

# Numerical Methods for Ordinary Differential Equations

Habib Ammari

Konstantinos Alexopoulos



# Contents

Chapter 1. Some basics	5
1.1. What is a differential equation?	5
1.2. Some methods of resolution	7
1.3. Important examples of ODEs	10
1.4. Problems	17
Chapter 2. Existence, uniqueness, and regularity in the Lipschitz case	19
2.1. Banach fixed point theorem	19
2.2. Gronwall's lemma	19
2.3. Cauchy-Lipschitz theorem	20
2.4. Stability	23
2.5. Regularity	25
2.6. Problems	26
Chapter 3. Linear systems	27
3.1. Exponential of a matrix	27
3.2. Linear systems with constant coefficients	28
3.3. Linear system with non-constant real coefficients	29
3.4. Second order linear equations	32
3.5. Linearization and stability for autonomous systems	36
3.6. Periodic linear systems	39
3.7. Problems	41
Chapter 4. Numerical solution of ordinary differential equations	43
4.1. Introduction	43
4.2. The general explicit one-step method	43
4.3. Example of linear systems	50
4.4. Runge-Kutta methods	51
4.5. Multi-step methods	57
4.6. Stiff equations and systems	62
4.7. Perturbation theories for differential equations	62
4.8. Problems	66
Chapter 5. Geometrical numerical integration methods for differential equations	67
5.1. Introduction	67
5.2. Structure preserving methods for Hamiltonian systems	67
5.3. Runge-Kutta methods	75
5.4. Long-time behaviour of numerical solutions	78
5.5. Problems	78
Chapter 6. Finite difference methods	81
6.1. Introduction	81
6.2. Numerical algorithms for the heat equation	81
6.3. Numerical algorithms for the wave equation	90

Bibliography

97

## CHAPTER 1

### Some basics

#### 1.1. What is a differential equation?

An ordinary differential equation (ODE) is an equation that contains one or more derivatives of an unknown function  $x(t)$ . The equation may also contain  $x$  itself and constants. We say that an ODE is of order  $n$  if the  $n$ -th derivative of the unknown function is the highest order derivative in the equation. The following equations are examples of ODEs:

**Membrane equation as a neuron model:**

$$C \frac{dx(t)}{dt} + gx(t) = f(t), \quad (1.1)$$

where  $x(t)$  is the membrane potential, *i.e.*, the voltage difference between the inside and the outside of the neuron,  $f(t)$  is the current flow due to excitation,  $C$  is the capacitance and  $g$  is the conductance (the inverse of the resistance) of the membrane.

Equation (1.1) is linear ODE of order 1.

**The theta model:** The theta model is a simple one-dimensional model for the spiking of a neuron. It takes the form

$$\frac{d\theta(t)}{dt} = 1 - \cos \theta(t) + (1 + \cos \theta(t))f(t), \quad (1.2)$$

where  $f(t)$  are the inputs to the model. The variable  $\theta$  lies on the unit circle and ranges between 0 and  $2\pi$ . When  $\theta = \pi$  the neuron spikes, that is, it produces an action potential. By the change of variables,  $x(t) = \tan(\theta(t)/2)$ , (1.2) leads to the quadratic model

$$\frac{dx(t)}{dt} = x^2(t) + f(t). \quad (1.3)$$

**Population growth under competition for resources:**

$$\frac{dx(t)}{dt} = rx(t) - \frac{r}{k}x^2(t), \quad (1.4)$$

where  $r$  and  $k$  are positive parameters. In (1.4),  $x(t)$  is the number of cells at time instant  $t$ ,  $rx(t)$  is the growth rate and  $-(r/k)x^2(t)$  is the death rate. Equations (1.2), (1.3), and (1.4) are nonlinear ODEs of order 1.

**FitzHugh-Nagumo model:**

$$\begin{cases} \frac{dV}{dt} = f(V) - W + I \\ \frac{dW}{dt} = a(V - bW), \end{cases} \quad (1.5)$$

where  $f(V)$  is a polynomial of third degree, and  $a$  and  $b$  are constant parameters. The FitzHugh-Nagumo model is a two-dimensional simplification of the Hodgkin-Huxley model of spike generation in squid giant axons. It aims at isolating the mathematical properties of excitation and propagation from the electrochemical properties of sodium and potassium ion flow. In (1.5),  $V$  is the membrane potential,  $W$  is a recovery variable, and  $I$  is the magnitude of stimulus current. Equation (1.5) is a system of nonlinear ODEs of order 1.

