

Problem 4.1 (Consistency of two-stage ERK, 4 Points)

A general explicit Runge-Kutta method with two stages has the following form

$$y_n = y_{n-1} + h_n F(h_n; t_{n-1}, y_{n-1})$$

$$F(h_n; t_{n-1}, y_{n-1}) = b_1 f_{n-1} + b_2 f(t_{n-1} + h_n c_2, y_{n-1} + h_n a_{21} f_{n-1}),$$

where $f_{n-1} := f(t_{n-1}, y_{n-1})$.

- a) Use Taylor-Expansion to find conditions on a_i, b_i, c_i such that they imply consistency of order 2.

Proof. We do a Taylor-Expansion of $u(t_1)$ around t_0 and obtain:

$$u(t_1) = u(t_0) + (t_1 - t_0) \frac{d}{dt} u(t) \Big|_{t=t_0} + \frac{(t_1 - t_0)^2}{2} \cdot \frac{d^2}{dt^2} u(t) \Big|_{t=t_0} + \mathcal{O}((t_1 - t_0)^3)$$

As we have seen in the proof of 2.3.9 the Taylor-expansion of the exact solⁿ up to $\mathcal{O}(h^2)$ is

$$u(t_1) = u(t_0) + h f(t_0, u_0) + \frac{h^2}{2} \left(\underbrace{\partial_t f(t_0, u_0)}_{:= f_t} + \underbrace{\partial_y f(t_0, u_0) f(t_0, u_0)}_{:= f_y} \right) + \mathcal{O}(h^3) \quad (*)$$

We now do a Taylor-Expansion of

$f(t_0 + h c_1, y_0 + h a_{21} f(t_0, u_0))$ around (t_0, y_0) . We get:

$$f(t_0 + h c_1, y_0 + h a_{21} f) = f + \partial_t f(t_0, u_0)(t_0 + h c_1 - t_0) + \partial_y f(t_0, y_0)(y_0 + h a_{21} f - y_0) + \mathcal{O}(h^3)$$

$$= f + h c_1 f_t + f_y (h a_{21} f) + f + h c_1 f_t + h a_{21} f_y f$$

Thus, we obtain:

$$y_1 = y_0 + h (b_1 f + b_2 (f + h c_1 f_t + h a_{21} f_y f) + \mathcal{O}(h^3))$$

$$= y_0 + h b_1 f + h b_2 f + h^2 c_1 b_2 f_t + h^2 a_{21} b_2 f_y f + \mathcal{O}(h^3) \quad (**)$$

Comparing $(**)$ w/ $(*)$ we get:

$$y_1 - u_1 = y_0 - u_0 + h f(b_1 + b_2 - 1) + h^2 f_t (c_1 b_2 - \frac{1}{2}) + h^2 f_y f (a_{21} b_2 - \frac{1}{2}) + \mathcal{O}(h^3)$$

For y_n to be of consistency order 2 it is sufficient therefore sufficient that:

$$b_1 + b_2 = 1, \quad c_1 b_2 = \frac{1}{2}, \quad a_{21} b_2 = \frac{1}{2}$$

□

- b) Thus, verify that both the *modified Euler method* and the *Heun method* satisfy your conditions.

• Modified Euler: We have: $b_1 = 0, b_2 = 1 \Rightarrow b_1 + b_2 = 0 + 1 = 1$. ✓

$$\text{and } c_1 = \frac{1}{2}, b_2 = 1 \Rightarrow c_1 b_2 = \frac{1}{2} \quad \checkmark$$

$$\text{and } a_{21} = \frac{1}{2}, b_2 = 1 \Rightarrow a_{21} b_2 = \frac{1}{2} \quad \checkmark$$

• Heun-Method: We have: $b_0 = b_1 = \frac{1}{2} \Rightarrow b_1 + b_2 = 1$

$$\text{and } c_2 = 1, b_2 = \frac{1}{2} \Rightarrow c_2 b_2 = \frac{1}{2}$$

$$\text{and } a_{21} = 1, b_2 = \frac{1}{2} \Rightarrow a_{21}b_2 = \frac{1}{2}$$

Thus both methods satisfy the derived conditions. \square

- c) Are there any other explicit two-stage Runge-Kutta methods with consistency of order 2?

