\mathbf{y}_{n+1} = \mathbf{p}(t_{n+1}) = \mathbf{p}(t_n) + h \sum_{j=1}^m \beta_j \mathbf{k}_j = \mathbf{y}_n + h \sum_{j=1}^m \beta_j \mathbf{k}_j$$

mit $\mathbf{k}_j = \mathbf{p}'(t^{(j)}) = \mathbf{f}(t^{(j)}, \mathbf{p}(t^{(j)})) = \mathbf{f}(t_n + \gamma_j h, \mathbf{y}_n + \sum_{\ell=1}^m \alpha_{j,\ell} \mathbf{k}_\ell)$.

Mit anderen Worten: Jedes Kollokationsverfahren ist ein (implizites) RKV. Implementiert wird es in der Form (RKV) bzw. (RKV*), d.h. zu gegebenen γ_j ($j = 1, 2, \dots, m$) bestimmt man zunächst

$$\alpha_{i,j} = \int_0^{\gamma_i} \ell_j(r) \, dr \quad \text{und} \quad \beta_j = \int_0^1 \ell_j(r) \, dr \quad (i, j = 1, 2, \dots, m).$$

Nicht jedes implizite RKV ist ein Kollokationsverfahren.

Beispiel.

0	1/4	-1/4
2/3	1/4	-5/12
	1/4	3/4

 repräsentiert kein Kollokationsverfahren.

Satz 4.3 (Konsistenzordnung bei Kollokationsverfahren)

Für ein m -stufiges Kollokationsverfahren mit dem Butcher-Tableau

$$\begin{array}{c|c} \mathbf{c} & A \\ \hline & \mathbf{b}^\top \end{array}$$

sind die folgenden drei Aussagen äquivalent:

(a) Das Verfahren besitzt die Konsistenzordnung $m + p$.

(b) $\int_0^1 \tau^j \prod_{j=1}^m (\tau - \gamma_j) d\tau = 0$ für $j = 0, 1, \dots, p - 1$.

(c) $\mathbf{b}^\top C^{\ell-1} \mathbf{e} = 1/\ell$ für $\ell = 1, 2, \dots, m + p$.

Dabei sind $C := \text{diag}(\gamma_1, \gamma_2, \dots, \gamma_m)$ und $\mathbf{e} := [1, 1, \dots, 1]^\top \in \mathbb{R}^m$.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Es gibt keine zufriedenstellende Definition der Bauart

„eine DG heißt steif, wenn ...“.

Wir beschreiben verschiedene Aspekte des Phänomens „Steifheit einer DG“ an Beispielen.

Beispiel 1. Die beiden AWP

$$\mathbf{y}' = \begin{bmatrix} -2 & 1 \\ 1 & -2 \end{bmatrix} \mathbf{y} + \begin{bmatrix} 2 \sin t \\ 2(\cos t - \sin t) \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad (\text{AWP}_1)$$

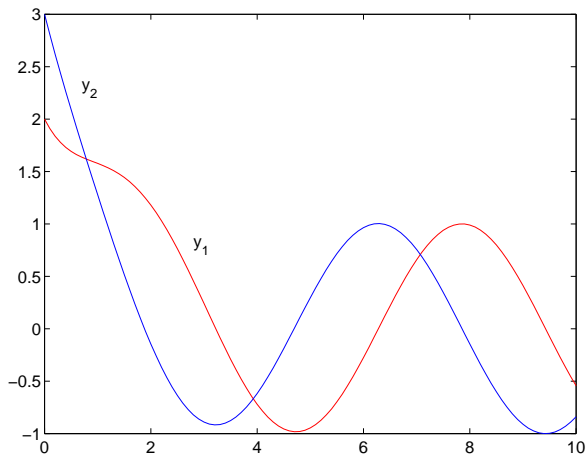
$$\mathbf{y}' = \begin{bmatrix} -2 & 1 \\ 998 & -999 \end{bmatrix} \mathbf{y} + \begin{bmatrix} 2 \sin t \\ 999(\cos t - \sin t) \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad (\text{AWP}_2)$$

besitzen dieselbe Lösung.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

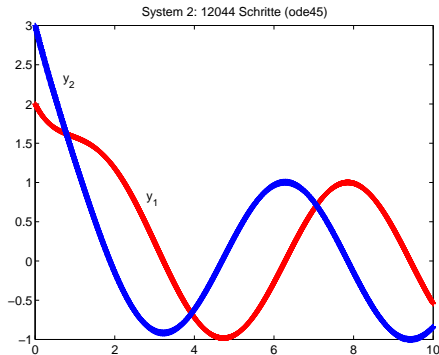
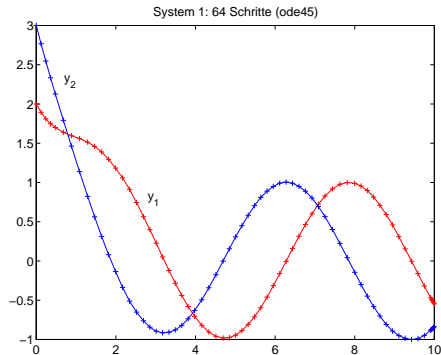
$$\mathbf{y}(t) = \exp(-t) \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} \sin t \\ \cos t \end{bmatrix}.$$



Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Wir lösen beide Probleme für $t \in [0, 10]$ mit der Matlab-Routine `ode45` (basierend auf dem Verfahren von Dormand-Prince, einem eingebettetes RKV, bei dem zwei RKV der Ordnungen 4 bzw. 5 kombiniert werden), wobei wir eine relative Toleranz von 0.01 vorgeben.



Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Obwohl die exakte Lösung in beiden Fällen dieselbe ist, erfordert die Lösung von (AWP_2) etwa 200-mal mehr Aufwand als die von (AWP_1):

AWP	h_{\min}	h_{\max}	h_{\emptyset}	Schritte
(AWP_1)	2.27e-2	2.03e-1	1.56e-1	64
(AWP_2)	5.98e-4	2.88e-3	8.30e-3	12044

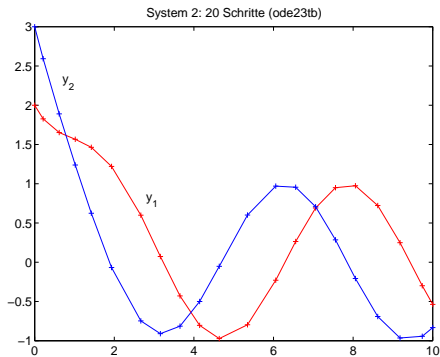
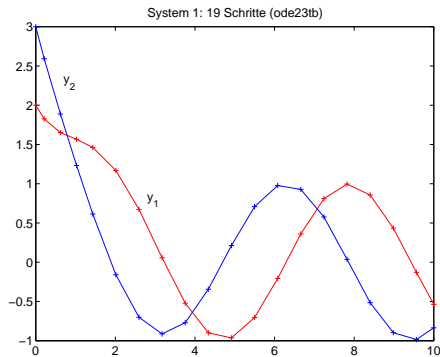
Lösen wir stattdessen beide Probleme mit der Matlab-Routine `ode23tb` (die die Trapezregel mit der BDF-Formel der Ordnung 2 kombiniert), so ergibt sich (relative Toleranz wie oben = 0.01):

AWP	h_{\min}	h_{\max}	h_{\emptyset}	Schritte
(AWP_1)	2.13e-1	5.81e-1	5.26e-1	19
(AWP_2)	2.13e-1	7.36e-1	5.00e-1	20

Beide Probleme werden ohne Schwierigkeiten gelöst.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?



Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Die Unterschiede zwischen den beiden AWPen werden sichtbar, wenn wir alle Lösungen der zugehörigen Systeme betrachten. Im ersten Fall ergibt sich

$$\mathbf{y} = \kappa_1 \exp(-t) \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \kappa_2 \exp(-3t) \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} \sin t \\ \cos t \end{bmatrix},$$

während im zweiten Fall

$$\mathbf{y} = \kappa_1 \exp(-t) \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \kappa_2 \exp(-1000t) \begin{bmatrix} 1 \\ -998 \end{bmatrix} + \begin{bmatrix} \sin t \\ \cos t \end{bmatrix}$$

die allgemeine Lösung ist. (Beide Systeme sind inhomogen und linear mit konstanten Koeffizienten, wobei die Eigenwerte der Koeffizientenmatrix

- im ersten Fall $\lambda_1 = -1$, $\lambda_2 = -3$
- und im zweiten Fall $\lambda_1 = -1$, $\lambda_2 = -1000$

sind.) Die Lösungen der konkreten AWPen sind jeweils durch $\kappa_1 = 2$ und $\kappa_2 = 0$ (!) gegeben.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

- In beiden Fällen ist die Lösung aus einem Gleichgewichtsanteil, nämlich $[\sin t, \cos t]^\top$, und einem transienten Anteil der Form $\exp(\lambda_1 t) \mathbf{a}_1 + \exp(\lambda_2 t) \mathbf{a}_2$ zusammengesetzt. Die Unterschiede liegen in der Geschwindigkeit, mit der der transiente Teil (für $t \rightarrow \infty$) verschwindet.
- Obwohl in der exakten Lösung der AWP_e gar nicht erkennbar ist, wie schnell der abklingende Anteil verschwindet ($\kappa_2 = 0$), bestimmt er die erforderliche Schrittweite.
- Würde man in (AWP₂) die AB in $\mathbf{y}(0) = [0, 1]^\top$ ändern, so würde $\kappa_1 = \kappa_2 = 0$ folgen, d.h. der transiente Teil ist in der Lösung $\mathbf{y}(t) = [\sin t, \cos t]^\top$ vollständig unsichtbar. Selbst dann wäre ode45 nur mit extrem kleiner Schrittweite in der Lage, die Aufgabe zu lösen.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

- Wir betrachten die Stabilitätsbereiche \mathcal{R}_A : Für ode45 kann man zeigen, dass $\mathcal{R}_A \cap \mathbb{R} \sim [-3, 0]$ (**Stabilitätsintervall**). Will man für (**AWP₁**) also $\lambda h \in \mathcal{R}_A$ (für alle Eigenwerte λ) garantieren, genügt es, $h < 1$ zu fordern (durchschn. Schrittweite $h_\emptyset \sim 0.156$ resultiert aus der geforderten Genauigkeit!).
- Dagegen muss für (**AWP₂**) $h < 0.003$ gefordert werden, damit $\lambda h \in \mathcal{R}_A$ für alle λ folgt und es sind Stabilitätsforderungen, die zu der kleinen durchschnittlichen Schrittweite von $h_\emptyset \sim 0.008$ führen.
- Das Verfahren ode23tb ist absolut stabil, so dass hier λh für alle $h > 0$ (in jedem der beiden Fälle (**AWP₁**) und (**AWP₂**)) im Stabilitätsbereich liegt.

Für $y' = Ay + b(t)$ heißt

$$\frac{\max_{\lambda \in \Lambda(A)} |\operatorname{Re} \lambda|}{\min_{\lambda \in \Lambda(A)} |\operatorname{Re} \lambda|}$$

Steifigkeitsquotient des linearen DG-Systems.

Ein lineares DG-System mit konstanten Koeffizienten heißt steif, wenn seine Eigenwerte alle negativen Realteil besitzen und sein Steifigkeitsquotient groß ist.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Um die Problematik dieser gebräuchlichen „Definition“ zu erläutern, betrachten wir ein weiteres AWP (mit derselben Lösung):

$$\mathbf{y}' = \begin{bmatrix} -2 & 1 \\ -1.999 & 0.999 \end{bmatrix} \mathbf{y} + \begin{bmatrix} 2 \sin t \\ -0.999(\cos t - \sin t) \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 2 \\ 3 \end{bmatrix}. \quad (\text{AWP}_3)$$

Die Matrix dieses Systems besitzt die Eigenwerte

$$\lambda_1 = -1, \quad \lambda_2 = -0.001$$

und sein Steifigkeitsquotient beträgt 1000 (wie bei (AWP₂)).

Trotzdem hat ode45 keine ernsten Probleme:

AWP	h_{\min}	h_{\max}	h_{\emptyset}	Schritte
(AWP ₃)	1.19e-1	2.50e-1	2.27e-1	44 (ode45)
(AWP ₃)	1.25e-1	4.92e-1	4.55e-1	22 (ode23tb)

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Ob ein System steif ist oder nicht kann also nicht immer aus dem Steifigkeitsquotienten abgelesen werden.

Auch „Definitionen“ wie

„ein System ist steif, wenn $\max |\operatorname{Re} \lambda|$ groß ist (etwa $\max |\operatorname{Re} \lambda| \gg 1$)“

sind natürlich wenig hilfreich (die Variablentransformation $t \mapsto 0.001t$ macht aus (AWP₂) ein Problem, das dieselbe Steifigkeit besitzt, bei dem aber $\max |\operatorname{Re} \lambda| = 1$ gilt).

Pragmatische „Definitionen“ für Steifigkeit sind:

Ein System ist steif, wenn Stabilitätsanforderungen — und nicht Genauigkeitsanforderungen — die Größe der Schrittweite bestimmen.

Ein System heißt steif, wenn gewisse Komponenten der Lösung sehr viel schneller abklingen als andere.

Ein System heißt steif, wenn explizite Verfahren nicht funktionieren.

Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Beispiel 2. Die Differentialgleichung

$$y'(t) = \lambda(y(t) - g(t)) + g'(t) \quad (*)$$

besitzt die allgemeine Lösung $y(t) = \gamma \exp(\lambda t) + g(t)$ ($\gamma \in \mathbb{R}$).

Wir wählen g als glatte Funktion, z.B. $g(t) = \arctan t$ und $\lambda = -10$. Auch hier setzt sich die Lösung aus einem glatten ($g(t)$ = Gleichgewichtsanteil) und einem schnell abklingenden Teil ($\gamma \exp(\lambda t)$ = transienter Anteil) zusammen.