**Langevin equation of motion for a single particle:**

$$\frac{dx(t)}{dt} = -ax(t) + \eta(t), \quad (1.6)$$

where  $x(t)$  is the position of the particle at time instant  $t$ ,  $a > 0$  is coefficient of friction, and  $\eta$  is a random variable that represents some uncertainties or stochastic effects perturbing the particle. Equation (1.6) represents diffusion-like motion from the probabilistic perspective of a single microscopic particle moving in a fluid medium. Equation (1.6) is a linear stochastic ODE of order 1.

**Vander der Pol equation:**

$$\frac{d^2x(t)}{dt^2} - a(1 - x^2(t))\frac{dx(t)}{dt} + x(t) = 0, \quad (1.7)$$

where  $a$  is a positive parameter, which controls the nonlinearity and the strength of the damping. Equation (1.7) is used to generate waveforms corresponding to electrocardiogram patterns. Equation (1.7) is a nonlinear ODE of order 2.

**1.1.1. Higher order ODEs.** Here we introduce higher order ODEs. Let  $\Omega \subset \mathbb{R}^{n+2}$  and  $n \in \mathbb{N}$ . Then an ODE of order  $n$  is an equation of the form:

$$F(t, x(t), \frac{dx}{dt}(t), \dots, \frac{d^n x}{dt^n}(t)) = 0,$$

where  $x$  is a real-valued unknown function and  $dx(t)/dt, \dots, d^n x(t)/dt^n$  are its derivatives. We say that  $\varphi \in \mathcal{C}^n(I)$  is a solution of the differential equation if  $I$  is an open interval,

$$(t, \varphi(t), \frac{d\varphi}{dt}(t), \dots, \frac{d^n \varphi}{dt^n}(t)) \in \Omega$$

for all  $t \in I$ , and

$$F(t, \varphi(t), \frac{d\varphi}{dt}(t), \dots, \frac{d^n \varphi}{dt^n}(t)) = 0$$

for all  $t \in I$ . When  $x$  is a vector valued function, *i.e.*,  $x(t) \in \mathbb{R}^d$ , then  $\Omega \subset \mathbb{R} \times \mathbb{R}^{(n+1)d}$ .

Next we consider the following form of  $n$ -th order ODE:

$$x^{(n)}(t) = f(t, x, \frac{dx}{dt}, \dots, \frac{d^{n-1}x}{dt^{n-1}}), \quad t \in I. \quad (1.8)$$

where  $x(t) \in \mathbb{R}^d$  and  $f : I \times \mathbb{R}^{nd} \rightarrow \mathbb{R}^d$ . To ensure uniqueness of the solution, (1.8) has to be augmented with the initial condition:

$$(x(t_0), x'(t_0), x''(t_0), \dots, x^{(n-1)}(t_0))^\top.$$

Here  $\top$  denotes the transpose.

We can reduce the high order ODE (1.8) into a first order ODE. Let us define

$$y(t) := (x(t), dx(t)/dt, \dots, d^{n-1}x(t)/dt^{n-1})^\top \in \mathbb{R}^{nd}$$

and

$$F(t, y) := (y_2, \dots, y_n, f(t, y_1, \dots, y_n))^\top$$

for  $y = (y_1, \dots, y_n)^\top \in \mathbb{R}^{nd}$  and  $y_i \in \mathbb{R}^d$  for  $i = 1, 2, \dots, n$ . Then the  $n$ -th order ODE (1.8) is equivalent to the following first order ODE:

$$\frac{dy}{dt} = F(t, y(t)).$$

**EXAMPLE 1.1.** Consider the second order ODE given by

$$\frac{d^2x}{dt^2} + p(t)\frac{dx}{dt} + q(t)x(t) = g(t).$$

Then we have

$$\frac{d}{dt} \begin{bmatrix} x \\ \frac{dx}{dt} \end{bmatrix} = \begin{bmatrix} \frac{dx}{dt} \\ -p(t)\frac{dx}{dt} - q(t)x(t) + g(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -q(t) & -p(t) \end{bmatrix} \begin{bmatrix} x \\ \frac{dx}{dt} \end{bmatrix} + \begin{bmatrix} 0 \\ g(t) \end{bmatrix}.$$

The main problems concerning ordinary differential equations are:

- (i) Existence of solutions;
- (ii) Uniqueness of solutions with suitable initial conditions;
- (iii) Regularity and stability of solutions (e.g. dependence on the initial conditions, large time stability, higher regularity);
- (iv) Computation of solutions.

The existence of solutions can be proved by fixed point theorems, by the implicit function theorem in Banach spaces, and by functional analysis techniques. The problem of uniqueness is typically more difficult. Only in a very few special cases is it possible to compute solutions in some explicit form.

### 1.2. Some methods of resolution

In the following subsections, we present several examples of exactly solvable ODEs and then explain how to solve them.

