

Overview of numerical methods for ODEs

In this course note we provide a brief overview of common numerical methods to approximate the solution of the initial value problem (IVP)

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

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. We have seen in AM 214 that if $\mathbf{f}(\mathbf{y}, t)$ is Lipschitz continuous in an open set $D \subseteq \mathbb{R}^n$ then the IVP (1) well-posed in D , i.e., the solution exists and is unique for all $\mathbf{y}_0 \in D$, at least for some time $\tau > 0$. The initial value problem (1) can be 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 some time interval $[0, t_1]$ by the *sequence of functions*

$$\mathbf{y}^{[0]}(t) \rightarrow \mathbf{y}^{[1]}(t) \rightarrow \cdots \rightarrow \mathbf{y}^{[k]}(t) \rightarrow \cdots \quad t \in [0, t_1] \quad (3)$$

generated by equation (2), i.e.,

$$\mathbf{y}^{[k+1]}(t) = \mathbf{y}_0 + \int_0^t \mathbf{f}(\mathbf{y}^{[k]}(s), s) ds, \quad k = 0, 1, 2, \dots, \quad t \in [0, t_1]. \quad (4)$$

Setting $\mathbf{y}^{[0]}(t) = \mathbf{y}_0$ in (4) yields $\mathbf{y}^{[1]}(t)$. With $\mathbf{y}^{[1]}(t)$ available we can compute $\mathbf{y}^{[2]}(t)$, etc. A necessary and sufficient condition for convergence of the sequence of functions $\mathbf{y}^{[k]}(t)$ is that the nonlinear operator at the right-hand-side of (4), i.e.,

$$\mathbf{N}(\mathbf{y})(t) = \mathbf{y}_0 + \int_0^t \mathbf{f}(\mathbf{y}(s), s) ds \quad (5)$$

is a *contraction*¹ (Banach fixed point theorem). This statement essentially sets an upper bound on t_1 (integration time) that depends on the Lipschitz constant of \mathbf{f} . The larger the Lipschitz constant, the smaller t_1 . In other words, Picard iterations converge if and only if t_1 is chosen small enough, where “how small” depends on the Lipschitz constant of \mathbf{f} . Clearly, once $\mathbf{y}(t)$ has been computed to the desired accuracy in $[0, t_1]$, we can then use $\mathbf{y}(t_1)$ as initial condition for the next Picard iterations, i.e.,

$$\mathbf{y}^{[k+1]}(t) = \mathbf{y}(t_1) + \int_{t_1}^t \mathbf{f}(\mathbf{y}^{[k]}(s), s) ds \quad t \in [t_1, t_2]. \quad (8)$$

The Picard iteration method is not practical as it involves the evaluation of a integral at each iteration k . Moreover, the integral has “ t ” as one of the endpoints.

¹Let X be a function space with norm $\|\cdot\|_X$. An operator $\mathbf{N} : X \mapsto X$ is called a (strict) *contraction* if

$$\|\mathbf{N}(\mathbf{y}^{[k]}) - \mathbf{N}(\mathbf{y}^{[j]})\|_X < \|\mathbf{y}^{[k]} - \mathbf{y}^{[j]}\|_X \quad (6)$$

In the case of (5) we have

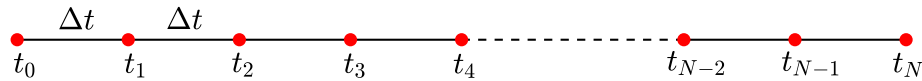
$$\|\mathbf{N}(\mathbf{y}^{[k]}) - \mathbf{N}(\mathbf{y}^{[j]})\|_X \leq \left\| \int_0^{t_1} [\mathbf{f}(\mathbf{y}^{[k]}(s), s) - \mathbf{f}(\mathbf{y}^{[j]}(s), s)] ds \right\|_X. \quad (7)$$

Elementary numerical methods

Let us partition the time interval $[0, T]$ into an evenly-spaced grid with N time instants

$$t_0 = 0, \quad t_N = T, \quad t_{i+1} = t_i + \Delta t \quad i = 0, \dots, N-1, \quad (9)$$

where $\Delta t = T/N$.



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

$$\mathbf{y}(t_{k+1}) = \mathbf{y}(t_k) + \int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds. \quad (10)$$

This formulation is quite convenient for developing numerical methods for ODEs based on *quadrature formulas*, i.e., numerical approximations of the one-dimensional temporal integral appearing at the right hand side of (10).

Euler and Crank-Nicolson methods

These are basic numerical schemes that can be obtained by approximating the integral at the right hand side of equation (10) by using the rectangle rule or the trapezoidal rule (see Figure 1). Specifically,

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}(\mathbf{y}(t_k), t_k), \quad (11)$$

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}), \quad (12)$$

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \frac{\Delta t}{2} [\mathbf{f}(\mathbf{y}(t_k), t_k) + \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1})]. \quad (13)$$

Note that all these approximations are obtained via interpolatory quadratures, i.e., by replacing $\mathbf{f}(\mathbf{y}(s), s)$ with a polynomial in $[t_k, t_{k+1}]$ and then integrating the polynomial. In particular, in (11) and (12) we use constant polynomials while in (13) we use a linear polynomial. A substitution of (11)-(13) into (10) yields, respectively, the following numerical schemes

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k) \quad \text{Euler forward method (explicit),} \quad (14)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) \quad \text{Euler backward method (implicit),} \quad (15)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_k, t_k) + \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1})] \quad \text{Crank-Nicolson method (implicit).} \quad (16)$$

In these methods \mathbf{u}_k represents an approximation of the exact solution $\mathbf{y}(t_k)$. In Figure 2 we sketch the difference between the analytical solution of (2) evaluated on a temporal grid, i.e., $\{\mathbf{y}(t_0), \mathbf{y}(t_1), \dots\}$ and the numerical solution \mathbf{u}_k obtained by iterating any of the schemes (14)-(16) (as a matter of fact any numerical scheme).

