

# NUMERICS OF DIFFERENTIAL EQUATIONS

MICHAEL INNERBERGER AND DIRK PRAETORIUS

## 2. EXPLICIT ONE-STEP METHODS

Throughout this section, we consider the following model problem: Let  $[t_0, T]$  be a given time-interval. For given  $n \in \mathbb{N}$ , let  $f \in C([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and suppose that  $f$  is Lipschitz continuous in  $y$ , i.e.,

$$\forall t \in [t_0, T] \forall y, \tilde{y} \in \mathbb{R}^n : \quad \|f(t, y) - f(t, \tilde{y})\| \leq L \|y - \tilde{y}\|, \quad (2.1)$$

where  $\|\cdot\|$  is an arbitrary norm on  $\mathbb{R}^n$  and  $L > 0$  is a fixed constant. Then, for any initial value  $y_0 \in \mathbb{R}^n$ , the Picard–Lindelöf theorem guarantees existence and uniqueness of  $y \in C^1([t_0, T]; \mathbb{R}^n)$  such that

$$y(t_0) = y_0 \quad \text{and} \quad y'(t) = f(t, y(t)) \quad \text{for all } t \in [t_0, T]. \quad (2.2)$$

---

**Remark 2.1.** If  $f(t, y) = g(t)$ , then the exact solution of (2.2) reads

$$y(t) = y_0 + \int_{t_0}^t g \, ds,$$

i.e.,  $y$  is the antiderivative of  $g$ . Therefore, the solution of an initial value problem (2.2) is also called **integration of the ODE** and the numerical solvers are also called (**numerical**) **integrators**.

---

**2.1. Notation.** In practice, the solution of the initial boundary value problem (2.2) cannot be computed in closed form. If quantitative results are required, one usually considers numerical methods, which provide approximations  $y_\ell \approx y(t_\ell)$  at certain time-steps  $t_\ell \in [t_0, T]$ . From now on, we shall implicitly use the following notation:

---

**Definition 2.2.** Suppose that we have given time-steps  $t_0 < t_1 < \dots < t_N = T$ . The set  $\Delta := \{t_0, \dots, t_N\}$  is called **mesh** (or **grid**) of the time interval  $[t_0, T]$ . The quantities  $h_\ell := t_{\ell+1} - t_\ell$  are called **step-sizes**. Moreover, the **maximum step-size** of  $\Delta$  is defined as  $h_\Delta := \max_{\ell=0, \dots, N-1} h_\ell$ .

---

**Definition 2.3.** Given an **incremental function**  $\Phi : [t_0, T] \times \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}^n$ , the **inductive procedure**

$$y_{\ell+1} := y_\ell + h_\ell \Phi(t_\ell, y_\ell, h_\ell) \quad \text{for all } \ell = 0, \dots, N-1 \quad (2.3)$$

---

**Note:** These a posteriori lecture notes provide the way in which I would have liked to present the course material. If you face any typos, please let me know: [dirk.praetorius@asc.tuwien.ac.at](mailto:dirk.praetorius@asc.tuwien.ac.at).

31 is called **explicit one-step method**. Given  $y_\ell \approx y(t_\ell)$ , we can compute  $y_{\ell+1} \approx y(t_{\ell+1})$ .  
 32 Note that this computation involves only the last time-step, but not the full history  
 33  $y_0, \dots, y_{\ell-1}$ . This will be different for multi-step methods considered below.

---

34 **Remark 2.4.** Later, we will also consider implicit one-step methods, where  $\Phi$  depends  
 35 additionally on the sought approximation  $y_{\ell+1}$ , i.e.,

$$36 \qquad y_{\ell+1} = y_\ell + h_\ell \Phi(t_\ell, y_\ell, y_{\ell+1}, h_\ell).$$

38 Then, each time-step of the implicit one-step method will require the solution of one  
 39 (possibly nonlinear) equation.

---

40 **Example 2.5 (Explicit Euler method / forward Euler method).** If  $f \in C^1([t_0, T] \times$   
 41  $\mathbb{R}^n; \mathbb{R}^n)$  and hence  $y \in C^2([t_0, T]; \mathbb{R}^n)$ , the Taylor theorem proves that

$$42 \qquad y(t+h) = y(t) + h y'(t) + \mathcal{O}(h^2) = y(t) + h f(t, y(t)) + \mathcal{O}(h^2). \quad (2.4)$$

44 With  $y(t_{\ell+1}) \approx y_{\ell+1}$  and  $y(t_\ell) \approx y_\ell$ , the explicit Euler method reads

$$45 \qquad y_{\ell+1} = y_\ell + h_\ell \Phi(t_\ell, y_\ell, h_\ell) \quad \text{with} \quad \Phi(t_\ell, y_\ell, h_\ell) := f(t_\ell, y_\ell). \quad (2.5)$$


---

47 **Example 2.6 (Implicit Euler method / backward Euler method).** If  $f \in$   
 48  $C^1([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and hence  $y \in C^2([t_0, T]; \mathbb{R}^n)$ , the Taylor theorem proves that

$$49 \qquad y(t) = y(t+h) - h y'(t+h) + \mathcal{O}(h^2) = y(t+h) - h f(t+h, y(t+h)) + \mathcal{O}(h^2). \quad (2.6)$$

51 With  $y(t_{\ell+1}) \approx y_{\ell+1}$  and  $y(t_\ell) \approx y_\ell$ , the implicit Euler method reads

$$52 \qquad y_{\ell+1} = y_\ell + h_\ell \Phi(t_\ell, y_\ell, y_{\ell+1}, h_\ell) \quad \text{with} \quad \Phi(t_\ell, y_\ell, y_{\ell+1}, h_\ell) := f(t_{\ell+1}, y_{\ell+1}), \quad (2.7)$$

54 where we note that  $t_{\ell+1} = t_\ell + h_\ell$ .

---

## 56 2.2. Consistency.

57 **Definition 2.7.** Let  $\Phi(t, y, h)$  be the incremental function of an explicit one-step method.  
 58 We say that the one-step method is **consistent**, if

$$59 \qquad \forall t \in [t_0, T) : \quad \lim_{h \rightarrow 0^+} \frac{\|y(t+h) - [y(t) + h \Phi(t, y(t), h)]\|}{h} = 0, \quad (2.8)$$

61 i.e., if we apply the one-step method to the exact solution, the discretization error of one  
 62 step vanishes as  $h \rightarrow 0$ . For  $p \geq 1$ , we say that the one-step method has **consistency**

---

**Leonhard Euler (1707–1783)** was a Swiss mathematician. From 1720–1723, he was studying at the University of Basel (attending also classes by Johann I Bernoulli). In 1727, he moved to the Imperial Russian Academy of Sciences in Saint Petersburg, where he became professor of mathematics in 1733. In 1741, he moved to the Prussian Academy of Sciences in Berlin, but returned to the Imperial Russian Academy of Sciences in Saint Petersburg in 1766.

63 **order**  $p$ , if for all  $f \in C^p([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and hence  $y \in C^{p+1}([t_0, T]; \mathbb{R}^n)$  for the solution  
 64 of (2.2), it holds that

$$\begin{aligned} & \exists C > 0 : \forall t \in [t_0, T] \forall h \in (0, T - t] : \\ & \|y(t+h) - [y(t) + h\Phi(t, y(t), h)]\| \leq Ch^{p+1}. \end{aligned} \quad (2.9)$$

67 **Remark 2.8.** Suppose that the incremental function  $\Phi(t, y(t), h)$  is continuous at  $h = 0$ .  
 68 Then,

$$\begin{aligned} \text{consistency (2.8)} & \iff \forall t \in [t_0, T] : \lim_{h \rightarrow 0^+} \left\| \frac{y(t+h) - y(t)}{h} - \Phi(t, y(t), h) \right\| = 0 \\ & \iff \forall t \in [t_0, T] : f(t, y(t)) = \Phi(t, y(t), 0), \end{aligned}$$

$$\text{since } f(t, y(t)) = y'(t) = \lim_{h \rightarrow 0^+} \frac{y(t+h) - y(t)}{h}.$$

73 **Example 2.9 (Explicit Euler has consistency order  $p = 1$ ).** The explicit Euler  
 74 method (2.5) is always consistent, since  $\Phi(t, y, h) = f(t, y)$ . According to (2.4) it has  
 75 consistency order  $p = 1$ . Moreover, the constant  $C$  in the consistency estimate (2.9)  
 76 takes the form  $C = \|y''\|_{\infty, [t_0, T]}/2$ , since the Taylor expansion yields that

$$y(t+h) = y(t) + hy'(t) + \frac{h^2}{2}y''(\xi) \stackrel{(2.5)}{=} y(t) + h\Phi(t, y(t), h) + \frac{h^2}{2}y''(\xi)$$

79 for some appropriate  $\xi \in [t, t+h]$ .

80 **Remark 2.10.** If the mesh  $\Delta$  has step-size  $h$ , then it takes  $N \sim \frac{T-t_0}{h} = \mathcal{O}(h^{-1})$  steps  
 81 of the numerical scheme to obtain the approximation  $y_N \approx y(T)$ . If each step has a  
 82 (cumulative consistency) error of order  $\mathcal{O}(h^{p+1})$ , one can hope for a total error estimate

$$\|y(T) - y_N\| \sim \sum_{j=0}^{N-1} \mathcal{O}(h^{p+1}) = \mathcal{O}(Nh^{p+1}) = \mathcal{O}(h^p).$$

85 In the following section, we aim to rigorously prove this expectation.

### 87 2.3. Convergence.

88 **Lemma 2.11 (Discrete Gronwall lemma).** Let  $A > 0$ ,  $B \geq 0$ ,  $h_\ell, a_\ell \geq 0$  such that

$$0 \leq a_{\ell+1} \leq (1 + h_\ell A)a_\ell + h_\ell B \quad \text{for all } \ell \in \mathbb{N}_0. \quad (2.10)$$

91 Then, it holds that

$$0 \leq a_\ell \leq \frac{B}{A} \left( \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) - 1 \right) + a_0 \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) \quad \text{for all } \ell \in \mathbb{N}_0. \quad (2.11)$$

94 *Proof.* It holds that

$$95 \quad 1 + x \leq \sum_{j=0}^{\infty} \frac{x^j}{j!} = \exp(x) \quad \text{for all } x \geq 0. \quad (2.12)$$

97 Note that (2.11) holds for  $\ell = 0$  (even with equality  $a_\ell = a_0$ ). Arguing by induction on  
98  $\ell$ , we may suppose that (2.11) holds up to some  $\ell \in \mathbb{N}_0$ . For  $\ell + 1$ , note that

$$\begin{aligned} 99 \quad a_{\ell+1} &\stackrel{(2.10)}{\leq} (1 + h_\ell A) a_\ell + h_\ell B \\ 100 \quad &\stackrel{(2.11)}{\leq} (1 + h_\ell A) \left[ \frac{B}{A} \left( \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) - 1 \right) + a_0 \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) \right] + h_\ell B \\ 101 \quad &\leq \frac{B}{A} \left( \underbrace{(1 + h_\ell A)}_{\stackrel{(2.12)}{\leq} \exp(h_\ell A)} \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) - (1 + h_\ell A) \right) + a_0 \underbrace{(1 + h_\ell A)}_{\stackrel{(2.12)}{\leq} \exp(h_\ell A)} \exp \left( A \sum_{j=0}^{\ell-1} h_j \right) + h_\ell B \\ 102 \quad &\leq \frac{B}{A} \left( \exp \left( A \sum_{j=0}^{\ell} h_j \right) - (1 + h_\ell A) \right) + a_0 \exp \left( A \sum_{j=0}^{\ell} h_j \right) + h_\ell B \\ 103 \quad &= \frac{B}{A} \left( \exp \left( A \sum_{j=0}^{\ell} h_j \right) - 1 \right) + a_0 \exp \left( A \sum_{j=0}^{\ell} h_j \right). \\ 104 \end{aligned}$$

105 This concludes the proof. □

---

106 **Theorem 2.12 (Consistency plus stability implies convergence).** Let  $\Phi(t, y, h)$  be  
107 the incremental function of an explicit one-step method. Suppose stability

$$108 \quad \forall t \in [t_0, T] \forall h \in (0, T - t] \forall y, \tilde{y} \in \mathbb{R}^n : \quad \|\Phi(t, y, h) - \Phi(t, \tilde{y}, h)\| \leq C_{\text{stab}} \|y - \tilde{y}\| \quad (2.13)$$

109 for some constant  $C_{\text{stab}} > 0$ . Let  $f \in C^p([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  with corresponding solution  
110  $y \in C^{p+1}([t_0, T]; \mathbb{R}^n)$ . Then, consistency order  $p \geq 1$  implies that

$$112 \quad \max_{\ell=1, \dots, N} \|y_\ell - y(t_\ell)\| \leq \frac{C_{\text{cons}}}{C_{\text{stab}}} [\exp(C_{\text{stab}}(T - t_0)) - 1] h_\Delta^p, \quad (2.14)$$

114 where  $C_{\text{cons}} > 0$  is the consistency constant; see (2.9) in Definition 2.7.

---

115 *Proof.* We define the consistency error

$$116 \quad \eta(t_{\ell+1}) := y(t_{\ell+1}) - [y(t_\ell) + h_\ell \Phi(t_\ell, y(t_\ell), h_\ell)] \stackrel{(2.9)}{=} \mathcal{O}(h_\ell^{p+1}).$$

118 With this notation, the error satisfies

$$\begin{aligned} 119 \quad E_{\ell+1} &:= y_{\ell+1} - y(t_{\ell+1}) \\ 120 \quad &= [y_\ell + h_\ell \Phi(t_\ell, y_\ell, h_\ell)] - [y(t_\ell) + h_\ell \Phi(t_\ell, y(t_\ell), h_\ell)] - \eta(t_{\ell+1}) \\ 121 \quad &= [y_\ell - y(t_\ell)] + h_\ell [\Phi(t_\ell, y_\ell, h_\ell) - \Phi(t_\ell, y(t_\ell), h_\ell)] - \eta(t_{\ell+1}). \end{aligned}$$

123 With  $a_\ell := \|y_\ell - y(t_\ell)\|$ , stability (2.13) and consistency (2.9) lead to

$$124 \quad 0 \leq a_{\ell+1} \leq a_\ell (1 + h_\ell C_{\text{stab}}) + C_{\text{cons}} h_\ell^{p+1} \leq a_\ell (1 + h_\ell C_{\text{stab}}) + C_{\text{cons}} h_\Delta^p h_\ell. \quad 125$$

126 With  $a_0 = 0$  and Lemma 2.11, we infer that

$$127 \quad \|y_\ell - y(t_\ell)\| = a_\ell \stackrel{(2.11)}{\leq} \frac{C_{\text{cons}}}{C_{\text{stab}}} \left( \exp \left( C_{\text{stab}} \sum_{j=0}^{\ell-1} h_j \right) - 1 \right) h_\Delta^p \quad \text{for all } \ell \in \mathbb{N}_0. \quad 128$$

129 The claim follows from  $\sum_{j=0}^{\ell-1} h_j \leq T - t_0$ . □

130 **Remark 2.13 (Stability of explicit Euler method).** *For the explicit Euler method,*  
 131 *it holds that  $\Phi(t, y, h) = f(t, y)$ . Hence, stability (2.13) with  $C_{\text{stab}} = L$  follows from*  
 132 *Lipschitz continuity of  $f$  in  $y$ .*

133 **Corollary 2.14 (One step methods are stable).** *Under the assumptions of Theo-*  
 134 *rem 2.12, we suppose that the computed approximations face rounding errors such that*

$$135 \quad \tilde{y}_0 = y_0 + \varepsilon_0 \quad \text{and} \quad \tilde{y}_{\ell+1} = \tilde{y}_\ell + h_\ell \Phi(t_\ell, \tilde{y}_\ell, h_\ell) + \delta_\ell \quad \text{for all } \ell = 0, \dots, N-1, \quad (2.15) \quad 136$$

137 where  $\|\varepsilon_0\| \leq \varepsilon$  and  $\|\delta_\ell\| \leq \delta h_\ell$ . Then,

$$138 \quad \max_{\ell=1, \dots, N} \|y_\ell - y(t_\ell)\| \leq C (h_\Delta^p + \delta) + \varepsilon \exp(C_{\text{stab}}(T - t_0)), \quad (2.16) \quad 139$$

140 where  $C = \frac{\max\{1, C_{\text{cons}}\}}{C_{\text{stab}}} [\exp(C_{\text{stab}}(T - t_0)) - 1]$ .

141 *Proof.* We argue as for the proof of Theorem 2.12. Define the consistency error

$$142 \quad \eta(t_{\ell+1}) := y(t_{\ell+1}) - [y(t_\ell) + h_\ell \Phi(t_\ell, y(t_\ell), h_\ell)] = \mathcal{O}(h_\ell^{p+1}). \quad 143$$

144 With this notation, the perturbed error satisfies

$$145 \quad \begin{aligned} \tilde{E}_{\ell+1} &:= \tilde{y}_{\ell+1} - y(t_{\ell+1}) \\ &= [\tilde{y}_\ell + h_\ell \Phi(t_\ell, \tilde{y}_\ell, h_\ell) + \delta_\ell] - [y(t_\ell) + h_\ell \Phi(t_\ell, y(t_\ell), h_\ell)] - \eta(t_{\ell+1}) \\ &= [\tilde{y}_\ell - y(t_\ell)] + h_\ell [\Phi(t_\ell, \tilde{y}_\ell, h_\ell) - \Phi(t_\ell, y(t_\ell), h_\ell)] - \eta(t_{\ell+1}) + \delta_\ell \end{aligned} \quad 146 \quad 147 \quad 148$$

149 With  $a_\ell := \|\tilde{y}_\ell - y(t_\ell)\|$ , stability (2.13) and consistency (2.9) lead to

$$150 \quad \begin{aligned} 0 \leq a_{\ell+1} &\leq a_\ell (1 + h_\ell C_{\text{stab}}) + C_{\text{cons}} h_\ell^{p+1} + \delta_\ell \\ &\leq a_\ell (1 + h_\ell C_{\text{stab}}) + \max\{1, C_{\text{cons}}\} (h_\Delta^p + \delta) h_\ell \end{aligned} \quad 151 \quad 152$$

153 With  $a_0 = \|\tilde{y}_0 - y_0\| \leq \varepsilon$  and Lemma 2.11, we infer that

$$154 \quad \begin{aligned} \|\tilde{y}_\ell - y(t_\ell)\| = a_\ell &\stackrel{(2.11)}{\leq} (h_\Delta^p + \delta) \frac{\max\{1, C_{\text{cons}}\}}{C_{\text{stab}}} \left( \exp \left( C_{\text{stab}} \sum_{j=0}^{\ell-1} h_j \right) - 1 \right) \\ &\quad + \varepsilon \exp \left( C_{\text{stab}} \sum_{j=0}^{\ell-1} h_j \right) \quad \text{for all } \ell \in \mathbb{N}_0. \end{aligned} \quad 155 \quad 156$$

157 With  $\sum_{j=0}^{\ell-1} h_j \leq T - t_0$ , this concludes the proof. □

158

159 **2.4. Examples.** We aim to construct explicit one-step methods with higher conver-  
 160 gence order. One natural approach is to start from the Taylor expansion of the exact  
 161 solution

$$162 \quad y(t+h) = y(t) + \sum_{k=1}^p \frac{y^{(k)}(t)}{k!} h^k + \mathcal{O}(h^{p+1}), \quad (2.17)$$

163  
 164 where we implicitly assume smoothness of  $f \in C^p([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and hence of  $y \in$   
 165  $C^{p+1}([t_0, T]; \mathbb{R}^n)$ . Since  $y$  is unknown (and so is  $y^{(k)}$ ), we need to express the derivatives  
 166 in terms of the given right-hand side  $f$ . First, recall that

$$167 \quad y'(t) = f(t, y(t)). \quad (2.18)$$

168  
 169 To write the second derivative of  $f$ , define  $g(t) := (t, y(t))$ . Then, the chain rule gives

$$170 \quad \begin{aligned} y''(t) &= d_t f(t, y(t)) = d_t [f(g(t))] = Df(t, y(t)) Dg(t) \\ &= (\partial_t f)(t, y(t)) + (\partial_y f)(t, y(t)) y'(t) \\ &= (\partial_t f)(t, y(t)) + (\partial_y f)(t, y(t)) f(t, y(t)). \end{aligned} \quad (2.19)$$

171  
 172 Proceeding with the chain rule, we can express all derivatives of  $y$  in terms of partial  
 173 derivatives of  $f$ . We leave this to the interested reader, but state the following example  
 174 for  $p = 2$ .

---

175 **Example 2.15 (Second-order one-step method based on Taylor expansion).**  
 176 Given  $f \in C^2([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$ , define the incremental function

$$177 \quad \Phi(t, y, h) := f(t, y) + \frac{h}{2} [(\partial_t f)(t, y) + (\partial_y f)(t, y) f(t, y)]. \quad (2.20)$$

178  
 179 Then, the corresponding explicit one-step method has consistency order  $p = 2$ .

---

180 One drawback of the latter construction is that it requires high regularity of  $f$  to write  
 181 down the incremental function. Moreover, one has to provide explicit formulas for the  
 182 derivatives of  $f$  and, more importantly, stability (2.13) of the one-step method requires  
 183 additional assumptions.

184 Alternatively, one can consider nested evaluations of  $f$  to avoid the evaluation of the  
 185 derivatives of  $f$ .

---

186 **Proposition 2.16.** Given  $a, b_1, b_2, c \in \mathbb{R}$ , define the incremental function

$$187 \quad \Phi(t, y, h) := b_1 f(t, y) + b_2 f(t + ch, y + ah f(t, y)). \quad (2.21)$$

188  
 189 Then, the following statements are equivalent:

- 190 (i) The corresponding one-step method has consistency order  $p = 2$ .  
 191 (ii)  $a = c$ ,  $b_1 + b_2 = 1$ , and  $b_2 a = 1/2$ .
- 

192 **Remark 2.17.** The conditions in Proposition 2.16(ii) show that one does not have four  
 193 degrees of freedom (i.e.,  $a, b_1, b_2, c \in \mathbb{R}$ ), but only one to ensure second-order consistency.  
 194 Fixing any of these constants, all other constants are determined as well. Moreover,  
 195 second-order consistency requires  $a \neq 0 \neq c$ . Finally, we note that a reasonable method

196 will additionally impose the restriction  $c \in (0, 1]$  to ensure that the evaluation point  $t + ch$   
 197 will belong to the time interval  $[t_0, T]$ .