**1.2.1. Separation of variables.** Let  $I$  and  $J$  be two open intervals and let  $f \in C^0(I)$  and  $g \in C^0(J)$  be two continuous functions. We look for solutions to the first order equation

$$\frac{dx}{dt} = f(t)g(x). \quad (1.9)$$

Let  $t_0 \in I$  and  $x_0 \in J$ . If  $g(x_0) = 0$  for some  $x_0 \in J$ , then the constant function  $x(t) = x_0$  for  $t \in I$  is a solution to (1.9). Suppose that  $g(x_0) \neq 0$ . Then  $g \neq 0$  in a neighborhood of  $x_0$  and we can divide (1.9) by  $g(x)$  and hence, **separate the variables**. We find

$$\frac{dx}{g(x)} = f(t)dt. \quad (1.10)$$

Integrating (1.10) gives

$$\int \frac{dx}{g(x)} = \int f(t)dt + c,$$

where the constant  $c$  is uniquely determined by the initial condition.

Let  $F$  and  $G$  be the primitives of  $f$  and  $1/g$ , respectively. The function  $G$  is strictly monotonic, because  $G'(x) \neq 0$ , and thus invertible. The solution of the differential equation (1.9) is then

$$x(t) = G^{-1}(F(t) + c).$$

This method of solving ODEs is called the **method of separation of variables** and (1.9) is called a **separable equation**.

EXAMPLE 1.2. Consider the following ODE:

$$\begin{cases} \frac{dx}{dt} = \frac{1+2t}{\cos x(t)}, \\ x(0) = \pi. \end{cases}$$

In this case, we have  $g(x) = 1/\cos x$  and  $f(t) = 1+2t$ . Note that  $g$  is defined for  $x \neq \pi/2 + k\pi, k \in \mathbb{Z}$ . By separating variables, we get

$$\cos x dx = 1 + 2t dt.$$

By integration, we have

$$\sin x(t) = t^2 + t + C,$$

for some constant  $C \in \mathbb{R}$ . Then, from the initial condition  $x(0) = \pi$ , we see that  $C = 0$ .

One might think that we can obtain the solution by taking the arcsin. But the function  $x(t) = \arcsin(t^2 + t)$  is not the solution because  $x(0) = \arcsin(0) = 0$ . In order to get the correct solution, we note that arcsin is the inverse of sin on  $[-\pi/2, \pi/2]$ , whereas  $x(t)$  takes the values in a neighborhood of  $\pi$ . Letting  $w(t) = x(t) - \pi$ , we have  $w(0) = x(0) - \pi = 0$ . So, we have  $w(t) = -\arcsin(t^2 + t)$ . Therefore, we get the following correct solution:

$$x(t) = \pi - \arcsin(t^2 + t).$$

**1.2.2. Change of variables.** There are a few important first-order equations that can be solved using some transformation.

1.2.2.1. *Homogeneous equation.* Consider the following ODE:

$$\frac{dx}{dt} = f\left(\frac{x(t)}{t}\right), \quad (1.11)$$

where  $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$  is a continuous function on some open interval  $I \subset \mathbb{R}$ . The ODE (1.11) is called **homogeneous**. By the **change of variables**  $x(t) = ty(t)$  where  $y(t)$  is the new unknown function, the above ODE can be changed to a separable equation. Since

$$\frac{dx}{dt} = y(t) + t \frac{dy}{dt} = f(y(t)),$$

we have a separable equation for  $y$ , which reads:

$$\frac{dy}{f(y) - y} = \frac{dt}{t}.$$

Therefore, (1.11) can be solved by the method of separation of variables.

EXAMPLE 1.3. *Consider*

$$\frac{dx}{dt} = \frac{t^2 + x^2}{xt}.$$

In this case,  $f(s) = s + 1/s$  with  $s = x/t$ . By letting  $y(t) = x(t)/t$ , we get  $ydy = dt/t$ . So, we have  $(1/2)y^2 = \ln t + C$ . Therefore, we obtain

$$x(t) = \pm t \sqrt{2(\ln t + C)}.$$

1.2.2.2. *Bernoulli equation.* A differential equation is of **Bernoulli** type if it is of the form

$$\frac{dx}{dt} = f(t)x + g(t)x^n, \quad n \neq 0, 1. \quad (1.12)$$

The transformation  $y = x^{1-n}$  gives the linear equation

$$\frac{dy}{dt} = (1-n)f(t)y + (1-n)g(t).$$

1.2.2.3. *Riccati equation.* A differential equation is of **Riccati** type if it is of the form

$$\frac{dx}{dt} = f(t)x + g(t)x^2 + h(t). \quad (1.13)$$

Assume that a particular solution  $x_p$  of (1.13) is known. Then the transformation  $y = 1/(x - x_p(t))$  yields the linear equation

$$\frac{dy}{dt} = -(f(t) + 2x_p(t)g(t))y - g(t).$$



**1.2.3. Method of integrating factors.** Consider

$$\frac{dx(t)}{dt} = f(t). \quad (1.14)$$

By integrating (1.14), it follows that the solution  $x(t)$  is given by

$$x(t) = x(0) + \int_0^t f(s) ds.$$

Consider

$$\frac{dx}{dt} + p(t)x(t) = g(t), \quad (1.15)$$

where  $p$  and  $g$  are functions of  $t$ .