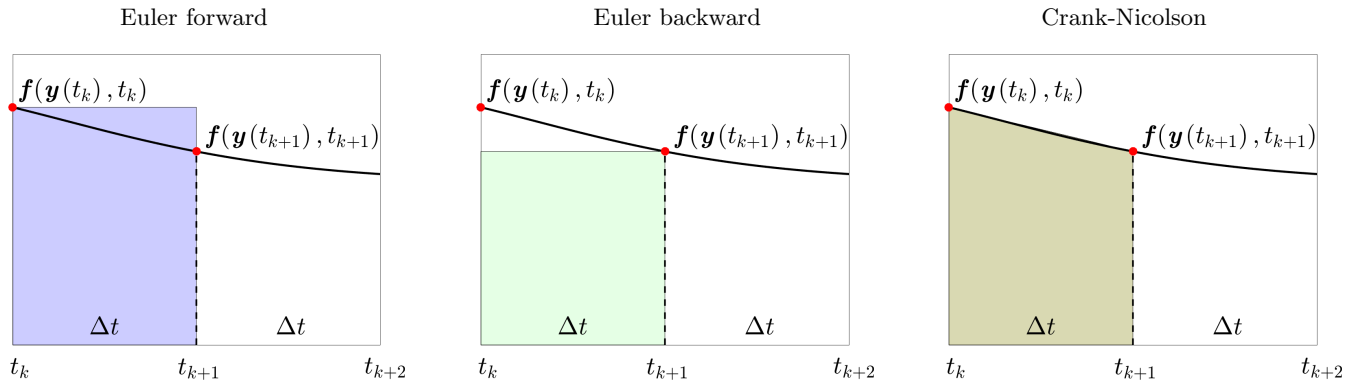


Figure 1: Elementary approximation of the integral $\int_{t_i}^{t_{i+1}} \mathbf{f}(\mathbf{y}(s), s) ds$ in equation (10) leading to Euler forward, Euler backward, and Crank-Nicolson methods.

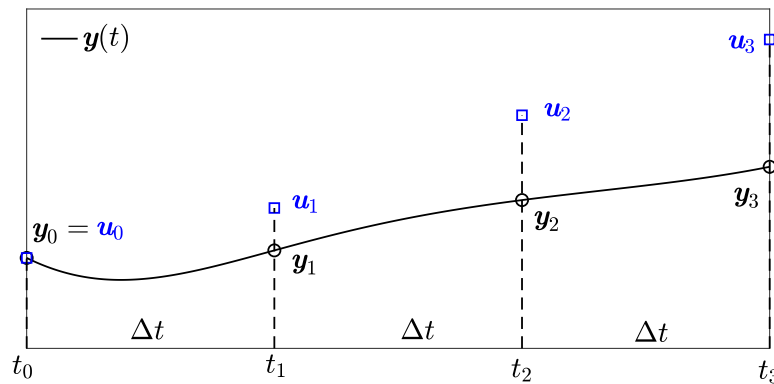


Figure 2: Analytical solution $\mathbf{y}(t)$ of the ODE (2) evaluated on a temporal grid $\{t_0, t_1, \dots\}$, i.e., $\mathbf{y}_i = \mathbf{y}(t_i)$, and numerical solution \mathbf{u}_i obtained by iterating any of the schemes (14)-(16). Here we assume that the initial condition used for the numerical scheme is exact, i.e. we set $\mathbf{u}_0 = \mathbf{y}_0$.

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 based on 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 are *implicit* methods. 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 (algebraic) equation. Specifically, for the Euler backward method we need to solve

$$\mathbf{u}_{k+1} = \mathbf{G}(\mathbf{u}_{k+1}) \quad \text{where} \quad \mathbf{G}(\mathbf{u}_{k+1}) = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \tag{17}$$

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

$$\mathbf{F}(\mathbf{u}_{k+1}) = \mathbf{u}_{k+1} - \mathbf{G}(\mathbf{u}_{k+1}) \tag{18}$$

we can write (17) as

$$\mathbf{F}(\mathbf{u}_{k+1}) = \mathbf{0}, \quad \Leftrightarrow \quad \mathbf{u}_{k+1} \text{ is a zero of } \mathbf{F}. \tag{19}$$

Equivalently, we can say that \mathbf{u}_{k+1} is a fixed point of the map \mathbf{G} defined in (17).

Newton’s method for nonlinear systems. To solve the nonlinear system of equations (19) one can use Newton’s method or any other iterative root finding method. 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 Theorem 1 below, and [4, Ch. 7]). The Newton's method for nonlinear systems can be written as²

$$\underbrace{\left[\mathbf{I} - \mathbf{J}_{\mathbf{G}} \left(\mathbf{u}_{k+1}^{[j]} \right) \right]}_{\text{matrix}} \underbrace{\left(\mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right)}_{\text{vector}} = \underbrace{\mathbf{G} \left(\mathbf{u}_{k+1}^{[j]} \right) - \mathbf{u}_{k+1}^{[j]}}_{\text{vector}} \quad j = 0, 1, \dots, \quad (22)$$

where $\mathbf{J}_{\mathbf{G}} \left(\mathbf{u}_{k+1}^{[j]} \right)$ is the Jacobian matrix of \mathbf{G} defined in equation (17), evaluated at $\mathbf{u}_{k+1}^{[j]}$. Equation (22) allows us to compute $\mathbf{u}_{k+1}^{[j+1]}$ given $\mathbf{u}_{k+1}^{[j]}$ by solving a linear system. In practice, 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 Δt sufficiently small the matrix at the left hand side of (22) is a perturbation of the identity (the norm of $\mathbf{J}_{\mathbf{G}}$ goes to zero linearly in Δt), and therefore $\mathbf{I} - \mathbf{J}_{\mathbf{G}}$ is always invertible for sufficiently small Δt . Indeed, $\mathbf{I} - \mathbf{J}_{\mathbf{G}}$ is diagonally dominant for small Δt . 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 (18) 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 Δt), and $\mathbf{J}_{\mathbf{F}}(\mathbf{x})$ is Lipschitz continuous in a neighborhood of \mathbf{x}^* , i.e.,

$$\|\mathbf{J}_{\mathbf{F}}(\mathbf{x}) - \mathbf{J}_{\mathbf{F}}(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad (23)$$

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

$$\mathbf{J}_{\mathbf{F}} \left(\mathbf{x}^{[k]} \right) \left(\mathbf{x}^{[k+1]} - \mathbf{x}^{[k]} \right) = -\mathbf{F} \left(\mathbf{x}^{[k]} \right) \quad (24)$$

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

$$\left\| \mathbf{x}^{[k+1]} - \mathbf{x}^* \right\| \leq C \left\| \mathbf{x}^{[k]} - \mathbf{x}^* \right\|^2 \quad \text{for all } k \geq m. \quad (25)$$

In equation (23) the norm at the left-hand-side is any matrix norm compatible with the vector norm at the right hand side (see Appendix A).

