

## Overview of numerical methods for ODEs

In this lecture we provide a brief overview of the most common numerical methods to approximate the solution of an initial value problem for systems of ODEs of the form

$$\begin{cases} \frac{d\mathbf{y}}{dt} = \mathbf{f}(\mathbf{y}, t) \\ \mathbf{y}(0) = \mathbf{y}_0 \end{cases} \quad (1)$$

where  $\mathbf{y}(t) = [y_1(t) \cdots y_n(t)]^T$  is a (column) vector of phase variables,  $\mathbf{f} : D \times [0, T] \rightarrow \mathbb{R}^n$ ,  $D$  is a subset of  $\mathbb{R}^n$ , and  $T$  is the integration period. Most of the material presented in this lecture can be found, e.g., in [2, 3, 4]. We assume that the initial value problem (1) well-posed, i.e., that it has a unique solution<sup>1</sup> (at least for some time  $t\tau > 0$ ). We have seen that this is equivalent to assume that  $\mathbf{f}(\mathbf{y}, t)$  is at least Lipschitz continuous in  $D$  and that  $\mathbf{y}_0 \in D$ .

The initial value problem (1) can be equivalently written as

$$\mathbf{y}(t) = \mathbf{y}(0) + \int_0^t \mathbf{f}(\mathbf{y}(s), s) ds. \quad (2)$$

i.e., as an integral equation for  $\mathbf{y}(t)$ .

**Picard iteration method.** The Picard iteration method is rarely used in practice to compute numerical solutions to ODEs, but rather to prove existence and uniqueness of solutions to ODEs. Picard's method is essentially a fixed point iteration in which the solution  $\mathbf{y}(t)$  is approximated within a fixed time interval  $[0, t]$  by a *sequence of functions*

$$\mathbf{y}^{[0]}(t) \rightarrow \mathbf{y}^{[1]}(t) \rightarrow \cdots \rightarrow \mathbf{y}^{[k]}(t) \rightarrow \cdots \quad (3)$$

generated by the fixed iteration rule

$$\mathbf{y}^{[k+1]}(t) = \mathbf{y}_0 + \int_0^t \mathbf{f}(\mathbf{y}^{[k]}(s), s) ds. \quad (4)$$

In (4) we can set  $\mathbf{y}^{[0]}(t) = \mathbf{y}_0$  which, in principle, allows us to compute  $\mathbf{y}^{[1]}(t)$ . With  $\mathbf{y}^{[1]}(t)$  available we can compute  $\mathbf{y}^{[2]}(t)$ , and so on and so forth. For the sequence of functions  $\mathbf{y}^{[k]}(t)$  to converge from an arbitrary  $\mathbf{y}^{[0]}(t)$ , we need to make sure that the Fréchet derivative of the nonlinear operator at the right hand side of (4) has a norm smaller than one, i.e., that the operator is a contraction. This translates to an upper bound for  $t$ , which depends on the Lipschitz constant of  $\mathbf{f}$ . The larger the Lipschitz constant, the smaller  $t$ . In other words, Picard iterations converge only if  $t$  is smaller than some  $t_{max}$  that depends on the Lipschitz constant of  $\mathbf{f}$ . Clearly, once  $\mathbf{y}(t)$  has been computed to the desired accuracy we can use  $\mathbf{y}(t)$  as initial condition for the next sequence of Picard iterations, e.g., the sequence that yields the solution within the time interval  $[t, 2t]$ . The Picard iteration method is clearly not practical, as it involves the evaluation of a time integral at each iteration. Moreover, the function  $\mathbf{y}^{[k]}(t)$  that is being integrated is defined by an integral with  $t$  as one of the endpoints of the integral (see equation (4)).

Next, consider a partition of time interval  $[0, T]$  into an evenly-spaced set of grid points:

$$t_0 = 0, \quad t_N = T, \quad t_{i+1} = t_i + \Delta t \quad i = 0, \dots, N-1, \quad (5)$$

<sup>1</sup>Of course, if the initial value problem (1) does not admit a unique solution then its numerical approximation might just pick one of the possible solutions.

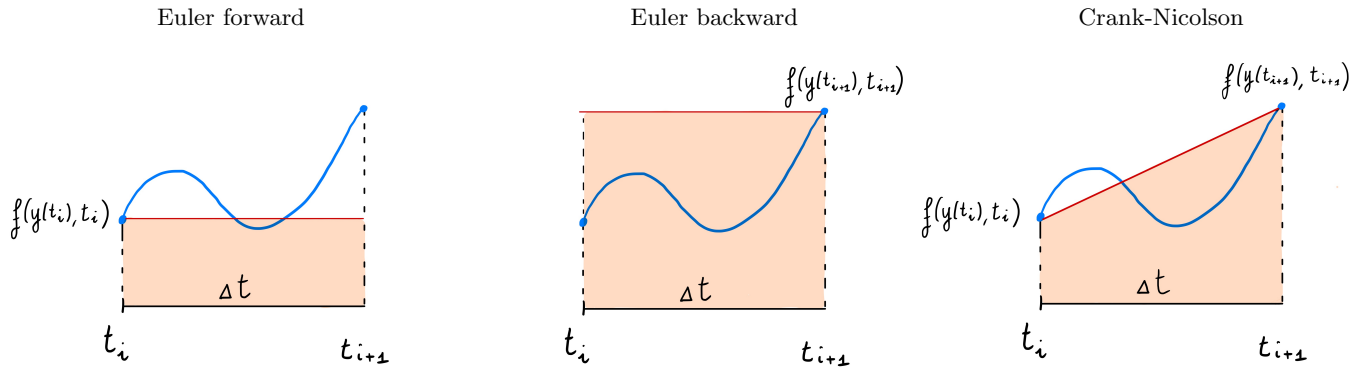


Figure 1: Approximations of the integral  $\int_{t_i}^{t_{i+1}} \mathbf{f}(\mathbf{y}(s), s) ds$  in (6) leading to well-known numerical schemes.

Figure 2: Show how you integrate rectangle and trapezoidal rule to obtain the three methods (Euler forward/backward and CN - 3 Figures).

where  $\Delta t = T/N$ .