If (1.15) were of the form (1.14), then we could immediately write down a solution in terms of integrals. By (1.15) being of the form (1.14), we mean that the left-hand side is expressed as the derivative of our unknown quantity. To make this happen, we can multiply (1.15) by a function,  $\mu(t)$ , and ask whether the resulting equation can be put in the form (1.14).

Let us look for  $\mu(t)$  such that

$$\mu(t) \frac{dx}{dt} + \mu(t)p(t)x(t) = \frac{d}{dt}(\mu(t)x(t)).$$

Taking derivatives, we have  $(1/\mu)d\mu/dt = p(t)$  or

$$\frac{d}{dt} \ln \mu(t) = p(t). \quad (1.16)$$

Integrating (1.16) gives

$$\mu(t) = \exp\left(\int_0^t p(s)ds\right),$$

up to a multiplicative constant. The equation (1.15) is transformed to

$$\frac{d}{dt}(\mu(t)x(t)) = \mu(t)g(t).$$

This equation is precisely of the form (1.14), so we can immediately conclude

$$x(t) = \frac{1}{\mu(t)} \left( \int_0^t \mu(s)g(s)ds \right) + \frac{C}{\mu(t)},$$

where the constant  $C$  can be determined from the initial condition  $x(0) = x_0$ . The function  $\mu(t)$  is called the **integrating factor**.

EXAMPLE 1.4. Consider

$$\begin{cases} \frac{dx}{dt} + \frac{1}{t+1}x(t) = (1+t)^2, & t \geq 0, \\ x(0) = 1. \end{cases}$$

In this case,  $p(t) = 1/(t+1)$  and  $g(t) = (1+t)^2$ . Then the integrating factor  $\mu$  is

$$\mu(t) = \exp\left(\int_0^t p(s)ds\right) = e^{\ln(t+1)} = t+1.$$

Therefore, we get

$$x(t) = \frac{1}{t+1} \int_0^t (s+1)^3 ds + \frac{C}{t+1} = \frac{(t+1)^3}{4} + \frac{C - \frac{1}{4}}{t+1}.$$

Then, from the initial condition  $x(0) = 1$ , we obtain  $C = 1$ .

EXAMPLE 1.5. (*Bernoulli's equation*) Consider

$$\frac{dx}{dt} + p(t)x(t) = g(t)x^\alpha(t). \quad (1.17)$$

Here  $\alpha$  is a real parameter satisfying  $\alpha \notin \{0, 1\}$ . Letting  $x = z^{\frac{1}{1-\alpha}}$ , we get

$$\frac{dx}{dt} = \frac{1}{1-\alpha} z^{\frac{\alpha}{1-\alpha}} \frac{dz}{dt}.$$

Then (1.17) can be reduced to the following linear equation:

$$\frac{dz}{dt} + (1-\alpha)p(t)z(t) = (1-\alpha)g(t),$$

which can be solved by the method of integrating factors.

### 1.3. Important examples of ODEs

#### 1.3.1. Autonomous ODEs.

DEFINITION 1.6. The equation

$$\frac{dx(t)}{dt} = f(t, x(t)) \quad (1.18)$$

is called **autonomous** if  $f$  is independent of  $t$ .

Any ODE can be rewritten as an autonomous ODE on a higher-dimensional space. Writing  $y = (t, x(t))$ , (1.18) is equivalent to the autonomous ODE

$$\frac{dy(t)}{dt} = F(y(t)),$$

where  $F(y) = \begin{pmatrix} 1 \\ f(t, x(t)) \end{pmatrix}$ .

**1.3.2. Exact equations.** Let  $\Omega = I \times \mathbb{R} \subset \mathbb{R}^2$  with  $I \subset \mathbb{R}$  being an open interval. Let  $f, g \in C^0(\Omega)$ . We look for a solution  $x \in C^1(I)$  of the differential equation

$$f(t, x(t)) + g(t, x(t)) \frac{dx}{dt} = 0 \quad (1.19)$$

satisfying the initial condition  $x(t_0) = x_0$  for some  $(t_0, x_0) \in \Omega$ .