If there are then the coefficients must obey:

$$\begin{aligned} b_1 + b_2 &= 1 \\ c_2 b_2 &= \frac{1}{2} \quad (\Rightarrow) \quad b_2 = \frac{1}{2} c_2 \\ a_{21} b_2 &= \frac{1}{2} \quad (\Rightarrow) \quad b_2 = \frac{1}{2} a_{21} \end{aligned}$$

$$\Rightarrow \frac{1}{2} c_2 = \frac{1}{2} a_{21} \quad (\Rightarrow) \quad c_2 = a_{21} \quad (1.1)$$

$$\Rightarrow b_1 + b_2 = \frac{1}{2} c_2 + b_1 = 1 \Rightarrow b_1 = 1 - \frac{1}{2} c_2, \text{ and thus } b_2 = 1 - (1 - \frac{1}{2} c_2) = \frac{1}{2} c_2$$

Thus our coefficients have become:

$$\begin{cases} b_1 = 1 - \frac{1}{2} c_2, \quad b_2 = \frac{1}{2} c_2 \\ a_{21} = c_2 \\ c_2 \in (0, 1] \text{ arbitrary.} \end{cases}$$

Thus we conclude that there are indeed infinitely many two-stage ERKs w/ consistency order two. \square

Problem 4.2 (Maximal consistency order of ERK, 4 Points)

Let $\lambda, u_0 \in \mathbb{R}$ and consider the initial value problem

$$\begin{aligned} u'(t) &= \lambda u(t) \text{ for all } t \in \mathbb{R}, \\ u(0) &= u_0. \end{aligned}$$

Now let s be a natural number and assume we apply an explicit s -stage Runge-Kutta method to this problem. (See Definition 2.3.2.)

a) Let k_i for $i = 1, \dots, s$ be defined as in Definition 2.3.2. Show that

$$hk_i = p_i(\lambda h)y_0$$

with polynomials p_i with $\deg(p_i) \leq i$.

Proof. As in 2.3.2 we have

$$k_i = f(t_0 + c_i h, y_i) \quad \text{where} \quad y_i = y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j, \quad \text{for } i=1, \dots, s. \quad \text{In this particular case}$$

$f(t, y) = \lambda y$, thus, $k_i = \lambda y_i$. We show the statement using induction on i for a fixed s .

$i=1$: Then $y_1 = y_0$ and $k_1 = f(t_0 + c_1 h, y_1) = \lambda y_0 \Rightarrow h k_1 = h \lambda y_0$, thus $p_1 = x$ does the job.

Let $i \geq 1$ and $i \leq s-1$. Then we obtain

$$hk_{i+1} = \lambda y_{i+1} = \lambda \left(y_0 + \sum_{j=1}^i a_{ij} k_j \right) = h \lambda \left(y_0 + \sum_{j=1}^i a_{ij} p_j(\lambda h) y_0 \right) = h y_0 \left(1 + \sum_{j=1}^i a_{ij} p_j(\lambda h) \right)$$

We define

$$p_{i+1}(T) = T \left(1 + \sum_{j=1}^i a_{ij} p_j(T) \right) \quad \text{since } \deg p_j \leq j, \forall j=1, \dots, i. \quad \text{We obtain}$$

$$\deg \left(1 + \sum_{j=1}^i a_{ij} p_j(T) \right) \leq i, \quad \text{thus} \quad \deg \left(T \left(1 + \sum_{j=1}^i a_{ij} p_j(T) \right) \right) = \underbrace{\deg(T)}_{\leq 1} + \underbrace{\deg \left(1 + \sum_{j=1}^i a_{ij} p_j(T) \right)}_{\leq i} \leq i+1.$$

Thus we obtain: $h k_{i+1} = y_0 \lambda h \left(1 + \sum_{j=1}^i a_{ij} p_j(\lambda h) \right) = y_0 p_{i+1}(\lambda h)$, with $\deg p_{i+1} \leq i+1$.

which finishes the proof. \square

b) Thus, show that an explicit s -stage Runge-Kutta method cannot be consistent of order $p > s$.