By using the semi-group property of (1) we can write (2) within each time interval  $[t_k, t_{k+1}]$  as

$$\mathbf{y}(t_{k+1}) = \mathbf{y}(t_k) + \int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds. \quad (6)$$

This formulation is quite convenient for developing numerical methods based on *numerical quadrature formulas*, i.e., numerical approximations of the one-dimensional temporal integral appearing at the right hand side of (6).

**Euler and Crank-Nicolson methods** These are elementary methods obtained by approximating the integral at the right hand side of (6) by using the rectangle rule or the trapezoidal rule (see Figure 1). Specifically, consider the approximations

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}(\mathbf{y}(t_k), t_k), \quad (7)$$

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}), \quad (8)$$

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \frac{\Delta t}{2} (\mathbf{f}[\mathbf{y}(t_k), t_k] + \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1})). \quad (9)$$

These quadrature rules yield, respectively, the following numerical schemes

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k) \quad \text{Euler forward method (explicit),} \quad (10)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) \quad \text{Euler backward method (implicit),} \quad (11)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_k, t_k) + \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1})] \quad \text{Crank-Nicolson method (implicit).} \quad (12)$$

In these methods  $\mathbf{u}_k$  represents an approximation of the exact solution  $\mathbf{y}(t_k)$ , and we set  $\mathbf{u}_0 = \mathbf{y}_0$  (initial condition). Both Euler and Crank-Nicolson methods are *one-step methods*. This means that the approximate solution at time  $t_{k+1}$ , i.e.,  $\mathbf{u}_{k+1}$ , can be computed by knowing only the solution (or its approximation)

at time  $t_k$ . The Euler forward method allows us to compute  $\mathbf{u}_{k+1}$  explicitly, given  $\mathbf{u}_k$  and  $\mathbf{f}(\mathbf{u}_k, t_k)$ . On the other hand, the Euler backward and Crank-Nicolson methods are “implicit”. This is because the approximate solution at time  $t_{k+1}$ , i.e.,  $\mathbf{u}_{k+1}$ , cannot be computed explicitly based on  $\mathbf{u}_k$ , but it requires solving a nonlinear equation. Specifically, in the case of the Euler backward method we need to solve the nonlinear equation

$$\mathbf{u}_{k+1} = \mathbf{G}(\mathbf{u}_{k+1}) \quad \text{where} \quad \mathbf{G}(\mathbf{u}_{k+1}) = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (13)$$

Nonlinear equations of this form are essentially fixed point problems which can be solved numerically using iterative methods such as the Newton’s method (provided  $\mathbf{f}$  is of class  $C^1$ ). Upon definition of

$$\mathbf{F}(\mathbf{u}_{k+1}) = \mathbf{u}_{k+1} - \mathbf{G}(\mathbf{u}_{k+1}) \quad (14)$$

we can equivalently write (13) as

$$\mathbf{F}(\mathbf{u}_{k+1}) = \mathbf{0}. \quad (15)$$

*The Newton’s method for nonlinear systems:* The solution to the system of nonlinear equations (15) can be approximated by using the Newton’s method or any other method for root finding. As is well known, the Newton’s method generates a sequence of vectors  $\mathbf{u}_{k+1}^{[j]}$  ( $j = 0, 1, \dots$ ) converging to  $\mathbf{u}_{k+1}$  under rather mild assumptions (see [4, Ch. 7]). The Newton’s method can be formulated as<sup>2</sup>

$$\underbrace{\left[ \mathbf{I} - \mathbf{J}_{\mathbf{G}} \left( \mathbf{u}_{k+1}^{[j]} \right) \right]}_{\text{matrix}} \underbrace{\left( \mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right)}_{\text{vector}} = \underbrace{\mathbf{G} \left( \mathbf{u}_{k+1}^{[j]} \right) - \mathbf{u}_{k+1}^{[j]}}_{\text{vector}} \quad j = 0, 1, \dots, \quad (18)$$

where  $\mathbf{J}_{\mathbf{G}} \left( \mathbf{u}_{k+1}^{[j]} \right)$  is the Jacobian matrix of  $\mathbf{G}$  defined in equation (13), evaluated at  $\mathbf{u}_{k+1}^{[j]}$ . It is convenient to set the initial guess  $\mathbf{u}_{k+1}^{[0]}$  for Newton’s iteration as  $\mathbf{u}_{k+1}^{[0]} = \mathbf{u}_k$ , i.e., the numerical solution at previous time step. For  $\Delta t$  sufficiently small this guarantees that  $\mathbf{u}_{k+1}^{[0]}$  is within the basin of attraction of  $\mathbf{u}_{k+1}$ . Note also that for  $\Delta t$  sufficiently small the matrix at the left hand side of (18) is a perturbation of the identity (the norm of  $\mathbf{J}_{\mathbf{G}}$  goes to zero linearly in  $\Delta t$ ), and therefore  $\mathbf{I} - \mathbf{J}_{\mathbf{G}}$  is always invertible for sufficiently small  $\Delta t$ . Indeed  $\mathbf{I} - \mathbf{J}_{\mathbf{G}}$  is a diagonally dominant matrix for small  $\Delta t$ . More rigorously, we have the following convergence result (see [4, Theorem 7.1]).

**Theorem 1** (Convergence of Newton’s method). Let  $\mathbf{F}(\mathbf{x})$  in equation (14) be of class  $C^1$  in a convex open set  $D \subseteq \mathbb{R}^n$  that contains a zero of  $\mathbf{F}$ , i.e., a point  $\mathbf{x}^* \in D$  such that  $\mathbf{F}(\mathbf{x}^*) = \mathbf{0}$ . If  $\mathbf{J}_{\mathbf{F}}(\mathbf{x})$  is invertible at  $\mathbf{x}^*$  (it always is for sufficiently small  $\Delta t$ ), and  $\mathbf{J}_{\mathbf{F}}(\mathbf{x})$  is Lipschitz continuous in a neighborhood of  $\mathbf{x}^*$ , i.e.<sup>3</sup>,

$$\|\mathbf{J}_{\mathbf{F}}(\mathbf{x}) - \mathbf{J}_{\mathbf{F}}(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad (19)$$

then there exists a neighborhood of  $\mathbf{x}^*$  such that for any initial guess  $\mathbf{x}^{[0]}$  in such a neighborhood we have that the sequence  $\mathbf{x}^{[k]}$  generated by

$$\mathbf{J}_{\mathbf{F}} \left( \mathbf{x}^{[k]} \right) \left( \mathbf{x}^{[k+1]} - \mathbf{x}^{[k]} \right) = -\mathbf{F} \left( \mathbf{x}^{[k]} \right) \quad (20)$$