198 *Proof of Proposition 2.16.* Let  $f \in C^2([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and hence  $y \in C^3([t_0, T]; \mathbb{R}^n)$  be  
 199 the exact solution. The Taylor expansion of  $\Phi$  around  $h = 0$  proves that

$$200 \quad \Phi(t, y, h) = \Phi(t, y, 0) + h \partial_h \Phi(t, y, 0) + \mathcal{O}(h^2).$$

201  
 202 Note that

$$203 \quad \Phi(t, y, 0) = b_1 f(t, y) + b_2 f(t, y) = (b_1 + b_2) f(t, y),$$

$$204 \quad \partial_h \Phi(t, y, 0) = b_2 c \partial_t f(t, y) + b_2 a \partial_y f(t, y) f(t, y).$$

205  
 206 Hence,

$$207 \quad \Phi(t, y, h) = (b_1 + b_2) f(t, y) + h [b_2 c \partial_t f(t, y) + b_2 a \partial_y f(t, y) f(t, y)] + \mathcal{O}(h^2). \quad (2.22)$$

208  
 209 Taylor expansion of  $y$  around  $t$  proves that

$$210 \quad y(t+h) = y(t) + h y'(t) + \frac{h^2}{2} y''(t) + \mathcal{O}(h^3) \quad (2.23)$$

$$211 \quad \stackrel{(2.19)}{=} y(t) + h f(t, y(t)) + \frac{h^2}{2} [(\partial_t f)(t, y(t)) + (\partial_y f)(t, y(t)) f(t, y(t))]$$

212 Combining (2.22)–(2.23), we can identify the consistency error. To ease the presentation,  
 213 we omit the arguments  $(t, y(t))$  for the  $f$ -terms. Then,

$$214 \quad y(t+h) - [y(t) + h \Phi(t, y(t), h)]$$

$$215 \quad = h f + \frac{h^2}{2} [\partial_t f + (\partial_y f) f] - h (b_1 + b_2) f - h^2 [b_2 c \partial_t f + b_2 a (\partial_y f) f] + \mathcal{O}(h^3)$$

$$216 \quad = h [1 - (b_1 + b_2)] f + h^2 [1/2 - b_2 c] \partial_t f + h^2 [1/2 - b_2 a] (\partial_y f) f + \mathcal{O}(h^3).$$

217  
 218 Consequently,

$$219 \quad y(t+h) - [y(t) + h \Phi(t, y(t), h)] = \mathcal{O}(h^3)$$

220 is equivalent to the fact that all lower-order powers of  $h$  vanish, i.e.,

$$221 \quad b_1 + b_2 = 1, \quad b_2 c = 1/2, \quad b_2 a = 1/2.$$

222 In particular, it follows that  $a \neq 0$ ,  $c \neq 0$ ,  $b_2 \neq 0$ , and  $a = c$ . This concludes the  
 223 proof.  $\square$

---

224 **Example 2.18 (Heun method (1900)).** Choose  $a = 1 = c$ . Then, Proposition 2.16(ii)  
 225 suggests the choice  $b_1 = 1/2 = b_2$ . Define

$$226 \quad \Phi(t, y, h) := \frac{1}{2} f(t, y) + \frac{1}{2} f(t+h, y+h f(t, y)). \quad (2.24)$$

---

**Karl Heun (1859–1929)** was a German mathematician. He took his PhD from the University of Göttingen in 1881 and completed his habilitation at LMU München. From 1890–1902, he worked as a teacher in Berlin, before he obtained the chair of theoretical mechanics at Technische Hochschule Darmstadt, where he retired in 1922.

230 According to Proposition 2.16, the resulting one-step method has consistency order  $p = 2$ .  
 231 Clearly, Lipschitz continuity of  $f$  in  $y$  yields stability (2.13). For the implementation, the  
 232 Heun method computes

$$\begin{aligned} 233 \quad k_1 &:= f(t, y), \\ 234 \quad k_2 &:= f(t + h, y + hk_1), \\ 235 \quad \Phi(t, y, h) &:= \frac{k_1 + k_2}{2}. \end{aligned}$$

237 **Example 2.19 (Modified Euler method (Runge 1895)).** Choose  $a = 1/2 = c$ .  
 238 Then, Proposition 2.16(ii) suggests the choice  $b_1 = 0$  and  $b_2 = 1$ . Define

$$239 \quad \Phi(t, y, h) := f\left(t + \frac{h}{2}, y + \frac{h}{2} f(t, y)\right). \quad (2.25)$$

241 According to Proposition 2.16, the resulting one-step method has consistency order  $p = 2$ .  
 242 Clearly, Lipschitz continuity of  $f$  in  $y$  yields stability (2.13). For the implementation,  
 243 Runge's modified Euler method computes

$$\begin{aligned} 244 \quad k_1 &:= f(t, y), \\ 245 \quad k_2 &:= f\left(t + \frac{h}{2}, y + \frac{h}{2} k_1\right), \\ 246 \quad \Phi(t, y, h) &:= k_2. \end{aligned}$$

248 **Remark 2.20.** As far as the effectivity of an integrator is concerned, one should plot the  
 249 error vs. the number of  $f$ -evaluations. Then, higher-order methods pay, if the solution is  
 250 smooth: For a uniform time-step size  $h$ , the number of time-steps satisfies  $N = (T - t_0)/h$ .

- 251 • The explicit Euler method has  $N$  evaluations of  $f$ , and the error decays like  $\mathcal{O}(h) =$   
 252  $\mathcal{O}(N^{-1})$ , i.e., we get an algebraic decay with rate 1.
- 253 • The modified Euler method has  $2N$  evaluations of  $f$ , and the error decays like  
 254  $\mathcal{O}(h^2) = \mathcal{O}(N^{-2})$ , i.e., we get an algebraic decay with rate 2.

255 Asymptotically, the modified Euler method will thus beat the explicit Euler method with  
 256 respect to accuracy and computational time.

257 **Remark 2.21.** Second-order consistency / convergence is only visible if  $f$  (and hence  $y$ )  
 258 are sufficiently smooth. If  $f$  is not smooth, then numerical experiments with uniform  
 259 meshes lead to order reductions.

**Carl Runge (1856–1927)** was a German mathematician. He studied at LMU München and Uni-  
 versity of Bremen, before he took his PhD in 1880 at HU Berlin supervised by Weierstrass. He did his  
 habilitation in 1883 at HU Berlin and became Professor at University of Hannover in 1886. In 1904,  
 he accepted the Chair of Applied Mathematics at University of Göttingen (the first chair for applied  
 mathematics in Germany). Runge is the father-in-law of Richard Courant.



260

261 **2.5. Explicit Runge–Kutta methods.** In this section, we generalize / formalize  
262 the idea of Proposition 2.16.

263 **Definition 2.22.** Let  $A \in \mathbb{R}^{m \times m}$  be strictly lower triangular (i.e.,  $A_{ij} = 0$  for  $i \leq j$ ),  
264  $b, c \in \mathbb{R}^m$  with  $0 \leq c_1 \leq c_2 \leq \dots \leq c_m \leq 1$ . Then, a one-step method with incremental  
265 function

266 
$$\Phi(t, y, h) := \sum_{j=1}^m b_j k_j, \tag{2.26}$$

267 where the so-called **increments** satisfy that

268 
$$k_i = f\left(t + c_i h, y + h \sum_{j=1}^{i-1} A_{ij} k_j\right) \quad \text{for all } i = 1, \dots, m, \tag{2.27}$$

269 is called **explicit m-stage Runge–Kutta method**. The y-arguments  $y + h \sum_{j=1}^{i-1} A_{ij} k_j$   
270 of the increments  $k_i$  are called **stages**. Usually, Runge–Kutta methods are denoted by  
271 their **Butcher tableau**  $\begin{array}{c|c} c & A \\ \hline & b^\top \end{array}$ . If the data are explicitly given, usually the zero entries  
272 of  $A$  are omitted (see examples below).  
273  
274

---

275 **Example 2.23.**

- 276 • The **explicit Euler method** has the Butcher tableau

277 
$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} = \begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array}.$$

- 278  
279 • The **modified Euler** has the Butcher tableau

280 
$$\begin{array}{c|cc} 0 & & \\ \hline 1/2 & 1/2 & \\ \hline & 0 & 1 \end{array} = \begin{array}{c|ccc} 0 & & 0 & 0 \\ \hline 1/2 & 1/2 & 0 & \\ \hline & 0 & 0 & 1 \end{array}.$$

- 281  
282 • The **Heun method** has the Butcher tableau

283 
$$\begin{array}{c|cc} 0 & & \\ \hline 1 & 1 & \\ \hline & 1/2 & 1/2 \end{array} = \begin{array}{c|ccc} 0 & & 0 & 0 \\ \hline 1 & 1 & 0 & \\ \hline & 1/2 & 1/2 & \end{array}.$$

284  
**Martin Wilhelm Kutta (1867–1944)** was a German mathematician. He studied mathematics at TU München. In 1901, we took his PhD with the thesis “*Beitrag zur näherungsweise Integration totaler Differentialgleichungen*”, where he was generalizing the ideas of Runge. Later, he was associate professor at University of Jena (1909–1910), RWTH Aachen (1910–1912), before he became full professor at University of Stuttgart in 1912. He retired in 1935.

**John Charles Butcher (born 1933)** is a mathematician from New Zealand. He took his PhD at University of Sydney in 1961. In 1966, he became professor at the University of Auckland. He retired in 1999. He has made fundamental contributions to the mathematical understanding of Runge–Kutta methods. In 2010, Butcher was awarded the Jones Medal from the Royal Society of New Zealand for his “exceptional lifetime work on numerical methods for the solution of differential equations”.

---

285 **Proposition 2.24 (Consistent Runge–Kutta methods provide quadrature rules).**

286 Let  $\frac{c}{b^\top} \Big| \frac{A}{b^\top}$  be an (explicit)  $m$ -stage Runge–Kutta method with consistency order  $p \geq 1$ .

287 Then, the vectors  $c, b \in \mathbb{R}^m$  provide a quadrature, which is exact for polynomials  $q \in \mathbb{P}_{p-1}$ ,  
 288 i.e.,

$$\sum_{j=1}^m b_j q(c_j) = \int_0^1 q(s) \, ds \quad \text{for all } q \in \mathbb{P}_{p-1}. \quad (2.28)$$

291 *Proof.* Consider the integration problem  $f(t, y) := t^p$  and  $y(0) = 0$ . The exact solution  
 292 is the antiderivative  $y(t) = \int_0^t t^p \, ds = \frac{t^{p+1}}{p+1}$ . On the one hand, it holds that

$$k_j \stackrel{(2.27)}{=} (t + c_j h)^p = \sum_{i=0}^p \binom{p}{i} t^{p-i} c_j^i h^i$$

295 and hence

$$\Phi(t, y, h) = \sum_{j=1}^m b_j k_j = \sum_{i=0}^p \binom{p}{i} t^{p-i} h^i \sum_{j=1}^m b_j c_j^i. \quad (2.29)$$

298 On the other hand, note that

$$y^{(i)}(t) = \frac{1}{p+1} \frac{(p+1)!}{(p+1-i)!} t^{p+1-i} = \frac{p!}{(p+1-i)!} \frac{(i-1)!}{(i-1)!} t^{p+1-i} = \binom{p}{i-1} (i-1)! t^{p+1-i}$$

301 and hence

$$y(t+h) - y(t) = \sum_{i=1}^{p+1} \frac{y^{(i)}(t)}{i!} h^i = \sum_{i=1}^{p+1} \binom{p}{i-1} \frac{1}{i} t^{p+1-i} h^i = \sum_{i=0}^p \binom{p}{i} \frac{1}{i+1} t^{p-i} h^{i+1}. \quad (2.30)$$

304 Combining (2.29)–(2.30), we obtain for the consistency error that

$$\mathcal{O}(h^{p+1}) = y(t+h) - [y(t) + h \Phi(t, y(t), h)] = \sum_{i=0}^p \binom{p}{i} t^{p-i} h^{i+1} \left[ \frac{1}{i+1} - \sum_{j=1}^m b_j c_j^i \right].$$

307 Consequently, the lower-order powers of  $h$  must vanish, i.e.,

$$\int_0^1 s^i \, ds = \frac{1}{i+1} = \sum_{j=1}^m b_j c_j^i \quad \text{for all } i = 0, \dots, p-1. \quad (2.31)$$

310 The identity (2.31) proves that the quadrature rule is exact for all monomials  $1, s, s^2, \dots, s^{p-1}$ .  
 311 Due to linearity, this proves (2.28).  $\square$

---

312 **Example 2.25 (Classical Runge–Kutta method RK4 (Runge 1901)).** Recall the  
 313 *Simpson rule*

$$314 \int_0^1 g \, ds \approx \frac{1}{6} [g(0) + 4g(1/2) + g(1)], \quad (2.32)$$

315  
 316 *which is exact for polynomials of degree 3, i.e.,*

$$317 \int_0^1 q \, ds = \frac{1}{6} [q(0) + 4q(1/2) + q(1)]. \quad (2.33)$$

319 *Due to Proposition 2.24, we aim to extend the Simpson rule to a Runge–Kutta method of*  
 320 *order  $p = 4$ . Runge introduced the method*

$$321 \begin{array}{c|ccc} 0 & & & \\ 1/2 & 1/2 & & \\ 1/2 & 0 & 1/2 & \\ 1 & 0 & 0 & 1 \\ \hline & 1/6 & 1/3 & 1/3 & 1/6 \end{array}. \quad (2.34)$$

322  
 323 *In explicit terms, the method reads*

$$324 \begin{aligned} k_1 &= f(t, y), \\ k_2 &= f\left(t + \frac{h}{2}, y + \frac{h}{2} k_1\right), \\ k_3 &= f\left(t + \frac{h}{2}, y + \frac{h}{2} k_2\right), \\ k_4 &= f(t + h, y + h k_3), \\ \Phi(t, y, h) &= \frac{1}{6} [k_1 + 2k_2 + 2k_3 + k_4]. \end{aligned}$$

328  
 329  
 330 *Note that (2.34) in fact induces the Simpson rule. By Taylor expansion, one can show*  
 331 *that this method has consistency order  $p = 4$ . Moreover, we will see below that it needs*  
 332 *(at least)  $m = 4$  stages to get consistency order  $p = 4$ . In this sense, RK4 is optimal.*

333 **Example 2.26 (Another extension of the Simpson rule).** Consider the Butcher  
 334 *tableau*

$$335 \begin{array}{c|cc} 0 & & \\ 1/2 & 1/2 & \\ 1 & -1 & 2 \\ \hline & 1/6 & 2/3 & 1/6 \end{array}. \quad (2.35)$$

336  
 337 *Note that (2.35) also induces the Simpson rule. Together with Runge’s method from*  
 338 *Example 2.25, we thus see that the extension of a quadrature rule to a Runge–Kutta*  
 339 *method is not unique. By Taylor expansion, one can show that the method (2.35) has*  
 340 *consistency order  $p = 3$ . Moreover, we will see below that this consistency order is*  
 341 *maximal (due to  $m = 3$  stages).*

342 **Theorem 2.27 (Stability of explicit Runge–Kutta methods).** Let  $\frac{c}{b^\top} \frac{A}{b^\top}$  be an  
 343 (explicit)  $m$ -stage Runge–Kutta method. Then, there exist  $C_{\text{stab}} > 0$  such that

$$344 \quad \forall t \in [t_0, T] \forall h \in (0, T - t] \forall y, \tilde{y} \in \mathbb{R}^n : \quad \|\Phi(t, y, h) - \Phi(t, \tilde{y}, h)\| \leq C_{\text{stab}} \|y - \tilde{y}\|. \quad (2.36)$$

345  
 346 Let  $L > 0$  be the Lipschitz constant of  $f$  in  $y$ . Then, it holds that

$$347 \quad C_{\text{stab}} \leq L p(hL), \quad \text{where } p(s) = \sum_{j=0}^{m-1} \mu_j s^j \in \mathbb{P}_{m-1} \quad (2.37)$$

348  
 349 with  $\mu_j \geq 0$  and  $\mu_0 = \sum_{j=1}^m |b_j|$

---

350 *Proof.* Recall the increments  $k_i = f(x + c_i h, y + h \sum_{j=1}^{i-1} A_{ij} k_j)$ . Let  $\tilde{k}_i$  denote the incre-  
 351 ments with respect to  $\tilde{y}$ . The proof is split into two steps.

352 **Step 1.** We prove by induction on  $i = 1, \dots, m$  that

$$353 \quad \|k_i - \tilde{k}_i\| \leq L q_i(hL) \|y - \tilde{y}\| \quad \text{with} \quad q_i(s) = \sum_{j=0}^{i-1} \lambda_j s^j \in \mathbb{P}_{i-1}, \quad \lambda_j \geq 0, \quad \lambda_0 = 1. \quad (2.38)$$

354  
 355 For  $i = 1$ , Lipschitz continuity of  $f$  in  $y$  proves that

$$356 \quad \|k_1 - \tilde{k}_1\| = \|f(t, y) - f(t, \tilde{y})\| \leq L \|y - \tilde{y}\|.$$

357  
 358 This proves (2.38) for  $i = 1$ . In the induction step  $(i - 1) \rightsquigarrow i$ , note that

$$359 \quad \|k_i - \tilde{k}_i\| \leq L \left\| y - \tilde{y} + h \sum_{j=1}^{i-1} A_{ij} (k_j - \tilde{k}_j) \right\|$$

$$360 \quad \leq L \left[ \|y - \tilde{y}\| + h \sum_{j=1}^{i-1} |A_{ij}| \|k_j - \tilde{k}_j\| \right]$$

$$361 \quad \leq L \|y - \tilde{y}\| \underbrace{\left[ 1 + hL \sum_{j=1}^{i-1} |A_{ij}| q_j(hL) \right]}_{=: q_i(hL)}.$$

362  
 363 Hence, (2.38) holds for all  $i = 1, \dots, m$ .

364 **Step 2.** Finally, it follows that

$$365 \quad \|\Phi(t, y, h) - \Phi(t, \tilde{y}, h)\| = \left\| \sum_{j=1}^m b_j (k_j - \tilde{k}_j) \right\| \stackrel{(2.38)}{\leq} L \underbrace{\sum_{j=1}^m |b_j| q_j(hL)}_{=: p(hL)} \|y - \tilde{y}\|.$$

366  
 367 This concludes the proof. □

---

368 **Theorem 2.28 (Necessary consistency conditions for Runge–Kutta methods).**

369 Let  $\frac{c}{b^\top} \left| \frac{A}{b^\top} \right|$  be an (explicit)  $m$ -stage Runge–Kutta method with consistency order  $p \geq 1$ .

370 Then, there hold the following statements:

371 (i)  $\sum_{j=1}^m b_j = 1$ .

372 (ii)  $\sum_{j=1}^m b_j c_j^\ell = \frac{1}{\ell + 1}$  for all  $\ell = 0, \dots, p - 1$ .

373 (iii)  $\sum_{j=1}^m b_j (A^\ell \mathbf{1})_j = \frac{1}{(\ell + 1)!}$  for all  $\ell = 0, \dots, p - 1$ , where  $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^m$ .

374 *Proof.* (i) follows already from (ii) and (iii) for  $\ell = 0$ . (ii) has already been proved in  
 375 Proposition 2.24; see the identity (2.31). Hence, it only remains to prove (iii).

376 To prove (iii), we consider the scalar problem  $y' = f(t, y) := y$  on  $[0, 1]$  with  $y(0) = 1$ .  
 377 The unique solution is  $y(t) = e^t$ . Let  $k := (k_1, \dots, k_m)^\top$  and note that

378  $k_i = f\left(t + c_i h, y + h \sum_{j=1}^{i-1} A_{ij} k_j\right) = y + h \sum_{j=1}^{i-1} A_{ij} k_j = y + h (Ak)_i$  for all  $i = 1, \dots, m$ .

379 In vector form, this identity reads

381  $(I - hA)k = k - hAk = y\mathbf{1}$ .

382 For small  $h$ , it holds that  $h\|A\| < 1$ . Hence, the so-called Neumann series proves that

384  $I - hA$  is invertible with  $(I - hA)^{-1} = \sum_{\ell=0}^{\infty} (hA)^\ell$ .

385 Therefore,

386  $\Phi(t, y, h) = \sum_{j=1}^m b_j k_j = b \cdot (I - hA)^{-1} (y\mathbf{1}) = b \cdot \sum_{\ell=0}^{\infty} (hA)^\ell (y\mathbf{1}) = y \sum_{\ell=0}^{\infty} h^\ell b \cdot (A^\ell \mathbf{1})$ .

387 On the other hand, Taylor expansion for  $y(t) = e^t$  gives

388  $y(t+h) - y(t) = \sum_{\ell=1}^{\infty} \frac{y^{(\ell)}(t)}{\ell!} h^\ell = y(t) \sum_{\ell=1}^{\infty} \frac{h^\ell}{\ell!} = y(t) \sum_{\ell=0}^{\infty} \frac{h^{\ell+1}}{(\ell+1)!}$ .

389 Combining these two identities, we obtain for the consistency error that

390  $\mathcal{O}(h^{p+1}) = y(t+h) - [y(t) + h\Phi(t, y(t), h)] = y(t) \sum_{\ell=0}^{\infty} h^{\ell+1} \left[ \frac{1}{(\ell+1)!} - b \cdot (A^\ell \mathbf{1}) \right]$ .

391 Hence, lower-order powers of  $h$  must vanish, i.e.,

392  $\frac{1}{(\ell+1)!} = b \cdot (A^\ell \mathbf{1})$  for all  $\ell = 0, \dots, p - 1$ .

393 This concludes the proof. □

---

399 **Corollary 2.29 (Naive Butcher barriers for explicit Runge–Kutta methods).**  
400 *For any explicit  $m$ -step Runge–Kutta method, the consistency order  $p \geq 1$  satisfies that*  
401  *$p \leq m$ .*

---

402 *Proof.* Let  $\left. \begin{array}{c} c \\ \hline b^\top \end{array} \right| A$  be the Butcher tableau of an explicit  $m$ -stage Runge–Kutta method.  
403 Recall that  $A = (A_{ij}) \in \mathbb{R}^{m \times m}$  is strictly lower triangular, i.e.,  $A_{ij} = 0$  for all  $i \leq j$ . Let  
404  $A^\ell = (A_{ij}^{(\ell)})$ . We show by induction on  $\ell = 1, \dots, m$  that

$$405 \quad A_{ij}^{(\ell)} = 0 \quad \text{for all } i \leq j + \ell - 1. \quad (2.39)$$

407 Obviously, the claim is OK for  $\ell = 1$ . For the induction step  $(\ell - 1) \rightsquigarrow \ell$ , note that

$$408 \quad A_{ij}^{(\ell)} = \sum_{k=1}^m A_{ik} A_{kj}^{(\ell-1)} = \sum_{k=1}^{i-1} A_{ik} A_{kj}^{(\ell-1)} \stackrel{(2.39)}{=} \sum_{k=j+\ell-1}^{i-1} A_{ik} A_{kj}^{(\ell-1)}.$$

410 The latter sum is empty (and hence 0) if  $i - 1 < j + \ell - 1$ , i.e.,  $i < j + \ell$ . This concludes  
411 the proof of (2.39).