Wir wählen die AB $y(0) = 0$ und approximieren die exakte Lösung

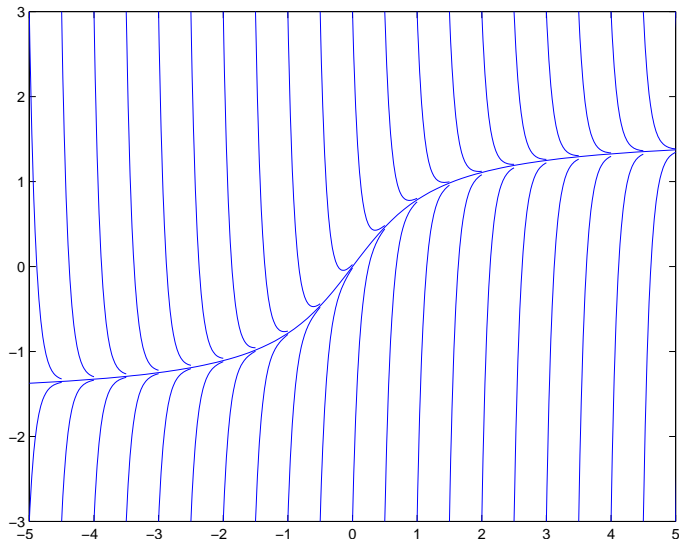
$$y(t) = \arctan t \quad \text{für} \quad t \in [0, 5]$$

mit dem

- expliziten Euler-Verfahren ($\mathcal{R}_A = \{\hat{h} : |\hat{h} + 1| < 1\}$) und dem
- impliziten Euler-Verfahren ($\mathcal{R}_A = \mathbb{C} \setminus \{\hat{h} : |\hat{h} - 1| \leq 1\}$).

Steife Differentialgleichungen

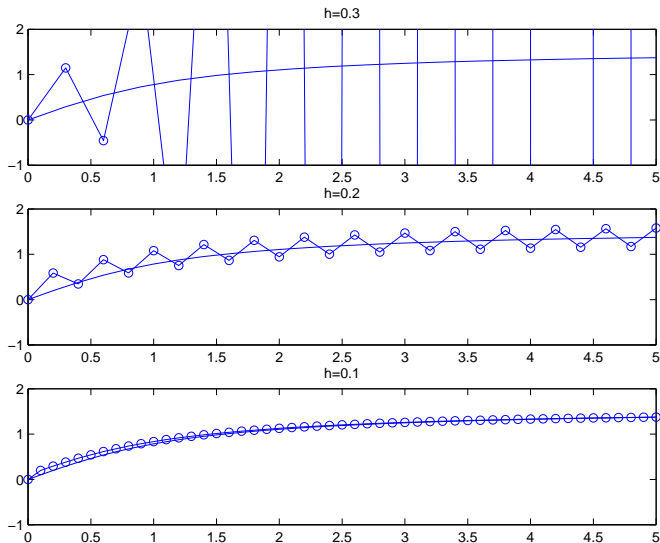
Was sind steife Differentialgleichungen?



Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

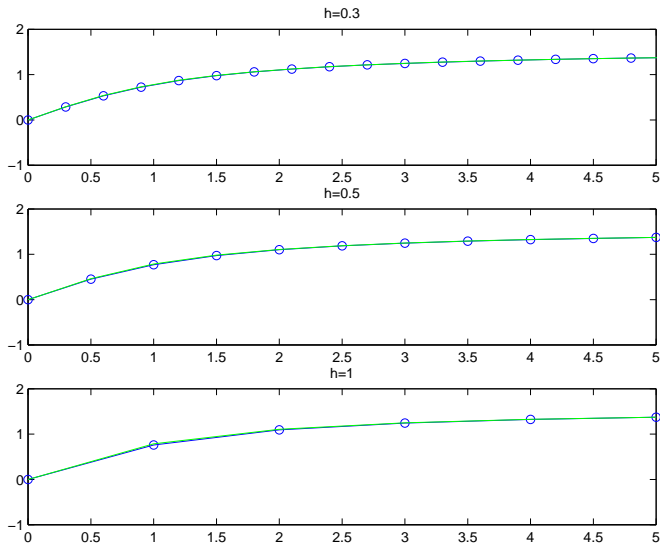
Explizit:



Steife Differentialgleichungen

Was sind steife Differentialgleichungen?

Implizit:



- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Erinnerung. Ein numerisches Verfahren zur Lösung von

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}), \quad \mathbf{y}(0) = \mathbf{y}_0,$$

heißt A -stabil, wenn folgendes gilt:

Wendet man das Verfahren auf ein lineares Problem $\mathbf{y}' = A\mathbf{y}$ ($t \in [0, \infty)$) an und liegen die Eigenwerte von A alle in der linken Halbebene $\{\zeta : \operatorname{Re} \zeta < 0\}$, dann soll für die Näherungslösung stets $\lim_{n \rightarrow \infty} \|\mathbf{y}_n\| = 0$ gelten (die exakte Lösung erfüllt $\lim_{t \rightarrow \infty} \|\mathbf{y}(t)\| = 0$).

In konkreten Fällen wird die äquivalente Bedingung

$$\mathcal{R}_A \supseteq \{\hat{h} : \operatorname{Re} \hat{h} < 0\}$$

überprüft. Zur Definition des Stabilitätsbereichs \mathcal{R}_A vergleiche Abschnitt 6 für lineare Mehrschrittverfahren und Abschnitt 3 für RKV (oder allgemeine Einschrittverfahren).

Abgeschwächte Stabilitätsbegriffe sind

- **$A(\alpha)$ -stabil** für $\alpha \in (0, \pi/2)$ $\Leftrightarrow \mathcal{R}_A \supseteq \{\hat{h} : -\alpha < \pi - \arg \hat{h} < \alpha\}$,
- **A_0 -stabil** $\Leftrightarrow \mathcal{R}_A \supseteq (-\infty, 0)$,
- **steif-stabil** $\Leftrightarrow \mathcal{R}_A \supseteq \mathcal{R}_1(\beta) \cup \mathcal{R}_2(\beta, \gamma)$ für positive β und γ , wobei

$$\mathcal{R}_1(\beta) := \{\hat{h} : \operatorname{Re} \hat{h} < -\beta\} \text{ und}$$

$$\mathcal{R}_2(\beta, \gamma) := \{\hat{h} : -\beta \leq \operatorname{Re} \hat{h} < 0, |\operatorname{Im} \hat{h}| \leq \gamma\}.$$

Außerdem heißt ein Einschrittverfahren

- **L -stabil**, wenn es A -stabil ist und zusätzlich gilt: Bei Anwendung auf die Testgleichung

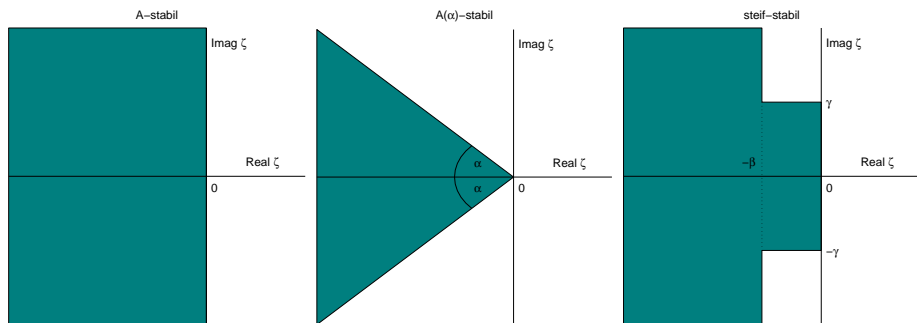
$$y' = \lambda y, \quad \operatorname{Re} \lambda < 0,$$

resultiert eine Näherung

$$y_{n+1} = R(h\lambda)y_n \quad \text{mit} \quad \lim_{h\lambda \rightarrow -\infty} |R(h\lambda)| = 0.$$

Steife Differentialgleichungen

Stabilitätsbegriffe



Offensichtlich gilt $(0 < \alpha \leq \arg(\beta + i\gamma))$

$$L\text{-stabil} \Rightarrow A\text{-stabil} \Rightarrow \text{steif-stabil} \Rightarrow A(\alpha)\text{-stabil} \Rightarrow A_0\text{-stabil}.$$

Steife Differentialgleichungen

Stabilitätsbegriffe

- Die Trapezregel

$$y_{n+1} = y_n + \frac{1}{2}h(f_{n+1} + f_n), \quad \text{d.h.} \quad R(\hat{h}) = \frac{1 + \hat{h}/2}{1 - \hat{h}/2},$$

ist A -stabil, aber nicht L -stabil.