<sup>2</sup>Consider the Taylor series

$$\mathbf{F} \left( \mathbf{u}_{k+1}^{[j+1]} \right) = \mathbf{F} \left( \mathbf{u}_{k+1}^{[j]} \right) + \mathbf{J}_{\mathbf{F}} \left( \mathbf{u}_{k+1}^{[j]} \right) \left( \mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right) + \dots \quad j = 0, 1, \dots \quad (16)$$

Setting  $\mathbf{F} \left( \mathbf{u}_{k+1}^{[j+1]} \right) = \mathbf{0}$  yields

$$\mathbf{J}_{\mathbf{F}} \left( \mathbf{u}_{k+1}^{[j]} \right) \left( \mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right) = -\mathbf{F} \left( \mathbf{u}_{k+1}^{[j]} \right) \quad j = 0, 1, \dots \quad (17)$$

which, upon substitution of (14), coincides with the Newton method (18).

<sup>3</sup>In equation (19) the matrix norm  $\|\cdot\|$  at the left hand side is induced by the vector norm at the right hand side.

converges to  $\mathbf{x}^*$  with order 2. In other words, there exists  $m \in \mathbb{N}$  such that

$$\|\mathbf{x}^{[k+1]} - \mathbf{x}^*\| \leq C \|\mathbf{x}^{[k]} - \mathbf{x}^*\|^2 \quad \text{for all } k \geq m. \quad (21)$$

If we replace  $\mathbf{u}_{k+1}$  at the right hand side of (12) with one step of the Euler forward scheme (10) we obtain the *Heun method*

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_k, t_k) + \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k), t_{k+1})] \quad \text{Heun method (explicit)} \quad (22)$$

The Heun method is a one-step explicit method that belongs to the class of two-stage explicit Runge-Kutta methods.

*Remark:* The Euler methods (10)-(11) can be derived also by replacing  $d\mathbf{y}/dt$  in (1) with the first-order forward and backward finite-differentiation formulas

$$\frac{d\mathbf{y}(t_k)}{dt} \simeq \frac{\mathbf{y}(t_{k+1}) - \mathbf{y}(t_k)}{\Delta t} = \mathbf{f}(\mathbf{y}(t_k), t_k), \quad (23)$$

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{\mathbf{y}(t_{k+1}) - \mathbf{y}(t_k)}{\Delta t} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}). \quad (24)$$

**The midpoint method.** By approximating the integral at the right hand side of (6) using the midpoint rule yields

$$\int_{t_i}^{t_{i+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}\left(\mathbf{y}\left(t_i + \frac{\Delta t}{2}\right), t_i + \frac{\Delta t}{2}\right). \quad (25)$$

At this point, we approximate  $\mathbf{y}(t_i + \Delta t/2)$  using the Euler forward method to obtain

$$\mathbf{y}\left(t_i + \frac{\Delta t}{2}\right) \simeq \mathbf{y}(t_i) + \frac{\Delta t}{2} \mathbf{f}(\mathbf{y}(t_i), t_i) \quad (26)$$

to obtain

$$\mathbf{u}_{i+1} = \mathbf{u}_i + \Delta t \mathbf{f}\left(\mathbf{u}_i + \frac{\Delta t}{2} \mathbf{f}(\mathbf{u}_i, t_i), t_i + \frac{\Delta t}{2}\right) \quad \text{explicit midpoint method} \quad (27)$$

where  $\mathbf{u}_i$  is an approximation of  $\mathbf{y}(t_i)$ . The explicit midpoint method is a one-step method that belongs to the class of two-stage explicit Runge-Kutta methods.

Approximating  $\mathbf{y}(t_i + \Delta t/2)$  by the average

$$\mathbf{y}\left(t_i + \frac{\Delta t}{2}\right) \simeq \frac{\mathbf{y}(t_i) + \mathbf{y}(t_{i+1})}{2} \quad (28)$$

yields

$$\mathbf{u}_{i+1} = \mathbf{u}_i + \Delta t \mathbf{f}\left(\frac{\mathbf{u}_i + \mathbf{u}_{i+1}}{2}, t_i + \frac{\Delta t}{2}\right) \quad \text{implicit midpoint method} \quad (29)$$

The implicit midpoint method is a *one-step symplectic integrator*, i.e., the scheme preserves the Hamiltonian when applied to Hamiltonian dynamical systems, e.g., pendulum or double-pendulum equations. There is a vast literature on *structure-preserving* integration methods for ordinary differential equations (see, e.g., [1] and the references therein).

*Exercise:* Solve the pendulum equations with the implicit midpoint method and the Euler forward method. Verify that the Hamiltonian is preserved in the case of the implicit midpoint method is used.

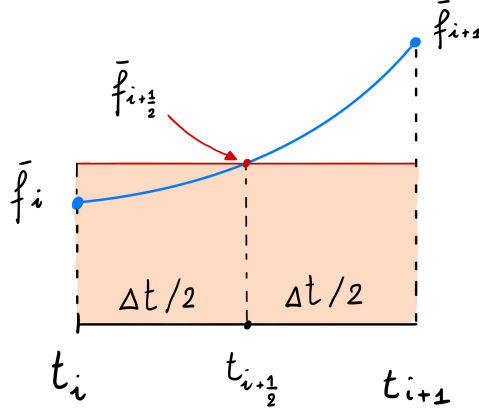


Figure 3: Approximation of the integral  $\int_{t_i}^{t_{i+1}} \mathbf{f}(\mathbf{y}(s), s) ds$  in (6) leading to midpoint method. In this figure we set  $\mathbf{f}_k = \mathbf{f}(\mathbf{y}(t_k), t_k)$ .