Remark: We emphasize that explicit and implicit Euler methods (14)-(15) can be also derived by replacing dy/dt in (1) with the first-order forward/backward finite-differentiation approximation

$$\frac{d\mathbf{y}(t_k)}{dt} \simeq \frac{\mathbf{y}(t_{k+1}) - \mathbf{y}(t_k)}{\Delta t} = \mathbf{f}(\mathbf{y}(t_k), t_k), \quad (26)$$

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{\mathbf{y}(t_{k+1}) - \mathbf{y}(t_k)}{\Delta t} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}). \quad (27)$$

²To derive the Newton's method for nonlinear systems of algebraic equations, consider the Taylor series

$$\mathbf{F} \left(\mathbf{u}_{k+1}^{[j+1]} \right) = \mathbf{F} \left(\mathbf{u}_{k+1}^{[j]} \right) + \mathbf{J}_{\mathbf{F}} \left(\mathbf{u}_{k+1}^{[j]} \right) \left(\mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right) + \dots \quad j = 0, 1, \dots \quad (20)$$

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

$$\mathbf{J}_{\mathbf{F}} \left(\mathbf{u}_{k+1}^{[j]} \right) \left(\mathbf{u}_{k+1}^{[j+1]} - \mathbf{u}_{k+1}^{[j]} \right) = -\mathbf{F} \left(\mathbf{u}_{k+1}^{[j]} \right) \quad j = 0, 1, \dots \quad (21)$$

which, upon substitution of (18), coincides with the Newton method (22).

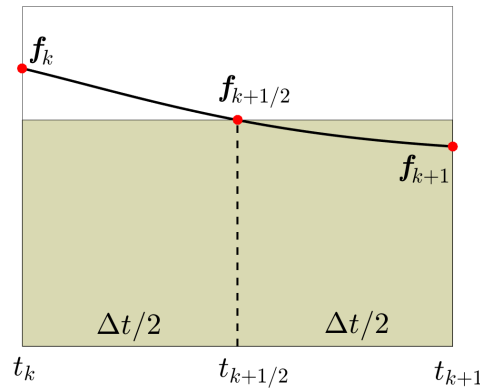


Figure 3: Approximation of the integral $\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds$ in (10) leading to midpoint method. Here we set $\mathbf{f}_k = \mathbf{f}(\mathbf{y}(t_k), t_k)$.

Heun's method

If we replace \mathbf{u}_{k+1} at the right hand side of (16) with one step of the Euler forward scheme (14) we obtain the *Heun method*

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_k, t_k) + \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k), t_{k+1})] \quad \text{Heun method (explicit)} \quad (28)$$

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

Midpoint method

We approximate the integral at the right hand side of (10) using the midpoint rule yields

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \mathbf{f}\left(\mathbf{y}\left(t_k + \frac{\Delta t}{2}, t_k + \frac{\Delta t}{2}\right)\right). \quad (29)$$

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

$$\mathbf{y}\left(t_k + \frac{\Delta t}{2}\right) \simeq \mathbf{y}(t_k) + \frac{\Delta t}{2} \mathbf{f}(\mathbf{y}(t_k), t_k) \quad (30)$$

to obtain

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{f}(\mathbf{u}_k, t_k), t_k + \frac{\Delta t}{2}\right) \quad \text{explicit midpoint method} \quad (31)$$

where \mathbf{u}_k is an approximation of $\mathbf{y}(t_k)$. The explicit midpoint method is a one-step method that belongs to the class of two-stage explicit Runge-Kutta methods. Note that the explicit midpoint method is different from the Heun method (28).

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

$$\mathbf{y}\left(t_k + \frac{\Delta t}{2}\right) \simeq \frac{\mathbf{y}(t_k) + \mathbf{y}(t_{k+1})}{2} \quad (32)$$

yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}\left(\frac{\mathbf{u}_k + \mathbf{u}_{k+1}}{2}, t_k + \frac{\Delta t}{2}\right) \quad \text{implicit midpoint method} \quad (33)$$

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 ODEs (see, e.g., [1]).

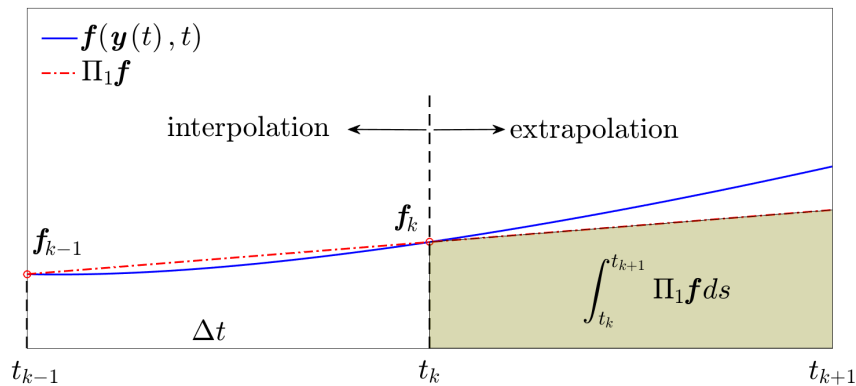


Figure 4: Derivation of the two-step Adams-Bashforth scheme (AB2). We first construct the polynomial (line) $\Pi_1 \mathbf{f}$ that interpolates $\mathbf{f}(\mathbf{y}(s), s)$ at t_{k-1} and t_k . Subsequently, we extrapolate $\Pi_1 \mathbf{f}$ to $[t_k, t_{k+1}]$ and compute the integral in equation (34).