Consider the differential form

$$\omega = f(t, x)dt + g(t, x)dx.$$

DEFINITION 1.7. The differential form is called **exact** if there exists  $F \in C^1(\Omega)$  such that

$$\omega = dF = \frac{\partial F}{\partial t} dt + \frac{\partial F}{\partial x} dx.$$

The function  $F$  is called a **potential** of  $\omega$ . In this case the differential equation (1.19) is called an **exact equation**.

THEOREM 1.8 (Implicit function theorem). Suppose that  $F(t, x)$  is continuously differentiable in a neighborhood of  $(t_0, x_0) \in \mathbb{R} \times \mathbb{R}^d$  and  $F(t_0, x_0) = 0$ . Suppose that  $\det \partial F / \partial x(t_0, x_0) \neq 0$ . Then there exist  $\delta > 0$  and  $\epsilon > 0$  such that for each  $t$  satisfying  $|t - t_0| < \delta$ , there exists a unique  $x$  such that  $|x - x_0| < \epsilon$  for which  $F(t, x) = 0$ . This correspondence defines a function  $x(t)$  continuously differentiable on  $\{|t - t_0| < \delta\}$  such that

$$F(t, x) = 0 \Leftrightarrow x = x(t).$$

**THEOREM 1.9.** *Suppose that  $\omega$  is an exact form with potential  $F$  such that*

$$\det \frac{\partial F}{\partial x}(t_0, x_0) \neq 0.$$

*Then the equation  $F(t, x) = F(t_0, x_0)$  implicitly defines a function  $x \in C^1(I)$  for some open interval  $I$  containing  $t_0$ , which solves (1.19) with the initial condition  $x(t_0) = x_0$ . This solution is unique on  $I$ .*

**PROOF.** Suppose without loss of generality that  $F(t_0, x_0) = 0$ . By the **implicit function theorem**, there exist  $\delta, \eta > 0$  and  $x \in C^1(t_0 - \delta, t_0 + \delta)$  such that

$$\{(t, x) \in \Omega : |t - t_0| < \delta, |x - x_0| < \eta, F(t, x) = 0\} = \{(t, x(t)) \in \Omega : |t - t_0| < \delta\}.$$

By differentiating the identity  $F(t, x(t)) = 0$ , we get

$$0 = \frac{d}{dt}F(t, x(t)) = \frac{\partial F}{\partial t}(t, x(t)) + \frac{\partial F}{\partial x}(t, x(t)) \frac{dx}{dt} = f(t, x(t)) + g(t, x(t)) \frac{dx}{dt},$$

and hence  $x(t)$  is a solution of the differential equation. Moreover,  $x(t_0) = x_0$ .

On the other hand, if  $z \in C^1(I)$  is a solution to (1.19) such that  $z(t_0) = x_0$ , then

$$\frac{d}{dt}F(t, z(t)) = 0 \implies F(t, z(t)) = F(t_0, z(t_0)) = 0 \implies z(t) = x(t).$$

□

**DEFINITION 1.10.** *Let  $f, g \in C^1(\Omega)$ . The differential form  $\omega = fdt + gdx$  is **closed** in  $\Omega$  if*

$$\frac{\partial f}{\partial x} = \frac{\partial g}{\partial t}$$

for all  $(t, x) \in \Omega$ .

**PROPOSITION 1.11.** *An exact differential form  $\omega = fdt + gdx$  with a potential  $F \in C^2$  is closed since by Schwarz's theorem*

$$\frac{\partial^2 F}{\partial t \partial x} = \frac{\partial^2 F}{\partial x \partial t}$$

for all  $(t, x) \in \Omega$ . The converse is also true if  $\Omega$  is simply connected: If  $\omega$  is closed then  $\omega$  is exact and is associated to a potential  $F \in C^2$ .

Closed forms always have a potential (at least locally).

**EXAMPLE 1.12.** *Consider the equation*

$$tx^2 + x - t \frac{dx}{dt} = 0. \tag{1.20}$$

Here,  $f(t, x) = tx^2 + x$  and  $g(t, x) = -t$ . Since

$$\frac{\partial f}{\partial x} = 2xt + 1 \neq \frac{\partial g}{\partial t} = -1,$$

equation (1.20) is not exact.

**EXAMPLE 1.13.** *The equation*

$$t + \frac{1}{x} - \frac{t}{x^2} \frac{dx}{dt} = 0$$

is exact with the potential function  $F$  given by

$$F(t, x) = \frac{t^2}{2} + \frac{t}{x} + C, \quad C \in \mathbb{R}.$$

The equation  $F(t, x) = 0$  implicitly defines the solutions (locally for  $t \neq 0$  and  $x \neq 0$  such that  $\partial F / \partial x(t, x) \neq 0$ ).