Proof. Let $p > s$. A method is consistent of order p iff (using that for ERKs we have $h_k = k$)

$$\max_{k=1, \dots, h} |I_k| \leq ch^p \quad \text{where} \quad I_k = \frac{u_k - u_{k-1}}{h} - F_h(t_{k-1}, y_{k-1}) \quad (*). \quad \text{We have}$$

$$F_h(t_{k-1}, y_{k-1}) = \sum_{i=1}^s b_i f(t_0 + c_i h, y_i) = \lambda \sum_{i=1}^s b_i g_i \quad \text{where} \quad g_i = y_{k-1} + \sum_{j=1}^{i-1} a_{ij} k_j \quad (*)$$

An exact soln of the IVP is given by $u(t) = u_0 e^{\lambda t}$ $(*)$, because $u(0) = u_0 e^{\lambda \cdot 0} = u_0$ and $u'(t) = \lambda u_0 e^{\lambda t} = \lambda u(t)$.

$$T_1 h^{-p} = \frac{u_1 - u_0}{h^{p+1}} - \frac{\lambda}{h^p} \sum_{i=1}^s b_i (y_0 + \sum_{j=1}^{i-1} a_{ij} \underbrace{p_j(h) y_0}_{h^{p+1}})$$

$$= \frac{u_1 - u_0}{h^{p+1}} - \lambda \sum_{i=1}^s \frac{b_i y_0}{h^p} + \sum_{j=1}^{i-1} a_{ij} \underbrace{p_j(h) y_0}_{h^{p+1}}$$

$p_j(h)$ are polys w/ degree $\leq j \leq s < p$, thus

$\underbrace{p_j(h)}_{h^{p+1}}$ diverges $\forall j = 1, \dots, s$.

Problem 4.3 (Local error control of ERK methods)

In Algorithm 2.4.2 of the lecture notes the optimal step size is determined by

$$h_{opt} = h_k \sqrt[p+1]{\frac{\varepsilon}{\|y_k - \hat{y}_k\|}}$$

where ε is a tolerance given by the user, y_k is the approximation for u_k obtained from a method with consistency of order p using the time step h_k and \hat{y}_k likewise, but obtained from a method with consistency of order $\geq p+1$.

- In the algorithm, if $\|y_k - \hat{y}_k\| > \varepsilon$, then y_k is rejected. Why?
- If y_k is rejected as above, step k is repeated with step size $h_k^{new} = h_{opt}$. Why is this a good choice?
- If the step is accepted, i.e.

$$\|y_k - \hat{y}_k\| \leq \varepsilon,$$

the step size for the next step is set to $h_{k+1} = h_{opt}$. Why is this a good choice?

Note: You will have to make reasonable assumptions for this algorithm to make sense. Extracting them from the paragraph preceding Algorithm 2.4.2 (or coming up with them yourself) is part of the exercise!

- a) If our threshold $\varepsilon > 0$ (for $\|y_k - \hat{y}_k\|$) is greater than or equals to $\|y_k - \hat{y}_k\|$, then we know that the allowed error (ε) was too big and we choose a smaller time-step (h_{opt}) and recompute y_k, \hat{y}_k . Thus we guarantee that $\|y_k - \hat{y}_k\| < \varepsilon$, for all k .
- b) If the time-step h_k and the computations for y_k, \hat{y}_k are rejected, we h_{opt} is a good choice for h_k^{new} , because $h_k^{new} < h_k$ and as the local error $|u_k - y_k| \leq C e^{LT} h^p$ (2.25) is bounded by (2.25), it is reasonable to choose a smaller h .

c) This is only a good choice if $h_{opt} \leq T - t_k$.
Thus we use a time-step s.t. $t_{k+1} \leq T$.
Since $\|y_k - \tilde{y}_k\| \leq \varepsilon$, and thus $h_{opt} \geq h_k$, it is reasonable to assume that h_{opt} is the maximal time-step we can make (using \tilde{y}_k as u_0) without making a local error greater than ε .