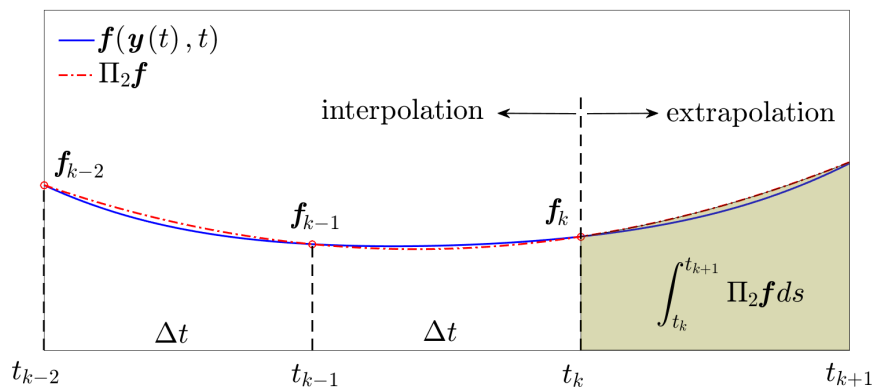


Figure 5: Derivation of the three-step Adams-Bashforth scheme (AB3). We first construct the polynomial $\Pi_2 \mathbf{f}$ that interpolates $\mathbf{f}(\mathbf{y}(s), s)$ at t_{k-2} , t_{k-1} and t_k . Subsequently, we extrapolate $\Pi_2 \mathbf{f}$ to $[t_k, t_{k+1}]$ and compute the integral in equation (34).

Adams-Bashforth methods

These are *explicit multistep methods* constructed by replacing the integral at the right hand side of (10) 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

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_q \mathbf{f}(\mathbf{y}(s), s) ds \quad (34)$$

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

$$\{\mathbf{f}_k, \mathbf{f}_{k-1}, \dots, \mathbf{f}_{k-q}\} \quad \text{at} \quad \{t_k, t_{k-1}, \dots, t_{k-q}\} \quad (35)$$

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):**

$$q = 0 \quad \Rightarrow \quad \Pi_0 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k \quad (36)$$

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

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_0 \mathbf{f}(\mathbf{y}(s), s) ds = \Delta t \mathbf{f}(\mathbf{y}(t_k), t_k). \quad (37)$$

This yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k) \quad \text{AB1 method.} \quad (38)$$

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

- **Two-step Adams-Bashforth method (AB2):** Here we have (see Figure 4)

$$q = 1 \quad \Rightarrow \quad \Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s) \quad (39)$$

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

$$l_k(s) = \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad l_{k-1}(s) = \frac{s - t_k}{t_{k-1} - t_k}. \quad (40)$$

The interpolant can be explicitly written as

$$\Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k \frac{s - t_{k-1}}{t_k - t_{k-1}} + \mathbf{f}_{k-1} \frac{s - t_k}{t_{k-1} - t_k}. \quad (41)$$

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

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_k), t_k) = \mathbf{f}_k, \quad (42)$$

$$\begin{aligned} \Pi_1 \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) &= \mathbf{f}_k \frac{t_{k+1} - t_{k-1}}{t_k - t_{k-1}} + \mathbf{f}_{k-1} \frac{t_{k+1} - t_k}{t_{k-1} - t_k} \\ &= 2\mathbf{f}_k - \mathbf{f}_{k-1}. \end{aligned} \quad (43)$$

which gives us

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_1 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{2} (3\mathbf{f}_k - \mathbf{f}_{k-1}). \quad (44)$$

Substituting this back into (10) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [3\mathbf{f}(\mathbf{u}_k, t_k) - \mathbf{f}(\mathbf{u}_{k-1}, t_{k-1})] \quad \text{AB2 method.} \quad (45)$$

- **Three-step Adams-Bashforth method (AB3):** With reference to Figure 5

$$q = 2 \quad \Rightarrow \quad \Pi_2 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s) + \mathbf{f}_{k-2} l_{k-2}(s), \quad (46)$$

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

$$l_k(s) = \frac{s - t_{k-1}}{t_k - t_{k-1}} \frac{s - t_{k-2}}{t_k - t_{k-2}}, \quad (47)$$

$$l_{k-1}(s) = \frac{s - t_k}{t_{k-1} - t_k} \frac{s - t_{k-2}}{t_{k-1} - t_{k-2}}, \quad (48)$$

$$l_{k-2}(s) = \frac{s - t_k}{t_{k-2} - t_k} \frac{s - t_{k-1}}{t_{k-2} - t_{k-1}}. \quad (49)$$

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

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_2 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{12} (23\mathbf{f}_k - 16\mathbf{f}_{k-1} + 5\mathbf{f}_{k-2}). \quad (50)$$

Substituting this back into (10) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{12} [23\mathbf{f}(\mathbf{u}_k, t_k) - 16\mathbf{f}(\mathbf{u}_{k-1}, t_{k-1}) + 5\mathbf{f}(\mathbf{u}_{k-2}, t_{k-2})] \quad \text{AB3 method.} \quad (51)$$

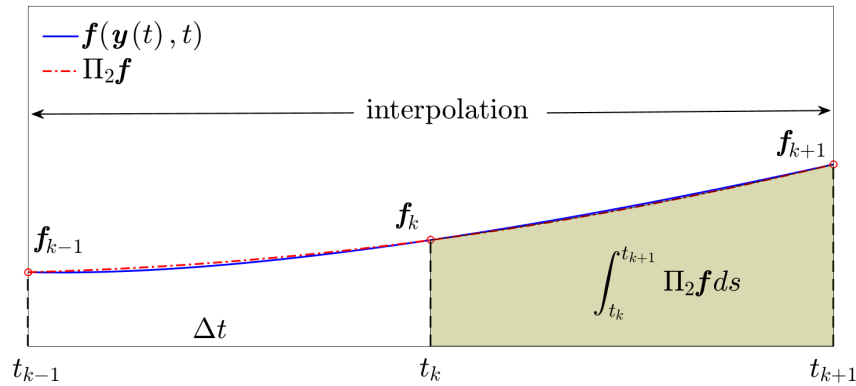


Figure 6: Derivation of the two-step Adams-Moulton scheme (AM2). We first construct the polynomial $\Pi_2 \mathbf{f}$ that interpolates $\mathbf{f}(\mathbf{y}(s), s)$ at t_{k-1} , t_k and t_{k+1} . Subsequently, we use $\Pi_2 \mathbf{f}$ polynomial to approximate the integral in equation (34).

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 (45)

Adams-Moulton methods

These are implicit multistep methods in which the time integral at the right hand-side of (10) 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

$$\{\mathbf{f}_{k+1}, \mathbf{f}_k, \dots, \mathbf{f}_{k-q+1}\} \quad \text{at} \quad \{t_{k+1}, t_{k-1}, \dots, t_{k-q+1}\} \quad (52)$$

The main difference with respect to Adams-Bashforth methods is that there is no extrapolation step, i.e., the point $(t_{k+1}, \mathbf{f}_{k+1})$ is included in the interpolation process (see Figure 6). Let us derive the Adams-Moulton schemes for $q = 0, 1, 2$.