EXAMPLE 1.14. Consider the equation

$$-2t^2 + 2x - x^2 + t(1-x)\frac{dx}{dt} = 0. \quad (1.21)$$

Here,  $f(t, x) = -2t^2 + 2x - x^2$  and  $g(t, x) = t(1-x)$ . Since

$$\frac{\partial f}{\partial x} = 2 - 2x \neq \frac{\partial g}{\partial t} = 1 - x,$$

equation (1.21) is not exact. However, multiplying (1.21) by  $t$  gives

$$-2t^3 + 2xt - tx^2 + t^2(1-x)\frac{dx}{dt} = 0.$$

We see from this that  $f(t, x) = -2t^3 + 2tx - tx^2$  and  $g(t, x) = t^2(1-x)$ . This leads to

$$\frac{\partial f}{\partial x} = 2t - 2tx, \quad \frac{\partial g}{\partial t} = 2t(1-x),$$

which satisfies the condition  $\frac{\partial f}{\partial x} = \frac{\partial g}{\partial t}$ . Thus, there must exist a function  $F(t, x)$  such that

$$\frac{\partial F}{\partial t} = f(t, x) \quad \text{and} \quad \frac{\partial F}{\partial x} = g(t, x). \quad (1.22)$$

Integrating equations (1.22) with respect to  $t$  and  $x$  and comparing the obtained formulas yields

$$F(t, x) = \frac{1}{2}t^4 - t^2x + \frac{1}{2}t^2x^2 + C,$$

for some constant  $C$ . Therefore, the differential equation (1.21) has the general solution  $F(t, x) = 0$  (locally for  $t \neq 0$  and  $x \neq 1$ ).

### 1.3.3. Hamiltonian systems.

DEFINITION 1.15. Let  $M$  be a subset of  $\mathbb{R}^d$  and let  $H : \mathbb{R}^d \times M \rightarrow \mathbb{R}$  be a  $C^1$  function.

The Hamiltonian system with Hamiltonian  $H$  is given by the first-order system of ODEs

$$\begin{cases} \frac{dp}{dt} = -\frac{\partial H}{\partial q}(p, q), \\ \frac{dq}{dt} = \frac{\partial H}{\partial p}(p, q). \end{cases} \quad (1.23)$$

EXAMPLE 1.16. An important basic example of a Hamiltonian system is the simple harmonic oscillator with Hamiltonian

$$H(p, q) = \frac{1}{2} \frac{p^2}{m} + \frac{1}{2} kq^2,$$

where  $m$  and  $k$  are positive constants. Given a potential  $V$ , Hamiltonian systems of the form

$$H(p, q) = \frac{1}{2} p^\top M^{-1} p + V(q),$$

where  $M$  is symmetric positive definite matrix and  $\top$  denotes the transpose, are widely used in **molecular and biological dynamics**.

We now introduce the notion of an invariant (also called **first integral**) for a system of ODEs.

DEFINITION 1.17. Let  $\Omega = I \times D$ , where  $I \subset \mathbb{R}$  and  $D \subset \mathbb{R}^d$ . Consider

$$\frac{dx}{dt} = f(t, x(t)), \quad (1.24)$$

where  $f : \Omega \rightarrow \mathbb{R}^d$ . We call  $F : D \rightarrow \mathbb{R}$  an **invariant** of (1.24) if  $F(x(t)) = \text{Constant}$ . A point  $(t, x) \in I \times D$  is called a **stationary point** if  $f(t, x) = 0$ .

EXAMPLE 1.18. Consider the system of **Lotka-Volterra's** ODEs given by

$$\begin{cases} \frac{du}{dt} = u(v - 2), \\ \frac{dv}{dt} = v(1 - u). \end{cases} \quad (1.25)$$

The system of ODEs (1.25) is used to describe the dynamics of biological systems in which two species interact, one as a predator and the other as prey.

Define

$$F(u, v) := \ln u - u + 2 \ln v - v.$$

$F(u, v)$  is an invariant of (1.25). In fact, by differentiating with respect to time, we have

$$\begin{aligned} \frac{d}{dt}F(u, v) &= \frac{1}{u} \frac{du}{dt} - \frac{du}{dt} + \frac{2}{v} \frac{dv}{dt} - \frac{dv}{dt} \\ &= v - 2 - \frac{du}{dt} + 2(1 - u) - \frac{dv}{dt} \\ &= (v - 2) - u(v - 2) + 2(1 - u) + v(1 - u) \\ &= (v - 2)(1 - u) + (2 - v)(1 - u) \\ &= 0. \end{aligned}$$

For the system (1.25),  $(u, v) = (1, 2)$  and  $(u, v) = (0, 0)$  are two stationary points.

