

# Zur Lösung differential-algebraischer Gleichungssysteme

Eine Einführung von M. Burkhardt und R. Seifried

## 1 Differential-algebraische Gleichungen

Die implizite Differentialgleichung

$$\mathbf{h}(t, \mathbf{x}, \dot{\mathbf{x}}) = \mathbf{0} \quad (1)$$

beschreibt die Änderung  $\dot{\mathbf{x}}$  einer  $n$ -dimensionalen Größe  $\mathbf{x}$  über der Zeit  $t$ . Ist die Jacobi-Matrix  $\frac{\partial \mathbf{h}}{\partial \dot{\mathbf{x}}}$  regulär, so handelt es sich bei Gleichung (1) um ein System gewöhnlicher Differentialgleichungen. Ist diese Jacobi-Matrix jedoch singular, so spricht man von einem differential-algebraischen System. Häufig ist es möglich, diese implizite Form (1) durch die semi-explizite Form

$$\mathbf{I} \cdot \dot{\mathbf{x}} = \mathbf{k}(t, \mathbf{x}) \quad (2)$$

darzustellen. Die darin enthaltene Matrix  $\mathbf{I}$  ist im Falle von differential-algebraischen Gleichungen singular und im Falle einer gewöhnlichen Differentialgleichung regulär. Diese Form eignet sich jedoch nicht zur numerischen Lösung. Durch eine Koordinatentransformation lässt sich Gleichung (2) durch

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}, \mathbf{u}) \quad (3a)$$

$$\mathbf{0} = \mathbf{c}(t, \mathbf{y}, \mathbf{u}) \quad (3b)$$

ausdrücken. Die Trennung in eine Differentialgleichung (3a) und eine algebraische Gleichung (3b) führt auf die explizite Darstellungsform, die zur numerischen Lösung herangezogen wird.

## 2 Runge-Kutta-Verfahren

Bei Runge-Kutta-Verfahren handelt es sich um Einschrittverfahren zur numerische Lösung von Differentialgleichungen der Form

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}), \quad (4)$$

die sich durch die Berechnungsvorschrift

$$\mathbf{g}_i = \mathbf{y}_k + h \sum_{j=1}^s a_{ij} \mathbf{f}(t_k + c_j h, \mathbf{g}_j) \quad i = 1 \dots s \quad (5a)$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h \sum_{j=1}^s b_j \mathbf{f}(t_k + c_j h, \mathbf{g}_j) \quad (5b)$$

ausdrücken lassen. Dies ist eine Approximation der Lösung der Differentialgleichung vom Zeitpunkt  $t_k$  bis zum Zeitpunkt  $t_{k+1} = t_k + h$ , wobei  $h$  die Schrittweite kennzeichnet. Dieses Intervall wird nun an den Stützstellen  $\mathbf{g}_i$  mit der Differentialgleichung in Einklang gebracht. Mit Hilfe des Butcher-Blocks

$$\begin{array}{c|c} \mathbf{c} & \mathbf{A} \\ \hline & \mathbf{b} \end{array} = \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & \ddots & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline & b_1 & \cdots & b_s \end{array} \quad (6)$$

lassen sich die Koeffizienten  $a_{ij}$ ,  $b_j$  und  $c_j$  eines  $s$ -stufigen Verfahrens kompakt darstellen. Runge-Kutta-Verfahren lassen sich in drei verschiedene Verfahrensklassen unterteilen. Ist  $\mathbf{A}$  eine strikte untere Dreiecksmatrix, so spricht man von einem expliziten Verfahren. Das heißt, dass die Berechnungsvorschriften der Stützstellen nur von bereits berechneten Stützstellen abhängen und so eine sequenzielle

Auswertung möglich ist. Ein semi-explizites Verfahren ergibt sich, wenn  $\mathbf{A}$  ein Dreiecksmatrix ist. Für jede andere Form der Matrix  $\mathbf{A}$  handelt es sich um ein implizites Verfahren. Hierbei kann Gleichung (5a) nicht getrennt für die einzelnen Stützstellen betrachtet werden, was zur Folge hat, dass zu jedem Zeitschritt ein nichtlineares Gleichungssystem gelöst werden muss. Für die Lösung differential-algebraischer Gleichungen eignen sich nun besonders die zuletzt genannten impliziten Verfahren.