- Das implizite Euler-Verfahren

$$y_{n+1} = y_n + hf_{n+1}, \quad \text{d.h.} \quad R(\hat{h}) = \frac{1}{1 - \hat{h}},$$

ist dagegen L -stabil.

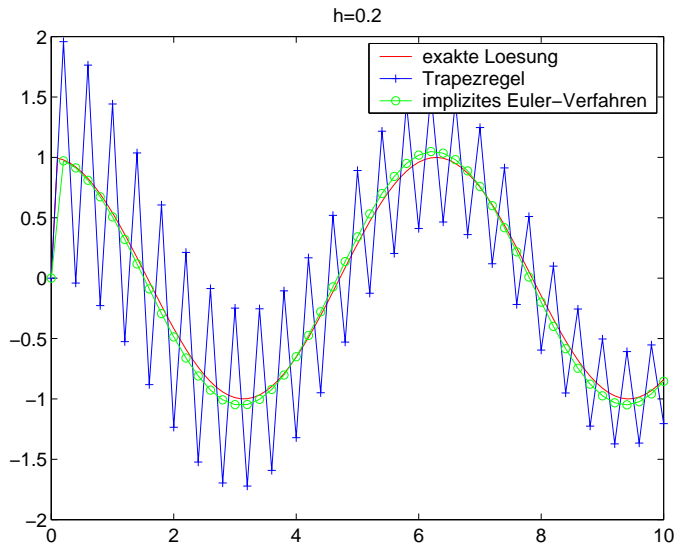
Wir illustrieren den Unterschied, indem wir mit beiden Verfahren die Lösung von (AWP₂) aus Abschnitt 1 approximieren. Um den transienten Anteil der Lösung besser sichtbar zu machen, werden die ABen $\mathbf{y}(0) = [0, 0]^\top$ vorgeschrieben. Die exakte Lösung ist dann

$$\mathbf{y}(t) = \frac{-1}{999} \exp(-t) \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \frac{1}{999} \exp(-1000t) \begin{bmatrix} 1 \\ -998 \end{bmatrix} + \begin{bmatrix} \sin t \\ \cos t \end{bmatrix}.$$

Die folgende Abbildung zeigt die zweite Komponente der (Näherungs)lösung für $h = 0.2$.

Steife Differentialgleichungen

Stabilitätsbegriffe



Das Beispiel

$$\mathbf{y}' = \begin{bmatrix} 42.2 & 50.1 & -42.1 \\ -66.1 & -58.0 & 58.1 \\ 26.1 & 42.1 & -34.0 \end{bmatrix} \mathbf{y}, \quad \mathbf{y}(0) = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix},$$

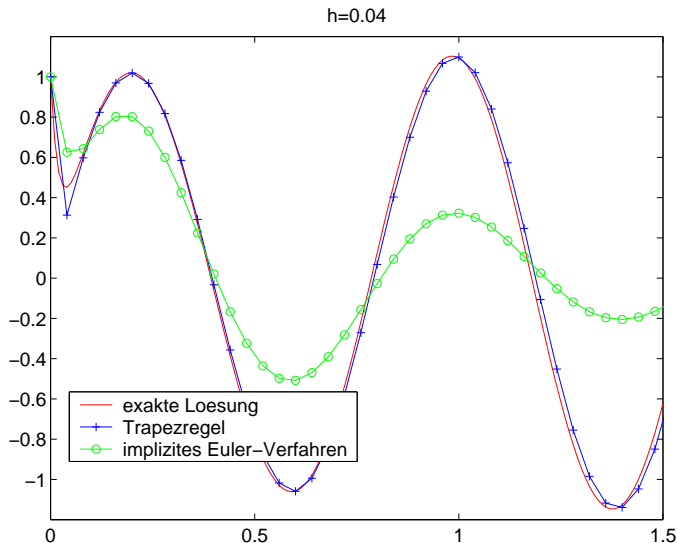
mit der exakten Lösung

$$\mathbf{y}(t) = \begin{bmatrix} \exp(0.8t) \sin(8t) + \exp(-50t) \\ \exp(0.1t) \cos(8t) - \exp(-50t) \\ \exp(0.1t) (\sin(8t) + \cos(8t)) + \exp(-50t) \end{bmatrix}$$

zeigt, dass L -stabile Verfahren nicht immer besser sind als A -stabile (die folgende Abbildung zeigt die erste Komponente der (Näherungs)lösung für $h = 0.04$ und $t \in [0, 1.5]$). Beim impliziten Euler-Verfahren müsste man $h < 0.0031 \dots$ wählen, um eine akzeptable Näherung zu erhalten (warum?).

Steife Differentialgleichungen

Stabilitätsbegriffe



- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungsterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Steife Differentialgleichungen

Ordnungssterne

Wendet man ein Einschrittverfahren (ESV) auf die Testgleichung $y' = \lambda y$ an, so gilt $y_{n+1} = R(\lambda h)y_n$. Dabei ist R ein Polynom (bei expliziten Verfahren) oder eine (gebrochen) rationale Funktion in $\hat{h} = \lambda h$.

Besitzt das ESV die Konsistenzordnung p , so ist

$$R(\hat{h}) - \exp(\hat{h}) = O(\hat{h}^{p+1}) \quad (\hat{h} \rightarrow 0).$$

Natürlicher Ansatz: Bestimme eine rationale Funktion R vom Typ (k, ℓ) (k = Zählergrad, ℓ = Nennergrad), so dass der Exponent μ in

$$R(\hat{h}) - \exp(\hat{h}) = O(\hat{h}^\mu)$$

maximal wird. Es zeigt sich, dass das Maximum $\mu = k + \ell + 1$ durch die in Abschnitt 4 eingeführten Padé-Approximationen $(k, \ell)_{\text{exp}}$ realisiert wird.