LEMMA 1.19. The Hamiltonian  $H$  is an invariant of the associated Hamiltonian system (1.23).

PROOF. We have

$$\begin{aligned} \frac{d}{dt}H(p(t), q(t)) &= \frac{\partial H}{\partial p}(p(t), q(t)) \frac{dp}{dt} + \frac{\partial H}{\partial q}(p(t), q(t)) \frac{dq}{dt} \\ &= -\frac{\partial H}{\partial p}(p(t), q(t)) \frac{\partial H}{\partial q}(p(t), q(t)) + \frac{\partial H}{\partial q}(p(t), q(t)) \frac{\partial H}{\partial p}(p(t), q(t)) = 0. \end{aligned}$$

Hence,  $H(p, q)$  is an invariant of the system of equations (1.23).  $\square$

EXAMPLE 1.20. Consider the system of equations

$$\begin{cases} \frac{dp}{dt} = -\sin q, \\ \frac{dq}{dt} = p. \end{cases}$$

Here,  $H(p, q) = \frac{1}{2}p^2 - \cos q$  is the Hamiltonian of the above system, because

$$\begin{cases} \frac{\partial H}{\partial q} = \sin q = -\frac{dp}{dt}, \\ \frac{\partial H}{\partial p} = p = \frac{dq}{dt}. \end{cases}$$

There is another equivalent expression for Hamiltonian systems. Let  $x = (p, q)^\top$  (note that  $p, q \in \mathbb{R}^d$ ), and let

$$J = \begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix}, \quad (1.26)$$

where  $I$  denotes the  $d \times d$  identity matrix. Note that

$$J^{-1} = J^\top.$$

We can rewrite the Hamiltonian system (1.23) in the form

$$\frac{dx}{dt} = J^{-1} \nabla H(x). \quad (1.27)$$

Here, we use the notation  $\nabla H(x) := (\frac{\partial H}{\partial x})^\top = (\frac{\partial H}{\partial x_1}, \dots, \frac{\partial H}{\partial x_{2d}})^\top$ . For a vector function  $f : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$ ,  $f(x) = (f_1(x), \dots, f_{2d}(x))$ , we define the Jacobian matrix  $f'$  of  $f$  by

$$f'(x) := \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_{2d}} \\ \vdots & \cdots & \vdots \\ \frac{\partial f_{2d}}{\partial x_1} & \cdots & \frac{\partial f_{2d}}{\partial x_{2d}} \end{pmatrix}.$$

DEFINITION 1.21 (Symplectic linear mapping). *A matrix  $A \in \mathbb{R}^{2d} \times \mathbb{R}^{2d}$  (which is also a linear mapping from  $\mathbb{R}^{2d}$  to  $\mathbb{R}^{2d}$ ) is called **symplectic** if  $A^\top J A = J$ .*

DEFINITION 1.22 (Symplectic mapping). *A differentiable map  $g : U \rightarrow \mathbb{R}^{2d}$  is called **symplectic** if the **Jacobian matrix**  $g'(p, q)$  is everywhere symplectic, i.e., if*

$$g'(p, q)^\top J g'(p, q) = J.$$

Taking the transpose of both sides of the above equation, we also have

$$g'(p, q)^\top J^\top g'(p, q) = J^\top,$$

or equivalently,

$$g'(p, q)^\top J^{-1} g'(p, q) = J^{-1}.$$

THEOREM 1.23. *If  $g$  is a symplectic mapping, then it preserves the Hamiltonian form of the equation.*

PROOF. Let  $x = (p, q)^\top$ ,  $y = g(p, q)^\top$  and let  $G(y) := H(x)$ . By using the chain rule, we have

$$\begin{aligned} \frac{\partial}{\partial x} H(x) &= \frac{\partial}{\partial x} G(y) \\ &= \frac{\partial}{\partial y} G(y) \frac{\partial}{\partial x} y(x) \\ &= \frac{\partial}{\partial y} G(y) g'^\top(p, q). \end{aligned}$$

Then,

$$\begin{aligned} \frac{dy}{dt} &= g'^\top(p, q) \frac{dx}{dt} \\ &= g'^\top(p, q) J^{-1} \left( \frac{\partial H(x)}{\partial x} \right)^\top \\ &= g'^\top J^{-1} g' \nabla_y G(y) \\ &= J^{-1} \nabla_y G(y), \end{aligned}$$

and therefore,

$$\frac{dy}{dt} = J^{-1} \nabla_y G(y).$$

□

DEFINITION 1.24 (Flow). *We define the **flow** by  $\phi_t(p_0, q_0) = (p(t, p_0, q_0), q(t, p_0, q_0))$ ,  $\phi_t : U \rightarrow \mathbb{R}^{2d}$ ,  $U \subset \mathbb{R}^{2d}$ , and  $p_0$  and  $q_0$  are the initial data at  $t = 0$ .*