### 3 Implizite Runge-Kutta-Verfahren

Eine Vielzahl an implizite Verfahren lassen sich als Kollokationsverfahren auffassen. Der Grundgedanke dieser Verfahren ist die stückweise Approximation der Lösung im Intervall  $t \in [t_k, t_{k+1}]$  durch Polynome, die an den Stützstellen die Differentialgleichung bzw. differential-algebraischen Gleichung erfüllen. Zur Bestimmung des Butcher-Blocks wird von verschobenen Legendre-Polynomen ausgegangen:

$$\frac{d^s}{dx^s} (x^s(x-1)^s) \quad (\text{Gauß-Verfahren}) \quad (7a)$$

$$\frac{d^{s-1}}{dx^{s-1}} (x^s(x-1)^{s-1}) \quad (\text{RadauIA-Verfahren}) \quad (7b)$$

$$\frac{d^{s-1}}{dx^{s-1}} (x^{s-1}(x-1)^s) \quad (\text{RadauIIA-Verfahren}) \quad (7c)$$

$$\frac{d^{s-2}}{dx^{s-2}} (x^{s-1}(x-1)^{s-1}) \quad (\text{Lobatto-Verfahren}) \quad (7d)$$

Diese Polynome zeichnen sich dadurch aus, dass alle Nullstellen im abgeschlossenen Intervall  $[0, 1]$  liegen. Diese Nullstellen ergeben die  $s$  Koeffizienten  $c_i$ , auf deren Basis die Koeffizienten  $a_{ij}$  und  $b_i$  berechnet werden können.

Im weiteren Verlauf werden anhand des dreistufigen *RadauIIA*-Verfahren<sup>1</sup> die einzelnen Schritte zur Lösung differential-algebraischer Gleichungen aufgezeigt. Dieses Verfahren bietet sich aufgrund seiner Struktur besonders für die Lösung steifer und differential-algebraischer Gleichungen an. Drei Eigenschaften sollten dabei Erwähnung finden:

1. Es gilt  $c_s = 1$  und somit unterliegen die algebraischen Nebenbedingungen zum Zeitpunkt  $t_{k+1}$  der Fehlerkontrolle. Das heißt, dass beim nächsten Zeitschritt konsistente Anfangsbedingungen vorliegen.
2. Die Matrix  $\mathbf{A}$  vollen Rang und ist somit invertierbar.
3. Der Vektor  $\mathbf{b}$  entspricht der letzten Zeile der Matrix  $\mathbf{A}$ , wodurch der neue Zustandsvektor der letzten Stützstelle entspricht und somit weitere Funktionsauswertungen entfallen.

Wie sich später zeigt, ermöglichen die beiden letzten Eigenschaften eine sowohl gut konditionierte als auch effiziente Berechnung.

Ausgehend von Gleichung (5) wird die neue Größe  $\mathbf{z}_i = \mathbf{g}_i - \mathbf{y}_k$ ,  $i = 1, \dots, s$  eingeführt. Damit lässt sich Gleichung (5a) zu

$$\mathbf{z}_i = h \sum_{j=1}^s a_{ij} \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j) \quad (8)$$

umformulieren. In Matrixschreibweise lautet diese Gleichung

$$\begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_s \end{bmatrix} = \mathbf{A} \otimes \mathbf{E}_n \begin{bmatrix} h\mathbf{f}(t_k + c_1 h, \mathbf{y}_k + \mathbf{z}_1) \\ \vdots \\ h\mathbf{f}(t_k + c_s h, \mathbf{y}_k + \mathbf{z}_s) \end{bmatrix}, \quad (9)$$

<sup>1</sup>Der entsprechende Butcher-Block ist im Anhang dargestellt.

wobei  $\mathbf{E}_n$  eine Einheitsmatrix der Größe  $n$  ist und  $\otimes$  das Kronecker-Produkt<sup>2</sup> vermittelt. Stellt man das Gleichungssystem nach dem rechten Vektor um und setzt diesen in Gleichung (5b) ein, ergibt sich

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \mathbf{b} \mathbf{A}^{-1} \otimes \mathbf{E}_n \begin{bmatrix} z_1 \\ \vdots \\ z_s \end{bmatrix} \quad (10)$$

Da der Vektor  $\mathbf{b}$  der letzten Zeile der Matrix  $\mathbf{A}$  entspricht, ergibt sich

$$\mathbf{y}_{k+1} = \mathbf{y}_k + z_s . \quad (11)$$

Damit ergibt sich die Berechnungsvorschrift des Zustandsvektors  $\mathbf{y}_{k+1}$  in Abhängigkeit der letzten Stützstelle  $z_s$  und des Startwerts  $\mathbf{y}_k$ .