**Adams-Bashforth methods.** These are *explicit multistep methods* constructed by replacing the integral at the right hand side of (6) with the integral of a polynomial interpolant of  $\mathbf{f}(\mathbf{y}(s), s)$  at  $\{t_k, t_{k-1}, \dots, t_{k-q}\}$  which is then extrapolated into  $[t_k, t_{k+1}]$  to compute the integral. In other words, we introduce the following approximation

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_q \mathbf{f}(\mathbf{y}(s), s) ds \quad (30)$$

where  $\Pi_q \mathbf{f}(\mathbf{y}(s), s)$  is a polynomial of degree  $q$  interpolating

$$\{\mathbf{f}_k, \mathbf{f}_{k-1}, \dots, \mathbf{f}_{k-q}\} \quad \text{at} \quad \{t_k, t_{k-1}, \dots, t_{k-q}\} \quad (31)$$

where  $\mathbf{f}_k = \mathbf{f}(\mathbf{y}(t_k), t_k)$ . It is very convenient to use Lagrangian interpolation to derive the polynomial  $\Pi_q \mathbf{f}$ . Hereafter we derive the Adams-Bashforth (AB) methods for  $q = 0, 1, 2$ .

- *One-step Adams-Bashforth method (AB1):*

$$q = 0 \quad \Rightarrow \quad \Pi_0 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k \quad (32)$$

where  $\mathbf{f}_k = \mathbf{f}(\mathbf{y}(t_k), t_k)$ . Hence,

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_0 \mathbf{f}(\mathbf{y}(s), s) ds = \Delta t \mathbf{f}(\mathbf{y}(t_k), t_k). \quad (33)$$

This yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k) \quad \text{AB1 method.} \quad (34)$$

Note that the AB1 method coincides with the Euler forward method.

- *Two-step Adams-Bashforth method (AB2):*

$$q = 1 \quad \Rightarrow \quad \Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s) \quad (35)$$

where  $l_k(s)$  and  $l_{k-1}(s)$  are Lagrange characteristic polynomials

$$l_k(s) = \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad l_{k-1}(s) = \frac{s - t_k}{t_{k-1} - t_k}. \quad (36)$$

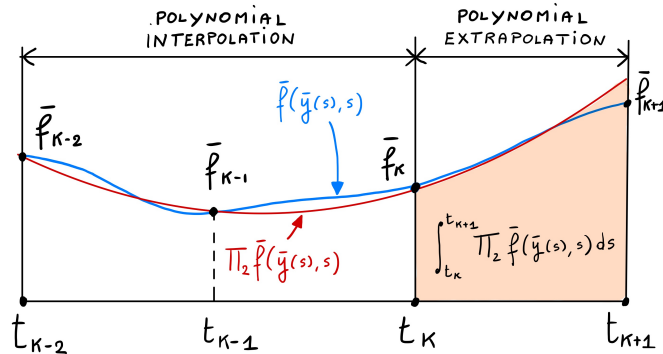


Figure 4: Derivation of the three-step Adams-Bashforth scheme (AB3). We first construct the polynomial  $\Pi_2 \mathbf{f}$  that interpolates  $\mathbf{f}(\mathbf{y}(s), s)$  at  $t_{k-2}$ ,  $t_{k-1}$  and  $t_k$ . Subsequently, we extrapolate  $\Pi_2 \mathbf{f}$  to  $[t_k, t_{k+1}]$  so that we can compute the integral in equation (30).

This yields

$$\Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k \frac{s - t_{k-1}}{t_k - t_{k-1}} + \mathbf{f}_{k-1} \frac{s - t_k}{t_{k-1} - t_k}. \quad (37)$$

The (linear) polynomial  $\Pi_1 \mathbf{f}(\mathbf{y}(s), s)$  is constructed in  $[t_{k-1}, t_k]$  and it can be integrated exactly in  $[t_k, t_{k+1}]$  using the trapezoidal rule. To this end, we notice that

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_k), t_k) = \mathbf{f}_k, \quad (38)$$

$$\begin{aligned} \Pi_1 \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) &= \mathbf{f}_k \frac{t_{k+1} - t_{k-1}}{t_k - t_{k-1}} + \mathbf{f}_{k-1} \frac{t_{k+1} - t_k}{t_{k-1} - t_k} \\ &= 2\mathbf{f}_k - \mathbf{f}_{k-1}. \end{aligned} \quad (39)$$

which gives us

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_1 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{2} (3\mathbf{f}_k - \mathbf{f}_{k-1}). \quad (40)$$

Substituting this back into (6) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [3\mathbf{f}(\mathbf{u}_k, t_k) - \mathbf{f}(\mathbf{u}_{k-1}, t_{k-1})] \quad \text{AB2 method} \quad (41)$$

- *Three-step Adams-Bashforth method (AB3)*: With reference to Figure (4)

$$q = 3 \quad \Rightarrow \quad \Pi_2 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s) + \mathbf{f}_{k-2} l_{k-2}(s) \quad (42)$$

where  $l_k(s)$ ,  $l_{k-1}(s)$  and  $l_{k-2}(s)$  are Lagrange characteristic polynomials

$$l_k(s) = \frac{s - t_{k-1}}{t_k - t_{k-1}} \frac{s - t_{k-2}}{t_k - t_{k-2}}, \quad (43)$$

$$l_{k-1}(s) = \frac{s - t_k}{t_{k-1} - t_k} \frac{s - t_{k-2}}{t_{k-1} - t_{k-2}}, \quad (44)$$

$$l_{k-2}(s) = \frac{s - t_k}{t_{k-2} - t_k} \frac{s - t_{k-1}}{t_{k-2} - t_{k-1}}. \quad (45)$$

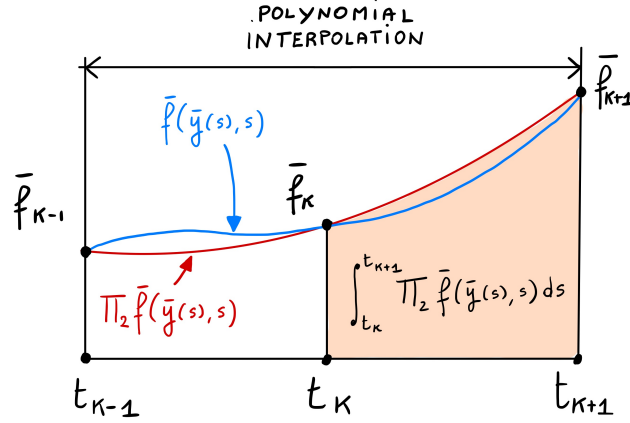


Figure 5: Derivation of the two-step Adams-Moulton scheme (AM2). We first construct the polynomial  $\Pi_2 \mathbf{f}$  that interpolates  $\mathbf{f}(\mathbf{y}(s), s)$  at  $t_{k-1}$ ,  $t_k$  and  $t_{k+1}$ . Subsequently, we use  $\Pi_2 \mathbf{f}$  polynomial to approximate the integral in equation (30).