412 For  $\ell = m$ , it follows from (2.39) that  $A^m = 0$ . Hence, Theorem 2.28(iii) cannot be  
413 satisfied for  $\ell = m$ . Therefore, we conclude that  $p \leq m$ .  $\square$

---

414 **Exercise 2.30.** Let  $\left. \begin{array}{c} c \\ \hline b^\top \end{array} \right| A$  be an (explicit)  $m$ -stage Runge–Kutta method. Show that the  
415 method has at least consistency order  $p = 1$ , if and only if  $\sum_{j=1}^m b_j = 1$ .

---

416 **Exercise 2.31.** Let  $\left. \begin{array}{c} c \\ \hline b^\top \end{array} \right| A$  be an (explicit)  $m$ -stage Runge–Kutta method. Suppose that

$$417 \quad \sum_{j=1}^m b_j = 1, \quad \sum_{j=1}^m b_j c_j = \frac{1}{2}, \quad \text{and} \quad \sum_{j=1}^{i-1} A_{ij} = c_i \quad \text{for all } i = 1, \dots, m. \quad (2.40)$$

419 Argue as for Proposition 2.16 to show that the considered Runge–Kutta method has at  
420 least consistency order  $p = 2$ .

---

421 **Remark 2.32.** In Exercise 2.31, we have seen that the consistency conditions of Theo-  
422 rem 2.28 are (essentially) sharp for  $m = 2 = p$ . However, it is known that the Butcher  
423 barriers  $p \leq m$  are not sharp in general. For  $m \geq 5$ , there are no explicit Runge–Kutta  
424 methods with consistency order  $p = m$ , i.e., it holds that  $m - p \geq 1$ ; see [But08, The-  
425 orem 324B]. Moreover  $m - p \geq 2$  for  $p \geq 7$  and  $m - p \geq 3$  for  $p \geq 8$ ; see [But08,  
426 page 188].

427 The precise growth of the maximal order  $p_{\max}(m)$  with respect to  $m$  as well as the the  
428 minimal stage number  $m_{\min}(p)$  with respect to  $p$  are still unknown, however  $p_{\max}(m) \rightarrow \infty$   
429 as  $m \rightarrow \infty$ ; see [But08, Theorem 324C].

430 The monograph [SWP12, Satz 2.4.6] refers to the original literature for the table

$$431 \frac{p}{m_{\min}(p)} \left| \begin{array}{cccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \end{array} \right| p \geq 9$$

$$432 \frac{p}{m_{\min}(p)} \left| \begin{array}{cccccccc} 1 & 2 & 3 & 4 & 6 & 7 & 9 & 11 \end{array} \right| m_{\min}(p) \geq p + 3$$

433 Conversely,

$$434 \frac{m}{p_{\max}(m)} \left| \begin{array}{cccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \end{array} \right| m \geq 9$$

$$435 \frac{m}{p_{\max}(m)} \left| \begin{array}{cccccccc} 1 & 2 & 3 & 4 & 4 & 5 & 6 & 6 & 7 \end{array} \right| p_{\max}(m) \leq m - 2$$

436 For  $p = 10$ , Hairer (1978) found an explicit Runge–Kutta method with  $m = 17$ . If  
437  $m_{\min}(10) < 17$  appears to be still open.

438 A nice motivation for RK methods is as follows: We employ quadrature on  $[0, 1]$  to  
439 approximate

$$440 y(t+h) = y(t) + \int_t^{t+h} y'(s) ds = y(t) + h \int_0^1 y'(t+sh) ds$$

$$441 = y(t) + h \int_0^1 f(t+sh, y(t+sh)) ds$$

$$442 \approx y(t) + h \sum_{j=1}^m b_j \underbrace{f(t+c_jh, y(t+c_jh))}_{\approx k_j}$$

$$443$$

$$444$$

444 This idea is also reflected by the following exercise.

445 **Exercise 2.33.** Let  $\frac{c}{b^\top}$  be an (explicit)  $m$ -stage Runge–Kutta method. Suppose that

$$446 \sum_{j=1}^m b_j = 1 \quad \text{and} \quad \sum_{j=1}^{i-1} A_{ij} = c_i \quad \text{for all } i = 1, \dots, m. \quad (2.41)$$

$$447$$

448 Let  $f \in C^1([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and suppose that  $f$  is Lipschitz continuous in  $y$ . Prove that

$$449 y'(t+c_jh) - k_j = f(t+c_jh, y(t+c_jh)) - k_j = \mathcal{O}(h^2). \quad (2.42)$$

$$450$$

451

452 **2.6. Adaptive time-step control.** In this section, we aim to (heuristically) design  
453 an algorithm of the following type:

454 **Algorithm 2.34 (Theoretical adaptive algorithm).**

455 **Input:** time interval  $[t_0, T]$ , initial value  $y_0$ , right-hand side  $f(\cdot, \cdot)$ , one-step method  
456 with incremental function  $\Phi(t, y, h)$ , tolerance  $\tau > 0$ , initial time-step size  $h_0$ , counter  
457  $\ell = 0$ .

458 REPEAT

**Ernst Hairer** (born 1949) is an Austrian mathematician. He took his PhD in 1972 at University of Innsbruck (supervised by Gerhard Wanner). In 1985, he became Associate Professor at University of Geneva, Switzerland. Since 1999, he is full professor at University of Geneva. He has coauthored the two-volume monograph “Solving Ordinary Differential Equations”, which is the main reference for research in the field of numerical integrators. Ernst Hairer is the father of the mathematician Martin Hairer, who won the Fields medal in 2014.

- 459 • Determine  $h > 0$  such that  $t + h \leq T$  and  $\|z(t_\ell + h) - [y_\ell + h \Phi(t_\ell, y_\ell, h)]\| \approx \tau h$ ,
- 460 where  $z$  solves  $z(t_\ell) = y_\ell$  and  $z' = f(t, z)$  in  $[t_\ell, T]$ .
- 461 • Define  $h_\ell := h$ ,  $t_{\ell+1} := t_\ell + h_\ell$ ,  $y_{\ell+1} := y_\ell + h_\ell \Phi(t_\ell, y_\ell, h_\ell)$ .
- 462 • Update counter  $\ell \mapsto \ell + 1$ .

463 UNTIL  $t_\ell = T$

464 **Output:** mesh  $\Delta = \{t_0, \dots, t_N = T\}$ , approximations  $y_\ell \approx y(t_\ell)$  for all  $\ell = 0, \dots, N$ .

---

465 **Remark 2.35 (Why adaptive time-step control?).**

- 466 • **Efficiency:** We want to compute / approximate  $y(T)$  as cheaply as possible,
  - 467 i.e., with as few evaluations of  $f$  as possible such that the error satisfies that
  - 468  $\|y(T) - y_N\| \approx \tau$ .
  - 469 • **Reliability:** Even if the solution is smooth, the error behavior  $\|y(T) - y_N\| =$
  - 470  $\mathcal{O}(h_\Delta^p)$  hides a multiplicative constant. Possibly, a chosen uniform step-size  $h =$
  - 471  $h_\Delta > 0$  is too coarse to ensure that  $\|y(T) - y_N\| \approx \tau$ .
  - 472 • If the solution  $y$  is non-smooth, then this will spoil the experimental convergence
  - 473 order. For uniform meshes, one usually obtains a convergence  $\mathcal{O}(h^q) = \mathcal{O}(N^{-q})$ ,
  - 474 where  $q > 0$  is (much) smaller than the consistency order  $p \geq 1$ . However, if the
  - 475 mesh is locally adapted to the points of reduced regularity (i.e., the singularities),
  - 476 in many cases one can recover  $\|y(T) - y_N\| = \mathcal{O}(N^{-p})$  for the error.
- 

477 **Remark 2.36 (What is the challenge?).** Since the exact solutions to ODEs are in

478 general unknown, the algorithm cannot evaluate the consistency error to ensure that

$$479 \quad \|z(t_\ell + h) - [y_\ell + h \Phi(t_\ell, y_\ell, h)]\| \approx \tau h,$$

481 where  $z$  solves  $z(t_\ell) = y_\ell$  and  $z' = f(t, z)$  in  $[t_\ell, T]$ . Therefore, we require appropriate

482 heuristics to circumvent this lack of knowledge.

---

483 **Heuristics 2.37 (Step-size control based on an  $p - (p + 1)$  strategy).** Let  $\Phi(t, y, h)$

484 be the incremental function of an explicit one-step method with consistency order  $p \geq 1$ .

485 Let  $\tilde{\Phi}(t, y, h)$  be the incremental function of an explicit one-step method with consistency

486 order  $p + 1$ . Let  $h > 0$  be given and  $\Phi(t, y, h)$  as well as  $\tilde{\Phi}(t, y, h)$  be computed. Then,

487 (up to higher-order terms)

$$488 \quad H := \left[ \frac{\tau}{\|\Phi(t, z(t), h) - \tilde{\Phi}(t, z(t), h)\|} \right]^{1/p} h \quad (2.43a)$$

490 would be OK to ensure that

$$491 \quad z(t + H) - [z(t) + H \Phi(t, z(t), H)] \approx H \tau. \quad (2.43b)$$

493 In explicit terms: If you compute with step-size  $h$ , you get a feedback on the appropriate

494 step-size  $H$ .

---

495 To see that Heuristics 2.37 makes sense, recall that consistency proofs are obtained by

496 means of the Taylor expansion. For instance, let us consider the explicit Euler method.



497 If the solution is smooth, then

$$\begin{aligned}
498 \quad z(t+h) &= z(t) + h z'(t) + \frac{h^2}{2} z''(t) + \mathcal{O}(h^3) \\
499 \quad &= z(t) + h \Phi(t, z(t), h) + \left[ \frac{h^2}{2} z''(t) + \mathcal{O}(h^3) \right]. \\
500
\end{aligned}$$

501 For a general one-step method with consistency order  $p \geq 1$ , one can similarly show that

$$502 \quad z(t+h) = z(t) + h \Phi(t, z(t), h) + h^{p+1} c(t) + \mathcal{O}(h^{p+2}), \quad (2.44)$$

503 where  $c(t) \sim z^{(p+1)}(t)$ . Define

$$\begin{aligned}
504 \quad z_1 &:= z(t) + h \Phi(t, z(t), h), \\
505 \quad \tilde{z}_1 &:= z(t) + h \tilde{\Phi}(t, z(t), h). \\
506 \quad & \\
507
\end{aligned}$$

508 Then, it holds that

$$\begin{aligned}
509 \quad z(t+h) - z_1 &= h^{p+1} c(t) + \mathcal{O}(h^{p+2}), \\
510 \quad z(t+h) - \tilde{z}_1 &= \mathcal{O}(h^{p+2}), \\
511
\end{aligned}$$

512 where we recall that  $\tilde{\Phi}(t, y, h)$  belongs to a method of order  $p+1$ . Hence,

$$513 \quad z_1 - \tilde{z}_1 = [z(t+h) - \tilde{z}_1] - [z(t+h) - z_1] = -c(t) h^{p+1} + \mathcal{O}(h^{p+2}).$$

514 We obtain that

$$515 \quad \|z_1 - \tilde{z}_1\| = \|c(t)\| h^{p+1} + \mathcal{O}(h^{p+2})$$

516 and consequently

$$517 \quad \|c(t)\| = \frac{\|z_1 - \tilde{z}_1\|}{h^{p+1}} + \mathcal{O}(h). \quad (2.45)$$

518 With the ansatz

$$\begin{aligned}
519 \quad \tau H &\stackrel{!}{=} \|z(t+H) - [z(t) + H \Phi(t, z(t), H)]\| \stackrel{(2.44)}{=} \|c(t)\| H^{p+1} + \mathcal{O}(H^{p+2}) \\
520 \quad &\stackrel{(2.45)}{=} \frac{\|z_1 - \tilde{z}_1\|}{h^{p+1}} H^{p+1} + \mathcal{O}(h H^{p+1}) + \mathcal{O}(H^{p+2}). \\
521
\end{aligned}$$

522 If we neglect the higher-order terms  $\mathcal{O}(h H^{p+1})$  and  $\mathcal{O}(H^{p+2})$ , we obtain that

$$523 \quad \tau H \approx \frac{\|z_1 - \tilde{z}_1\|}{h^{p+1}} H^{p+1}.$$

524 Rearranging this estimate, we are led to

$$525 \quad H^p \approx \frac{\tau h^{p+1}}{\|z_1 - \tilde{z}_1\|} = \frac{\tau h^p}{\|\Phi(t, z(t), h) - \tilde{\Phi}(t, z(t), h)\|},$$

526 where we recall that  $z_1 - \tilde{z}_1 = h [\Phi(t, z(t), h) - \tilde{\Phi}(t, z(t), h)]$ . This proves (2.43).

---

527 **Algorithm 2.38 (Practical adaptive algorithm).**

528 *Input:* time interval  $[t_0, T]$ , initial value  $y_0$ , right-hand side  $f(\cdot, \cdot)$ , one-step method  
529  $\Phi(t, y, h)$  of order  $p \geq 1$ , auxiliary one-step method  $\tilde{\Phi}(t, y, h)$  of order  $p+1$ , tolerance  
530  $\tau > 0$ , initial time-step size  $h > 0$ , minimal time-step size  $h_{\min} > 0$ , conformity factor  
531  $\lambda \geq 1$ , safety factor  $0 < \varrho \leq 1$ , counter  $\ell = 0$ .

537 REPEAT

538 [1]  $h := \min \{T - t_\ell, \max\{h_{\min}, h\}\}$

539 [2]  $F := \tilde{\Phi}(t_\ell, y_\ell, h)$

540 [3]  $H := \varrho \left[ \frac{\tau}{\|\Phi(t_\ell, y_\ell, h) - F\|} \right]^{1/p} h$

541 [4] IF  $h \leq H$  OR  $h \leq h_{\min}$

542 [5]  $t_{\ell+1} := t_\ell + h$

543 [6]  $y_{\ell+1} := y_\ell + hF$

544 [7] IF  $t_{\ell+1} < T$

545 [8]  $h := \min\{H, \lambda h\}$

546 [9] Update counter  $\ell \mapsto \ell + 1$

547 [10] END IF

548 [11] ELSE

549 [12]  $h := \min\{H, h/\lambda\}$

550 [13] END IF

551 UNTIL  $t_{\ell+1} = T$

552 **Output:** mesh  $\Delta = \{t_0, \dots, t_N = T\}$  and corresponding approximations  $y_\ell \approx y(t_\ell)$  for  
553 all  $\ell = 0, \dots, N$ .

---

554 **Remark 2.39 (Comments on Algorithm 2.38).**

555 **Line [1] and [4]:** To ensure that the algorithm is finite, we have to guarantee that  
556  $h \not\rightarrow 0$ . In fact, the algorithm ensures that  $h \geq h_{\min}$  up to the final time-step (where  
557 possibly  $T - t < h_{\min}$ ).

558 **Line [3] (and [2]):** Up to the safety factor  $0 < \varrho \leq 1$ , we use the formula for  $H$   
559 derived in (2.43). The safety factor tries to cover the fact that (2.43) was derived by  
560 neglecting higher-order terms (which could have an impact in our computation).

561 **Line [4]:** In order to avoid  $f$ -evaluations, we accept the time-step, if  $h \leq H$ , i.e., if  
562 the current step-size is smaller than the step-size allowed by our heuristics (2.43) (i.e., a  
563 longer time-step would have been OK).

564 **Line [6] (and [2]):** Usually  $\tilde{\Phi}(t_\ell, y_\ell, h)$  is more accurate. Therefore, we use  $\tilde{\Phi}(t_\ell, y_\ell, h)$   
565 instead of  $\Phi(t_\ell, y_\ell, h)$  for the next time step.

566 **Line [7] and [8]:** If the time-step was accepted, but the final time  $T$  has not been  
567 reached, the algorithm requires a guess for the step-size in the next time-step. To this  
568 end, we aim to choose  $h = H$ , but we enforce that the growth of the step-size is not too  
569 big (i.e., the ratio is bounded by  $\lambda$  if the step-size is increased). Again, we intend to avoid  
570  $f$ -evaluations. Therefore, the step-size guess should be accepted in the next time-step.

571 **Line [11] and [12]:** If  $H < h$ , then we cannot accept the time-step  $t_\ell + h$  and we have  
572 to recompute with a smaller  $h$ . To this end, we aim to choose  $h = H$ , but enforce at least  
573 a uniform reduction of the current step-size (by the factor  $\lambda^{-1}$ ).

574 **Attention:** Since we have made certain simplifications to get the error estimate (2.43)  
575 for the consistency error, we cannot rigorously guarantee that the adaptive algorithm is  
576 mathematically reliable in the sense that  $\|y(T) - y_N\| \leq \tau$  is guaranteed.

---

577 **Remark 2.40.** *Practical choices of the parameters are the following:*

- 578 •  $h \sim \tau^{1/p}$ , e.g.,  $h := \tau^{1/p}/10$ , since  $\tau \approx \|y(T) - y_N\| = \mathcal{O}(h^p)$  for smooth  $y$ .
- 579 •  $h_{\min} = \tau$ .
- 580 •  $\lambda = 2$ .
- 581 •  $\varrho = 0.8$ .

582 **Remark 2.41.** *For the implementation of Algorithm 2.38, one aims to choose  $\Phi(t, y, h)$*   
 583 *and  $\tilde{\Phi}(t, y, h)$  in a way that minimizes the number of  $f$ -evaluations. Practically relevant*  
 584 *are so-called **embedded Runge–Kutta methods**, which use the same increments  $k_j$*   
 585 *(and  $c_j$ ), but differ only for the vector  $b \in \mathbb{R}^m$ . They are usually denoted by*

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

586 where  $b \in \mathbb{R}^m$  belongs to the higher-order method and  $\beta \in \mathbb{R}^m$  belongs to the lower-order  
 589 method, i.e.,

$$\Phi(t, y, h) = \sum_{j=1}^m \beta_j k_j \quad \text{and} \quad \tilde{\Phi}(t, y, h) = \sum_{j=1}^m b_j k_j.$$

592 **Example 2.42 (Bogacki–Shampine pair RK3(2)).** *In 1989, Bogacki and Shampine*  
 593 *proposed the embedded 4-stage Runge–Kutta method*

$$\begin{array}{c|ccc} 0 & & & \\ 1/2 & 1/2 & & \\ 3/4 & 0 & 3/4 & \\ 1 & 2/9 & 1/3 & 4/9 \\ \hline & 2/9 & 1/3 & 4/9 & 0 \\ & 7/24 & 1/4 & 1/3 & 1/8 \end{array}$$

596 The vector  $b = (2/9, 1/3, 4/9, 0)^\top$  belongs to a third-order method, while the vector  $\beta =$   
 597  $(7/24, 1/4, 1/3, 1/8)^\top$  belongs to a second-order method.

598 For one time-step, the method needs 4  $f$ -evaluations instead of  $3 + 2 = 5$  for non-  
 599 embedded methods. Moreover, after the first time-step, RK3(2) has only 3  $f$ -evaluations  
 600 because of the **FSAL property** (i.e., first same as last): Since the last row of  $A$  coincides  
 601 with  $b$  (together with  $c_1 = 0$  and  $c_4 = 1$ ), the last increment of the current time-step  
 602 coincides with the first increment of the next time-step (i.e.,  $k_1(t_{\ell+1}) = k_4(t_\ell)$ ).

603 An adaptive Bogacki–Shampine method is provided by MATLAB function `ode23`.

**Przemyslaw Bogacki**, Professor at Old Dominion University, Virginia, USA; see [webpage]  
**Lawrence F. Shampine**, Professor Emeritus at Southern Methodist University, Texas, USA; see [webpage]

604 **Example 2.43 (Dormand–Prince pair RK5(4)).** In 1980, Dormand and Prince  
 605 proposed the embedded 7-stage Runge–Kutta method

606	0							
	1/5	1/5						
	3/10	3/40	9/40					
	4/5	44/45	−56/15	32/9				
606	8/9	19372/6561	−25360/2187	64448/6561	−212/729			
	1	9017/3168	−355/33	46732/5247	49/176	−5103/18656		
	1	35/384	0	500/1113	125/192	−2187/6784	11/84	
		35/384	0	500/1113	125/192	−2187/6784	11/84	0
607		5179/57600	0	7571/16695	393/640	−92097/339200	187/210	1/40

608 The vector  $b = (35/384, 0, 500/1113, 125/192, -2187/6784, 11/84, 0)^\top \in \mathbb{R}^7$  belongs to the fifth-order  
 609 method, while  $\beta = (5179/57600, 0, 7571/16695, 393/640, -92097/339200, 187/210, 1/40)^\top$  belongs to the  
 610 fourth-order method.

611 For one time-step, the method needs 7  $f$ -evaluations instead of  $5 + 4 = 9$  for non-  
 612 embedded methods. As the Bogacki–Shampine pair RK3(2) from Example 2.42, RK5(4)  
 613 has the FASL property and hence requires only 6  $f$ -evaluations after the first time-step.

614 An adaptive Dormand–Prince method is provided by MATLAB function `ode45`.

---

615

616 **2.7. Extrapolation.** If no auxiliary one-step method of order  $p + 1$  is at hand, one  
 617 can use the Richardson extrapolation to obtain the higher-order method.

618 **Heuristics 2.44 (Two-step Richardson extrapolation of one-step method).** Let  
 619  $\Phi(t, y, h)$  be the incremental function of an explicit one-step method with consistency order  
 620  $p \geq 1$ . Let

621 
$$z_1 := z(t) + h \Phi(t, z(h), h) \tag{2.46}$$

622 be one step of this method with step-size  $h$ . Moreover, let

624 
$$\widehat{z}_1 := \widehat{z}_{1/2} + \frac{h}{2} \Phi(t + h/2, \widehat{z}_{1/2}, h/2) \quad \text{with} \quad \widehat{z}_{1/2} := z(t) + \frac{h}{2} \Phi(t, z(t), h/2) \tag{2.47}$$

625 be the result of two successive steps with step-size  $h/2$ . Then, it holds that

627 
$$z(t + h) - \frac{2^p \widehat{z}_1 - z_1}{2^p - 1} = \mathcal{O}(h^{p+2}). \tag{2.48}$$

628 This procedure thus gives rise to an auxiliary one-step method of order  $p + 1$ .

---

**John R. Dormand:** According to his book on “Numerical methods for differential equations”, Dormand took his PhD in Physics at the University of York and then used to be a senior lecturer at the Department of Mathematics and Statistics, Teesside Polytechnique, Middlesbrough, Cleveland, UK. According to Scopus, his 1980 paper on RK5(4) has been cited more than 1400 times. His last paper appeared in 2003.