### Beispiel: Euler rückwärts

Das implizite Euler-Verfahren ist ein einstufiges Runge-Kutta-Verfahren, dass durch den Butcher-Block

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array} \quad (12)$$

beschrieben wird. Die Berechnungsvorschrift ergibt sich demnach zu

$$\mathbf{g}_1 = \mathbf{y}_k + h \mathbf{f}(t_k + h, \mathbf{g}_1) \quad (13a)$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h \mathbf{f}(t_k + h, \mathbf{g}_1) \quad (13b)$$

bzw. mit den hier möglichen Vereinfachungen  $\mathbf{g}_1 = \mathbf{y}_{k+1}$  und  $t_k + h = t_{k+1}$  zu

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h \mathbf{f}(t_{k+1}, \mathbf{y}_{k+1}) . \quad (14)$$

Dieses Verfahren nähert die Lösung der Differentialgleichung durch ein Polynom

$$\mathbf{p}(t) \text{ für } t \in [t_k, t_{k+1}] \quad (15)$$

ersten Grades an. Dabei entspricht die erste zeitliche Ableitung des Polynoms der an der Stützstelle ausgewerteten Differentialgleichung.

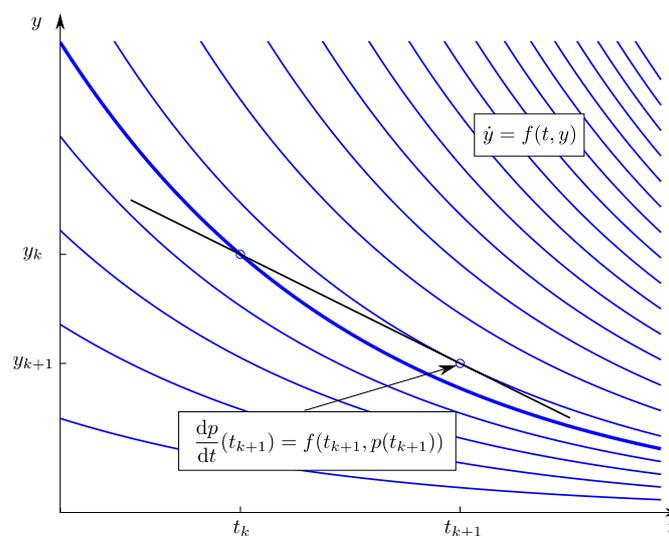


Abbildung 1: Grafische Interpretation des Euler-rückwärts-Verfahrens.

<sup>2</sup>s. Anhang

Abbildung 1 zeigt die grundlegenden Eigenschaften dieses Verfahrens. Die Steigung der Gerade entspricht der an der Stützstelle ausgewerteten Differentialgleichung. Um diese Stützstelle zu finden, bedient man sich eines vereinfachten Newton-Verfahrens.

## 4 Vereinfachtes Newton-Verfahren

Um Gleichung (9) nach den Unbekannten  $z_1, \dots, z_s$  zu lösen, wird die Gleichung zunächst in ihr implizites Äquivalent gemäß

$$\begin{bmatrix} z_1 \\ \vdots \\ z_s \end{bmatrix} - \mathbf{A} \otimes \mathbf{E}_n \begin{bmatrix} h\mathbf{f}(t_n + c_1 h, \mathbf{y}_n + \mathbf{z}_1) \\ \vdots \\ h\mathbf{f}(t_k + c_s h, \mathbf{y}_k + \mathbf{z}_s) \end{bmatrix} = \mathbf{0} \quad (16)$$

umgeformt. Zur Lösung dieses nichtlinearen Gleichungssystems wird ein vereinfachtes Newton-Verfahren verwendet. Es wird jedoch statt der analytischen Jacobi-Matrix die approximierte Jacobi-Matrix

$$\mathbf{J}_{global} = \mathbf{E} - h\mathbf{A} \otimes \mathbf{J}_{approx} \quad (17)$$