Steife Differentialgleichungen

Ordnungssterne

$(k, \ell)_{\text{exp}}(\zeta)$ für $0 \leq k, \ell \leq 3$:

k, ℓ	0	1	2	3
0	1	$\frac{1}{1-\zeta}$	$\frac{1}{1-\zeta+\frac{1}{2}\zeta^2}$	$\frac{1}{1-\zeta+\frac{1}{2}\zeta^2-\frac{1}{6}\zeta^3}$
1	$1+\zeta$	$\frac{1+\frac{1}{2}\zeta}{1-\frac{1}{2}\zeta}$	$\frac{1+\frac{1}{3}\zeta}{1-\frac{2}{3}\zeta+\frac{1}{6}\zeta^2}$	$\frac{1+\frac{1}{4}\zeta}{1-\frac{3}{4}\zeta+\frac{1}{4}\zeta^2-\frac{1}{24}\zeta^3}$
2	$1+\zeta+\frac{1}{2}\zeta^2$	$\frac{1+\frac{2}{3}\zeta+\frac{1}{6}\zeta^2}{1-\frac{1}{3}\zeta}$	$\frac{1+\frac{1}{2}\zeta+\frac{1}{12}\zeta^2}{1-\frac{1}{2}\zeta+\frac{1}{12}\zeta^2}$	$\frac{1+\frac{2}{5}\zeta+\frac{1}{20}\zeta^2}{1-\frac{3}{5}\zeta+\frac{3}{20}\zeta^2-\frac{1}{60}\zeta^3}$
3	$1+\zeta+\frac{1}{2}\zeta^2+\frac{1}{6}\zeta^3$	$\frac{1+\frac{3}{4}\zeta+\frac{1}{4}\zeta^2+\frac{1}{24}\zeta^3}{1-\frac{1}{4}\zeta}$	$\frac{1+\frac{3}{5}\zeta+\frac{3}{20}\zeta^2+\frac{1}{60}\zeta^3}{1-\frac{2}{5}\zeta+\frac{1}{20}\zeta^2}$	$\frac{1+\frac{1}{2}\zeta+\frac{1}{10}\zeta^2+\frac{1}{120}\zeta^3}{1-\frac{1}{2}\zeta+\frac{1}{10}\zeta^2-\frac{1}{120}\zeta^3}$

Eine rationale Approximation $R(\zeta)$ an $\exp(\zeta)$ heißt

- **A-akzeptabel**, wenn $|R(\zeta)| < 1$ für alle $\zeta \in \mathbb{C}$ mit $\operatorname{Re} \zeta < 0$,
- **A₀-akzeptabel**, wenn $|R(\zeta)| < 1$ für alle $\zeta \in \mathbb{R}$ mit $\zeta < 0$,
- **L-akzeptabel**, wenn R A-akzeptabel ist und $R(\zeta) \rightarrow 0$ für $\zeta \rightarrow -\infty$ gilt.

Satz 5.1

Es gilt

1. $(k, k)_{\exp}$ ist A-akzeptabel (G. Birkhoff, R.S. Varga 1962).
2. $(k, \ell)_{\exp}$ ist A₀-akzeptabel für $k \leq \ell$ (R.S. Varga 1961).
3. $(k, \ell)_{\exp}$ ist L-akzeptabel für $k \in \{\ell - 1, \ell - 2\}$ (B.L. Ehle, 1969).
4. $(k, \ell)_{\exp}$ ist genau dann A-akzeptabel, wenn $k \in \{\ell, \ell - 1, \ell - 2\}$
(Ehle-Vermutung; G. Wanner, E. Hairer, S.P. Nørsett 1978).

Steife Differentialgleichungen

Ordnungssterne

Der Beweis der vierten Aussage basiert auf der Theorie der Ordnungssterne. Ist $R(\zeta)$ eine rationale Funktion, so heißt

$$\mathcal{S}_R := \{\zeta \in \mathbb{C} : |R(\zeta)| > |\exp(\zeta)|\}$$

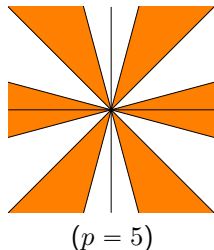
der zugehörige **Ordnungsstern**. Sei $\mathcal{T}_R := \mathbb{C} \setminus \mathcal{S}_R$ sein Komplement.

Vier Eigenschaften.

- (1) R ist genau dann eine Approximation der Ordnung p an die Exponentialfunktion, d.h.

$$R(\zeta) - \exp(\zeta) = O(\zeta^{p+1}) \quad (\zeta \rightarrow 0),$$

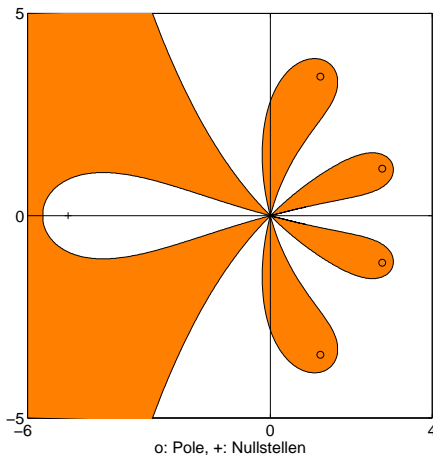
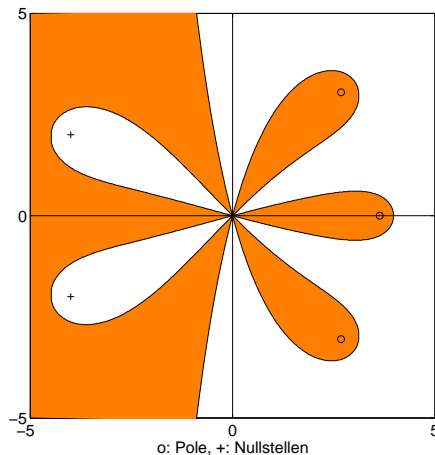
wenn Folgendes gilt: Für $\zeta \rightarrow 0$ besteht sowohl \mathcal{S}_R als auch \mathcal{T}_R aus genau $p + 1$ Sektoren mit Winkel jeweils $\pi/(p + 1)$. Die Sektoren von \mathcal{S}_R und \mathcal{T}_R trennen sich.



- (2) Der Rand von \mathcal{S}_R enthält genau zwei unbeschränkte Zweige.
- (3) Komponenten von \mathcal{S}_R , die k Sektoren enthalten, nennt man **Finger der Ordnung k** . Komponenten von \mathcal{I}_R , die k Sektoren enthalten, heißen **duale Finger der Ordnung k** .
Es gilt: Jeder beschränkte Finger der Ordnung k enthält mindestens k Pole von R (Vielfachheiten mitzählen). Jeder beschränkte duale Finger der Ordnung k enthält mindestens k Nullstellen von R (Vielfachheiten mitzählen).
- (4) R ist genau dann A -akzeptabel, wenn \mathcal{S}_R keine Punkte der imaginären Achse enthält und R keine Pole in der linken Halbebene $\operatorname{Re} \zeta < 0$ besitzt.

Steife Differentialgleichungen

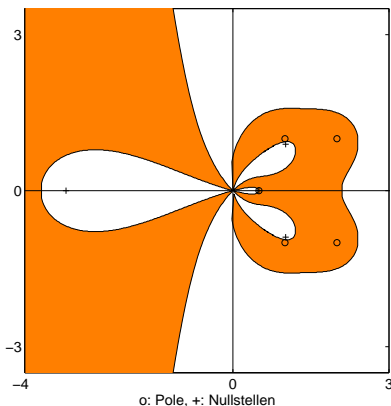
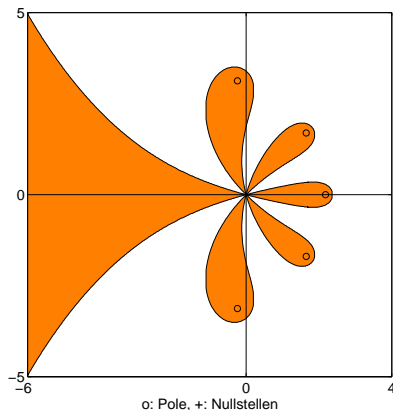
Ordnungssterne



Ordnungssterne der Padé-Approximationen $(2,3)_{\text{exp}}$ (links) und $(1,4)_{\text{exp}}$ (rechts).