**Peter J. Prince** used to work at Department of Mathematics and Statistics, Teesside Polytechnique, Middlesbrough, Cleveland, UK. From 1978–1999, he was publishing papers on Runge–Kutta methods (mainly together with John R. Dormand).

630 To see that Heuristics 2.45 makes sense, we argue as before: Recall that consistency  
631 proofs are obtained by means of the Taylor expansion. For instance, let us consider the  
632 explicit Euler method. If the solution is smooth, then

$$633 \quad z(t+h) = z(t) + h \Phi(t, z(t), h) + \left[ \frac{h^2}{2} z''(t) + \mathcal{O}(h^3) \right].$$

635 For a general one-step method with consistency order  $p \geq 1$ , one can similarly show that

$$636 \quad z(t+h) = z(t) + h \Phi(t, z(t), h) + h^{p+1} c(t) + \mathcal{O}(h^{p+2}), \quad (2.49)$$

638 where  $c(t) \sim z^{(p+1)}(t)$  is at least  $C^1$ . Define

$$639 \quad z_1^* := z(t+h/2) + \frac{h}{2} \Phi(t+h/2, z(t+h/2), h/2).$$

641 Then,

$$642 \quad z(t+h) - \widehat{z}_1 = [z(t+h) - z_1^*] + [z_1^* - \widehat{z}_1]$$

$$643 \quad \stackrel{(2.49)}{=} [(h/2)^{p+1} c(t+h/2) + \mathcal{O}(h^{p+2})]$$

$$644 \quad + \left[ z(t+h/2) - \widehat{z}_{1/2} + \frac{h}{2} \left( \Phi(t+h/2, z(t+h/2), h/2) - \Phi(t+h/2, \widehat{z}_{1/2}, h/2) \right) \right].$$

646 If  $\Phi$  is stable in the sense of (2.13), consistency order  $p$  guarantees that

$$647 \quad \left\| \Phi(t+h/2, z(t+h/2), h/2) - \Phi(t+h/2, \widehat{z}_{1/2}, h/2) \right\| \leq L \|z(t+h/2) - \widehat{z}_{1/2}\|$$

$$648 \quad = \mathcal{O}(h^{p+1}).$$

650 Hence, we get that

$$z(t+h) - \widehat{z}_1 = (h/2)^{p+1} c(t+h/2) + [z(t+h/2) - \widehat{z}_{1/2}] + \mathcal{O}(h^{p+2})$$

$$651 \quad \stackrel{(2.49)}{=} (h/2)^{p+1} [c(t+h/2) + c(t)] + \mathcal{O}(h^{p+2}) \quad (2.50)$$

$$652 \quad = 2(h/2)^{p+1} c(t) + \mathcal{O}(h^{p+2}),$$

653 where we have finally used the Taylor expansion  $c(t+h/2) = c(t) + \mathcal{O}(h)$ . Together with

$$654 \quad z(t+h) - z_1 \stackrel{(2.49)}{=} c(t) h^{p+1} + \mathcal{O}(h^{p+2}),$$

656 we are led to

$$657 \quad \widehat{z}_1 - z_1 = [z(t+h) - z_1] - [z(t+h) - \widehat{z}_1] = c(t) h^{p+1} [1 - 2^{-p}] + \mathcal{O}(h^{p+2}).$$

659 This yields that

$$660 \quad c(t) = \frac{\widehat{z}_1 - z_1}{1 - 2^{-p}} h^{-(p+1)} + \mathcal{O}(h) \quad (2.51)$$

661

662 and hence

$$\begin{aligned}
z(t+h) - \widehat{z}_1 &\stackrel{(2.50)}{=} 2(h/2)^{p+1} c(t) + \mathcal{O}(h^{p+2}) \\
&\stackrel{(2.51)}{=} 2(h/2)^{p+1} \frac{\widehat{z}_1 - z_1}{1 - 2^{-p}} h^{-(p+1)} + \mathcal{O}(h^{p+2}) \\
&= 2^{-p} \frac{\widehat{z}_1 - z_1}{1 - 2^{-p}} + \mathcal{O}(h^{p+2}) \\
&= \frac{\widehat{z}_1 - z_1}{2^p - 1} + \mathcal{O}(h^{p+2}).
\end{aligned} \tag{2.52}$$

664 Since

$$\widehat{z}_1 + \frac{\widehat{z}_1 - z_1}{2^p - 1} = \frac{(2^p - 1)\widehat{z}_1 + (\widehat{z}_1 - z_1)}{2^p - 1} = \frac{2^p \widehat{z}_1 - z_1}{2^p - 1},$$

668 we conclude that

$$z(t+h) - \frac{2^p \widehat{z}_1 - z_1}{2^p - 1} = z(t+h) - \left[ \widehat{z}_1 + \frac{\widehat{z}_1 - z_1}{2^p - 1} \right] \stackrel{(2.52)}{=} \mathcal{O}(h^{p+2}),$$

671 i.e., the extrapolated method has consistency order  $p + 1$ .

---

672 **Heuristics 2.45 (Step-size control based on an  $h - h/2$  strategy).** We use the  
673 notation of Heuristics 2.44. For given  $h > 0$ , let  $z_1$  and  $\widehat{z}_1$  be computed by (2.46)–(2.47).  
674 Then, (up to higher-order terms)

$$H := \left[ \frac{2^p - 1}{2^p} \frac{\tau h^{p+1}}{\|z_1 - \widehat{z}_1\|} \right]^{1/p} \tag{2.53a}$$

677 would be OK to ensure that

$$z(t+H) - [z(t) + H\Phi(t, z(t), H)] \approx H\tau. \tag{2.53b}$$

680 *In explicit terms: If you compute with step-size  $h$ , you get a feedback on the appropriate*  
681 *step-size  $H$ . In particular, we can also employ (2.53) to steer Algorithm 2.38.*

---

682 We build on Heuristics 2.37 and use the auxiliary method provided by Heuristics 2.44:  
683 Note that

$$\begin{aligned}
\Phi(t, z(t), h) - \widetilde{\Phi}(t, z(t), h) &= z_1 - \frac{2^p \widehat{z}_1 - z_1}{2^p - 1} = \frac{(2^p - 1)z_1 - (2^p \widehat{z}_1 - z_1)}{2^p - 1} \\
&= \frac{2^p}{2^p - 1} (z_1 - \widehat{z}_1).
\end{aligned}$$

687 Therefore, (2.43) implies (2.53).

### 1 3. IMPLICIT ONE-STEP METHODS

2 Throughout this section, we consider the following model problem: Let  $[t_0, T]$  be a given  
3 time-interval. For given  $n \in \mathbb{N}$ , let  $f \in C([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and suppose that  $f$  is Lipschitz  
4 continuous in  $y$ , i.e.,

$$\forall t \in [t_0, T] \forall y, \widetilde{y} \in \mathbb{R}^n : \quad \|f(t, y) - f(t, \widetilde{y})\| \leq L \|y - \widetilde{y}\|, \tag{3.1}$$

7 where  $\|\cdot\|$  is an arbitrary norm on  $\mathbb{R}^n$  and  $L > 0$  is a fixed constant. Then, for any  
 8 initial value  $y_0 \in \mathbb{R}^n$ , the Picard–Lindelöf theorem guarantees existence and uniqueness  
 9 of  $y \in C^1([t_0, T]; \mathbb{R}^n)$  such that

$$10 \quad y(t_0) = y_0 \quad \text{and} \quad y'(t) = f(t, y(t)) \quad \text{for all } t \in [t_0, T]. \quad (3.2)$$

12 Let  $\Delta = \{t_0 < t_1 < \dots < t_N = T\}$  be a given mesh with local mesh-sizes  $h_\ell := t_{\ell+1} - t_\ell >$   
 13  $0$  for all  $\ell = 0, \dots, N-1$  and maximum mesh-size  $h_\Delta := \max_{\ell=0, \dots, N-1} h_\ell$ . As before, our  
 14 task is to compute approximations

$$15 \quad y_\ell \approx y(t_\ell) \quad \text{for all } \ell = 1, \dots, N. \quad (3.3)$$

17

18 **3.1. Fundamentals.** This section aims to briefly transfer concepts and results from  
 19 *explicit* one-step methods to *implicit* one-step methods.

20 **Definition 3.1.** For a given **incremental function**  $\Phi : [t_0, T] \times \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}^n$ ,  
 21 the inductive procedure

$$22 \quad y_{\ell+1} := y_\ell + h_\ell \Phi(t_\ell, y_\ell, y_{\ell+1}, h_\ell) \quad \text{for all } \ell = 0, \dots, N-1 \quad (3.4)$$

24 is called **implicit one-step method**.

25 **Exercise 3.2.** Suppose that the incremental function  $\Phi(t, y, z, h)$  is Lipschitz continuous  
 26 in  $z$ , i.e.,

$$27 \quad \exists L > 0 \forall t \in [t_0, T] \forall y \in \mathbb{R}^n \forall z, \tilde{z} \in \mathbb{R}^n \forall h \in (0, T-t) : \quad (3.5)$$

$$28 \quad \|\Phi(t, y, z, h) - \Phi(t, y, \tilde{z}, h)\| \leq L \|z - \tilde{z}\|.$$

29 Then, it holds that

$$30 \quad \exists H > 0 \forall t \in [t_0, T] \forall y \in \mathbb{R}^n \forall 0 < h \leq \min\{T-t, H\} \exists! z \in \mathbb{R}^n : \quad (3.6)$$

$$31 \quad z = y + h \Phi(t, y, z, h).$$

32 In particular, for sufficiently small  $h_\Delta > 0$ , the implicit one-step method (3.4) is well-  
 33 defined.

34 **Remark 3.3.** (i) While the proof of the last exercise is based on the Banach fixpoint the-  
 35 orem, one rather uses the Newton method than the Banach fixpoint iteration to compute  
 36  $y_{\ell+1}$  from (3.4) in practice. Usually, the fixpoint iteration requires the restrictive con-  
 37 dition  $h_\Delta L < 1$ , while the Newton method converges (in practice!) under much weaker  
 38 conditions.

39 (ii) If  $\Phi(t, y, z, h) = M(t, y, h)z + b(t, y, h)$  is affine in  $z$ , then one has to solve only one  
 40 linear system

$$41 \quad (I - h_\ell M(t_\ell, y_\ell, h_\ell))y_{\ell+1} = y_\ell + h_\ell b(t_\ell, y_\ell, h_\ell)$$

42 to compute the solution  $y_{\ell+1}$  of (3.4).  
 43

44 **Definition 3.4.** *The implicit one-step method corresponding to the incremental function*  
 45  *$\Phi(t, y, z, h)$  has consistency order  $p \geq 1$ , if the following is satisfied: Given  $f \in$*   
 46  *$C^p([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$ , the exact solution  $y \in C^{p+1}([t_0, T]; \mathbb{R}^n)$  of (3.2) satisfies that*

$$\exists C > 0 \forall t \in [t_0, T] \forall h \in (0, T - t] : \quad \|y(t+h) - [y(t) + h \Phi(t, y(t), y(t+h), h)]\| \leq C h^{p+1}. \quad (3.7)$$

---

49 **Example 3.5 (Implicit Euler method).** *Recall that*

$$y_{\ell+1} := y_\ell + h_\ell f(t_{\ell+1}, y_{\ell+1}), \quad \text{i.e.,} \quad \Phi(t, y, z, h) := f(t+h, z).$$

50  
 51  
 52 *According to the Taylor theorem for  $t = (t+h) - h$ , it holds that*

$$y(t) = y(t+h) - h y'(t+h) + \mathcal{O}(h^2).$$

53  
 54  
 55 *Hence, we see that*

$$\begin{aligned} y(t+h) - [y(t) + h \Phi(t, y(t), y(t+h), h)] &= y(t+h) - [y(t) + h f(t+h, y(t+h))] \\ &= y(t+h) - [y(t) + h y'(t+h)] = \mathcal{O}(h^2). \end{aligned}$$

56  
 57  
 58  
 59 *Therefore, the implicit Euler method has consistency order  $p = 1$ .*

---

60 **Example 3.6 (Implicit midpoint rule).** *The implicit midpoint rule is defined as*

$$y_{\ell+1} := y_\ell + h_\ell f\left(t_\ell + \frac{h_\ell}{2}, \frac{y_\ell + y_{\ell+1}}{2}\right), \quad \text{i.e.,} \quad \Phi(t, y, z, h) := f\left(t + \frac{h}{2}, \frac{y+z}{2}\right).$$

61  
 62  
 63 *According to the Taylor theorem for  $(t+h/2) \pm h/2$ , it holds that*

$$y(t) = y(t+h/2) - \frac{h}{2} y'(t+h/2) + \frac{h^2}{8} y''(t+h/2) + \mathcal{O}(h^3),$$

$$y(t+h) = y(t+h/2) + \frac{h}{2} y'(t+h/2) + \frac{h^2}{8} y''(t+h/2) + \mathcal{O}(h^3).$$

64  
 65  
 66  
 67 *On the one hand, this shows that*

$$y(t+h) - y(t) = h y'(t+h/2) + \mathcal{O}(h^3) = h f(t+h/2, y(t+h/2)) + \mathcal{O}(h^3).$$

68  
 69  
 70 *On the other hand, this shows that*

$$\frac{y(t+h) + y(t)}{2} = y(t+h/2) + \mathcal{O}(h^2)$$

71  
 72  
 73 *and hence*

$$f\left(t+h/2, y(t+h/2)\right) = f\left(t+h/2, \frac{y(t+h) + y(t)}{2}\right) + \mathcal{O}(h^2).$$

74  
 75  
 76 *Overall, this results in*

$$y(t+h) - \left[ y(t) + h f\left(t + \frac{h}{2}, \frac{y(t+h) + y(t)}{2}\right) \right] = \mathcal{O}(h^2).$$

77  
 78  
 79 *Therefore, the implicit midpoint rule has consistency order  $p = 2$ .*

---



80 **Exercise 3.7 (Trapezoidal rule).** *The trapezoidal rule is defined as*

$$81 \quad y_{\ell+1} := y_\ell + h_\ell \frac{f(t_\ell, y_\ell) + f(t_{\ell+1}, y_{\ell+1})}{2}, \quad \text{i.e.,} \quad \Phi(t, y, z, h) := \frac{f(t, y) + f(t+h, z)}{2}. \quad 82$$

83 *Show that the trapezoidal rule has consistency order  $p = 2$ .*

---

84 **Exercise 3.8 (Stability plus consistency implies convergence).** *Let  $\Phi(t, y, z, h)$*   
 85 *be the incremental function of an implicit one-step method with consistency order  $p \geq 1$ .*  
 86 *Suppose that  $\Phi(t, y, z, h)$  is stable, i.e.,*

$$87 \quad \exists L > 0 \forall t \in [t_0, T) \forall y, \tilde{y} \in \mathbb{R}^n \forall z, \tilde{z} \in \mathbb{R}^n \forall h \in (0, T-t] : \quad 88$$

$$\|\Phi(t, y, z, h) - \Phi(t, \tilde{y}, \tilde{z}, h)\| \leq L (\|y - \tilde{y}\| + \|z - \tilde{z}\|). \quad (3.8)$$

89 *Then, the solution  $y \in C^{p+1}([t_0, T]; \mathbb{R}^n)$  satisfies that*

$$90 \quad \max_{\ell=1, \dots, N} \|y(t_\ell) - y_\ell\| = \mathcal{O}(h_\Delta^p), \quad 91 \quad (3.9)$$

92 *whenever the discrete solutions exist (e.g.,  $h_\Delta \leq H$ ).*

---

93 The following proposition shows that all implicit one-step methods are locally explicit.  
 94 As far as, e.g., the adaptive step-size control is concerned, we can thus simply employ  
 95 the ideas developed for explicit one-step methods (since our ideas above have built only  
 96 on local Taylor expansions).

97 **Proposition 3.9 (Implicit methods are locally explicit).** *Let  $\Phi(t, y, z, h)$  be the in-*  
 98 *cremental function of an implicit one-step method with consistency order  $p \geq 1$ . Suppose*  
 99 *that  $\Phi(t, y, z, h)$  is Lipschitz continuous in  $z$  (see (3.5)) and continuously differentiable*  
 100 *in  $(t, y, z)$ . Let  $h > 0$  and  $(t, y, z)$  with*

$$101 \quad z = y + h \Phi(t, y, z, h) \quad \text{and} \quad h \|D_z \Phi(t, y, z, h)\| < 1. \quad 102 \quad (3.10)$$

103 *Then, there exist open sets  $U \subset [0, T] \times \mathbb{R}^n$  and  $V \subset \mathbb{R}^n$  with  $(t, y) \in U$  and  $z \in V$  as*  
 104 *well as a function  $g \in C^1(U; V)$  such that*

$$105 \quad \forall (\tilde{t}, \tilde{y}) \in U \forall \tilde{z} \in V : \quad (\tilde{z} = \tilde{y} + h \Phi(\tilde{t}, \tilde{y}, \tilde{z}, h) \iff \tilde{z} = \tilde{g}(\tilde{t}, \tilde{y})). \quad 106 \quad (3.11)$$

107 *In particular, one step of the one-step method (3.4) with step-size  $h$  is well-defined and*  
 108 *even explicit, since  $\Phi(t, y, z, h) = \Phi(t, y, g(t, y), h)$ .*

---

109 *Proof.* Consider  $F((\tilde{t}, \tilde{y}), \tilde{z}) := \tilde{z} - [\tilde{q} + h \Phi(\tilde{t}, \tilde{y}, \tilde{z}, h)]$ . Then,  $F \in C^1([t_0, T] \times \mathbb{R}^n \times \mathbb{R}^n; \mathbb{R}^n)$   
 110 and  $F((t, y), z) = 0$  by assumption. Moreover,  $\kappa := h \|D_z \Phi(t, y, z, h)\| < 1$  allows to  
 111 employ the Neumann series to see that

$$112 \quad D_z F((t, y), z) = I - h D_z \Phi(t, y, z, h) \quad 113$$

114 is regular with

$$115 \quad [D_z F((t, y), z)]^{-1} = \sum_{k=0}^{\infty} (h D_z \Phi(t, y, z, h))^k. \quad 116$$

117 Therefore, the claim follows from the implicit function theorem. □

118

119 **3.2. Implicit Runge–Kutta methods.**

120 **Definition 3.10.** Let  $A \in \mathbb{R}^{m \times m}$ ,  $b, c \in \mathbb{R}^m$  with  $0 \leq c_1 \leq c_2 \leq \dots \leq c_m \leq 1$ . Then, a  
 121 one-step method with incremental function

$$122 \quad \Phi(t, y, h) := \sum_{j=1}^m b_j k_j, \quad (3.12)$$

123 where the so-called **stages** satisfy the implicit conditions

$$125 \quad k_j = f\left(t + c_j h, y + h \sum_{\ell=1}^m A_{j\ell} k_\ell\right) \quad \text{for all } j = 1, \dots, m, \quad (3.13)$$

126 is called ***m*-stage Runge–Kutta method**. A method is called **implicit *m*-stage Runge–**  
 127 **Kutta method**, if the matrix  $A$  is not strictly lower triangular. Usually, Runge–Kutta

128 methods are denoted by their **Butcher tableau**  $\frac{c}{b^\top} \left| \begin{array}{c} A \\ b^\top \end{array} \right.$ .

129 We stress that well-posedness of an implicit Runge–Kutta method is not obvious, since  
 130 the equations for the stages (3.13) are (nonlinearly) coupled and implicit. However, well-  
 131 posedness is shown in the following proposition, if the step-size  $h$  is sufficiently small.

132 **Proposition 3.11.** Let  $\frac{c}{b^\top} \left| \begin{array}{c} A \\ b^\top \end{array} \right.$  be an *m*-stage Runge–Kutta method. Then, there exists  
 133  $H > 0$  such that the stages (3.13) are well-defined, i.e.,

134  $\forall t \in [t_0, T] \forall y \in \mathbb{R}^n \forall 0 < h < \min\{T - t, H\} \exists! k_1, \dots, k_m \in \mathbb{R}^n \forall j = 1, \dots, m :$

$$136 \quad k_j = f\left(t + c_j h, y + h \sum_{\ell=1}^m A_{j\ell} k_\ell\right) \quad (3.14)$$

137 *Proof.* We write  $K := (k_1, \dots, k_m) \in \mathbb{R}^{n \times m}$ . Consider

$$139 \quad \Psi : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \times m}, \quad (\Psi(K))_j := f\left(t + c_j h, y + h \sum_{\ell=1}^m A_{j\ell} k_\ell\right) \in \mathbb{R}^n \quad \text{for all } j = 1, \dots, m.$$

140 Note that the stage conditions (3.13) are equivalent to  $\Psi(K) = K$ . Hence, we only need  
 141 to show that  $\Psi$  has a unique fixpoint. To this end, we aim to apply the Banach fixpoint  
 142 theorem. Note that

$$144 \quad \begin{aligned} \|\Psi(K) - \Psi(\tilde{K})\|_\infty &:= \max_{j=1, \dots, m} \left\| f\left(t + c_j h, y + h \sum_{\ell=1}^m A_{j\ell} k_\ell\right) - f\left(t + c_j h, y + h \sum_{\ell=1}^m A_{j\ell} \tilde{k}_\ell\right) \right\| \\ &\leq hL \max_{j=1, \dots, m} \left\| \sum_{\ell=1}^m A_{j\ell} (k_\ell - \tilde{k}_\ell) \right\| \leq hL \max_{j=1, \dots, m} \sum_{\ell=1}^m |A_{j\ell}| \max_{i=1, \dots, m} \|k_i - \tilde{k}_i\| \\ &= \left( hL \max_{j=1, \dots, m} \sum_{\ell=1}^m |A_{j\ell}| \right) \|K\|_\infty. \end{aligned}$$

148 For sufficiently small  $h > 0$ , the mapping  $\Psi$  is hence a contraction. Therefore, the Banach  
 149 fixed point theorem concludes the proof.  $\square$

150 **Example 3.12.** We consider the Runge–Kutta method  $\frac{1}{1} \mid \frac{1}{1}$ . By definition, one step of  
 151 this method leads to  $y_{\ell+1} = y_{\ell} + h k_1$ , where

$$152 \quad k_1 = f(t_{\ell} + h_{\ell}, y_{\ell} + h_{\ell} k_1) = f(t_{\ell+1}, y_{\ell} + h_{\ell} k_1) = f(t_{\ell+1}, y_{\ell+1}).$$

154 We hence obtain that

$$155 \quad y_{\ell+1} = y_{\ell} + h f(t_{\ell+1}, y_{\ell+1}),$$

157 which is the implicit Euler method.