By integrating  $\Pi_2 \mathbf{f}(\mathbf{y}(s), s)$  from  $t_k$  to  $t_{k+1}$  we obtain

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_2 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{12} (23\mathbf{f}_k - 16\mathbf{f}_{k-1} + 5\mathbf{f}_{k-2}). \quad (46)$$

Substituting this back into (6) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{12} [23\mathbf{f}(\mathbf{u}_k, t_k) - 16\mathbf{f}(\mathbf{u}_{k-1}, t_{k-1}) + 5\mathbf{f}(\mathbf{u}_{k-2}, t_{k-2})] \quad \text{AB3 method} \quad (47)$$

Higher-order Adams-Bashforth schemes can be obtained similarly. Note that in order to start-up a linear multistep scheme we need to compute the solution at the intermediate steps using different methods. For example, we could start-up the AB2 method with one step of the Heun method (to compute  $\mathbf{u}_1$ ) and then carry on the integration using the scheme (41)

**Adams-Moulton methods.** These are implicit multistep methods in which the time integral at the right hand-side of (6) is approximated by replacing  $\mathbf{f}(\mathbf{y}(s), s)$  by a polynomial  $\Pi_q \mathbf{f}(\mathbf{y}(s), s)$  of degree  $q$  interpolating

$$\{\mathbf{f}_{k+1}, \mathbf{f}_k, \dots, \mathbf{f}_{k-q+1}\} \quad \text{at} \quad \{t_{k+1}, t_{k-1}, \dots, t_{k-q+1}\} \quad (48)$$

The main difference with respect to the Adams-Bashforth method is that there is no extrapolation step, i.e., the point  $(t_{k+1}, \mathbf{f}_{k+1})$  is included in the interpolation (see Figure 5). Let us derive the Adams-Moulton schemes for  $q = 0, 1, 2$ .

- *One-step Adams-Moulton method (AM0):*

$$q = 0 \quad \Rightarrow \quad \Pi_0 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} \quad (49)$$

where  $\mathbf{f}_{k+1} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1})$ . Hence,

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_0 \mathbf{f}(\mathbf{y}(s), s) ds = \Delta t \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}). \quad (50)$$

This yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) \quad \text{AM0 method.} \quad (51)$$

Note that the AM0 method coincides with the Euler backward method.

- *One-step Adams-Moulton method (AM1):*

$$q = 1 \quad \Rightarrow \quad \Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} l_{k+1}(s) + \mathbf{f}_k l_k(s) \quad (52)$$

where  $l_k(s)$  and  $l_{k+1}(s)$  are Lagrange characteristic polynomials

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k}, \quad l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}}. \quad (53)$$

This yields

$$\Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} \frac{s - t_k}{t_{k+1} - t_k} + \mathbf{f}_k \frac{s - t_{k+1}}{t_k - t_{k+1}}. \quad (54)$$

The (linear) polynomial  $\Pi_1 \mathbf{f}(\mathbf{y}(s), s)$  is constructed in  $[t_k, t_{k+1}]$  and it can be integrated exactly in the same interval. To this end, we notice that

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_k), t_k) = \mathbf{f}_k, \quad (55)$$

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) = \mathbf{f}_{k+1} \quad (56)$$

which imply

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_1 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{2} (\mathbf{f}_{k+1} + \mathbf{f}_k). \quad (57)$$

Substituting this back into (6) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) + \mathbf{f}(\mathbf{u}_k, t_k)] \quad \text{AM1 method} \quad (58)$$

Hence, the AM1 scheme coincides with the Crank-Nicolson scheme.

- *Two-step Adams-Moulton method (AM2):*

$$q = 2 \quad \Rightarrow \quad \Pi_2 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} l_{k+1}(s) + \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s). \quad (59)$$

where  $l_{k+1}(s)$ ,  $l_k(s)$  and  $l_{k-1}(s)$  are Lagrange characteristic polynomials

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k} \frac{s - t_{k-1}}{t_{k+1} - t_{k-1}}, \quad (60)$$

$$l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad (61)$$

$$l_{k-1}(s) = \frac{s - t_{k+1}}{t_{k-1} - t_{k+1}} \frac{s - t_k}{t_{k-1} - t_k}. \quad (62)$$

By integrating  $\Pi_2 \mathbf{f}(\mathbf{y}(s), s)$  from  $t_k$  to  $t_{k+1}$  we obtain

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_2 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{12} (5\mathbf{f}_{k+1} + 8\mathbf{f}_k - \mathbf{f}_{k-1}). \quad (63)$$

Substituting this back into (6) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{12} [5\mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) + 8\mathbf{f}(\mathbf{u}_k, t_k) - \mathbf{f}(\mathbf{u}_{k-1}, t_{k-1})] \quad \text{AM2 method.} \quad (64)$$

**Backward differentiation formulas (BDF) methods.** These are linear implicit multistep methods that perform well for stiff problems. These methods are obtained by approximating  $d\mathbf{y}(t)/dt$  in (1) using a backward finite-difference formula. These formulas are obtained by interpolating  $\mathbf{y}(s)$  at  $\{t_{k+1}, t_k, \dots, t_{k-q+1}\}$  with a polynomial of degree  $q$ , differentiating the polynomial and evaluating the derivative at  $t = t_{k+1}$  (see Figure 6). Let us derive the first few BDF methods.