Steife Differentialgleichungen

Ordnungssterne



Ordnungssterne der Padé-Approximation $(0, 5)_{\text{exp}}$ (links) und einer (nicht A -akzeptablen) $(4, 5)$ -Approximation R der Ordnung 5 (rechts) ($R(\zeta) = [\frac{41}{24}\zeta^4 + \frac{7}{6}\zeta^3 - 9\zeta^2 + 14\zeta - 5]/[\zeta^5 + \frac{13}{2}\zeta^4 + 18\zeta^3 - \frac{51}{2}\zeta^2 + 19\zeta - 5]$).

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Steife Differentialgleichungen

Lineare MSV für steife Probleme

Satz 5.2 (Zweite Dahlquist-Barriere)

Es gilt:

- Kein explizites lineares MSV ist A -stabil.
- Ein A -stabiles lineares MSV besitzt höchstens die Ordnung 2.
- Das A -stabile lineare MSV mit Konsistenzordnung 2 und kleinster Fehlerkonstante ist die Trapezregel.

Satz 5.3 (C.W. Cryer, 1973)

Es gilt:

- Kein explizites lineares MSV ist A_0 -stabil.
- Es gibt A_0 -stabile lineare MSV beliebig hoher Konsistenzordnung.

Beispielsweise sind die BDF-Verfahren aus Abschnitt 7 alle A_0 -stabil.

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Entscheidend (vgl. Abschnitt 3)

$$R(\hat{h}) = 1 + \hat{h} \mathbf{b}^\top (I - \hat{h} A)^{-1} \mathbf{e} = \frac{\det(I - \hat{h}(A - \mathbf{e} \mathbf{b}^\top))}{\det(I - \hat{h} A)}.$$

	Ordnung	Typ von R	Stabilität
Gauß	$2m$	(m, m)	A -stabil
Gauß-Radau	$2m - 1$	$(m - 1, m)$	L -stabil
Gauß-Lobatto	$2m - 2$	$(m - 1, m - 1)$	A -stabil

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ④ Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
 - 5.1 Was sind steife Differentialgleichungen?
 - 5.2 Stabilitätsbegriffe
 - 5.3 Ordnungssterne
 - 5.4 Lineare MSV für steife Probleme
 - 5.5 RKV für steife Probleme
 - 5.6 Nichtlineare Stabilitätstheorie

⑥ Ausblick

Steife Differentialgleichungen

Nichtlineare Stabilitätstheorie

Alle bisher getroffenen Stabilitätsaussagen basieren auf der Untersuchung linearer Systeme mit konstanten Koeffizienten (**lineare Stabilitätstheorie**). Oft wird versucht, diese Aussagen mit folgender Argumentation zu verallgemeinern:

Die Lösung \mathbf{y} des AWP's $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$, $\mathbf{y}(t_0) = \mathbf{y}_0$ kann durch die Lösung der linearen **Variationsgleichung**

$$\mathbf{y}' = \mathbf{f}(t, \mathbf{y}^*(t)) + \mathbf{f}_y(t, \mathbf{y}^*(t))[\mathbf{y} - \mathbf{y}^*(t)] \quad (*)$$

approximiert werden. Lokal ist $\mathbf{f}_y(t, \mathbf{y}^*)$ nahezu zeitunabhängig, so dass $(*)$ näherungsweise die Form $\mathbf{y}' = A\mathbf{y} + \mathbf{g}(t)$ besitzt.

Zwei Beispiele sollen zeigen, dass Skepsis gegenüber diesem “frozen Jacobian argument” angebracht ist.

Beispiel 1.

$$\mathbf{y}' = \begin{bmatrix} -1 - 9 \cos^2(6t) + 6 \sin(12t) & 12 \cos^2(6t) + 4.5 \sin(12t) \\ -12 \sin^2(6t) + 4.5 \sin(12t) & -1 - 9 \sin^2(6t) - 6 \sin(12t) \end{bmatrix} \mathbf{y}.$$

Die Systemmatrix $A(t)$ besitzt die (t -unabhängigen!) Eigenwerte

$$\lambda_1 = -1, \quad \lambda_2 = -10.$$

Obwohl die Eigenwerte negativ sind, ist die allgemeine Lösung

$$\mathbf{y}(t) = \kappa_1 \exp(2t) \begin{bmatrix} \cos(6t) + 2 \sin(6t) \\ 2 \cos(6t) - \sin(6t) \end{bmatrix} + \kappa_2 \exp(-13t) \begin{bmatrix} \sin(6t) - 2 \cos(6t) \\ 2 \sin(6t) + \cos(6t) \end{bmatrix}$$

sicher nicht monoton fallend für $t \geq 0$.

Beispiel 2.

$$\mathbf{y}' = \begin{bmatrix} \frac{-1}{2t} & \frac{2}{t^3} \\ \frac{-t}{2} & \frac{-1}{2t} \end{bmatrix} \mathbf{y}, \quad t \geq 1. \quad (\diamond)$$

Die Systemmatrix $A(t)$ besitzt die Eigenwerte

$$\lambda_{1,2} = \frac{-1 \pm 2i}{2t},$$

die für $t \geq 1$ negativen Realteil besitzen. Die allgemeine Lösung ist

$$\mathbf{y}(t) = \kappa_1 \begin{bmatrix} t^{-3/2} \\ -\frac{1}{2}t^{1/2} \end{bmatrix} + \kappa_2 \begin{bmatrix} 2t^{-3/2} \log t \\ t^{1/2}(1 - \log t) \end{bmatrix}.$$

Für $\kappa_1 = 1$, $\kappa_2 = 0$ ist

$$\|\mathbf{y}\|_2^2 = t^{-3} + \frac{1}{4}t$$

streng monoton wachsend (für $t \geq 12^{1/4} = 1.86 \dots$).

Definition 5.4

Das System $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ heißt **dissipativ** in $[t_0, t_{\text{end}}]$, wenn

$$\langle \mathbf{f}(t, \mathbf{y}) - \mathbf{f}(t, \tilde{\mathbf{y}}), \mathbf{y} - \tilde{\mathbf{y}} \rangle \leq 0$$

für alle $t \in [t_0, t_{\text{end}}]$ und alle $\mathbf{y}, \tilde{\mathbf{y}}$ aus dem Definitionsbereich von \mathbf{f} gilt.

$\langle \cdot, \cdot \rangle$ bezeichnet ein (z.B. das Euklidische) Innenprodukt im \mathbb{R}^n .

Satz 5.5

Ist das System $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ dissipativ in $[t_0, t_{\text{end}}]$, dann ist es auch **kontraktiv** in $[t_0, t_{\text{end}}]$. Dies bedeutet: Für je zwei Lösungen $\mathbf{y}, \tilde{\mathbf{y}}$ von $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ mit Anfangsbedingungen $\mathbf{y}(t_0) = \mathbf{y}_0, \tilde{\mathbf{y}}(t_0) = \tilde{\mathbf{y}}_0, \mathbf{y}_0 \neq \tilde{\mathbf{y}}_0$, gilt

$$\|\mathbf{y}(t_2) - \tilde{\mathbf{y}}(t_2)\| \leq \|\mathbf{y}(t_1) - \tilde{\mathbf{y}}(t_1)\|$$

für alle t_1, t_2 mit $t_0 \leq t_1 \leq t_2 \leq t_{\text{end}}$. ($\|\cdot\| := \langle \cdot, \cdot \rangle^{1/2}$)

Ob ein System dissipativ (kontraktiv) ist, kann mit Hilfe **logarithmischer Normen** entschieden werden: Sei $\|\cdot\|$ eine (von einer Vektornorm induzierte) Matrixnorm im $\mathbb{R}^{n \times n}$. Dann heißt

$$\mu(A) = \mu_{\|\cdot\|}(A) := \lim_{h \rightarrow 0+} \frac{\|I_n + hA\| - 1}{h} \quad (A \in \mathbb{R}^{n \times n})$$

die zu $\|\cdot\|$ gehörige **logarithmische Norm**. Ist $\|\cdot\| = \|\cdot\|_p$, so schreibt man $\mu_p(\cdot)$ für die zugehörige logarithmische Norm.