158 **Example 3.13.** We consider the Runge–Kutta method  $\frac{0}{1} \mid \begin{array}{cc} 0 & 0 \\ 1/2 & 1/2 \end{array}$ . Then,

$$159 \quad k_1 = f(t_{\ell}, y_{\ell}),$$

$$160 \quad k_2 = f\left(t_{\ell} + h_{\ell}, y_{\ell} + h_{\ell} \frac{k_1 + k_2}{2}\right),$$

$$161 \quad y_{\ell+1} = y_{\ell} + h_{\ell} \frac{k_1 + k_2}{2}.$$

163 From the last equation, we obtain that  $k_2 = f(t_{\ell+1}, y_{\ell+1})$ . Overall, we see that

$$164 \quad y_{\ell+1} = y_{\ell} + h_{\ell} \frac{f(t_{\ell}, y_{\ell}) + f(t_{\ell+1}, y_{\ell+1})}{2},$$

165 which is the implicit trapezoidal rule.

167 **Example 3.14.** We consider the Runge–Kutta method  $\frac{1/2}{1} \mid \frac{1/2}{1}$ . Then,

$$168 \quad k_1 = f\left(t_{\ell} + \frac{h_{\ell}}{2}, y_{\ell} + \frac{h_{\ell}}{2} k_1\right),$$

$$169 \quad y_{\ell+1} = y_{\ell} + h_{\ell} k_1.$$

171 From the last equation, we obtain that  $h_{\ell} k_1 = y_{\ell+1} - y_{\ell}$  and hence  $k_1 = f\left(t_{\ell} + \frac{h_{\ell}}{2}, \frac{y_{\ell} + y_{\ell+1}}{2}\right)$ .

172 Overall, we see that

$$173 \quad y_{\ell+1} = y_{\ell} + h_{\ell} f\left(t_{\ell} + \frac{h_{\ell}}{2}, \frac{y_{\ell} + y_{\ell+1}}{2}\right),$$

174 which is the implicit midpoint rule.

176 **Exercise 3.15.** Consider a general  $m$ -stage Runge–Kutta method  $\frac{c}{b^T} \mid \frac{A}{1}$ . Note that the  
 177 corresponding stages  $k_j = k_j(t, y, h)$  depend on  $(t, y, h)$  and that  $\Phi(t, y, h) = \sum_{j=1}^m b_j k_j$ .

178 By exploiting the Lipschitz continuity of  $f$  in  $y$ , show that the Runge–Kutta method is  
 179 stable in the sense of

$$\begin{aligned} & \exists H > 0 \exists C > 0 \forall t \in [t_0, T] \forall y, \tilde{y} \in \mathbb{R}^n \forall 0 < h < \min\{T - t, H\} : \\ & \|\Phi(t, y, h) - \Phi(t, \tilde{y}, h)\| \leq C \|y - \tilde{y}\|. \end{aligned} \tag{3.15}$$

182 Using stability and consistency, formulate and prove a convergence theorem for general  
 183 Runge–Kutta methods.

184 **Remark 3.16.** We have already proved that an  $m$ -stage Runge–Kutta method with con-  
 185 sistency order  $p \geq 1$  satisfies the following statements (i)–(iii):

- 186 (i)  $\sum_{j=1}^m b_j = 1.$
- 187 (ii)  $\sum_{j=1}^m b_j c_j^\ell = \frac{1}{\ell + 1}$  for all  $\ell = 0, \dots, p - 1.$
- 188 (iii)  $\sum_{j=1}^m b_j (A^\ell \mathbf{1})_j = \frac{1}{(\ell + 1)!}$  for all  $\ell = 0, \dots, p - 1$ , where  $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^m.$

189 **Proposition 3.17 (Consistency  $\leq 2m$ ).** Consider a general  $m$ -stage Runge–Kutta  
 190 method  $\frac{c}{b^\top} \Big| \frac{A}{b^\top}$  and suppose consistency order  $p \geq 1$ . Then, it holds that  $p \leq 2m$ .  
 191 Moreover, if  $p = 2m$ , then the vectors  $b, c \in \mathbb{R}^m$  are unique: The coefficients of  $c$  are  
 192 the nodes of the Gaussian quadrature rule on  $[0, 1]$  and that of  $b$  are the corresponding  
 193 weights.

194 *Proof.* Consistency order  $p$  yields that

$$\sum_{j=1}^m b_j c_j^\ell \stackrel{!}{=} \frac{1}{\ell + 1} = \int_0^1 t^\ell dt \quad \text{for all } \ell = 0, \dots, p - 1.$$

197 Hence, this quadrature is exact for polynomials of degree  $p - 1$ , i.e.,

$$\sum_{j=1}^m b_j q(c_j) = \int_0^1 q(t) dt \quad \text{for all } q \in \mathbb{P}_{p-1}.$$

200 The claim follows from the basic lecture on numerical analysis: First, the maximum  
 201 exactness for a quadrature rule with  $m$  nodes is  $2m - 1$  and hence  $p \leq 2m$ . Second, the  
 202 Gaussian quadrature is the unique quadrature rule with exactness  $2m - 1$ .  $\square$

203 **Remark 3.18.** Runge–Kutta methods can be interpreted as the application of appropriate  
 204 quadrature schemes employed to the integral representation of the exact solution. A simple

205 substitution shows that

$$\begin{aligned}
 206 \quad y(t+h) - y(t) &= \int_t^{t+h} y'(s) \, ds = h \int_0^1 y'(t+sh) \, ds = h \int_0^1 f(t+sh, y(t+sh)) \, ds \\
 207 \quad &\approx h \sum_{j=1}^m b_j f(t+c_j h, y(t+c_j h)) \approx h \sum_{j=1}^m b_j k_j.
 \end{aligned}$$

208  
209 To derive a formula for  $k_j \approx f(t+c_j h, y(t+c_j h))$ , we proceed analogously:

$$\begin{aligned}
 210 \quad y(t+c_j h) - y(t) &= \int_t^{t+c_j h} y'(s) \, ds = h \int_0^{c_j} y'(t+sh) \, ds = h \int_0^{c_j} f(t+sh, y(t+sh)) \, ds \\
 211 \quad &\approx h \sum_{\ell=1}^m A_{j\ell} f(t+c_\ell h, y(t+c_\ell h)) \approx h \sum_{\ell=1}^m A_{j\ell} k_\ell.
 \end{aligned}$$

---

213 With the foregoing interpretation, the following proposition gives a sufficient criterion  
214 for consistency of Runge–Kutta methods in terms of quadrature.

215 **Proposition 3.19 (Consistency in terms of quadrature rules).** *Let  $p \geq 1$ . Consider*

216 *a general Runge–Kutta method  $\begin{array}{c|c} c & A \\ \hline & b^\top \end{array}$  and suppose that*

$$217 \quad \sum_{j=1}^m b_j q(c_j) = \int_0^1 q \, dt \quad \text{for all } q \in \mathbb{P}_{p-1} \tag{3.16a}$$

218  
219 *as well as*

$$220 \quad \sum_{j=1}^m A_{\ell j} q(c_j) = \int_0^{c_\ell} q \, dt \quad \text{for all } q \in \mathbb{P}_{p-2} \text{ and all } \ell = 1, \dots, m, \tag{3.16b}$$

222 *where  $\mathbb{P}_s$  denotes the space of all polynomials of degree  $\leq s$ . Then, the Runge–Kutta*  
223 *method has at least consistency order  $p$ .*

---

224 **Lemma 3.20 (Lagrange interpolation).** *Let  $g \in C^s[a, b]$  and  $a \leq t_1 < \dots < t_s \leq b$ .*  
225 *Then, there exists a unique polynomial*

$$226 \quad q \in \mathbb{P}_{s-1} \quad \text{such that} \quad q(t_j) = g(t_j) \quad \text{for all } j = 1, \dots, s, \tag{3.17}$$

227  
228 *and it holds that*

$$229 \quad q(t) = \sum_{j=1}^s g(t_j) L_j(t) \quad \text{with the **Lagrange polynomials** } L_j(t) := \prod_{\substack{k=1 \\ k \neq j}}^s \frac{t - t_k}{t_j - t_k}. \tag{3.18}$$

230  
231 *Moreover, for all  $k = 0, \dots, s-1$  and all  $t \in [a, b]$ , there holds the error identity*

$$232 \quad g^{(k)}(t) - q^{(k)}(t) = \frac{g^{(s)}(\xi)}{(s-k)!} \prod_{\ell=1}^{s-k} (t - \zeta_\ell) \tag{3.19}$$

234 with appropriate scalars  $\xi = \xi(k, t), \zeta_\ell = \zeta_\ell(k) \in [a, b]$ . In particular, it follows that

$$235 \quad \|g^{(k)} - q^{(k)}\|_{\infty, [a, b]} \leq \frac{\|g^{(s)}\|_{\infty, [a, b]}}{(s-k)!} |b-a|^{s-k}. \quad (3.20)$$

236

237 *Proof.* The proof is split into four steps.

238 **Step 1.** Consider the operator

$$239 \quad \mathbb{T} : \mathbb{P}_{s-1} \rightarrow \mathbb{R}^s, \quad \mathbb{T}p := (p(t_1), \dots, p(t_s)).$$

240

241 Clearly,  $\mathbb{T}$  is linear and  $\dim \mathbb{P}_{s-1} = s$ . Hence,  $\mathbb{T}$  is bijective if and only if it is surjective (or  
 242 injective). Given  $z \in \mathbb{R}^s$ , define  $p := \sum_{j=1}^s z_j L_j$ . Note that  $L_j(t_j) = 1$ , while  $L_j(t_k) = 0$   
 243 for all  $j, k = 1, \dots, n$  with  $j \neq k$ . Therefore, we obtain that  $\mathbb{T}p = z$ , i.e.,  $\mathbb{T}$  is surjective  
 244 and hence bijective. In explicit terms, we have thus shown that  $p = \sum_{j=1}^s z_j L_j$  is the  
 245 unique polynomial in  $\mathbb{P}_{s-1}$  such that

$$246 \quad \forall j = 1, \dots, s : \quad p(t_j) = z_j.$$

247

248 This concludes the proof of (3.17).

249 **Step 2.** The error  $e := g - q \in C^s[a, b]$  has at least  $s$  zeros in  $[a, b]$  (at the  $t_j$ ).  
 250 According to the mean value theorem, between two zeros of  $e$ , one has one zero of  $e'$ .  
 251 Hence,  $e' \in C^{s-1}[a, b]$  has at least  $s-1$  zeros. Inductively, we obtain that  $e^{(k)} \in C^{s-k}[a, b]$   
 252 has at least  $s-k$  zeros  $\zeta_1, \dots, \zeta_{s-k} \in [a, b]$ . For  $t = \zeta_\ell$ , the error identity (3.19) is trivial.  
 253 Without loss of generality, we can thus assume that  $t \notin \{\zeta_1, \dots, \zeta_{s-k}\}$ .

254 **Step 3.** We consider the function

$$255 \quad G(x) := e^{(k)}(t) \omega(x) - \omega(t) e^{(k)}(x), \quad \text{where} \quad \omega(x) := \prod_{\ell=1}^{s-k} (x - \zeta_\ell).$$

256

257 Clearly,  $G \in C^{s-k}[a, b]$  has at least  $s-k+1$  zeros in  $[a, b]$ . Inductively, the mean value  
 258 theorem shows that  $G^{(s-k)}$  has at least one zero  $\xi \in [a, b]$ . Hence,

$$259 \quad 0 = G^{(s-k)}(\xi) = e^{(k)}(t) \omega^{(s-k)}(\xi) - \omega(t) e^{(s)}(\xi) = e^{(k)}(t) (s-k)! - \omega(t) g^{(s)}(\xi).$$

260

261 Rearranging this estimate, we prove the error identity (3.19).

262 **Step 4.** The final estimate (3.20) follows from  $|t - \zeta_\ell| \leq |b - a|$  and hence

$$263 \quad |g^{(k)}(t) - q^{(k)}(t)| \stackrel{(3.19)}{\leq} \frac{\|g^{(s)}\|_{\infty, [a, b]}}{(s-k)!} |b-a|^{s-k}.$$

264

265 Taking the supremum over all  $t \in [a, b]$ , we conclude the proof. □

266 **Lemma 3.21 (Interpolatory quadrature).** Let  $g \in C^s[a, b]$  and  $a \leq t_1 < \dots < t_s \leq b$ .  
 267 Then,

$$268 \quad \left| \int_a^b g \, dt - \sum_{j=1}^s g(t_j) \int_a^b L_j \, dt \right| \leq \frac{\|f^{(s)}\|_{\infty, [a, b]}}{s!} |b-a|^{s+1}. \quad (3.21)$$

269

270 *Proof.* Let  $q := \sum_{j=1}^s g(t_j)L_j \in \mathbb{P}_{s-1}$  and recall the Lagrange interpolation (3.17)–(3.18).  
 271 Then,

$$272 \quad \sum_{j=1}^s g(t_j) \int_a^b L_j dt = \int_a^b q dt$$

273 and hence

$$275 \quad \left| \int_a^b g dt - \sum_{j=1}^s g(t_j) \int_a^b L_j dt \right| = \left| \int_a^b (g - q) dt \right| \stackrel{(3.20)}{\leq} \frac{\|g^{(s)}\|_{\infty, [a, b]}}{s!} |b - a|^{s+1}.$$

276 This concludes the proof. □

277 *Proof of Proposition 3.19.* The proof is split into two steps.

278 **Step 1.** First, note that

$$280 \quad y(t + c_\ell h) - y(t) = \int_t^{t+c_\ell h} y'(s) ds = h \int_0^{c_\ell} y'(t + sh) ds$$

$$281 \quad = h \int_0^{c_\ell} f(t + sh, y(t + sh)) ds \stackrel{(3.16b)}{=} h \left( \sum_{j=1}^m A_{\ell j} f(t + c_j h, y(t + c_j h)) + \mathcal{O}(h^{p-1}) \right).$$

282 Recall that  $k_j = f(t + c_j h, y(t) + h \sum_{i=1}^m A_{ji} k_i)$ . Using the last identity, we get that

$$284 \quad r_\ell := \left\| y(t + c_\ell h) - \left[ y(t) + h \sum_{j=1}^m A_{\ell j} k_j \right] \right\|$$

$$285 \quad = h \left\| \sum_{j=1}^m A_{\ell j} \left[ f(t + c_j h, y(t + c_j h)) - f(t + c_j h, y(t) + h \sum_{i=1}^m A_{ji} k_i) \right] \right\| + \mathcal{O}(h^p)$$

$$286 \quad \leq hL \sum_{j=1}^m |A_{\ell j}| \left\| y(t + c_j h) - \left[ y(t) + h \sum_{i=1}^m A_{ji} k_i \right] \right\| + \mathcal{O}(h^p)$$

$$287 \quad \leq hL \max_{i=1, \dots, m} \sum_{j=1}^m |A_{ij}| \|r_j\| + \mathcal{O}(h^p).$$

288 Note that the right-hand side is independent of  $\ell$ . Therefore, we get that

$$290 \quad R := \max_{\ell=1, \dots, m} r_\ell \leq \left( hL \max_{i=1, \dots, m} \sum_{j=1}^m |A_{ij}| \right) R + \mathcal{O}(h^p). \quad (3.22)$$

291 For sufficiently small  $h > 0$ , we thus obtain that  $R = \mathcal{O}(h^p)$ .

292 **Step 2.** Similarly, we see that

$$294 \quad y(t + h) - y(t) = \int_t^{t+h} y'(s) ds = h \int_0^1 y'(t + sh) ds$$

$$295 \quad = h \int_0^1 f(t + sh, y(t + sh)) ds \stackrel{(3.16a)}{=} h \left( \sum_{j=1}^m b_j f(t + c_j h, y(t + c_j h)) + \mathcal{O}(h^p) \right).$$

297 Finally, we can estimate the consistency error by

$$\begin{aligned}
298 \quad & \|y(t+h) - [y(t) + h\Phi(t, y(t), h)]\| = \left\| y(t+h) - \left[ y(t) + h \sum_{j=1}^m b_j k_j \right] \right\| \\
299 \quad & = h \left\| \sum_{j=1}^m b_j \left[ f(t + c_j h, y(t + c_j h)) - f(t + c_j h, y(t) + h \sum_{i=1}^m A_{ji} k_i) \right] \right\| + \mathcal{O}(h^{p+1}) \\
300 \quad & \leq hL \sum_{j=1}^m |b_j| \left\| y(t + c_j h) - \left[ y(t) + h \sum_{i=1}^m A_{ji} k_i \right] \right\| + \mathcal{O}(h^{p+1}) \\
301 \quad & \leq hL \sum_{j=1}^m |b_j| R + \mathcal{O}(h^{p+1}) \stackrel{(3.22)}{=} \mathcal{O}(h^{p+1}).
\end{aligned}$$

303 This concludes the proof. □

304

### 305 3.3. Collocation methods.

306 **Definition 3.22.** Let  $0 \leq c_1 < \dots < c_m \leq 1$ . Then, the following inductive procedure is  
307 called **collocation method**: Given a time-step  $t_\ell$  and the corresponding approximation  
308  $y_\ell \in \mathbb{R}^n$ , let  $q_\ell \in \mathbb{P}_m$  satisfy

$$309 \quad q_\ell(t_\ell) = y_\ell \quad \text{and} \quad q'_\ell(t_\ell + c_j h_\ell) = f(t_\ell + c_j h_\ell, q_\ell(t_\ell + c_j h_\ell)) \quad \text{for all } j = 1, \dots, m. \tag{3.23}$$

311 Then define  $y_{\ell+1} := q_\ell(t_{\ell+1})$ .

312 Note that a collocation method is a very natural strategy to solve (3.2). A collocation  
313 method provides a continuous piecewise polynomial, which satisfies the given ODE  
314 pointwise at finitely many collocation nodes  $t_\ell + c_j h_\ell$ . We stress that well-posedness of a  
315 collocation method is not obvious, since the interpolation conditions for  $q'_\ell$  in (3.23) are  
316 implicit and (thus possibly) nonlinear. However, well-posedness is shown in the following  
317 proposition, if the step-size  $h$  is sufficiently small.

318 **Proposition 3.23.** Let  $0 \leq c_1 < \dots < c_m \leq 1$ . Then, it holds that

$$319 \quad \exists H > 0 \forall t \in [t_0, T] \forall y \in \mathbb{R}^n \forall 0 < h \leq \min\{T - t, H\} \exists! q \in \mathbb{P}_m : \tag{3.24} \\
320 \quad \left( q(t) = y \quad \text{and} \quad \forall j = 1, \dots, m : \quad q'(t + c_j h) = f(t + c_j h, q(t + c_j h)) \right).$$

321 In particular, there exists a unique polynomial  $q \in \mathbb{P}_m$  which satisfies (3.23), if the step-  
322 size  $h_\ell$  is sufficiently small.

323 *Proof.* The proof is split into two steps.

324 **Step 1.** We consider the operator

$$325 \quad \mathbb{T} : \mathbb{P}_m \rightarrow \mathbb{R}^{m+1}, \quad \mathbb{T}q := (q(t), q'(t + c_1 h), \dots, q'(t + c_m h)).$$

327 Clearly,  $\mathbb{T}$  is linear and  $\dim \mathbb{P}_m = m + 1$ . Hence,  $\mathbb{T}$  is bijective if and only if  $\mathbb{T}$  is injective  
328 (or surjective). Let  $q \in \mathbb{P}_m$  with  $\mathbb{T}q = 0$ . Then,  $q' \in \mathbb{P}_{m-1}$  has at least  $m$  zeros (at  
329  $t + c_j h$ ). Hence, the fundamental theorem of algebra shows that  $q' = 0$  and hence  $q$  is



330 constant. Since  $q(t) = 0$ , it follows that  $q = 0$ . Overall, we have shown that  $T$  is injective  
 331 and hence bijective. In explicit terms, it thus holds that

$$332 \quad \forall z_0, z_1, \dots, z_m \in \mathbb{R} \exists! q \in \mathbb{P}_m : \left( q(t) = z_0 \quad \text{and} \quad \forall j = 1, \dots, m : \quad q'(t + c_j h) = z_j \right).$$

334 **Step 2.** By abuse of notation, we abbreviate  $(\mathbb{P}_m)^n$  by simply  $\mathbb{P}_m$  in the following. Let  
 335  $q \in \mathbb{P}_m$ . According to Step 1, there exists a unique  $\Psi q \in \mathbb{P}_m$  such that

$$336 \quad (\Psi q)(t) = y \quad \text{and} \quad \forall j = 1, \dots, m : \quad (\Psi q)'(t + c_j h) = f(t + c_j h, q(t + c_j h))$$

338 This gives rise to an operator  $\Psi : \mathbb{P}_m \rightarrow \mathbb{P}_m$ . In view of (3.24), we have to show that the  
 339 mapping  $\Psi$  has as unique fixpoint, provided that  $h$  is sufficiently small.

340 To this end, note that  $(\Psi q)' \in \mathbb{P}_{m-1}$  and Lemma 3.20 yields that

$$341 \quad (\Psi q)'(s) \stackrel{(3.18)}{=} \sum_{j=1}^m f(t + c_j h, q(t + c_j h)) L_j(s).$$

342 For  $0 \leq \lambda \leq 1$ , it follows that

$$343 \quad (\Psi q)(t + \lambda h) = (\Psi q)(t) + \int_t^{t+\lambda h} (\Psi q)'(s) ds = y + h \int_0^\lambda (\Psi q)'(s) ds$$

$$344 \quad = y + h \sum_{j=1}^m f(t + c_j h, q(t + c_j h)) \int_0^\lambda L_j(s) ds.$$

345 For  $q, \tilde{q} \in \mathbb{P}_m$ , it follows that

$$346 \quad \|(\Psi q)(t + \lambda h) - (\Psi \tilde{q})(t + \lambda h)\|$$

$$347 \quad = h \left\| \sum_{j=1}^m [f(t + c_j h, q(t + c_j h)) - f(t + c_j h, \tilde{q}(t + c_j h))] \int_0^\lambda L_j(s) ds \right\|$$

$$348 \quad \leq h \sum_{j=1}^m \left| \int_0^\lambda L_j(s) ds \right| \|f(t + c_j h, q(t + c_j h)) - f(t + c_j h, \tilde{q}(t + c_j h))\|$$

$$349 \quad \leq hL \sum_{j=1}^m \left( \int_0^\lambda |L_j| ds \right) \|q(t + c_j h) - \tilde{q}(t + c_j h)\|$$

$$350 \quad \leq hL \sum_{j=1}^m \left( \int_0^1 |L_j| ds \right) \|q - \tilde{q}\|_{\infty, [t, t+h]}.$$

351 Note that the right-hand side is independent of  $\lambda$ . Thus, we are led to

$$352 \quad \|\Psi q - \Psi \tilde{q}\|_{\infty, [t, t+h]} \leq hL \sum_{j=1}^m \left( \int_0^1 |L_j| ds \right) \|q - \tilde{q}\|_{\infty, [t, t+h]}.$$

353 This shows that  $\Psi$  is a contraction, if  $h$  is sufficiently small. Therefore, the Banach  
 354 fixpoint theorem concludes the proof.  $\square$

---

355 **Example 3.24 (Collocation methods with  $m = 1$ ).** Let  $m = 1$ . Then, the collocation  
 356 polynomial  $q_\ell$  from (3.23) is linear and hence its derivative is constant  $q'_\ell(t) = \frac{y_{\ell+1} - y_\ell}{h_\ell}$ .