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

$$q = 0 \quad \Rightarrow \quad \Pi_0 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} \quad (53)$$

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

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_0 \mathbf{f}(\mathbf{y}(s), s) ds = \Delta t \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}). \quad (54)$$

This yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) \quad \text{AM0 method.} \quad (55)$$

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

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

$$q = 1 \quad \Rightarrow \quad \Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} l_{k+1}(s) + \mathbf{f}_k l_k(s) \quad (56)$$

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

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k}, \quad l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}}. \quad (57)$$

This yields

$$\Pi_1 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} \frac{s - t_k}{t_{k+1} - t_k} + \mathbf{f}_k \frac{s - t_{k+1}}{t_k - t_{k+1}}. \quad (58)$$

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

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_k), t_k) = \mathbf{f}_k, \quad (59)$$

$$\Pi_1 \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) = \mathbf{f}_{k+1} \quad (60)$$

which imply

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_1 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{2} (\mathbf{f}_{k+1} + \mathbf{f}_k). \quad (61)$$

Substituting this back into (10) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) + \mathbf{f}(\mathbf{u}_k, t_k)] \quad \text{AM1 method} \quad (62)$$

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

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

$$q = 2 \quad \Rightarrow \quad \Pi_2 \mathbf{f}(\mathbf{y}(s), s) = \mathbf{f}_{k+1} l_{k+1}(s) + \mathbf{f}_k l_k(s) + \mathbf{f}_{k-1} l_{k-1}(s). \quad (63)$$

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

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k} \frac{s - t_{k-1}}{t_{k+1} - t_{k-1}}, \quad (64)$$

$$l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad (65)$$

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

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

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \int_{t_k}^{t_{k+1}} \Pi_2 \mathbf{f}(\mathbf{y}(s), s) ds = \frac{\Delta t}{12} (5\mathbf{f}_{k+1} + 8\mathbf{f}_k - \mathbf{f}_{k-1}). \quad (67)$$

Substituting this back into (10) yields the scheme

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{12} [5\mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) + 8\mathbf{f}(\mathbf{u}_k, t_k) - \mathbf{f}(\mathbf{u}_{k-1}, t_{k-1})] \quad \text{AM2 method.} \quad (68)$$

Backward differentiation formulas (BDF) methods

These are 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 (BFD) 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 7). Let us derive the first few BDF methods.

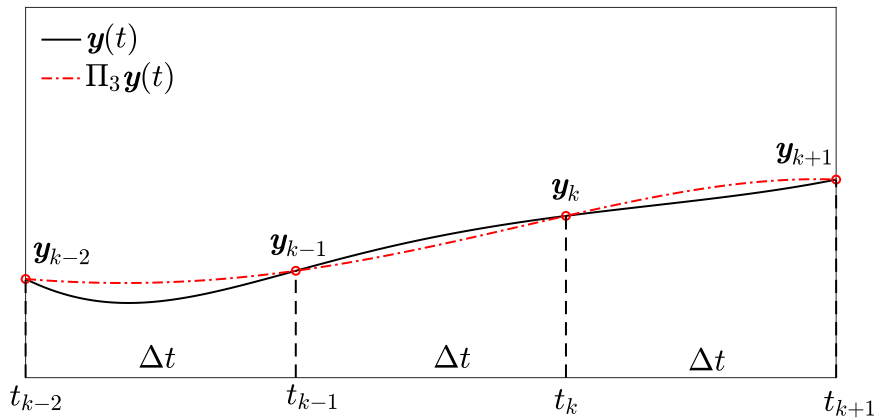


Figure 7: Derivation of the three-step backward differentiation formula (BDF3) method. We first construct the polynomial $\Pi_3 \mathbf{f}$ that interpolates $\mathbf{y}(t)$ at t_{k-2} , t_{k-1} , t_k , and t_{k+1} . Subsequently, approximate the derivative of $\mathbf{y}(t)$ at t_{k+1} with the derivative of the polynomial $\Pi_3 \mathbf{y}(t)$ at t_{k+1} .

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

$$\Pi_1 \mathbf{y}(s) = \mathbf{y}_{k+1} l_{k+1}(s) + \mathbf{y}_k l_k(s) \quad (69)$$

where the Lagrange characteristic polynomials are given by

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k}, \quad l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \quad (70)$$

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.,

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{d\Pi_1 \mathbf{y}(t_{k+1})}{dt} = \frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{\Delta t}. \quad (71)$$

Substituting this approximation into the exact equation

$$\frac{d\mathbf{y}(t_{k+1})}{dt} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) \quad (72)$$

yields the scheme

$$\mathbf{u}_{k+1} - \mathbf{u}_k = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (73)$$

Note that BDF1 coincides with the Euler backward scheme.

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

$$\Pi_2 \mathbf{y}(s) = \mathbf{y}_{k+1} l_{k+1}(s) + \mathbf{y}_k l_k(s) + \mathbf{y}_{k-1} l_{k-1}(s) \quad (74)$$

where the Lagrange characteristic polynomials are given by

$$l_{k+1}(s) = \frac{s - t_k}{t_{k+1} - t_k} \frac{s - t_{k-1}}{t_{k+1} - t_{k-1}}, \quad (75)$$

$$l_k(s) = \frac{s - t_{k+1}}{t_k - t_{k+1}} \frac{s - t_{k-1}}{t_k - t_{k-1}}, \quad (76)$$

$$l_{k-1}(s) = \frac{s - t_{k+1}}{t_{k-1} - t_{k+1}} \frac{s - t_k}{t_{k-1} - t_k}. \quad (77)$$

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.,

$$\frac{d\mathbf{y}(t_{k+1})}{dt} \simeq \frac{d\Pi_2\mathbf{y}(t_{k+1})}{dt} = \frac{3\mathbf{y}_{k+1} - 4\mathbf{y}_k + \mathbf{y}_{k-1}}{2\Delta t}. \quad (78)$$

Substituting this approximation into the exact equation

$$\frac{d\mathbf{y}(t_{k+1})}{dt} = \mathbf{f}(\mathbf{y}(t_{k+1}), t_{k+1}) \quad (79)$$

yields the scheme

$$\frac{3}{2}\mathbf{u}_{k+1} - 2\mathbf{u}_k + \frac{1}{2}\mathbf{u}_{k-1} = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (80)$$