Eigenschaften.

- (1) Die logarithmische Norm ist – trotz ihres Namens – keine Norm.
- (2) Wird $\|\cdot\|$ von einem Innenprodukt $\langle \cdot, \cdot \rangle$ induziert, so gilt

$$\mu(A) = \max_{\|x\|=1} \langle Ax, x \rangle.$$

- (3) $\mu_2(A) = \rho(\frac{1}{2}(A + A^T))$ (ρ = Spektralradius).
- (4) Ist $f(t, y) = A(t)y$, so ist

$$\langle f(t, y) - f(t, \tilde{y}), y - \tilde{y} \rangle = \langle A(t)(y - \tilde{y}), y - \tilde{y} \rangle \leq \mu(A(t)) \|y - \tilde{y}\|^2.$$

Satz 5.6

Seien $\|\cdot\|$ eine Norm im \mathbb{R}^n und $\nu(t) : \mathbb{R} \rightarrow \mathbb{R}$ eine stückweise stetige Funktion mit

$$\mu(\mathbf{f}_{\mathbf{y}}(t, \mathbf{y})) \leq \nu(t)$$

für alle $t \in [t_0, t_{\text{end}}]$ und alle \mathbf{y} aus dem Definitionsbereich von \mathbf{f} . Dann gilt für je zwei Lösungen $\mathbf{y}, \tilde{\mathbf{y}}$ von $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ mit Anfangsbedingungen $\mathbf{y}(t_0) = \mathbf{y}_0$, $\tilde{\mathbf{y}}(t_0) = \tilde{\mathbf{y}}_0$, $\mathbf{y}_0 \neq \tilde{\mathbf{y}}_0$,

$$\|\mathbf{y}(t_2) - \tilde{\mathbf{y}}(t_2)\| \leq \exp\left(\int_{t_1}^{t_2} \nu(s) \, ds\right) \|\mathbf{y}(t_1) - \tilde{\mathbf{y}}(t_1)\|$$

für alle t_1, t_2 mit $t_0 \leq t_1 \leq t_2 \leq t_{\text{end}}$.

Beispiel.

Für das System (\diamond) von Seite 279 gilt:

$$\mathbf{f}_y(t, \mathbf{y}) = A(t) = \begin{bmatrix} \frac{-1}{2t} & \frac{2}{t^3} \\ \frac{-t}{2} & \frac{-1}{2t} \end{bmatrix}.$$

Es folgt

$$\mu_2(A(t)) = \rho\left(\frac{1}{2}(A(t) + A(t)^\top)\right) = \max\left\{\left|\frac{-1}{2t} \pm \left(\frac{1}{t^3} - \frac{t}{4}\right)\right|\right\}$$

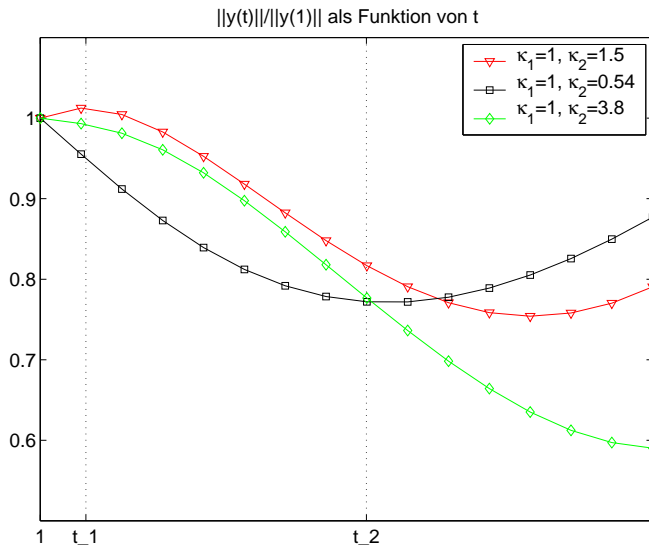
und damit $\mu_2(A(t)) \leq 0$ genau dann, wenn

$$t \in [t_1, t_2] \approx [1.1, 1.8] \quad \text{mit} \quad t_{1,2} = (\sqrt{5} \mp 1)^{1/2}.$$

Das System ist daher in $[t_1, t_2]$ dissipativ (und kontraktiv). Folglich ist $\|\mathbf{y}(t)\|_2$ für jede Lösung \mathbf{y} in $[t_1, t_2]$ monoton fallend.

Steife Differentialgleichungen

Nichtlineare Stabilitätstheorie



Definition 5.7

Ein Runge-Kutta-Verfahren gegeben durch das Butcher-Tableau $\begin{array}{c|c} c & A \\ \hline & b^\top \end{array}$ heißt **algebraisch stabil**, wenn die symmetrischen Matrizen

$$B := \text{diag}(\beta_1, \beta_2, \dots, \beta_m) \quad \text{und} \quad M := BA + A^\top B - \mathbf{b}\mathbf{b}^\top$$

beide positiv semidefinit sind.

Satz 5.8

Seien $\mathbf{y}' = \mathbf{f}(t, \mathbf{y})$ ein dissipatives System und $\{\mathbf{y}_n\}, \{\tilde{\mathbf{y}}_n\}, \mathbf{y}_0 \neq \tilde{\mathbf{y}}_0$, Näherungslösungen, die aus der Anwendung eines Runge-Kutta-Verfahrens resultieren. Dann gilt: Ist das Verfahren algebraisch stabil, so sind die Näherungen kontraktiv, d.h.

$$\|\mathbf{y}_{n+1} - \tilde{\mathbf{y}}_{n+1}\| \leq \|\mathbf{y}_n - \tilde{\mathbf{y}}_n\| \quad (n = 0, 1, \dots).$$

Beispiele.

- Die m -stufigen Gauß-Formeln ($m \in \{1, 2, 3\}$) aus Abschnitt 5 sind alle algebraisch stabil (hier ist M die Nullmatrix!).
- Die Gauß-Radau-Formel

$$\begin{array}{c|cc} \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ 1 & \frac{3}{4} & \frac{1}{4} \\ \hline & \frac{3}{4} & \frac{1}{4} \end{array}$$

ist algebraisch stabil, denn:

$$M = \frac{1}{16} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}, \text{ d.h. } \Lambda(M) = \{0, 1/8\}.$$

- Die Trapezregel (aufgefasst als implizites zweistufiges RKV) ist nicht algebraisch stabil ($\Lambda(M) = \{\pm 1/2\}$).