- 361 • If  $c_1 = 0$ , then  $\frac{y_{\ell+1} - y_\ell}{h_\ell} = q'_\ell(t_\ell) = f(t_\ell, q_\ell(t_\ell)) = f(t_\ell, y_\ell)$ . Hence, we obtain the  
 362 explicit Euler method.  
 363 • If  $c_1 = 1$ , then  $\frac{y_{\ell+1} - y_\ell}{h_\ell} = q'_\ell(t_{\ell+1}) = f(t_{\ell+1}, q_\ell(t_{\ell+1})) = f(t_{\ell+1}, y_{\ell+1})$ . Hence, we  
 364 obtain the implicit Euler method.  
 365 • If  $c = 1/2$ , then  $t_\ell + \frac{h_\ell}{2} = \frac{t_\ell + t_{\ell+1}}{2}$  and  $q_\ell(t_\ell + \frac{h_\ell}{2}) = \frac{q_\ell(t_\ell) + q_\ell(t_{\ell+1})}{2} = \frac{y_\ell + y_{\ell+1}}{2}$ . Hence,

366 
$$\frac{y_{\ell+1} - y_\ell}{h_\ell} = q'_\ell\left(\frac{t_\ell + t_{\ell+1}}{2}\right) = f\left(t_\ell + \frac{h_\ell}{2}, q_\ell\left(t_\ell + \frac{h_\ell}{2}\right)\right) = f\left(\frac{t_\ell + t_{\ell+1}}{2}, \frac{y_\ell + y_{\ell+1}}{2}\right),$$

368 and we obtain the implicit midpoint scheme.

369 **Theorem 3.25 (Collocation methods are Runge–Kutta methods).** For any nodes  
 370  $0 \leq c_1 < \dots < c_m \leq 1$ , the corresponding collocation scheme is an (possibly implicit)  
 371  $m$ -stage Runge–Kutta method with consistency order  $p \geq m$ , where

372 
$$A_{ij} := \int_0^{c_i} L_j dt \quad \text{and} \quad b_j := \int_0^1 L_j dt \quad \text{for all } i, j = 1, \dots, m \quad (3.25)$$

374 with the Lagrange polynomials  $L_j(t) := \prod_{\substack{k=1 \\ k \neq j}}^m \frac{t - c_k}{c_j - c_k}$ .

375 *Proof.* The proof is split into two steps.

376 **Step 1.** Let  $q_\ell \in \mathbb{P}_m$  be the collocation polynomial (3.23). Then,  $q'_\ell \in \mathbb{P}_{m-1}$  and hence

377 
$$q'_\ell(s) \stackrel{(3.18)}{=} \sum_{j=1}^m f(t + c_j h, q_\ell(t + c_j h)) L_j(s).$$

379 Therefore, we see that

380 
$$\begin{aligned} q_\ell(t_\ell + c_i h_\ell) &= q_\ell(t_\ell) + \int_{t_\ell}^{t_\ell + c_i h_\ell} q'_\ell(t) dt = y_\ell + h_\ell \int_0^{c_i} q'_\ell(t_\ell + s h_\ell) ds \\ &= y_\ell + h_\ell \sum_{j=1}^m f(t + c_j h, q_\ell(t + c_j h)) \int_0^{c_i} L_j(s) ds \\ &\stackrel{(3.25)}{=} y_\ell + h_\ell \sum_{j=1}^m A_{ij} f(t + c_j h, q_\ell(t + c_j h)). \end{aligned}$$

384 With  $k_i := f(t + c_i h, q_\ell(t + c_i h))$ , the last identity proves that

385 
$$k_i = f(t_\ell + c_i h_\ell, q_\ell(t_\ell + c_i h_\ell)) = f\left(t_\ell + c_i h_\ell, y_\ell + h_\ell \sum_{j=1}^m A_{ij} k_j\right).$$

386

387 The same argument proves that

$$\begin{aligned}
 388 \quad y_{\ell+1} &= q_\ell(t_{\ell+1}) = q_\ell(t_\ell) + \int_{t_\ell}^{t_\ell+h_\ell} q'_\ell(t) \, dt = y_\ell + h_\ell \int_0^1 q'_\ell(t_\ell + sh_\ell) \, ds \\
 389 \quad &= y_\ell + h_\ell \sum_{j=1}^m f(t + c_j h, q_\ell(t + c_j h)) \int_0^1 L_j(s) \, ds \stackrel{(3.25)}{=} y_\ell + h_\ell \sum_{j=1}^m b_j k_j. \\
 390
 \end{aligned}$$

391 This shows that the collocation method is indeed a Runge–Kutta method.

392 **Step 2.** For each polynomial  $q \in \mathbb{P}_{m-1}$ , it holds that

$$393 \quad q(s) \stackrel{(3.18)}{=} \sum_{j=1}^m q(c_j) L_j(s).$$

394 Hence, it follows that

$$396 \quad \sum_{j=1}^m b_j q(c_j) \stackrel{(3.25)}{=} \sum_{j=1}^m q(c_j) \int_0^1 L_j \, dt = \int_0^1 q \, dt$$

397 as well as

$$399 \quad \sum_{j=1}^m A_{ij} q(c_j) \stackrel{(3.25)}{=} \sum_{j=1}^m q(c_j) \int_0^{c_i} L_j \, dt = \int_0^{c_i} q \, dt.$$

400 Therefore, Proposition 3.19 proves that the consistency order is at least  $m$ .  $\square$

401 One major advantage of collocation methods is that they provide an approximation of  
 402 the sought solution  $y$  in the full time interval  $[t_0, T]$  (and not only at the time-steps  $t_j$ ).  
 403 Moreover, we even get natural approximations of the derivatives of  $y$ .  
 404

---

405 **Theorem 3.26 (Global convergence of collocation methods).** *Let  $0 \leq c_1 < \dots <$   
 406  $c_m \leq 1$  be given nodes of a collocation method with consistency order  $p$ . Define the spline  
 407  $q : [t_0, T] \rightarrow \mathbb{R}^n$  by  $q|_{[t_\ell, t_{\ell+1}]} := q_\ell$ , where  $q_\ell \in \mathbb{P}_m$  are the collocation polynomials (3.23) for  
 408 all  $\ell = 0, \dots, N-1$ . Then,  $q \in C([t_0, T]; \mathbb{R}^n)$ . Let  $r := \min\{p, m+1\}$  and note that  $m \leq$   
 409  $r \leq 2m$ . Suppose that  $f \in C^r([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$ . Then, the solution  $y \in C^{r+1}([t_0, T]; \mathbb{R}^n)$   
 410 of (3.2) satisfies that*

$$411 \quad \|y^{(k)} - q^{(k)}\|_{\infty, [t_0, T]} = \mathcal{O}(h_\Delta^{r-k}) \quad \text{for all } 0 \leq k \leq r, \quad (3.26)$$

412 where  $q^{(k)}$  is understood elementwise on  $[t_\ell, t_{\ell+1}]$  for all  $\ell = 0, \dots, N-1$ .

---

414 **Lemma 3.27 (Polynomial approximation).** *Let  $z \in C^{r+1}[t, t+h]$ ,  $q \in \mathbb{P}_r$ , and  
 415  $k \in \{1, \dots, r\}$ . Then, it holds that*

$$416 \quad \|z^{(k)} - q^{(k)}\|_{\infty, [t, t+h]} \leq C (h^{r+1-k} \|z^{(r+1)}\|_{\infty, [t, t+h]} + h^{-k} \|z - q\|_{\infty, [t, t+h]}), \quad (3.27)$$

417 where  $C > 0$  depends only on  $r$ , but is independent of  $z$ ,  $q$ , and  $h$ .

---

419 *Proof.* The proof follows from a so-called scaling argument and is split in three steps.

420 **Step 1.** Define  $\widehat{z}(s) := z(t + sh)$  and  $\widehat{q}(s) := q(t + sh)$ . Note that, e.g.,  $\widehat{z}^{(k)}(s) =$   
 421  $h^k z^{(k)}(t + sh)$ . In particular, one sees that

$$422 \quad \|\widehat{z}^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]} = h^k \|z^{(k)} - q^{(k)}\|_{\infty, [t, t+h]},$$

$$423 \quad \|\widehat{z}^{(r+1)}\|_{\infty, [0,1]} = h^{r+1} \|z^{(r+1)}\|_{\infty, [t, t+h]}.$$

425 **Step 2.** Let  $0 \leq s_0 < \dots < s_r \leq 1$  be interpolation nodes. Define

$$426 \quad (I\widehat{z})(s) := \sum_{j=0}^r \widehat{z}(s_j) L_j(s), \quad \text{where} \quad L_j(s) := \prod_{\substack{k=0 \\ k \neq j}}^r \frac{s - s_k}{s_j - s_k}.$$

427  
 428 The error estimate (3.20) for the Lagrange interpolation yields that

$$429 \quad \|\widehat{z} - I\widehat{z}\|_{\infty, [0,1]} \leq \frac{1}{(r+1)!} \|\widehat{z}^{(r+1)}\|_{\infty, [0,1]},$$

$$430 \quad \|\widehat{z}^{(k)} - (I\widehat{z})^{(k)}\|_{\infty, [0,1]} \leq \frac{1}{(r+1-k)!} \|\widehat{z}^{(r+1)}\|_{\infty, [0,1]}.$$

432 Moreover, norm equivalence on the finite-dimensional space  $\mathbb{P}_r$  provides a constant  $C_k > 0$   
 433 such that

$$434 \quad \|(I\widehat{z})^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]} \leq \|I\widehat{z} - \widehat{q}\|_{\infty, [0,1]} + \|(I\widehat{z})^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]}$$

$$435 \quad \leq C_k \max_{j=0, \dots, r} |(I\widehat{z})(s_j) - \widehat{q}(s_j)|,$$

437 where we note that  $\|p\| := \max_{j=0, \dots, r} |p(s_j)|$  is a norm on  $\mathbb{P}_r$  due to Lemma 3.20. Since  
 438  $(I\widehat{z})(s_j) = z(s_j)$  for all  $j = 0, \dots, r$ , we obtain that

$$439 \quad \|(I\widehat{z})^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]} \leq C_k \|\widehat{z} - \widehat{q}\|_{\infty, [0,1]}.$$

441 With the triangle inequality and  $C := \max_{k=1, \dots, r} \max\{C_k, 1/(r+1-k)!\}$ , we conclude that

$$442 \quad \|\widehat{z}^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]} \leq \|z^{(k)} - (I\widehat{z})^{(k)}\|_{\infty, [0,1]} + \|(I\widehat{z})^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]}$$

$$443 \quad \leq \frac{1}{(r+1-k)!} \|\widehat{z}^{(r+1)}\|_{\infty, [0,1]} + C_k \|\widehat{z} - \widehat{q}\|_{\infty, [0,1]}$$

$$444 \quad \leq C (\|\widehat{z}^{(r+1)}\|_{\infty, [0,1]} + \|\widehat{z} - \widehat{q}\|_{\infty, [0,1]}).$$

446 We stress that this already proves the claim (3.27) on the reference interval  $[0, 1]$ .

447 **Step 3.** It holds that

$$448 \quad \|z^{(k)} - q^{(k)}\|_{\infty, [t, t+h]} = h^{-k} \|\widehat{z}^{(k)} - \widehat{q}^{(k)}\|_{\infty, [0,1]}$$

$$449 \quad \leq C h^{-k} (\|\widehat{z}^{(r+1)}\|_{\infty, [0,1]} + \|\widehat{z} - \widehat{q}\|_{\infty, [0,1]})$$

$$450 \quad \leq C h^{-k} (h^{r+1} \|z^{(r+1)}\|_{\infty, [t, t+h]} + \|z - q\|_{\infty, [t, t+h]}).$$

452 This concludes the proof. We stress that the norm equivalence argument has to be used  
 453 on  $[0, 1]$  instead of  $[t, t+h]$  to ensure that  $C$  is independent of  $h$ .  $\square$

454 *Proof of Theorem 3.26.* Note that

$$455 \quad \|g\|_{\infty, [t_0, T]} = \max_{\ell=0, \dots, N-1} \|g\|_{\infty, [t_\ell, t_{\ell+1}]}.$$

456

457 Therefore, it suffices to prove (3.26) for one time interval  $[t_\ell, t_{\ell+1}]$ . Moreover, recall that  
 458  $q|_{[t_\ell, t_{\ell+1}]} = q_\ell \in \mathbb{P}_m$  and  $r \geq m$ . Therefore, Lemma 3.27 proves that

$$459 \quad \|y^{(k)} - q^{(k)}\|_{\infty, [t_\ell, t_{\ell+1}]} \leq C (h_\ell^{r+1-k} \|y^{(r+1)}\|_{\infty, [t_\ell, t_{\ell+1}]} + h_\ell^{-k} \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]}).$$

461 Hence, it suffices to consider the case  $k = 0$ .

462 **Step 1.** Let  $t \in [t_0, T]$  and  $h > 0$ . Let  $Iy \in \mathbb{P}_m$  be the (unique) polynomial such that

$$463 \quad Iy(t) = y(t) \quad \text{and} \quad \forall j = 1, \dots, m: \quad (Iy)'(t + c_j h) = y'(t + c_j h).$$

465 Let  $0 \leq \lambda \leq 1$ . With  $L_j(s) = \prod_{\substack{k=1 \\ k \neq j}}^m \frac{s - c_k}{c_j - c_k}$ , Lemma 3.21 yields that

$$466 \quad y(t + \lambda h) - Iy(t + \lambda h) = \int_t^{t+\lambda h} (y - Iy)'(s) ds = h \int_0^\lambda (y - Iy)'(t + sh) ds$$

$$467 \quad \stackrel{(3.21)}{=} h \left[ \sum_{j=1}^m (y - Iy)'(t + c_j h) \int_0^\lambda L_j(s) ds + \mathcal{O}(h^m) \right] = \mathcal{O}(h^{m+1}).$$

469 Since the right-hand side is independent of  $\lambda$ , this proves that

$$470 \quad \|y - Iy\|_{\infty, [t, t+h]} = \mathcal{O}(h^{m+1}).$$

472 **Step 2.** Employ the notation  $Iy \in \mathbb{P}_m$  from Step 1 for the interval  $[t_\ell, t_\ell + h_\ell] = [t_\ell, t_{\ell+1}]$ .  
 473 The triangle inequality and Step 1 prove that

$$474 \quad \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} \leq \|y - Iy\|_{\infty, [t_\ell, t_{\ell+1}]} + \|Iy - q\|_{\infty, [t_\ell, t_{\ell+1}]} \quad (3.28)$$

$$475 \quad = \mathcal{O}(h_\ell^{m+1}) + \|Iy - q\|_{\infty, [t_\ell, t_{\ell+1}]}.$$

476 Let  $0 \leq \lambda \leq 1$ . We argue as before:

$$477 \quad Iy(t_\ell + \lambda h_\ell) - q(t_\ell + \lambda h_\ell) = y(t_\ell) - q(t_\ell) + \int_{t_\ell}^{t_\ell + \lambda h_\ell} (Iy - q)'(s) ds$$

$$478 \quad = y(t_\ell) - q(t_\ell) + h_\ell \int_0^\lambda (Iy - q)'(t_\ell + sh_\ell) ds$$

$$479 \quad = y(t_\ell) - q(t_\ell) + h_\ell \sum_{j=1}^m (Iy - q)'(t_\ell + c_j h_\ell) \int_0^\lambda L_j(s) ds.$$

481 Recall that collocation methods are implicit Runge–Kutta schemes and hence stable.  
 482 From consistency, we thus get convergence

$$483 \quad y(t_\ell) - q(t_\ell) = y(t_\ell) - y_\ell = \mathcal{O}(h_\ell^p).$$

485 Note that

$$486 \quad (Iy - q)'(t_\ell + c_j h_\ell) = f(t_\ell + c_j h_\ell, y(t_\ell + c_j h_\ell)) - f(t_\ell + c_j h_\ell, q(t_\ell + c_j h_\ell)).$$

488 Combining the last three identities with the Lipschitz continuity of  $f$ , we obtain that

$$\begin{aligned}
 489 \quad \|(Iy - q)(t_\ell + \lambda h_\ell)\| &\leq h_\ell L \sum_{j=1}^m \left| \int_0^\lambda L_j(s) \, ds \right| \|(y - q)(t_\ell + c_j h_\ell)\| + \mathcal{O}(h_\Delta^p) \\
 490 &\leq \left( h_\ell L \sum_{j=1}^m \int_0^1 |L_j(s)| \, ds \right) \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} + \mathcal{O}(h_\Delta^p). \\
 491
 \end{aligned}$$

492 Note that the right-hand side is independent of  $0 \leq \lambda \leq 1$ . Hence, we are led to

$$493 \quad \|Iy - q\|_{\infty, [t_\ell, t_{\ell+1}]} \leq \left( h_\ell L \sum_{j=1}^m \int_0^1 |L_j(s)| \, ds \right) \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} + \mathcal{O}(h_\Delta^p).$$

494 In combination with (3.28), we see that

$$496 \quad \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} \leq \left( h_\ell L \sum_{j=1}^m \int_0^1 |L_j(s)| \, ds \right) \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} + \mathcal{O}(h_\Delta^p) + \mathcal{O}(h_\ell^{m+1}).$$

497 For sufficiently small  $h_\ell \leq h_\Delta$ , we thus infer that

$$499 \quad \|y - q\|_{\infty, [t_\ell, t_{\ell+1}]} = \mathcal{O}(h_\Delta^p) + \mathcal{O}(h_\Delta^{m+1}) = \mathcal{O}(h_\Delta^r).$$

500

□

501  
 502 The next goal is the theorem that a collocation method has consistency order  $p$  if and  
 503 only if the induced quadrature on  $[0, 1]$  (see Proposition 3.17) is exact for all  $q \in \mathbb{P}_{p-1}$ . In  
 504 particular, the  $m$ -stage Gaussian collocation method is the only collocation scheme with  
 505 consistency order  $p = 2m$ .

---

506 **Remark 3.28.** *Let  $f \in C^p([t_0, T] \times \mathbb{R}^n; \mathbb{R}^n)$  and suppose that  $f$  is Lipschitz in  $y$ . Let*  
 507  *$(t_0, y_0) \in \mathbb{R} \times \mathbb{R}^n$ . Then,*

$$508 \quad Y(t, t_0, y_0) := y_0 + \int_{t_0}^t f(s, Y(s, t_0, y_0)) \, ds \tag{3.29}$$

509  
 510 *is the unique solution of*

$$511 \quad y'(t) = f(t, y(t)) \text{ in } [t_0, T] \quad \text{subject to} \quad y(t_0) = y_0.$$

512  
 513 *Clearly,  $y(t) = Y(t, t_0, y_0)$  satisfies  $y \in C^{p+1}([t_0, T], \mathbb{R}^n)$ . Moreover, one can show that*  
 514  *$Y(t, t_0, y_0)$  is also  $C^p$  with respect to  $t_0$  and  $y_0$ , and all derivatives depend only on the*  
 515 *derivatives of  $f$ . We refer to [Wal00, Part III, Section §13, Subsection XI, Corollar].*

---

516 **Lemma 3.29 (Gröbner & Alekseev).** *Let  $f, \varepsilon \in C^1([t, t+h] \times \mathbb{R}^n; \mathbb{R}^n)$  be Lipschitz*  
 517 *in  $y$ . Let  $y, z \in C^2([t, t+h], \mathbb{R}^n)$  such that*

$$\begin{aligned}
 518 \quad y(t) &= y_0, \quad y'(t) = f(t, y(t)), \\
 519 \quad z(t) &= y_0, \quad z'(t) = f(t, z(t)) + \varepsilon(t, z(t)).
 \end{aligned} \tag{3.30}$$

520 *Then, it follows that*

$$521 \quad z(t+h) - y(t+h) = \int_t^{t+h} \partial_3 Y(t+h, s, z(s)) \varepsilon(s, z(s)) \, ds \tag{3.31}$$

522

523 where  $Y(t, s, z(s))$  is given by (3.29).

524 *Proof.* Let  $N \in \mathbb{N}$ . For  $\ell = 0, \dots, N$ , let  $t_\ell := t + \ell H$  with  $H := h/N$ . Note that  $t_0 = t$   
525 and  $t_N = t + h$ . The Taylor expansion shows that

$$\begin{aligned}
526 \quad d_i &:= z(t_{i+1}) - Y(t_{i+1}, t_i, z(t_i)) \\
527 \quad &= [z(t_i) + H \{f(t_i, z(t_i)) + \varepsilon(t_i, z(t_i))\} + \mathcal{O}(H^2)] - [z(t_i) + H f(t_i, z(t_i)) + \mathcal{O}(H^2)] \\
528 \quad &= H \varepsilon(t_i, z(t_i)) + \mathcal{O}(H^2) \\
529 \quad &= H \varepsilon(t_{i+1}, z(t_{i+1})) + \mathcal{O}(H^2). \\
530
\end{aligned}$$

531 By definition of  $d_i$ , the uniqueness of ODE solutions proves that

$$532 \quad Y(t + h, t_i, z(t_i)) = Y(t + h, t_{i+1}, z(t_{i+1}) - d_i). \\ 533$$

534 Therefore, the Taylor expansion shows that

$$\begin{aligned}
535 \quad D_i &:= Y(t + h, t_{i+1}, z(t_{i+1})) - Y(t + h, t_i, z(t_i)) \\
536 \quad &= Y(t + h, t_{i+1}, z(t_{i+1})) - Y(t + h, t_{i+1}, z(t_{i+1}) - d_i) \\
537 \quad &= \partial_3 Y(t + h, t_{i+1}, z(t_{i+1})) d_i + \mathcal{O}(|d_i|^2) \\
538 \quad &= H \underbrace{\partial_3 Y(t + h, t_{i+1}, z(t_{i+1})) \varepsilon(t_{i+1}, z(t_{i+1}))}_{=: g(t_{i+1})} + \mathcal{O}(H^2). \\
539
\end{aligned}$$

540 Note that  $\mathcal{O}(H^2) = \mathcal{O}(N^{-2})$ . Together with Riemann sum theory, it follows that

$$541 \quad \sum_{i=0}^{N-1} D_i = \mathcal{O}(N^{-1}) + \sum_{i=0}^{N-1} \left( (t_{i+1} - t_i) g(t_{i+1}) \right) \xrightarrow{N \rightarrow \infty} \int_t^{t+h} g(s) ds, \\ 542$$

543 since  $g$  is (at least) continuous. Moreover, the telescopic series proves that

$$\begin{aligned}
544 \quad \sum_{i=0}^{N-1} D_i &= Y(t + h, t_N, z(t_N)) - Y(t + h, t_0, z(t_0)) \\
545 \quad &= Y(t + h, t + h, z(t + h)) - Y(t + h, t, y(t)) = z(t + h) - y(t + h). \\
546
\end{aligned}$$

547 Combining the latter two identities, we prove (3.31). □

---

548 **Theorem 3.30 (Quadrature vs. collocation).** Let  $0 \leq c_1 < \dots < c_m \leq 1$  and  
549  $b_j := \int_0^1 L_j dt$ , where  $L_j \in \mathbb{P}_{m-1}$  are the Lagrange polynomials. Let  $p \in \mathbb{N}$ . Then, the  
550 following statements (i)–(ii) are equivalent:

- 551 (i)  $\sum_{j=1}^m b_j q(c_j) = \int_0^1 q dt$  for all  $q \in \mathbb{P}_{p-1}$ , i.e., the quadrature has exactness  $p - 1$ .
- 552 (ii) The corresponding collocation method has consistency order  $p$ .