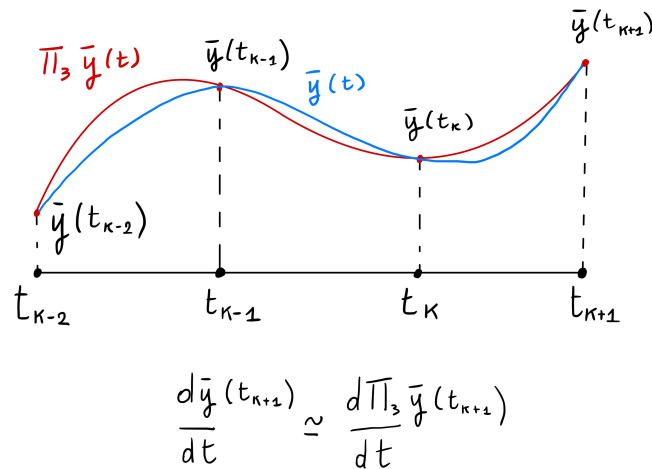


Figure 6: Derivation of the three-step backward differentiation formula (BDF3) method. We first construct the polynomial  $\Pi_3 \mathbf{f}$  that interpolates  $\mathbf{y}(t)$  at  $t_{k-2}$ ,  $t_{k-1}$ ,  $t_k$ , and  $t_{k+1}$ . Subsequently, approximate the derivative of  $\mathbf{y}(t)$  at  $t_{k+1}$  with the derivative of the polynomial  $\Pi_3 \mathbf{y}(t)$  at  $t_{k+1}$ .

- *One-step BDF method (BDF1):*

$$\Pi_1 \mathbf{y}(s) = \mathbf{y}_{k+1} l_{k+1}(s) + \mathbf{y}_k l_k(s) \quad (65)$$

where the Lagrange characteristic polynomials are given by

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k}, \quad l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \quad (66)$$

We approximate the derivative of  $\mathbf{y}(s)$  at  $t_{k+1}$  with the derivative of the polynomial  $\Pi_1 \mathbf{y}(s)$  at  $t_{k+1}$ , i.e.,

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{d\Pi_1 \mathbf{y}(t_{k+1})}{dt} = \frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{\Delta t}. \quad (67)$$

Substituting this approximation into the exact equation

$$\frac{d\mathbf{y}(t_{k+1})}{dt} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) \quad (68)$$

yields the scheme

$$\mathbf{u}_{k+1} - \mathbf{u}_k = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (69)$$

Note that BDF1 coincides with the Euler backward scheme.

- *Two-step BDF method (BDF2):*

$$\Pi_2 \mathbf{y}(s) = \mathbf{y}_{k+1} l_{k+1}(s) + \mathbf{y}_k l_k(s) + \mathbf{y}_{k-1} l_{k-1}(s) \quad (70)$$

where the Lagrange characteristic polynomials are given by

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k} \frac{s - t_{k-1}}{t_{k+1} - t_{k-1}}, \quad (71)$$

$$l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad (72)$$

$$l_{k-1}(s) = \frac{s - t_{k+1}}{t_{k-1} - t_{k+1}} \frac{s - t_k}{t_{k-1} - t_k}. \quad (73)$$

We approximate the derivative of  $\mathbf{y}(s)$  at  $t_{k+1}$  with the derivative of the polynomial  $\Pi_1\mathbf{y}(s)$  at  $t_{k+1}$ , i.e.,

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{d\Pi_2\mathbf{y}(t_{k+1})}{dt} = \frac{3\mathbf{y}_{k+1} - 4\mathbf{y}_k + \mathbf{y}_{k-1}}{2\Delta t}. \quad (74)$$

Substituting this approximation into the exact equation

$$\frac{d\mathbf{y}(t_{k+1})}{dt} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) \quad (75)$$

yields the scheme

$$\frac{3}{2}\mathbf{u}_{k+1} - 2\mathbf{u}_k - \frac{1}{2}\mathbf{u}_{k-1} = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (76)$$

- *Three-step BDF method (BDF3)*: By following a similar procedure as in BDF2 it is straightforward to show that

$$\frac{11}{6}\mathbf{u}_{k+1} - 3\mathbf{u}_k + \frac{3}{2}\mathbf{u}_{k-1} - \frac{1}{3}\mathbf{u}_{k-2} = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (77)$$

**General form of a linear multistep method (LMM).** The general form of a linear multistep method is

$$\sum_{j=0}^q \alpha_j \mathbf{u}_{k+j} = \Delta t \sum_{j=0}^q \beta_j \underbrace{\mathbf{f}(\mathbf{u}_{k+j}, t_{k+j})}_{\mathbf{f}_{k+j}}. \quad (78)$$

Note that Adams-Bashforth, Adams-Moulton and BDF methods are all in the form (78). To avoid non-uniqueness of coefficients due to rescaling we set  $\alpha_q = 1$ . Clearly, if  $\beta_q = 0$  and then the method is *explicit*. On the other hand, if  $\beta_q \neq 0$  then the method is *implicit*. Let us provide a few examples.

*Example:* The AB3 method

$$\mathbf{u}_{k+3} = \mathbf{u}_{k+2} + \frac{\Delta t}{12} (23\mathbf{f}_{k+2} - 16\mathbf{f}_{k+1} + 5\mathbf{f}_k) \quad (79)$$

can be written in the form (78) by setting

$$\begin{aligned} \alpha_3 &= 1 & \alpha_2 &= -1, & \alpha_1 &= 0, & \alpha_0 &= 0, \\ \beta_3 &= 0, & \beta_2 &= \frac{23}{12}, & \beta_1 &= -\frac{16}{12}, & \beta_0 &= \frac{5}{12}. \end{aligned}$$

Note that  $(\alpha_3, \beta_3) = (1, 0)$  (the method is explicit) and

$$\sum_{j=0}^3 \beta_j = 1. \quad (80)$$

*Example:* The BDF2 method

$$\mathbf{u}_{k+2} - \frac{4}{3}\mathbf{u}_{k+1} + \frac{1}{3}\mathbf{u}_k = \frac{2}{3}\Delta t \mathbf{f}_{k+2} \quad (81)$$

can be written in the form (78) by setting

$$\alpha_2 = 1, \quad \alpha_1 = -\frac{4}{3}, \quad \alpha_0 = \frac{1}{3}, \quad \beta_2 = \frac{2}{3}, \quad \beta_1 = 0, \quad \beta_0 = 0. \quad (82)$$

Note that  $(\alpha_2, \beta_2) = (1, 1/3)$  (the method is implicit) and

$$\sum_{j=0}^2 \alpha_j = 0, \quad \sum_{j=0}^2 (\beta_j - j\alpha_j) = 0. \quad (83)$$