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

$$\frac{11}{6}\mathbf{u}_{k+1} - 3\mathbf{u}_k + \frac{3}{2}\mathbf{u}_{k-1} - \frac{1}{3}\mathbf{u}_{k-2} = \Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (81)$$

Linear multistep methods (LMM)

The general form of a linear multistep method is

$$\sum_{j=0}^q \alpha_j \mathbf{u}_{k+j} = \Delta t \sum_{j=0}^q \beta_j \underbrace{\mathbf{f}(\mathbf{u}_{k+j}, t_{k+j})}_{\mathbf{f}_{k+j}}. \quad (82)$$

Note that Adams-Bashforth, Adams-Moulton and BDF methods are all in the form (82). 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

$$\mathbf{u}_{k+3} = \mathbf{u}_{k+2} + \frac{\Delta t}{12} (23\mathbf{f}_{k+2} - 16\mathbf{f}_{k+1} + 5\mathbf{f}_k) \quad (83)$$

can be written in the form (82) by setting

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

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

$$\sum_{j=0}^3 \beta_j = 1. \quad (84)$$

Example: The BDF2 method

$$\mathbf{u}_{k+2} - \frac{4}{3}\mathbf{u}_{k+1} + \frac{1}{3}\mathbf{u}_k = \frac{2}{3}\Delta t \mathbf{f}_{k+2} \quad (85)$$

can be written in the form (82) by setting

$$\alpha_2 = 1, \quad \alpha_1 = -\frac{4}{3}, \quad \alpha_0 = \frac{1}{3}, \quad \beta_2 = \frac{2}{3}, \quad \beta_1 = 0, \quad \beta_0 = 0. \quad (86)$$

Note that $(\alpha_2, \beta_2) = (1, 2/3)$ (the method is implicit) and

$$\sum_{j=0}^2 \alpha_j = 0, \quad \sum_{j=0}^2 (\beta_j - j\alpha_j) = 0. \quad (87)$$

As we will see the conditions (84) and (87) guarantee that AB3 and BDF2 are *consistent* methods, i.e., that the truncation error of these methods goes to zero as we send Δt to zero. Roughly speaking this means that the numerical schemes (83) and (85) 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

$$\sum_{j=0}^q \alpha_j = 0, \quad \sum_{j=0}^q (\beta_j - j\alpha_j) = 0. \quad (88)$$

Remark: The linear multistep methods we have seen (Adams and BDF) rely on a polynomial interpolation process on an evenly-spaced temporal grid. As we know, polynomial interpolation on evenly spaced grids is, in general, ill-conditioned and can undergo a severe Runge's phenomenon depending on the function we are interpolating. However, the interpolation process of a function with a polynomial of fixed degree, say $q + 1$, within a small a time interval (equal to $q\Delta t$ for a q -step method) is convergent as Δt goes to zero. To show this, recall the following error estimate we obtained in course note 1 (equation (23))

$$\|g - \Pi_q g\|_\infty \leq \Delta t^{q+1} \frac{\|g^{(q+1)}\|_\infty}{(q+1)}, \quad (89)$$

where $g(t)$ 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 $\{t, t - \Delta t, \dots, t - q\Delta t\}$. The upper bound in (89) goes to zero as we send Δt to zero.

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

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \sum_{i=1}^s b_i \mathbf{K}_i \quad (90)$$

where

$$\mathbf{K}_i = \mathbf{f} \left(\mathbf{u}_k + \Delta t \sum_{j=1}^s a_{ij} \mathbf{K}_j, t_k + c_i \Delta t \right). \quad (91)$$

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

$$\begin{array}{c|cccc} c_1 & a_{11} & a_{12} & \cdots & a_{1s} \\ c_2 & a_{21} & a_{22} & \cdots & a_{2s} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_s & a_{s1} & a_{s2} & \cdots & a_{ss} \\ \hline & b_1 & b_2 & \cdots & b_s \end{array}$$

a_{ij} in the Butcher array can be positive or negative, while c_i and b_i are real numbers in $[0, 1]$. The coefficients $c_i \in [0, 1]$ are sometimes referred to as “nodes” as they identify appropriate time instants between t_k and

t_{k+1} in which \mathbf{f} is evaluated. On the other hand, $b_i \in [0, 1]$ are sometimes referred to as “weights” since they appear in the weighted sum (90).

We will see that consistency of RK methods, i.e., the fact that (90) converges to (1) as $\Delta t \rightarrow 0$ implies the following condition

$$\sum_{j=1}^s b_j = 1. \quad (92)$$

Moreover, it is usually assumed that

$$c_i = \sum_{j=1}^s a_{ij}. \quad (93)$$

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

If $a_{ij} = 0$ for $j \geq i$ 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 of RK methods.

- **Euler forward method:** The Euler forward method (14) 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 (28) 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

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), \quad (94)$$

where

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (95)$$

$$\mathbf{K}_2 = \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{f}(\mathbf{u}_k, t_k), t_k + \Delta t). \quad (96)$$

- **Crank-Nicolson method:** The Crank-Nicolson (CN) method (16) 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

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), \quad (97)$$

where

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

$$\mathbf{K}_2 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2), t_{k+1}\right). \quad (99)$$

From equation (97) we see that

$$\frac{\Delta t}{2} (\mathbf{K}_1 + \mathbf{K}_2) = \mathbf{u}_{k+1} - \mathbf{u}_k. \quad (100)$$

Substituting this expression into (99) yields

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad \mathbf{K}_2 = \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}). \quad (101)$$

- **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

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

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

$$\mathbf{K}_3 = \mathbf{f}(\mathbf{u}_k - \Delta t \mathbf{K}_1 + 2\Delta t \mathbf{K}_2, t_k + \Delta t), \quad (104)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{6} (\mathbf{K}_1 + 4\mathbf{K}_2 + \mathbf{K}_3). \quad (105)$$

- **Runge-Kutta method (RK4):** The most famous RK method is perhaps the one proposed in the original paper by Runge and Kutta. The scheme is 4-stage explicit and can be written as