553 In any case, it holds that  $m \leq p \leq 2m$ .

---

554 *Proof.* The implication (ii)  $\implies$  (i) as well as the bound  $m \leq p \leq 2m$  are already known;  
555 see Proposition 3.17. It thus only remains to prove that (i)  $\implies$  (ii).

556 **Step 1.** Let  $g \in C^p[t, t+h]$ . Choose  $\tilde{c}_{m+1}, \dots, \tilde{c}_p \in [0, 1]$  such that  $c_1, \dots, c_m, \tilde{c}_{m+1}, \dots, \tilde{c}_p$   
 557 are pairwise distinct. According to Lemma 3.20, there exists a unique polynomial  $q \in \mathbb{P}_{p-1}$   
 558 such that

$$559 \quad q(t + c_j h) = g(t + c_j h) \quad \text{for all } j = 1, \dots, m,$$

$$560 \quad q(t + \tilde{c}_j h) = g(t + \tilde{c}_j h) \quad \text{for all } j = m + 1, \dots, p.$$

562 Then, it holds that

$$563 \quad \|g - q\|_{\infty, [t, t+h]} \stackrel{(3.20)}{\leq} \frac{\|g^{(p)}\|_{\infty, [t, t+h]}}{p!} h^p. \quad (3.32)$$

565 By assumption (i), it holds that

$$566 \quad \int_t^{t+h} q(s) ds = h \int_0^1 q(t + sh) ds \stackrel{(i)}{=} h \sum_{j=1}^m b_j q(t + c_j h) = h \sum_{j=1}^m b_j g(t + c_j h).$$

568 Therefore, it follows that

$$569 \quad \left| \int_t^{t+h} g ds - h \sum_{j=1}^m b_j g(t + c_j h) \right| = \left| \int_t^{t+h} (g - q) ds \right| \leq h \|g - q\|_{\infty, [t, t+h]} \quad (3.33)$$

$$\stackrel{(3.32)}{\leq} \frac{\|g^{(p)}\|_{\infty, [t, t+h]}}{p!} h^{p+1}.$$

571 **Step 2.** Let  $q \in \mathbb{P}_m$  be the collocation polynomial on  $[t, t+h]$ , i.e.,

$$572 \quad q(t) = y(t) \quad \text{and} \quad q(t + c_j h) = f(t + c_j h, q(t + c_j h)) \quad \text{for all } j = 1, \dots, m.$$

574 We apply Lemma 3.29 for the collocation polynomial and

$$575 \quad q'(t) = f(t, q(t)) + [q'(t) - f(t, q(t))], \quad \text{i.e.,} \quad \varepsilon(t, q(t)) = q'(t) - f(t, q(t)).$$

577 This leads to

$$578 \quad q(t+h) - y(t+h) = \int_t^{t+h} \underbrace{\partial_3 Y(t+h, s, q(s)) [q'(s) - f(s, q(s))]}_{=:g(s)} ds,$$

580 where we also define  $g$  for which we aim to apply Step 1. Note that  $g(t + c_j h) = 0$   
 581 for all  $j = 1, \dots, m$  by the collocation conditions. From Step 1, we thus infer that the  
 582 consistency error satisfies

$$583 \quad |y(t+h) - q(t+h)| = \left| \int_t^{t+h} g ds - \underbrace{h \sum_{j=1}^m b_j g(t + c_j h)}_{=0} \right| \stackrel{(3.33)}{\leq} \frac{\|g^{(p)}\|_{\infty, [t, t+h]}}{p!} h^{p+1}.$$

585 It only remains to argue that  $\|g^{(p)}\|_{\infty, [t, t+h]}$  is uniformly bounded in terms of  $y$  and  $f$ .  
 586 Note that

$$587 \quad g(s) = \partial_3 Y(t+h, s, q(s)) [q'(s) - f(s, q(s))],$$

589 where  $Y(t+h, s, q(s))$  solves

$$590 \quad \partial_1 Y(t, s, q(s)) = f(t, Y(t, s, q(s))) \quad \text{in } [t, t+h] \quad \text{subject to} \quad Y(s, s, q(s)) = q(s).$$



592 According to the chain rule (and Remark 3.28),  $g^{(p)}$  depends only on partial derivatives  
 593 of  $f$  and derivatives of  $q$ . For  $0 \leq k \leq m$ , it holds that

$$595 \quad \|q^{(k)}\|_{\infty,[t,t+h]} \leq \|y^{(k)} - q^{(k)}\|_{\infty,[t,t+h]} + \|y^{(k)}\|_{\infty,[t,t+h]} \stackrel{(3.26)}{=} \mathcal{O}(h^{m-k}) + \|y^{(k)}\|_{\infty,[t,t+h]} = \mathcal{O}(1).$$

596 For  $k > m$ , it holds that  $q^{(k)} = 0$ . Overall, we thus obtain that  $\|g^{(p)}\|_{\infty,[t,t+h]} = \mathcal{O}(1)$ .  
 597 This concludes the proof.  $\square$

598

1

## 4. STIFF ODES

2

3 **4.1. Introduction.** In 1952, it was observed by Curtiss and Hirschfelder that explicit  
 4 methods fail for the numerical integration of certain ODEs which model chemical reac-  
 5 tions. They said that the ODE is stiff, if certain components of the solution arrive in  
 6 a very short time in their equilibrium (i.e., the fast reacting components), while **other**  
 7 slowly changing components are more or less fixed (i.e., stiff).

8 Mathematically, this is, e.g., the case for

$$9 \quad y'(t) = My(t) + f(t) \tag{4.1}$$

11 if the eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  of  $M \in \mathbb{R}^{n \times n}$ , sorted by  $\operatorname{Re} \lambda_1 \geq \dots \geq \operatorname{Re} \lambda_n$ , satisfy  
 12 that

$$13 \quad |\operatorname{Re} \lambda_1| \sim 1, \quad \text{but} \quad \operatorname{Re} \lambda_n \ll 0. \tag{4.2}$$

---

15 **Example 4.1.** Consider the symmetric matrix  $M = \begin{pmatrix} a & b \\ b & a \end{pmatrix} \in \mathbb{R}^{2 \times 2}$  with  $a, b \in \mathbb{R}$ . The  
 16 characteristic polynomial is  $p(\lambda) = \det(A - \lambda I) = (a - \lambda)^2 - b^2 = \lambda^2 - 2a\lambda + (a^2 - b^2)$ .  
 17 The eigenvalues of  $M$  are the zeros of  $p$  and hence  $M$  has the eigenvalues  $\lambda = a \pm b$ .

**Joseph-Louis de Lagrange (1736–1813)** was an Italian mathematician (born in Turino as *Giuseppe Lodovico Lagrangia*). In 1755, he got the chair for mathematics at the military academy in Turino. In 1766, he became the successor of Euler at the Prussian Academy of Sciences in Berlin. In 1786, he moved to Paris to the French Academy of Sciences. From 1797, he was professor at the Ecole Polytechnique in Paris, where he was promoting the young Augustin-Louis Cauchy.

**Charles Francis Curtiss (1921–2007)** was a American chemist. He lived his academic live at the University of Wisconsin, where he did his bachelor in 1942 and his PhD (supervized by Hirschfelder) in 1948. He became assistant professor in 1949, associate professor in 1954, and full professor in 1960. He retired in 1989.

**Joseph Oakland Hirschfelder (1911–1990)** was an American chemist in physicist. He studied natural sciences at University of Minnesota (1927–1929) and Yale (1929–1931) and finished his PhD in chemistry and physics at Princeton in 1936. From 1936–1937, he was postdoc with John von Neumann at the Princeton Institute for Advanced Studies. In 1937, he went to the University of Wisconsin, where he became assistant professor in chemistry in 1941. 1944/45, he was group leader at the Manhattan project at Los Alamos and later leading researcher in the American nuclear program. In 1946, he became full professor at the University of Wisconsin. He retired in 1981.

Charles Francis Curtiss, Joseph Oakland Hirschfelder: *Integration of stiff equations*, Proceedings of the National Academy of Sciences of the USA; 38 (1952), 235–243. [This work introduces the BDF multistep methods for the solution of stiff ODEs]

18 Even for simple  $a = -50$ ,  $b = -51$ , we get eigenvalues  $\lambda_1 = 1$  and  $\lambda_2 = -101$ , so that  
 19 the corresponding ODE would be stiff.

---

20 **Proposition 4.2.** Given  $y_0 \in \mathbb{R}^n$  and  $M \in \mathbb{R}^{n \times n}$ , consider the initial value problem

$$21 \quad y'(t) = My(t) \text{ in } [t_0, T] \quad \text{subject to} \quad y(t_0) = y_0. \quad (4.3)$$

23 Then, the unique solution  $y \in C^1([t_0, T]; \mathbb{R}^n)$  reads

$$24 \quad y(t) = e^{M(t-t_0)} y_0 \quad (4.4)$$

26 with the matrix exponential function

$$27 \quad e^X := \sum_{k=0}^{\infty} \frac{1}{k!} X^k \quad \text{for } X \in \mathbb{R}^{n \times n}. \quad (4.5)$$

29 Suppose that  $M$  is diagonalizable, i.e., there exist  $\{v_1, \dots, v_n\} \subset \mathbb{C}^n$  linearly independent  
 30 eigenvectors for eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  such that  $Mv_j = \lambda_j v_j$  for all  $j = 1, \dots, n$ .  
 31 Then,

$$32 \quad y_0 = \sum_{j=1}^n \alpha_j v_j \quad \implies \quad y(t) = \sum_{j=1}^n (\alpha_j e^{\lambda_j(t-t_0)}) v_j, \quad (4.6)$$

34 i.e.,  $y = \sum_{j=1}^n y_j v_j$ , where the scalar functions  $y_j \in C^1[t_0, T]$  solve the ODEs

$$35 \quad y'_j(t) = \lambda_j y_j(t) \text{ in } [t_0, T] \quad \text{subject to} \quad y_j(t_0) = \alpha_j. \quad (4.7)$$

37 *Proof.* Existence and uniqueness of the solution of (4.3) (resp. (4.7)) follows from the  
 38 Picard–Lindelöf theorem. Consider  $y$  given by (4.4). Then,  $y(t_0) = y_0$  and

$$39 \quad y'(t) = d_t \left( \sum_{k=0}^{\infty} \frac{1}{k!} M^k (t-t_0)^k y_0 \right)$$

$$40 \quad = \sum_{k=1}^{\infty} \frac{1}{k!} k M^k (t-t_0)^{k-1} y_0$$

$$41 \quad = M \left( \sum_{k=1}^{\infty} \frac{1}{(k-1)!} M^{k-1} (t-t_0)^{k-1} y_0 \right)$$

$$42 \quad = M \left( \sum_{k=0}^{\infty} \frac{1}{k!} M^k (t-t_0)^k y_0 \right) = My(t).$$

44 If  $M$  is diagonalizable and  $y_0 = \sum_{j=1}^n \alpha_j v_j$ , then

$$\begin{aligned}
 45 \quad y(t) &= \sum_{k=0}^{\infty} \frac{1}{k!} M^k (t - t_0)^k y_0 \\
 46 \quad &= \sum_{j=1}^n \left( \sum_{k=0}^{\infty} \frac{1}{k!} M^k (t - t_0)^k \alpha_j v_j \right) \\
 47 \quad &= \sum_{j=1}^n \left( \alpha_j \sum_{k=0}^{\infty} \frac{1}{k!} \lambda_j^k (t - t_0)^k v_j \right) \\
 48 \quad &= \sum_{j=1}^n \left( \alpha_j e^{\lambda_j (t - t_0)} \right) v_j.
 \end{aligned}$$

50 Since  $y_j = e^{\lambda_j (t - t_0)} \alpha_j$ , this concludes the proof.  $\square$

51 More generally, one can even prove the following result. We note that the solution  
 52  $z \in C^1([t_0, T]; \mathbb{R}^n)$  of (4.9) can be obtained by solving  $n$  scalar ODEs (if  $\tilde{g}(t) := V^{-1}g(t)$   
 53 is explicitly known).

---

54 **Proposition 4.3.** *Given  $g \in C([t_0, T]; \mathbb{R}^n)$ ,  $M \in \mathbb{R}^{n \times n}$ , and  $y_0 \in \mathbb{R}^n$ , consider the*  
 55 *inhomogeneous initial value problem*

$$56 \quad y'(t) = My(t) + g(t) \text{ in } [t_0, T], \quad y(t_0) = y_0. \quad (4.8)$$

58 *Suppose that the matrix  $M \in \mathbb{R}^n$  is diagonalizable, i.e., there exist  $\{v_1, \dots, v_n\} \subset \mathbb{C}^n$*   
 59 *linearly independent eigenvectors for eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  such that  $Mv_j = \lambda_j v_j$*   
 60 *for all  $j = 1, \dots, n$ . Define  $\Lambda := \text{diag}(\lambda_1, \dots, \lambda_n) \in \mathbb{R}^{n \times n}$  and  $V := (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$*   
 61 *and consider the problem*

$$62 \quad z'(t) = \Lambda z(t) + V^{-1}g(t) \text{ in } [t_0, T], \quad z(t_0) = V^{-1}y_0. \quad (4.9)$$

64 *Then, (4.8)–(4.9) have unique solutions and it holds that  $y(t) = Vz(t)$ .*

---

65 *Proof.* Existence and uniqueness of the solutions of (4.8)–(4.9) follows from the Picard–  
 66 Lindelöf theorem. Clearly,  $V$  is invertible and independent of  $t$ . Therefore, (4.9) implies  
 67 that

$$68 \quad (Vz)'(t) = V\Lambda z(t) + g(t).$$

70 Note that  $V\Lambda = MV$ , since the columns of  $V$  are the eigenvectors of  $M$ . Hence,  $\tilde{y} := Vz$   
 71 solves that

$$72 \quad \tilde{y}'(t) = M\tilde{y}(t) + g(t) \quad \text{together with} \quad \tilde{y}(t_0) = Vz(t_0) = y_0.$$

74 From the uniqueness of solutions, we thus conclude that  $\tilde{y} = y$ .  $\square$

---

75 **Exercise 4.4.** *Consider the setting of Proposition 4.3. Let  $\frac{c}{b^\top} \Big| \frac{A}{b^\top}$  be an explicit  $m$ -stage*  
 76 *Runge–Kutta method. Let  $y_\ell, z_\ell \in \mathbb{R}^n$  be the resulting Runge–Kutta iterates of (4.8)–(4.9).*  
 77 *Prove that  $y_\ell = Vz_\ell$  for all  $\ell = 0, \dots, N$ , i.e., explicit Runge–Kutta methods commute*  
 78 *with diagonalization of the ODE system.*

79 If we consider the simple homogeneous model problem from Proposition 4.2, it follows  
 80 that, for a stiff ODE, we should employ Runge–Kutta methods, which appropriately  
 81 integrate the scalar ODEs

$$82 \quad y'_j(t) = \lambda_j y(t) \quad \text{in } [t_0, T]$$

83 simultaneously for all eigenvalues  $\lambda_j \in \mathbb{C}$  of  $M \in \mathbb{R}^{n \times n}$ .

---

85 **Remark 4.5.** For given  $\lambda \in \mathbb{C} \setminus \{0\}$ , let us consider the scalar ODE

$$86 \quad y'(t) = \lambda y(t) \text{ in } \mathbb{R}_{\geq 0}, \quad y(0) = y_0. \quad (4.10)$$

88 The unique solution is  $y(t) = e^{\lambda t} y_0$ . If we apply the explicit Euler method with constant  
 89 time-step size  $h > 0$ , we get that

$$90 \quad y_{\ell+1} = y_\ell + h f(t_\ell, y_\ell) = (1 + h\lambda) y_\ell \quad \text{for all } \ell \geq 0$$

91 and hence

$$93 \quad y_\ell = (1 + h\lambda)^\ell y_0 \quad \text{for all } \ell \geq 0. \quad (4.11)$$

95 If we apply the implicit Euler method, we get that

$$96 \quad y_{\ell+1} = y_\ell + h f(t_{\ell+1}, y_{\ell+1}) = y_\ell + h\lambda y_{\ell+1} \quad \text{for all } \ell \geq 0$$

98 and hence (if  $h\lambda \neq 1$ )

$$99 \quad y_\ell = (1 - h\lambda)^{-\ell} y_0 \quad \text{for all } \ell \geq 0. \quad (4.12)$$

101 We make the following observation: If  $\operatorname{Re} \lambda \gg 0$ , then both methods will need small  $h$ ,  
 102 since the exact solution grows exponentially. If  $\operatorname{Re} \lambda < 0$ , then the explicit Euler method  
 103 needs small  $h$  so that (4.11) leads to  $y_\ell \rightarrow 0$  as  $\ell \rightarrow \infty$ . If  $|1 + h\lambda| > 1$ , then (4.11)  
 104 leads to blow-up and oscillations. On the other hand, the implicit Euler method (4.12)  
 105 guarantees always  $y_\ell \rightarrow 0$  as  $\ell \rightarrow \infty$ , if  $\operatorname{Re} \lambda < 0$ .

---

106 **Example 4.6.** We consider the ODE

$$107 \quad y' = My \text{ in } \mathbb{R}_{\geq 0} \quad \text{with} \quad M = \begin{pmatrix} a & b \\ b & a \end{pmatrix}.$$

109 Obviously, the eigenvalues of  $M$  are  $\lambda_{1/2} = a \pm b$  with eigenvectors  $v_{1/2} = (1, \pm 1)^\top$ .  
 110 According to Proposition 4.2, the unique solution satisfies

$$111 \quad y(t) = \alpha_1 e^{(a+b)t} v_1 + \alpha_2 e^{(a-b)t} v_2 \quad \text{provided that} \quad y(0) = \alpha_1 v_1 + \alpha_2 v_2.$$

113 From the initial condition, we get that

$$114 \quad \left. \begin{aligned} y_1(0) &= \alpha_1 + \alpha_2 \\ y_2(0) &= \alpha_1 - \alpha_2 \end{aligned} \right\} \quad \text{and hence} \quad \begin{cases} \alpha_1 &= \frac{y_1(0) + y_2(0)}{2} \\ \alpha_2 &= \frac{y_1(0) - y_2(0)}{2}. \end{cases}$$

116 Altogether, the solution reads

$$117 \quad y_1(t) = \frac{y_1(0) + y_2(0)}{2} e^{(a+b)t} + \frac{y_1(0) - y_2(0)}{2} e^{(a-b)t},$$

$$118 \quad y_2(t) = \frac{y_1(0) + y_2(0)}{2} e^{(a+b)t} - \frac{y_1(0) - y_2(0)}{2} e^{(a-b)t}.$$

120 Let  $a = -51$ ,  $b = -50$ , so that  $a + b = -101$  and  $a - b = -1$ . For large  $t > 0$ , a good  
 121 approximation of  $y(t)$  is obtained, if we neglect the first summands with  $e^{(a+b)t}$ . But the  
 122 ODE is stiff. The explicit Euler method requires  $|1 + h(a \pm b)| < 1$  to avoid oscillations  
 123 (in both components) of  $y$ . However, note that the “component”  $y(t) \cdot (1, -1)^\top$  (which  
 124 corresponds to the “good” eigenvalue  $\lambda_2 = a - b = -1$  will nevertheless be approximated  
 125 well for any  $0 < h < 1$ .

---

126

127 **4.2. Stability domains.** The idea of all notions of stability is that the discrete solu-  
 128 tion should reflect certain properties of the continuous solution, at least for model prob-  
 129 lems, where the solution (or its qualitative behavior) is known. For stability domains (as  
 130 well as  $A$ -stability and  $L$ -stability), one considers the scalar model problem

$$131 \quad y'(t) = \lambda y(t) \text{ in } \mathbb{R}_{\geq 0}, \quad y(0) = y_0, \quad (4.13)$$

132 where  $\lambda \in \mathbb{C}$  with  $\operatorname{Re} \lambda < 0$ .

133 **Definition 4.7.** Let  $\Phi(t, y, z, h)$  be the incremental function of a one-step method. A  
 134 function  $R : \mathbb{C} \rightarrow \mathbb{C}$  is called **stability function** of the method, if

$$135 \quad y_{\ell+1} := y_\ell + h_\ell \Phi(t_\ell, y_\ell, y_{\ell+1}, h_\ell) \stackrel{!}{=} R(\lambda h_\ell) y_\ell. \quad (4.14)$$

136

137 **Example 4.8.** We rephrase Remark 4.5: The stability function of the explicit Euler  
 138 method is  $R(z) = 1 + z$ , and it is well-defined for all  $z \in \mathbb{C}$ . The stability function of the  
 139 implicit Euler method is  $R(z) = 1/(1 - z)$ , and it is well-defined for all  $z \in \mathbb{C} \setminus \{1\}$ .