As we will see the conditions (80) and (83) guarantee that AB3 and BDF2 are *consistent* methods, i.e., that the truncation error of these methods goes to zero as we send  $\Delta t$  to zero. Roughly speaking this means that the numerical schemes (79) and (81) converge to the ODE (1) as  $\Delta t \rightarrow 0$ . This is necessary but not sufficient for the numerical solution generated by the scheme to converge to the analytical solution. The other element that is needed for convergence is *zero-stability*. The consistency conditions for a general linear  $q$ -step method are

$$\sum_{j=0}^q \alpha_j = 0, \quad \sum_{j=0}^q (\beta_j - j\alpha_j) = 0. \quad (84)$$

*Remark:* All linear multistep methods rely on a polynomial interpolation process on an evenly-spaced temporal grid. As is well-known, polynomial interpolation on evenly spaced grids is, in general, ill-conditioned and can undergo a severe Gibbs phenomenon depending on the function. However, the process of interpolating a function with a polynomial of degree  $q + 1$  within a very small a time interval (equal to  $q\Delta t$  for a  $q$ -step method) is not ill-conditioned. The reason can be traced back to the fact that we are interpolating on a small time interval  $[t - q\Delta t, t]$  in which the function behaves more or less almost like a line. More rigorously, if  $g(t)$  is any function of class  $C^{q+1}$  in  $[t - q\Delta t, t]$  and  $\Pi_q g(t)$  is a polynomial of degree  $q$  that interpolates  $g(t)$  at  $\{t, t - \Delta t, \dots, t - q\Delta t\}$  then we have the error estimate

$$|g(t) - \Pi_q g(t)| = \frac{1}{(q+1)!} \left| \frac{d^{q+1}f(\xi)}{dt^{q+1}} \right| \prod_{j=0}^q |(t - t_j)| \leq \frac{(q\Delta t)^{q+1}}{(q+1)!} \left| \frac{d^{q+1}f(\xi)}{dt^{q+1}} \right| \quad \xi \in [t - q\Delta t, t], \quad (85)$$

which clearly goes to zero for sufficiently small  $\Delta t$ .

**Runge-Kutta methods.** Runge-Kutta (RK) methods are one-step methods (implicit or explicit) that aim at increasing accuracy by increasing the number of function evaluations within each time step. The general form of a RK method with  $s$  stages is

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \sum_{i=1}^s b_i \mathbf{K}_i \quad (86)$$

where

$$\mathbf{K}_i = \mathbf{f} \left( \mathbf{u}_k + \Delta t \sum_{j=1}^s a_{ij} \mathbf{K}_j, t_k + c_i \Delta t \right). \quad (87)$$

The coefficients of the RK method are usually collected in a table called *Butcher array*

$c_1$	$a_{11}$	$a_{12}$	$\cdots$	$a_{1s}$
$c_2$	$a_{21}$	$a_{22}$	$\cdots$	$a_{2s}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$c_s$	$a_{s1}$	$a_{s2}$	$\cdots$	$a_{ss}$
	$b_1$	$b_2$	$\cdots$	$b_s$

The elements  $a_{ij}$  in the Butcher table can be positive or negative. Consistency of RK methods, i.e., the fact that (86) converge to (1) as  $\Delta t \rightarrow 0$  implies the following conditions

$$\sum_{j=1}^s b_j = 1. \quad (88)$$

Moreover, we assume that

$$c_i = \sum_{j=1}^s a_{ij}. \quad (89)$$

Such a “row-sum condition” is not needed for a consistent method.

If  $a_{ij} = 0$  for  $i \geq j$  then each  $\mathbf{K}_i$  can be computed recursively from the previous ones and the RK method is *explicit*. Otherwise the RK method is implicit. Let us provide a few examples.

- *Euler forward method*: The Euler forward method can be seen as a one-step explicit RK method. The Butcher array corresponding to such explicit RK1 method is

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

- *Heun method*: The Heun method (22) is a two-stage explicit Runge-Kutta method. The Butcher array for the Heun method is:

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

This table corresponds to the following RK2 method

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), \quad (90)$$

where

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (91)$$

$$\mathbf{K}_2 = \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k), t_k + \Delta t). \quad (92)$$

- *Crank-Nicolson method*: The Crank-Nicolson (CN) method (12) is a two-stage implicit Runge-Kutta method. The Butcher array corresponding to such method is:

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

This table corresponds to the following RK2 method

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), \quad (93)$$

where

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (94)$$

$$\mathbf{K}_2 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), t_{k+1}\right). \quad (95)$$

From equation (93) we see that

$$\frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2) = \mathbf{u}_{k+1} - \mathbf{u}_k. \quad (96)$$

Substituting this expression into (95) yields

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad \mathbf{K}_2 = \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (97)$$

- *Kutta's method*: This method is an explicit 3-stage method corresponding to the Butcher array:

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 \\ 1 & -1 & 2 & 0 \\ \hline & 1/6 & 2/3 & 1/6 \end{array}$$

The method can be written as

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (98)$$

$$\mathbf{K}_2 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{K}_1, t_k + \frac{\Delta t}{2}\right), \quad (99)$$

$$\mathbf{K}_3 = \mathbf{f}(\mathbf{u}_k - \Delta t \mathbf{K}_1 + 2\Delta t \mathbf{K}_2, t_k + \Delta t), \quad (100)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{6} (\mathbf{K}_1 + 4\mathbf{K}_2 + \mathbf{K}_3). \quad (101)$$

- *Runge-Kutta method (RK4)*: The most famous RK method is perhaps the one proposed in the original paper by Runge and Kutta. Such a method is an explicit 4-stage method corresponding to the Butcher array:

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 0 & 1/2 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & 1/6 & 1/3 & 1/3 & 1/6 \end{array}$$

The method can be written as

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (102)$$

$$\mathbf{K}_2 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{K}_1, t_k + \frac{\Delta t}{2}\right), \quad (103)$$

$$\mathbf{K}_3 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{K}_2, t_k + \frac{\Delta t}{2}\right), \quad (104)$$

$$\mathbf{K}_4 = \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{K}_3, t_k + \Delta t), \quad (105)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{6} (\mathbf{K}_1 + 2\mathbf{K}_2 + 2\mathbf{K}_3 + \mathbf{K}_4). \quad (106)$$

$$\begin{array}{c|cc} 0 & 0 & 0 \\ c & c & 0 \\ \hline & b_1 & b_2 \end{array}$$

**Derivation of explicit RK2 methods.** Let us show how to derive an arbitrary explicit two-stage RK method. To this end we first write the Butcher array:

where  $b_1 + b_2 = 1$ . The RK2 method corresponding to this table can be written explicitly as