herangezogen. Dabei wird  $\mathbf{J}_{approx} = \frac{\partial \mathbf{f}}{\partial \mathbf{y}_k}(t_k, \mathbf{y}_k)$  anstatt der partiellen Ableitungen  $\mathbf{J}_{exakt} = \frac{\partial \mathbf{f}}{\partial z_i}(t_k + c_i h, \mathbf{y}_k + \mathbf{z}_i)$  verwendet. Dadurch wird zwar die quadratische Konvergenz des Newton-Verfahrens aufgegeben, jedoch werden  $(s-1)n^2$  Funktionsauswertungen eingespart. Mit Hilfe dieser Jacobi-Matrix kann nun nach den Stützstellen  $z_i$  iterativ aufgelöst werden.

Um ein differential-algebraisches Gleichungssystem zu lösen, wird gefordert, dass die Nebenbedingungen an allen Stützstellen gemäß

$$\begin{bmatrix} \mathbf{c}(t_k + c_1 h, \mathbf{y}_k + \mathbf{z}_1, \mathbf{u}_k + \mathbf{l}_1) \\ \vdots \\ \mathbf{c}(t_k + c_s h, \mathbf{y}_k + \mathbf{z}_s, \mathbf{u}_k + \mathbf{l}_s) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} \quad (18)$$

erfüllt werden. Die gesamte rechte Seite ergibt sich damit aus Gl. (16) und (18) zu

$$\begin{bmatrix} z_1 - h \sum_{j=1}^s a_{1j} \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j, \mathbf{u}_k + \mathbf{l}_j) \\ \vdots \\ z_s - h \sum_{j=1}^s a_{sj} \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j, \mathbf{u}_k + \mathbf{l}_j) \\ \mathbf{c}(t_k + c_1 h, \mathbf{y}_k + \mathbf{z}_1, \mathbf{u}_k + \mathbf{l}_1) \\ \vdots \\ \mathbf{c}(t_k + c_s h, \mathbf{y}_k + \mathbf{z}_s, \mathbf{u}_k + \mathbf{l}_s) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} . \quad (19)$$

Die globale Jacobi-Matrix ist dementsprechend auch um die partiellen Ableitungen nach  $\mathbf{u}$ , sowie um die partiellen Ableitungen der Nebenbedingungen gemäß

$$\mathbf{J}_{global} = \begin{bmatrix} \mathbf{E} - h\mathbf{A} \otimes \frac{\partial \mathbf{f}}{\partial \mathbf{y}_k}(t_k, \mathbf{y}_k, \mathbf{u}_k) & -h\mathbf{A} \otimes \frac{\partial \mathbf{f}}{\partial \mathbf{u}_k}(t_k, \mathbf{y}_k, \mathbf{u}_k) \\ \mathbf{E} \otimes \frac{\partial \mathbf{c}}{\partial \mathbf{y}_k}(t_k, \mathbf{y}_k, \mathbf{u}_k) & \mathbf{E} \otimes \frac{\partial \mathbf{c}}{\partial \mathbf{u}_k}(t_k, \mathbf{y}_k, \mathbf{u}_k) \end{bmatrix} \quad (20)$$

zu erweitern. Das Newton-Verfahren ergibt sich damit zu

$$\mathbf{J}_{global} \cdot \begin{bmatrix} \Delta \mathbf{z}_1 \\ \vdots \\ \Delta \mathbf{z}_s \\ \Delta \mathbf{l}_1 \\ \vdots \\ \Delta \mathbf{l}_s \end{bmatrix} = - \begin{bmatrix} \mathbf{z}_1 - h \sum_{j=1}^s a_{1j} \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j, \mathbf{u}_k + \mathbf{l}_j) \\ \vdots \\ \mathbf{z}_s - h \sum_{j=1}^s a_{sj} \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j, \mathbf{u}_k + \mathbf{l}_j) \\ \mathbf{c}(t_k + c_1 h, \mathbf{y}_k + \mathbf{z}_1, \mathbf{u}_k + \mathbf{l}_1) \\ \vdots \\ \mathbf{c}(t_k + c_s h, \mathbf{y}_k + \mathbf{z}_s, \mathbf{u}_k + \mathbf{l}_s) \end{bmatrix}. \quad (21)$$

Für die Lösung des linearen Gleichungssystems (21) bietet es sich an, die algebraischen Nebenbedingungen mit der Schrittweite  $h$  zu skalieren. Dadurch wird die Pivot-Suche dahingehend beeinflusst, dass die Nebenbedingungen nicht in den oberen Teil des Gleichungssystems wandern. Dies erhöht sowohl die Stabilität als auch die Genauigkeit.