- ① Einleitung
- ② Numerische Methoden für Anfangswertprobleme
- ③ Lineare Mehrschrittverfahren
- ① Runge-Kutta-Verfahren
- ⑤ Steife Differentialgleichungen
- ⑥ Ausblick

Diskretisierungen von **Reaktions-Diffusions-Transport Problemen** (z.B. reaktive Strömung) führen oft auf **steife** nichtlineare Anfangswertaufgaben der Form

$$\mathbf{u}'(t) = \mathbf{F}(\mathbf{u}(t)). \quad (*)$$

Nach Linearisierung um $\mathbf{u}(t_n)$

$$\mathbf{F}(\mathbf{u}(t)) = \mathbf{A}_n[\mathbf{u}(t) - \mathbf{u}(t_n)] + \mathbf{F}(\mathbf{u}(t_n)) + \mathbf{R}(\mathbf{u}(t_n)) \quad \mathbf{A}_n := \frac{\partial \mathbf{F}}{\partial \mathbf{u}}(\mathbf{u}(t_n))$$

lautet (*)

$$\mathbf{u}'(t) = \mathbf{A}_n[\mathbf{u}(t) - \mathbf{u}(t_n)] + \mathbf{F}(\mathbf{u}(t_n)) + \mathbf{R}(\mathbf{u}(t)) \quad (**)$$

- Zeitschrittverfahren behandeln typischerweise den **linearen Anteil** implizit (z.B. SIMPLE-Verfahren)
- Erfordert Lösung eines LGS in jedem Zeitschritt.

Alternative: Lösung von $(**)$ lautet

$$\begin{aligned} \mathbf{u}(t_n + \Delta t) = \mathbf{u}(t_n) + & \underbrace{(e^{\Delta t \mathbf{A}_n} - \mathbf{I}) \mathbf{A}_n^{-1} \mathbf{F}(\mathbf{u}(t_n))}_{=: \phi(\mathbf{A}) \mathbf{b}} \\ & + \int_{t_n}^{t_n + \Delta t} e^{(t_n + \Delta t - s) \mathbf{A}_n} \mathbf{R}(\mathbf{u}(s)) ds \end{aligned}$$

mit

$$\phi(z) = \frac{e^{\Delta t z} - 1}{z}, \quad \mathbf{b} = \mathbf{F}(\mathbf{u}(t_n)).$$

Krylov-Approximation von $\phi(\mathbf{A}_n) \mathbf{b}$ in Kombination mit verschiedenen Approximationen des Integraltermes führen auf sog. **exponentielle Integratoren** zur **expliziten** Zeitintegration steifer Systeme.

Hamiltonsche Systeme

$$\dot{p}_i = -\frac{\partial H}{\partial q_i}(p, q), \quad \dot{q}_i = \frac{\partial H}{\partial p_i}(p, q),$$

besitzen zwei entscheidende Eigenschaften:

- (a) Die Lösungen erhalten den Wert der Hamiltonschen Funktion $H(p, q)$.
- (b) Der zugehörige Fluß ist **symplektisch**, d.h. er erhält die 2-Form

$$\omega^2 = \sum_{i=1}^n dp_i \wedge dq_i$$

Wir betrachten den **harmonischen Oszillator**

$$H(p, q) = \frac{1}{2}(p^2 + k^2 q^2), \text{ d.h. } \dot{p} = -k^2 q, \quad \dot{q} = p.$$

Anwendung von Euler, implizitem Euler, verbessertem Euler und (impliziter) Mittelpunktsregel:

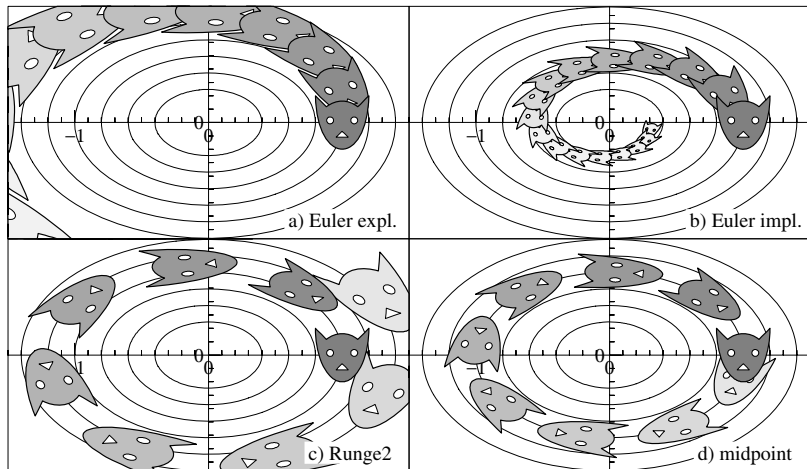


Fig. 16.1. Destruction of symplecticity of a Hamiltonian flow, $k = (\sqrt{5} + 1)/2$

Simulation verschiedener mathematischer Modelle (CME, SDE, ODE) der autokatalytischen chemischen Reaktion $S \rightarrow \emptyset$. ausgehend von 100 Molekülen.

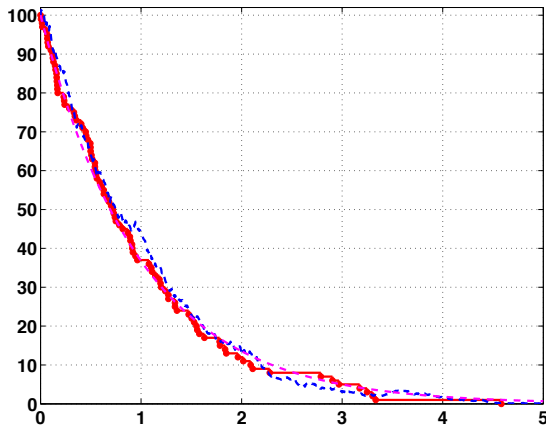


Bild: D. Higham

CME $p_i'(t) = \lambda(i+1)p_{i+1}(t) - \lambda ip_i(t)$
ergibt $\mathbf{E}[\mathbf{X}(T)] = Ne^{-\lambda t}$, $\mathbf{Var}[\mathbf{X}(t)] = Ne^{-\lambda t}(1 - e^{-\lambda t})$

SDE $d\mathbf{Y}(t) = -\lambda \mathbf{Y}(t)dt - \sqrt{\lambda \mathbf{Y}(t)} dW(t)$
ergibt $\mathbf{E}[\mathbf{Y}(T)] = Ne^{-\lambda t}$, $\mathbf{Var}[\mathbf{Y}(t)] = Ne^{-\lambda t}(1 - e^{-\lambda t})$

ODE $z'(t) = -\lambda z(t)$, $z(0) = N$
ergibt $z(t) = Ne^{-\lambda t}$.

Dieselben Simulationen ausgehend von 10 Molekülen.

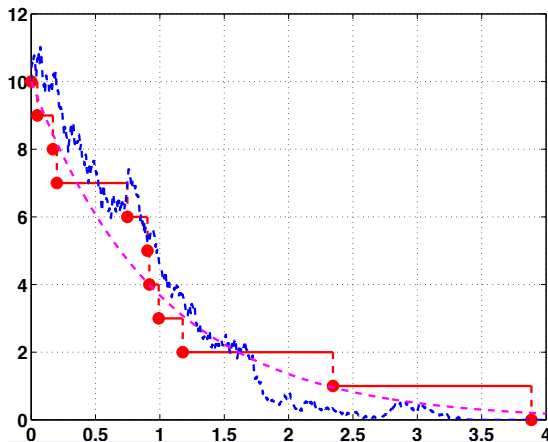


Bild: D. Higham