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad (106)$$

$$\mathbf{K}_2 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{K}_1, t_k + \frac{\Delta t}{2}\right), \quad (107)$$

$$\mathbf{K}_3 = \mathbf{f}\left(\mathbf{u}_k + \frac{\Delta t}{2} \mathbf{K}_2, t_k + \frac{\Delta t}{2}\right), \quad (108)$$

$$\mathbf{K}_4 = \mathbf{f}(\mathbf{u}_k + \Delta t \mathbf{K}_3, t_k + \Delta t), \quad (109)$$

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \frac{\Delta t}{6} (\mathbf{K}_1 + 2\mathbf{K}_2 + 2\mathbf{K}_3 + \mathbf{K}_4). \quad (110)$$

The Butcher array for this scheme is

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

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:

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

The RK2 method corresponding to this table can be written explicitly as

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t(b_1 \mathbf{K}_1 + b_2 \mathbf{K}_2) \quad (111)$$

where

$$\mathbf{K}_1 = \mathbf{f}(\mathbf{u}_k, t_k), \quad \mathbf{K}_2 = \mathbf{f}(\mathbf{u}_k + c\Delta t \mathbf{K}_1, t_k + c\Delta t). \quad (112)$$

Expand \mathbf{K}_2 in a Taylor series in Δt to obtain

$$\begin{aligned} \mathbf{K}_2 &= \mathbf{K}_1 + c\Delta t \left(\sum_{j=1}^n \frac{\partial \mathbf{f}}{\partial y_j} K_{1j} + \frac{\partial \mathbf{f}}{\partial t} \right) + \dots \\ &= \mathbf{f}(\mathbf{u}_k, t_k) + c\Delta t \left(\sum_{j=1}^n \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial y_j} f_j(\mathbf{u}_k, t_k) + \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial t} \right) + \dots \end{aligned} \quad (113)$$

A substitution of this expression into (111) yields

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t(b_1 + b_2) \mathbf{f}(\mathbf{u}_k, t_k) + b_2 c (\Delta t)^2 \left(\sum_{j=1}^n \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial y_j} f_j(\mathbf{u}_k, t_k) + \frac{\partial \mathbf{f}(\mathbf{u}_k, t_k)}{\partial t} \right) + \dots \quad (114)$$

Next, we expand the solution to the ODE (1) in a Taylor series

$$\mathbf{y}(t_{k+1}) = \mathbf{y}(t_k) + \frac{d\mathbf{y}(t_k)}{dt} \Delta t + \frac{(\Delta t)^2}{2} \frac{d^2 \mathbf{y}(t_k)}{dt^2} + \dots \quad (115)$$

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

$$\frac{d\mathbf{y}(t_k)}{dt} = \mathbf{f}(\mathbf{y}_k, t_k), \quad \frac{d^2 \mathbf{y}(t_k)}{dt^2} = \sum_{j=1}^n \frac{\partial \mathbf{f}(\mathbf{y}_k, t_k)}{\partial y_j} f_j(\mathbf{y}_k, t_k) + \frac{\partial \mathbf{f}(\mathbf{y}_k, t_k)}{\partial t}. \quad (116)$$

Assuming that $\mathbf{y}_k = \mathbf{u}_k$ and matching the terms multiplying the same powers of Δt in (114) and (115) we obtain³

$$b_1 + b_2 = 1 \quad \text{and} \quad cb_2 = \frac{1}{2}. \quad (117)$$

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

$$c = 1, \quad b_1 = b_2 = \frac{1}{2} \quad \text{yields the Heun method (28)}. \quad (118)$$

On the other hand, setting

$$c = \frac{1}{2}, \quad b_1 = 0 \quad b_2 = 1 \quad \text{yields the explicit midpoint method (31)}. \quad (119)$$

³The condition $b_1 + b_2 = 1$ is a consistency condition which guarantees that the scheme (111) converges to the ODE (1) in the limit $\Delta t \rightarrow 0$.

Derivation of higher order explicit RK methods. By using the Taylor series approach discussed in previous section, it is possible to derive conditions on the entries of the Butcher array for explicit RK methods with an arbitrary number of stages. Essentially, one can perform a Taylor series of the RK method in Δt 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., (117) are called *stage-order* conditions (see, e.g., [2, §5.9]). For example, it can be shown that the stage-order conditions for a three stage explicit RK method corresponding to the arbitrary Butcher array

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

are

$$\begin{cases} b_1 + b_2 + b_3 = 1 \\ b_2 c_2 + b_3 c_3 = \frac{1}{2} \\ b_2 c_2^2 + b_3 c_3^2 = \frac{1}{3} \\ b_3 a_{32} c_2 = \frac{1}{6} \end{cases} \quad (120)$$

This system admits one two-parameter family of solutions and two one-parameter families of solutions (see [2, p. 178]). Hereafter we list a few three-stage explicit RK methods satisfying (120):

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 \\ 2/3 & 0 & 2/3 & 0 \\ \hline & 1/4 & 0 & 3/4 \end{array} \quad \text{Heun's method}$$

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 8/15 & 8/15 & 0 & 0 \\ 2/3 & 1/4 & 5/12 & 0 \\ \hline & 1/4 & 0 & 3/4 \end{array} \quad \text{Van der Houwen's/Wray method}$$

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 \\ 3/4 & 0 & 3/4 & 0 \\ \hline & 2/9 & 1/3 & 4/9 \end{array} \quad \text{Ralson's method}$$

The Taylor series approach rapidly become intractable as the number of stages increases, and so does the corresponding stage-order conditions. Fortunately, there is a more efficient way to derive conditions such as (117) using *rooted trees* (see [2, §5.6]).

Derivation of implicit RK methods. Implicit RK methods are very expensive to implement and cannot rival explicit RK methods in terms of efficiency (for non-stiff problems). R can be derived from the integral formulation (10) of the Cauchy problem. In particular, consider a quadrature formula with s

nodes in $[t_k, t_{k+1}]$ to approximate the integral at the right hand side of (10),

$$\int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{y}(s), s) ds \simeq \Delta t \sum_{j=1}^s b_j \mathbf{f}(\mathbf{y}(t_k + c_j \Delta t), t_k + c_j \Delta t). \quad (121)$$

Here, we denoted by b_j the quadrature weights and by $t_k + c_j \Delta t$ the quadrature nodes. The coefficients a_{ij} of an implicit RK method can be derived based on the weights b_j and the nodes c_j of the quadrature formula (121) and appropriated stage-order conditions (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 (i.e., the midpoint of the interval $[t_k, t_{k+1}]$). This method corresponds to the Butcher array and can attain order 2.