---

140 **Theorem 4.9 (Stability function of Runge–Kutta methods).** Let  $\begin{array}{c|c} c & A \\ \hline & b^\top \end{array}$  be an  
 141  $m$ -step Runge–Kutta method. Then, the stability function reads

$$142 \quad R(z) = 1 + z b^\top (I - zA)^{-1} \mathbf{1}, \quad \text{where } \mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^m. \quad (4.15)$$

143 It is well-defined for all  $z \in \mathbb{C}$  such that  $1/z \notin \sigma(A) := \{\lambda \in \mathbb{C} : \lambda \text{ is eigenvalue of } A\}$ .  
 144 Moreover, there hold the following statements:

- 145 (i) If the method is explicit, then  $R \in \mathbb{P}_m$ .  
 146 (ii) If the method is implicit, then  $R = P/Q$  with  $P, Q \in \mathbb{P}_m$ .
- 

147 *Proof.* Recall the Runge–Kutta increments

$$148 \quad k_i = f\left(t + c_i h, y_\ell + h \sum_{j=1}^m A_{ij} k_j\right) \stackrel{!}{=} \lambda \left(y_\ell + h \sum_{j=1}^m A_{ij} k_j\right) \quad \text{for all } i = 1, \dots, m.$$

149 With the vector  $k = (k_1, \dots, k_m)^\top$ , it follows that  $k = \lambda y_\ell \mathbf{1} + \lambda h A k$  and hence

$$150 \quad k = \lambda y_\ell (I - \lambda h A)^{-1} \mathbf{1}, \quad \text{provided that } 1/(\lambda h) \notin \sigma(A).$$

151 By definition of a Runge–Kutta method, it follows that

$$152 \quad y_{\ell+1} = y_\ell + h \sum_{i=1}^m b_i k_i = y_\ell + h b^\top k = y_\ell [1 + \lambda h b^\top (I - \lambda h A)^{-1} \mathbf{1}] = y_\ell R(\lambda h).$$

153

158 This proves (4.15).

159 For the remainder of the proof, recall **Cramer's rule**: Let  $V = (v_1, \dots, v_n) \in \mathbb{R}^n$  be  
 160 invertible with columns  $v_j \in \mathbb{R}^n$ . Let  $w \in \mathbb{R}^n$ . Define  $V_i := (v_1, \dots, v_{i-1}, w, v_{i+1}, \dots, v_n) \in$   
 161  $\mathbb{R}^{n \times n}$ . Then, the vector  $x := V^{-1}w \in \mathbb{R}^n$  satisfies that

$$162 \quad x_i = \frac{\det(V_i)}{\det(V)} \quad \text{for all } i = 1, \dots, m. \quad (4.16)$$

164 To apply Cramer's rule, we write

$$165 \quad R(z) = 1 + zb^\top(I - zA)^{-1}\mathbf{1} = 1 + zb^\top\gamma \quad \text{where } (I - zA)\gamma = \mathbf{1}.$$

167 This is rewritten as

$$168 \quad \underbrace{\begin{pmatrix} I - zA & 0 \\ -zb^\top & 1 \end{pmatrix}}_{=:V} \begin{pmatrix} \gamma \\ R(z) \end{pmatrix} = \begin{pmatrix} \mathbf{1} \\ 1 \end{pmatrix}.$$

169 According to the Laplace cofactor expansion (for the  $(m + 1)$ -th column), it holds that

$$170 \quad \det(V) = (-1)^{(m+1)+(m+1)} \det(I - zA) = \det(I - zA) \neq 0 \quad \text{provided that } 1/z \notin \sigma(A).$$

171 Hence,  $V$  is invertible and Cramer's rule yields (for  $R(m)$  being the  $(m + 1)$ -th coefficient  
 172 of the solution) that

$$173 \quad R(z) = \frac{\det \begin{pmatrix} I - zA & \mathbf{1} \\ -zb^\top & 1 \end{pmatrix}}{\det \begin{pmatrix} I - zA & 0 \\ -zb^\top & 1 \end{pmatrix}} = \frac{\det \begin{pmatrix} I - zA & \mathbf{1} \\ -zb^\top & 1 \end{pmatrix}}{\det(I - zA)} =: \frac{P(z)}{Q(z)}.$$

174 Clearly, the characteristic polynomial  $Q(z) = \det(I - zA)$  satisfies  $Q \in \mathbb{P}_m$ . Moreover,  
 175 if the method is explicit, then  $A$  is strictly lower triangular and hence  $\det(I - zA) = 1$ .  
 176 Hence, it only remains to show that the numerator satisfies that  $P \in \mathbb{P}_m$ . With the  
 177 Laplace cofactor expansion (for the  $(m + 1)$ -th column), we write

$$178 \quad P(z) = \det \begin{pmatrix} I - zA & \mathbf{1} \\ -zb^\top & 1 \end{pmatrix} =: \det(S - zT \quad \mathbf{1}) = \sum_{i=1}^{m+1} (-1)^{i+(m+1)} \det(S_i - zT_i),$$

179 where  $S_i, T_i \in \mathbb{R}^{m \times m}$  are obtained by canceling the  $i$ -th row of  $S, T \in \mathbb{R}^{(m+1) \times m}$ . A simple  
 180 induction on the dimension  $m$  (together with the Laplace cofactor expansion) proves that  
 181  $\det(S_i - zT_i) \in \mathbb{P}_m$ . Therefore, we conclude that  $P \in \mathbb{P}_m$ .  $\square$

---

186 **Corollary 4.10 (Stability function is an approximation of exp).** *Let  $R(z)$  be the*  
 187 *stability function of a Runge–Kutta method with consistency order  $p \geq 1$ . Then,*

$$188 \quad R(z) = \exp(z) + \mathcal{O}(z^{p+1}). \quad (4.17)$$

---

**Gabriel Cramer** (1704–1752) was a Swiss Mathematician. He studied Mathematics at the University of Geneva and took his PhD in 1722. In 1724, he became professor at the University of Geneva. He had personal contact with his colleagues Johann Bernoulli, Leonhard Euler, and James Stirling.

**Pierre-Simon (Marquis de) Laplace** (1749–1827) was a French Mathematician and, from 1768–1771, a student of d'Alembert in Paris. In 1771, he became teacher for geometry, trigonometry, analysis, and statistics at the military academy in Paris. In 1773, he became associate member of the French Academy of Sciences.

190 Moreover, if the method is explicit, then it is a polynomial with

$$191 \quad R(z) = \sum_{j=1}^p \frac{z^j}{j!} + \mathcal{O}(z^{p+1}), \quad (4.18)$$

192 and  $\mathcal{O}(z^{p+1})$  vanishes for an explicit  $p$ -stage method of order  $p$ .

194 *Proof.* Consider the first step of the Runge–Kutta method for  $\lambda = -1$  and  $y_0 = 1$ . Then,

$$195 \quad \exp(-h) = y(h) \approx y_1 = R(-h)y_0 = R(-h).$$

197 Moreover, consistency order  $p$  thus implies that

$$198 \quad \exp(-h) - R(-h) = \mathcal{O}(h^{p+1}). \quad (4.19)$$

200 Note that  $R(z) = 1 + zb^\top(I - zA)^{-1}\mathbf{1}$  is smooth locally around  $z = 0$  (i.e., analytic).

201 Hence, it holds that

$$202 \quad f(z) := R(z) - \exp(z) = \sum_{n=0}^p \frac{f^{(n)}(0)}{n!} z^n + \mathcal{O}(z^{p+1}).$$

204 Due to (4.19), we see that  $f^{(n)}(0) = 0$  for all  $n = 0, \dots, p$ . Therefore, it follows that

$$205 \quad R(z) = \exp(z) + (R(z) - \exp(z)) = \exp(z) + \mathcal{O}(z^{p+1}).$$

207 If the method is explicit, then  $R \in \mathbb{P}_m$ , i.e.,  $R(z) = \sum_{j=0}^m a_j z^j$ . Then,

$$208 \quad \sum_{j=0}^p \left( a_j - \frac{1}{j!} \right) z^j = R(z) - \exp(z) + \mathcal{O}(z^{p+1}) \stackrel{(4.17)}{=} \mathcal{O}(z^{p+1}).$$

210 Consequently, the lower-order powers of  $z$  vanish, i.e.,  $a_j = 1/j!$  for  $j = 0, \dots, p$ . This  
211 concludes the proof.  $\square$

---

212 **Example 4.11.** The stability function of the Heun method (Example 2.18) and the mod-  
213 ified Euler method (Example 2.19) is  $R(z) = 1 + z + \frac{z^2}{2}$ , since both methods are explicit  
214 2-stage methods of order  $p = 2$ . The stability function of the classical RK4 method (Ex-  
215 ample 2.25) is  $R(z) = 1 + z + \frac{z^2}{2} + \frac{z^3}{6} + \frac{z^4}{24}$ , since RK4 is a 4-step method of order  
216  $p = 4$ .

---

217 **Definition 4.12.** Let  $R(z)$  be the stability function of a one-step method. Then,

$$218 \quad S := \{z \in \mathbb{C} : |R(z)| \leq 1\} \quad (4.20)$$

219 denotes the corresponding **stability domain**.

---

221 **Remark 4.13.** If  $\lambda h = z \in S$ , then the approximations  $y_\ell \approx y(t_\ell)$  of the model prob-  
222 lem (4.13) remain bounded. Note that  $y(t) = e^{\lambda t} y_0 \rightarrow 0$  as  $t \rightarrow \infty$ , since  $\operatorname{Re} \lambda < 0$ .

---

223 **Example 4.14.** The stability domain of the explicit Euler method is

$$224 \quad S = \{z \in \mathbb{C} : |1 + z| \leq 1\} = \overline{U_1(-1)}.$$

226 The stability domain of the implicit Euler method is

$$227 \quad S = \left\{ z \in \mathbb{C} : \left| \frac{1}{1-z} \right| \leq 1 \right\} = \{z \in \mathbb{C} : 1 \leq |1-z|\} = \mathbb{C} \setminus U_1(1). \\ 228$$

229 **Corollary 4.15.** For all Runge–Kutta methods with consistency order  $p \geq 1$ , it holds  
230 that  $0 \in \partial S$ .

231 *Proof.* According to Corollary 4.10, it holds that

$$232 \quad R(z) = \exp(z) + \mathcal{O}(z^{p+1}) = 1 + z + \mathcal{O}(z^2). \\ 233$$

234 For all sufficiently small  $z = \pm h$ , it hence follows that

- 235 •  $R(+h) = 1 + h + \mathcal{O}(h^2) \geq 1 + h/2 > 1$ , i.e.,  $+h \notin S$ ;
- 236 •  $R(-h) = 1 - h + \mathcal{O}(h^2) \leq 1 - h/2 < 1$ , i.e.,  $-h \in S$ .

237 Hence, each neighborhood of  $z = 0$  contains one point  $-h \in S$  and  $+h \notin S$ . This proves  
238 that  $0 \in \partial S$ . □

239

### 240 4.3. A-stability and L-stability.

241 **Definition 4.16.** Let  $R(z)$  be the stability function of a one-step method. The one-step  
242 method is

- 243 • **A-stable**, if  $\sup_{\operatorname{Re} z \leq 0} |R(z)| \leq 1$ ;
- 244 • **L-stable**, if  $\lim_{\operatorname{Re} z \rightarrow -\infty} |R(z)| = 0$ .

245 **Remark 4.17.** Due to the definition of the stability domain  $S$  in (4.20), A-stability is  
246 equivalent to  $\mathbb{C}^- := \{z \in \mathbb{C} : \operatorname{Re} z \leq 0\} \subseteq S$ . Hence, the explicit Euler method is not  
247 A-stable (see Example 4.14).

248 **Remark 4.18.** If a method is A-stable, then

$$249 \quad |y_{\ell+1}| = |R(\lambda h)| |y_\ell| \leq |y_\ell| \quad \text{for all } h > 0 \text{ and } \operatorname{Re} \lambda \leq 0, \\ 250$$

251 i.e., the discrete solutions of the model problem (4.13) are non-expansive and, in partic-  
252 ular, remain at least bounded.

253 If a method is L-stable, then

$$254 \quad |y_{\ell+1}| = |R(\lambda h)| |y_\ell| \rightarrow 0 \quad \text{as } h \rightarrow \infty \text{ for all } \operatorname{Re} \lambda \leq 0, \\ 255$$

256 i.e., the discrete solutions of the model problem (4.13) decay with larger time-steps.

257 **Example 4.19.** The implicit Euler method satisfies that  $R(z) = \frac{1}{1-z}$  and  $S = \mathbb{C} \setminus U_1(1) \supseteq$   
258  $\mathbb{C}^-$ . Hence, the implicit Euler method is L-stable and A-stable.



259 **Exercise 4.20.** Show that the stability function of the implicit midpoint rule is  $R(z) =$   
 260  $\frac{1+z/2}{1-z/2}$ . Determine the stability domain of the implicit midpoint rule! Is the implicit mid-  
 261 point rule  $A$ -stable and/or  $L$ -stable?

---

262 **Theorem 4.21 (Explicit Runge–Kutta methods fail).** No explicit Runge–Kutta  
 263 method is  $L$ -stable. No consistent explicit Runge–Kutta method is  $A$ -stable. In particular,  
 264 any  $A$ -stable and/or  $L$ -stable Runge–Kutta method is implicit.

---

265 *Proof ( $L$ -stability fails).* Recall that the stability function of an  $m$ -stage Runge–Kutta  
 266 method satisfies that  $R = P/Q$  with  $P, Q \in \mathbb{P}_m$ . Obviously,  $L$ -stability thus is equivalent  
 267 to the fact that  $P \in \mathbb{P}_{\nu-1}$  and  $Q \in \mathbb{P}_\nu \setminus \mathbb{P}_{\nu-1}$  for some  $1 \leq \nu \leq m$  to ensure that  
 268  $P/Q = \mathcal{O}(1/z)$  as  $|z| \rightarrow \infty$ .  $\square$

269 *Proof ( $A$ -stability fails).* The stability function of an explicit  $m$ -stage Runge–Kutta method  
 270  $\frac{c}{b^\top} \left| \begin{array}{c} A \\ \hline \end{array} \right.$  satisfies that  $R \in \mathbb{P}_m$ . Consistency implies that  $\sum_{j=1}^m b_j = 1$  and hence the  
 271 method has consistency order  $p \geq 1$ . Together with Corollary 4.10, this implies that

$$272 \quad R(z) = \sum_{j=0}^{\nu} a_j z^j \quad \text{for some } 1 \leq \nu \leq m \text{ with } a_\nu \neq 0.$$

273  
 274 The triangle inequality thus proves that

$$275 \quad |R(z)| \geq |a_\nu| |z|^\nu - \sum_{j=0}^{\nu-1} |a_j| |z|^j = |z|^\nu \left( |a_\nu| - \sum_{j=0}^{\nu-1} |a_j| \frac{1}{|z|^{\nu-j}} \right) \xrightarrow{\operatorname{Re} z \rightarrow -\infty} \infty.$$

276  
 277 Hence,  $\sup_{\operatorname{Re} z \leq 0} |R(z)| = \infty$  and the method is *not*  $A$ -stable.  $\square$

---

278 **Theorem 4.22.** Let  $\frac{c}{b^\top} \left| \begin{array}{c} A \\ \hline \end{array} \right.$  be an implicit Runge–Kutta method such that  $A$  is invertible.

279 Then, the following statements (i)–(ii) are equivalent:

280 (i)  $b^\top A^{-1} \mathbf{1} = 1$  with  $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^m$ .

281 (ii) The method is  $L$ -stable.

---

282 *Proof.* If  $\|\frac{1}{z} A^{-1}\| < 1$  (e.g.,  $|z| > 2\|A^{-1}\|$ ), then the Neumann series proves that

$$283 \quad \left( I - \frac{1}{z} A^{-1} \right)^{-1} = \sum_{j=0}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j = I + \sum_{j=1}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j = I + \frac{1}{z} A^{-1} \left( I + \sum_{j=1}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j \right).$$

284  
 285 Moreover,

$$286 \quad I - zA = (A^{-1} - zI)A = -z \left( I - \frac{1}{z} A^{-1} \right) A.$$

287  
 288 This leads to

$$289 \quad (I - zA)^{-1} = -\frac{1}{z} A^{-1} \left( I + \frac{1}{z} A^{-1} \right)^{-1} = -\frac{1}{z} A^{-1} \left[ I + \frac{1}{z} A^{-1} \left( I + \sum_{j=1}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j \right) \right].$$

290

291 With Theorem 4.9, it follows that

$$\begin{aligned}
 292 \quad R(z) &\stackrel{(4.15)}{=} 1 + zb^\top(I - zA)^{-1}\mathbf{1} = 1 - b^\top A^{-1} \left[ I + \frac{1}{z} A^{-1} \left( I + \sum_{j=1}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j \right) \right] \mathbf{1} \\
 293 \quad &= 1 - b^\top A^{-1} \mathbf{1} - \frac{1}{z} \left[ b^\top A^{-1} A^{-1} \left( I + \sum_{j=1}^{\infty} \left( \frac{1}{z} A^{-1} \right)^j \right) \right] \mathbf{1}.
 \end{aligned}$$

294  
295 Hence,  $L$ -stability is equivalent to  $1 - b^\top A^{-1} \mathbf{1} = 0$ . This concludes the proof.  $\square$

296 **Remark 4.23.** Suppose the assumptions of Theorem 4.22. If the  $j$ -th column of  $A$  is  
297 constant  $b_j \neq 0$ , then the  $j$ -th unit vector  $e_j \in \mathbb{R}^m$  leads to  $Ae_j = b_j \mathbf{1}$ . This implies that  
298  $A^{-1} \mathbf{1} = \frac{1}{b_j} e_j$  and results in

$$299 \quad b^\top A^{-1} \mathbf{1} = \frac{1}{b_j} b^\top e_j = 1.$$

300  
301 According to Theorem 4.22, the method is  $L$ -stable.

302 **Remark 4.24 (Radau-IIA methods).** The Radau-IIA methods are collocation meth-  
303 ods with  $c_m = 1$ . The remaining nodes  $c_j$  of  $c \in \mathbb{R}^m$  are chosen such that the induced  
304 quadrature rule has maximal exactness. Arguing as for the Gaussian quadrature rule, one  
305 obtains that the nodes  $0 < c_1 < c_2 < \dots < c_m = 1$  are unique. The quadrature rule  
306 has exactness  $2m - 2$  (i.e., one order less than the Gaussian quadrature). Hence, the  
307 Radau-IIA methods have consistency order  $p = 2m - 1$ . Recall from Theorem 3.25 that  
308 the Runge–Kutta data of a collocation method read

$$309 \quad A_{ij} = \int_0^{c_i} L_j dt \quad \text{and} \quad b_j = \int_0^1 L_j dt \quad \text{for all } i, j = 1, \dots, m.$$

311 For Radau-IIA methods, the last row of  $A \in \mathbb{R}^{m \times m}$  thus coincides with  $b \in \mathbb{R}^m$  (since  
312  $c_m = 1$ ). Moreover,  $A$  is invertible, since  $c_j > 0$ . With the  $m$ -th unit vector  $e_m \in \mathbb{R}^m$ ,  
313 this means that  $b^\top = e_m^\top A$  and hence  $b^\top A^{-1} = e_m^\top$ . Consequently,

$$314 \quad b^\top A^{-1} \mathbf{1} = e_m^\top \mathbf{1} = 1.$$

315  
316 According to Theorem 4.22, the Radau-IIA methods are  $L$ -stable.

317 **Exercise 4.25.** Show that the implicit Euler method is the unique Radau-IIA method for  
318  $m = 1$ .

319 **Remark 4.26.** The  $A$ -stability of collocation methods is fully understood mathematically;  
320 see [But08, Theorem 358A]. One can show that all Radau-IIA methods are  $A$ -stable and  
321  $L$ -stable. Moreover, one can show that all Gauss methods are  $A$ -stable, but fail to be  
322  $L$ -stable.

**Jean Charles Rodolphe Radau** (1835–1911) was a German-French mathematician. Born in Ger-  
many, he studied mathematics and astronomy at the University of Königsberg. In 1858, he moved to  
Paris and later became French citizen in 1873. In 1897, he became member of the French Academy of  
Sciences.



---

**Theorem A.1 (Banach fixpoint theorem).** *Let  $M$  be a complete metric space and  $A : M \rightarrow M$  be a contraction, i.e., there exists a contraction constant  $0 < \kappa < 1$  such that*

$$\forall x, y \in M : \quad d(Ax, Ay) \leq \kappa d(x, y). \quad (\text{A.1})$$

*Then,  $A$  has a unique fixpoint  $w \in M$ , i.e.,  $Aw = w$ . Moreover, for any initial value  $w_0 \in M$ , the sequence of Banach iterates  $w_{n+1} := Aw_n$  converges to  $w$  and it holds that*

$$d(w, w_n) \leq \frac{\kappa}{1 - \kappa} d(w_{n+1}, w_n) \leq \frac{\kappa^n}{1 - \kappa} d(w_1, w_0) \quad \text{for all } n \in \mathbb{N}_0. \quad (\text{A.2})$$


---

**Remark A.2.** *Often, the Banach fixpoint theorem is applied in the following setting: Let  $X$  be a Banach space and  $M \subseteq X$  be a closed subspace. Suppose that  $A : M \rightarrow M$  is a contraction with respect to the natural metric, i.e., there exists a contraction constant  $0 < \kappa < 1$  such that*

$$\forall x, y \in M : \quad \|Ax - Ay\|_X \leq \kappa \|x - y\|_X. \quad (\text{A.3})$$

*Then, the Banach fixpoint theorem applies.*

---

**Theorem A.3 (Implicit function theorem).** *Let  $f \in C^1(\mathbb{R}^m \times \mathbb{R}^n; \mathbb{R}^n)$ . Let  $(x, y) \in \mathbb{R}^m \times \mathbb{R}^n$  with  $f(x, y) = 0$ . Suppose that the Jacobi matrix  $D_y f(x, y) \in \mathbb{R}^{n \times n}$  is invertible. Then, there exist open sets  $U \subset \mathbb{R}^m$  and  $V \subset \mathbb{R}^n$  with  $(x, y) \in U \times V$  as well as a function  $g \in C^1(U; V)$  such that*

$$\forall (\tilde{x}, \tilde{y}) \in U \times V : \quad [f(\tilde{x}, \tilde{y}) = 0 \iff \tilde{y} = g(\tilde{x})] \quad (\text{A.4})$$


---

## REFERENCES

- [But08] JOHN C. BUTCHER: Numerical methods for ordinary differential equations, Wiley, Chichester, second edition, 2008.
- [SWP12] KARL STREHMEL, RÜDIGER WEINER, HELMUT PODHAISKY: Numerik gewöhnlicher Differentialgleichungen, Springer, Berlin, second edition, 2012 [in German].
- [Wal00] WOLFGANG WALTER: Gewöhnliche Differentialgleichungen, Springer, Berlin, second edition, 2000 [in German].

---

**Stefan Banach (1892–1945)** was a Polish mathematician. Being the founder of modern functional analysis, he is amongst the most important mathematicians of the 20th century. His major work was the 1932 book “*Théorie des opérations linéaires*”, the first monograph on the general theory of functional analysis. Banach finished his PhD in mathematics at Lviv University in 1922 (the thesis also contained the Banach fixpoint theorem). In the same year, he completed the habilitation and became associate professor at Lviv University. In 1926, he was promoted to full professor. In August 1945, he died of lung cancer.