## 5 Schrittweitensteuerung

Um die Stabilität und Konvergenz des Einschrittverfahrens zu garantieren, muss eine Anpassung der Schrittweite basierend auf einer Fehlerabschätzung erfolgen. Dabei kommt eine Schrittweitensteuerung, die die Lösungen von Verfahren verschiedener Ordnung miteinander vergleicht, zum Einsatz. Die Lösung des eingebetteten Verfahrens ergibt sich aus

$$\hat{\mathbf{y}}_{k+1} = \mathbf{y}_k + h \left( \hat{b}_0 \mathbf{f}(t_k, \mathbf{y}_k, \mathbf{u}_k) + \sum_{j=1}^s \hat{b}_j \mathbf{f}(t_k + c_j h, \mathbf{y}_k + \mathbf{z}_j, \mathbf{u}_k + \mathbf{l}_j) \right), \quad (22)$$

wobei der Faktor  $\hat{b}_0$  der Kehrwert des reellen Eigenwertes<sup>3</sup> der Matrix  $\mathbf{A}$  ist. Aus der Differenz der beiden Lösungen folgt

$$\hat{\mathbf{y}}_{k+1} - \mathbf{y}_{k+1} = \hat{b}_0 h \mathbf{f}(t_k, \mathbf{y}_k, \mathbf{u}_k) + \mathbf{e} \mathbf{z}. \quad (23)$$

Der Vektor  $\mathbf{e}$  ergibt sich dabei aus

$$\mathbf{e} = (\hat{\mathbf{b}} - \mathbf{b}) \mathbf{A}^{-1}, \quad (24)$$

wobei der Vektor  $\hat{\mathbf{b}}$  aus

$$\hat{\mathbf{b}} = \begin{bmatrix} 1 - \hat{b}_0 & \frac{1}{2} & \cdots & \frac{1}{s} \end{bmatrix} \begin{bmatrix} 1 & \cdots & c_i^{s-1} \\ \vdots & \ddots & \vdots \\ 1 & \cdots & c_s^{s-1} \end{bmatrix}^{-1} \quad (25)$$

folgt. Da die Differenz  $\hat{\mathbf{y}}_{k+1} - \mathbf{y}_{k+1}$  für  $h\lambda \rightarrow \infty$  gegen  $\hat{b}_0 h \lambda y_k$  strebt und damit unbeschränkt ist, schlug [?] vor, den Fehler  $\epsilon$  als

$$\epsilon = (\mathbf{E} - h \hat{b}_0 \mathbf{J}_{approx})^{-1} (\hat{\mathbf{y}}_{k+1} - \mathbf{y}_{k+1}) \quad (26)$$

zu definieren. Basierend auf diesem Fehler kann nun eine optimale Schrittweite gefunden werden. Dabei eignen sich vor allem die Fehlberg- ([?]) sowie die Gustafsson-Schrittweitensteuerung ([?]).

<sup>3</sup>Die Matrix  $\mathbf{A}$  hat bei ungerader Stufenzahl einen reellen Eigenwert und  $s - 1$  komplexe Eigenwerte

## Anhang

### Kronecker-Produkt

Das Kronecker-Produkt zweier Matrizen  $\mathbf{A} \in \mathcal{R}^{m \times n}$  und  $\mathbf{B} \in \mathcal{R}^{o \times p}$  ist definiert als

$$\mathbf{C} = \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}, \quad (27)$$

wobei die Matrix  $\mathbf{C}$  der Dimension  $mo \times np$  ist. Falls die Matrizen  $\mathbf{A}$  und  $\mathbf{B}$  invertierbar sind, gilt für die inverse Matrix

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}. \quad (28)$$

Diese Beziehung gilt ohne Forderung an die Matrizen  $\mathbf{A}$  und  $\mathbf{B}$  analog auch für die transponierte, konjugiert komplexe oder adjungierte Matrix.