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

Similarly, the following (implicit) Gauss-Legendre-RK method with two stages

$$\begin{array}{c|cc} 1/2 - \sqrt{3}/6 & 1/4 & 1/4 - \sqrt{3}/6 \\ 1/2 + \sqrt{3}/6 & 1/4 + \sqrt{3}/6 & 1/4 \\ \hline & 1/2 & 1/2 \end{array}$$

can attain order 4. In general, an implicit s -stage Gauss-RK method can attain order $2s$, while Gauss-Radau-RK and Gauss-Lobatto-RK methods can achieve order $2s - 1$ and $2s - 2$, respectively.

Appendix A: Vector norms and matrix norms

In this appendix we briefly review the most common vector norms in \mathbb{R}^n and their equivalence. We also discuss common matrix norms induced by vector norms.

Vector norms and their equivalence

As is well known, all norms defined in a finite-dimensional vector space (such as \mathbb{R}^n) are *equivalent*. This means that if we pick two arbitrary norms in \mathbb{R}^n , say $\|\cdot\|_a$ and $\|\cdot\|_b$, then there exist two numbers c_1 and c_2 such that

$$c_1 \|\mathbf{y}\|_a \leq \|\mathbf{y}\|_b \leq c_2 \|\mathbf{y}\|_a \quad \text{for all } \mathbf{y} \in \mathbb{R}^n. \quad (122)$$

The most common norms in \mathbb{R}^n are

$$\|\mathbf{y}\|_\infty = \max_{k=1,\dots,n} |y_k|, \quad (123)$$

$$\|\mathbf{y}\|_1 = \sum_{k=1}^n |y_k|, \quad (124)$$

$$\|\mathbf{y}\|_2 = \left(\sum_{k=1}^n |y_k|^2 \right)^{1/2}, \quad (125)$$

$$\vdots \quad (126)$$

$$\|\mathbf{y}\|_p = \left(\sum_{k=1}^n |y_k|^p \right)^{1/p} \quad p \in \mathbb{N} \setminus \{\infty\}. \quad (127)$$

Based on these definitions it is easy to show, e.g., that

$$\|\mathbf{y}\|_\infty \leq \|\mathbf{y}\|_1 \leq n \|\mathbf{y}\|_\infty, \quad (128)$$

$$\|\mathbf{y}\|_2 \leq \|\mathbf{y}\|_1 \leq \sqrt{n} \|\mathbf{y}\|_2, \quad (129)$$

$$\|\mathbf{y}\|_\infty \leq \|\mathbf{y}\|_2 \leq \sqrt{n} \|\mathbf{y}\|_\infty. \quad (130)$$

As an example, if $\mathbf{f}(\mathbf{y}, t)$ in (1) is Lipschitz continuous in D with respect to the 1-norm, i.e.,

$$\|\mathbf{f}(\mathbf{y}_1, t) - \mathbf{f}(\mathbf{y}_2, t)\|_1 \leq L_1 \|\mathbf{y}_1 - \mathbf{y}_2\|_1 \quad \text{for all } \mathbf{y}_1, \mathbf{y}_2 \in D, \quad \text{for all } t \geq 0 \quad (131)$$

then it is also Lipschitz continuous with respect to the uniform norm. In fact, by using (128) we have

$$\|\mathbf{f}(\mathbf{y}_1, t) - \mathbf{f}(\mathbf{y}_2, t)\|_\infty \leq \underbrace{L_1 n}_{L_\infty} \|\mathbf{y}_1 - \mathbf{y}_2\|_\infty. \quad (132)$$

Of course, $\mathbf{f}(\mathbf{y}, t)$ is also Lipschitz continuous with respect to the 2-norm.

Matrix norms induced by vector norms

Let us define the following matrix norm

$$\|\mathbf{A}\| = \sup_{\mathbf{y} \neq \mathbf{0}_{\mathbb{R}^n}} \frac{\|\mathbf{A}\mathbf{y}\|}{\|\mathbf{y}\|} = \sup_{\|\mathbf{y}\|=1} \|\mathbf{A}\mathbf{y}\|. \quad (133)$$

where the norm at the right hand side is any vector norm. Clearly, $\|\mathbf{A}\|$ is matrix norm (prove it as exercise), which satisfies, by definition, the following inequality

$$\|\mathbf{A}\| \geq \frac{\|\mathbf{A}\mathbf{y}\|}{\|\mathbf{y}\|} \quad \text{i.e.} \quad \|\mathbf{A}\mathbf{y}\| \leq \|\mathbf{A}\| \|\mathbf{y}\|. \quad (134)$$

It is straightforward to show that

$$\|\mathbf{A}\|_\infty = \max_{i=1,\dots,n} \left(\sum_{j=1}^n |A_{ij}| \right), \quad (135)$$

$$\|\mathbf{A}\|_1 = \max_{j=1,\dots,n} \left(\sum_{i=1}^n |A_{ij}| \right), \quad (136)$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} = \sigma_{\max}(\mathbf{A}), \quad (137)$$

where $\sigma_{\max}(\mathbf{A})$ is the largest singular value of the matrix \mathbf{A} . For example,

$$\|\mathbf{A}\mathbf{y}\|_\infty = \max_{i=1,\dots,n} \left| \sum_{j=1}^n A_{ij} y_j \right| \leq \max_{i=1,\dots,n} \left(\sum_{j=1}^n |A_{ij}| |y_j| \right) \leq \|\mathbf{y}\|_\infty \max_{i=1,\dots,n} \left(\sum_{j=1}^n |A_{ij}| \right) \quad (138)$$

which implies that

$$\frac{\|\mathbf{A}\mathbf{y}\|_\infty}{\|\mathbf{y}\|_\infty} \leq \max_{i=1,\dots,n} \left(\sum_{j=1}^n |A_{ij}| \right) \quad \text{for all } \mathbf{y} \neq \mathbf{0}_{\mathbb{R}^n}, \quad (139)$$

i.e.,

$$\sup_{\mathbf{y} \neq \mathbf{0}_{\mathbb{R}^n}} \frac{\|\mathbf{A}\mathbf{y}\|_\infty}{\|\mathbf{y}\|_\infty} = \max_{i=1,\dots,n} \left(\sum_{j=1}^n |A_{ij}| \right) = \|\mathbf{A}\|_\infty. \quad (140)$$

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.