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t(b_1 \mathbf{K}_1 + b_2 \mathbf{K}_2) \quad (107)$$

where

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k) \quad (108)$$

$$\mathbf{K}_2 = \mathbf{f}(\mathbf{u}_k + c\Delta t \mathbf{K}_1, t_k + c\Delta t). \quad (109)$$

Expand  $\mathbf{K}_2$  in a Taylor series in  $\Delta t$  to obtain

$$\mathbf{K}_2 = \mathbf{K}_1 + c\Delta t \left( \sum_{j=1}^n \frac{\partial \mathbf{f}}{\partial y_j} K_{1j} + \frac{\partial \mathbf{f}}{\partial t} \right) + \dots \quad (110)$$

A substitution of this expression into (107) yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t(b_1 + b_2)\mathbf{f}(\mathbf{u}_k, t_k) + b_2 c (\Delta t)^2 \left( \sum_{j=1}^n \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial y_j} f_j(\mathbf{u}_k, t_k) + \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial t} \right). \quad (111)$$

Next, we expand the solution to the ODE (1) in a Taylor series at time  $t_{k+1} = t_k + \Delta t$ , assuming that  $\mathbf{y}(t_k) = \mathbf{u}_k$ . This yields

$$\mathbf{y}(t_{k+1}) = \mathbf{u}(t_k) + \frac{d\mathbf{y}(t_k)}{dt} \Delta t + \frac{(\Delta t)^2}{2} \frac{d^2 \mathbf{y}(t_k)}{dt^2} + \dots \quad (112)$$

By using (1) and the chain rule we obtain:

$$\frac{d\mathbf{y}(t_k)}{dt} = \mathbf{f}(\mathbf{y}_k, t_k), \quad \frac{d^2 \mathbf{y}(t_k)}{dt^2} = \sum_{j=1}^n \frac{\partial \mathbf{f}(\mathbf{y}_k, t_k)}{\partial y_j} f_j(\mathbf{y}_k, t_k) + \frac{\partial \mathbf{f}(\mathbf{y}_k, t_k)}{\partial t}. \quad (113)$$

Assuming that  $\mathbf{y}_k = \mathbf{u}_k$  and matching the terms multiplying the same powers of  $\Delta t$  in (111) and (112) we obtain<sup>4</sup>

$$b_1 + b_2 = 1 \quad \text{and} \quad cb_2 = \frac{1}{2}. \quad (114)$$

This is a system of 2 equations in 3 unknowns. Thus, there is an *infinite number* of explicit (and consistent) RK2 methods. For example, setting

$$c = 1, \quad b_1 = b_2 = \frac{1}{2} \quad \text{yields the Heun method (22)}. \quad (115)$$

On the other hand, setting

$$c = \frac{1}{2}, \quad b_1 = 0 \quad b_2 = 1 \quad \text{yields the explicit midpoint method (27)}. \quad (116)$$

<sup>4</sup>The condition  $b_1 + b_2 = 1$  is a consistency condition which guarantees that the scheme (107) converges to the ODE (1) in the limit  $\Delta t \rightarrow 0$ .

By using the Taylor series approach discussed in this section, it is possible to derive conditions on the entries of the Butcher array for RK methods with an arbitrary number of stages. Essentially, we perform a Taylor series of the RK method in  $\Delta t$  and then match it with the Taylor series expansion of the solution up to a given order. The corresponding equations for the coefficients, e.g., (114) are called “stage-order” conditions. This is discussed, e.g., in [2, §5.9]. For example, it can be shown that for a three stage RK method

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ c_2 & a_{21} & 0 & 0 \\ c_3 & a_{31} & a_{32} & 0 \\ \hline & b_1 & b_2 & b_3 \end{array}$$

we obtain the order conditions

$$\begin{cases} b_1 + b_2 + b_3 = 1 \\ b_2 c_2 + b_3 c_3 = \frac{1}{2} \\ b_2 c_2^2 + b_3 c_3^2 = \frac{1}{3} \\ b_3 a_{32} c_2 = \frac{1}{6} \end{cases} \quad (117)$$

which yields to one two-parameter family of solutions and two one-parameter families of solutions [2, p. 178]. As easily seen, the Taylor series approach rapidly become intractable as the number of stages increases, and so does the corresponding set of order conditions. Fortunately, there is a more efficient way to derive conditions such as (114) by using rooted trees (see [2, §5.6]).

**Derivation of implicit RK methods.** Implicit RK methods can be derived from the integral formulation (6) of the Cauchy problem. In fact, if a Gauss quadrature formula with  $s$  nodes in  $[t_k, t_{k+1}]$  is employed to approximate the integral at the right hand side of (6), we obtain

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \sum_{j=1}^s b_j \mathbf{f}(\mathbf{y}(t_k + c_j \Delta t), t_k + c_j \Delta t). \quad (118)$$

Here, we denoted by  $b_j$  the quadrature weights and by  $t_k + c_j \Delta t$  the quadrature nodes. It can be proved that for any RK formula (86)-(87), there exists a correspondence between the coefficients  $b_j$ ,  $c_j$  of the formula and the weights and nodes of a Gauss quadrature rule (see, [2, §5.11] for details). For instance, the implicit midpoint method can be seen as a one step implicit RK method based on a 1 point Gauss-Legendre quadrature rule. This corresponds to the Butcher array

$$\begin{array}{c|c} 1/2 & 1/2 \\ \hline & 1 \end{array}$$

and can attain order 2. In general, an implicit  $s$ -stage Gauss RK method can achieve order  $2s$ . Similarly, Gauss-Radau and Gauss-Lobatto RK methods can achieve order  $2s - 1$  and  $2s - 2$ , respectively.

## References

- [1] E. Hairer, C. Lubich, and G. Wanner. *Geometric numerical integration: structure-preserving algorithms for ordinary differential equations*. Springer, 2006.
- [2] J. D. Lambert. *Numerical methods for ordinary differential systems: the initial value problem*. Wiley, 1991.
- [3] R. LeVeque. *Finite difference methods for ordinary and partial differential equations*. SIAM, 2007.
- [4] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Springer, 2007.