### Das drei-stufige RadauIIA-Verfahren

Eine beliebte Methode aus der Familie der impliziten Runge-Kutta-Verfahren ist das drei-stufige RadauIIA-Verfahren. Diese Methode eignet sich besonders für steife und differential-algebraische Probleme. Der Butcher-Block dieses Verfahrens ergibt sich zu

$$\begin{array}{c|ccc} \frac{4 - \sqrt{6}}{10} & \frac{88 - 7\sqrt{6}}{360} & \frac{296 - 169\sqrt{6}}{1800} & \frac{-2 + 3\sqrt{6}}{225} \\ \frac{4 + \sqrt{6}}{10} & \frac{296 + 169\sqrt{6}}{1800} & \frac{88 + 7\sqrt{6}}{360} & \frac{-2 - 3\sqrt{6}}{225} \\ 1 & \frac{16 - \sqrt{6}}{36} & \frac{16 + \sqrt{6}}{36} & \frac{1}{9} \\ \hline & \frac{16 - \sqrt{6}}{36} & \frac{16 + \sqrt{6}}{36} & \frac{1}{9} \end{array}. \quad (29)$$

Der Vektor  $\mathbf{b}$  entspricht der letzten Zeile von  $\mathbf{A}$ , was bedeutet, dass die letzte Stützstelle dem neuen Funktionswert entspricht.

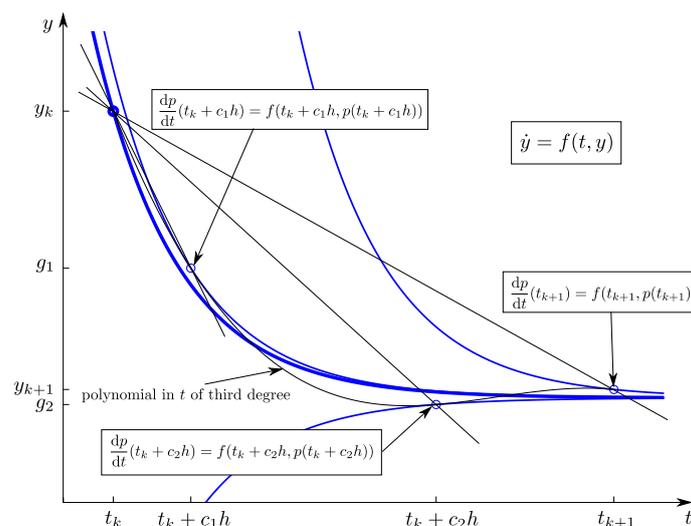


Abbildung 2: Grafische Interpretation des drei-stufigen RadauIIA-Verfahrens mit einer sehr großen Schrittweite  $h$ .



Abbildung 2 zeigt die typischen Eigenschaften des RadauIIA-Verfahrens. Die erste zeitliche Ableitung des Polynom stimmt an den drei Stützstellen  $g_1$ ,  $g_2$  und  $y_{k+1}$  mit der zugrunde liegenden Differentialgleichung überein.



## Literatur

- [Campbell07] Campbell, S.L.: High-index differential algebraic equations. *Mechanics of Structures and Machines*, Bd. 23, Nr. 2, S. 37–41, 2007.
- [Fehlberg69] Fehlberg, E.: Low-order classical Runge-Kutta formulas with step size control and their application to some heat transfer problems. *Techn. Bericht. 315*, National Aeronautics and Space Administration, 1969.
- [FhrerLeimkuhler91] Fhrer, C.; Leimkuhler, B.: Numerical solution of differential-algebraic equations for constrained mechanical motion. *Numerische Mathematik*, Bd. 59, S. 55–69, 1991.
- [Gustafsson94] Gustafsson, K.: Control-theoretic techniques for stepsize selection in implicit Runge-Kutta methods. *ACM Trans. Math. Soft.*, Bd. 20, S. 496–517, 1994.
- [HairerLubichRoche89] Hairer, E.; Lubich, C.; Roche, M.: *The Numerical Solution of Differential-Algebraic Systems by Runge-Kutta Methods*, Bd. 1409 von *Lecture Notes in Mathematics*. Springer-Verlag, 1989.
- [HairerWanner10] Hairer, E.; Wanner, G.: *Solving Ordinary Differential Equations II - Stiff and Differential-Algebraic Problems*. Berlin: Springer, 2. Aufl., 2010.
- [Shampine80] Shampine, L.: Implementation of implicit formulas for the solution of ODEs. *SIAM Journal of Stat. Computations*, Bd. 2, S. 44–53, 1980.