THEOREM 1.25 (Poincaré's theorem). *Suppose that  $H$  is twice differentiable. Then the flow  $\phi_t$  is a symplectic transformation whenever it is defined.*

PROOF. Let  $y_0 = (p_0, q_0)$ . Note that

$$\frac{d}{dt} \left( \frac{\partial \phi_t}{\partial y_0} \right) = J^{-1} \nabla^2 H(\phi_t(y_0)) \frac{\partial \phi_t}{\partial y_0}.$$

Then we have

$$\begin{aligned}
& \frac{d}{dt} \left( \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top J \left( \frac{\partial \phi_t}{\partial y_0} \right) \right) \\
&= \left( \frac{\partial \phi_t}{\partial y_0} \right)'^\top J \left( \frac{\partial \phi_t}{\partial y_0} \right) + \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top J \left( \frac{\partial \phi_t}{\partial y_0} \right)' \\
&= \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top \nabla^2 H J^{-\top} J \left( \frac{\partial \phi_t}{\partial y_0} \right) + \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top J J^{-1} \nabla^2 H \left( \frac{\partial \phi_t}{\partial y_0} \right) \\
&= 0,
\end{aligned}$$

where  $\nabla^2 H$  is the Hessian matrix of  $H(p, q)$  (and is symmetric). Moreover, since  $\partial \phi_t / \partial y_0$  at  $t = 0$  is the identity map, the identity

$$\left( \frac{\partial \phi_t}{\partial y_0} \right)^\top J \left( \frac{\partial \phi_t}{\partial y_0} \right) = J$$

is satisfied for all  $t$  and all  $(p_0, q_0)$  as long as the solution remains in the domain of definition of  $H$ .  $\square$

The following result shows that the symplecticity of the flow is a characteristic property of the Hamiltonian system.

**THEOREM 1.26.** *Let  $f : U \rightarrow \mathbb{R}^{2d}$  be continuously differentiable. Then  $\frac{dx}{dt} = f(x)$  is locally Hamiltonian if and only if  $\phi_t(x)$  is symplectic for all  $x \in U$  and for all sufficiently small  $t$ .*

**PROOF.** The necessity follows from Theorem 1.25. We therefore suppose that  $\phi_t$  is symplectic, and we have to prove the local existence of a Hamiltonian  $H$  such that  $f(x) = J^{-1} \nabla H(s)$ . Using the fact that  $\frac{\partial \phi_t}{\partial y_0}$  is a solution of

$$\frac{dy}{dt} = f'(\phi_t(y_0))y,$$

we obtain

$$\frac{d}{dt} \left( \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top J \left( \frac{\partial \phi_t}{\partial y_0} \right) \right) = \left( \frac{\partial \phi_t}{\partial y_0} \right)^\top [f'(\phi_t(y_0))^\top J + Jf'] \left( \frac{\partial \phi_t}{\partial y_0} \right) = 0.$$

Putting  $t = 0$ , it follows from  $J = -J^\top$  that  $Jf'(y_0)$  is a symmetric matrix for all  $y_0$ . The **integrability lemma** below shows that  $Jf(y)$  can be written as the gradient of a function  $H$ .  $\square$

**LEMMA 1.27 (Integrability lemma).** *Let  $D \subset \mathbb{R}^{2d}$  be an open set and let  $g : D \rightarrow \mathbb{R}^{2d}$  be of class  $C^1$ . Suppose that the Jacobian  $g'(y)$  is symmetric for all  $y \in D$ . Then, for every  $y_0 \in D$ , there exists a neighborhood of  $y_0$  and a function  $H(y)$  such that*

$$g(y) = \nabla H(y)$$

on this neighborhood.

**PROOF.** Suppose that  $y_0 = 0$ , and consider a ball around  $y_0$  which is contained in  $D$ . On this ball we define

$$H(y) = \int_0^1 y^\top g(ty) dt.$$

Differentiating with respect to  $y_k$ , and using the symmetry assumption

$$\frac{\partial g_i}{\partial y_k} = \frac{\partial g_k}{\partial y_i}$$

yields

$$\begin{aligned}\frac{\partial H}{\partial y_k} &= \int_0^1 (g_k(ty) + y^\top \frac{\partial g}{\partial y_k}(ty)t) dt \\ &= \int_0^1 \frac{d}{dt}(tg_k(ty)) dt = g_k(y),\end{aligned}$$

which proves that

$$\nabla H = g.$$

□

**1.3.4. Gradient systems.** Finally, consider the gradient systems.

DEFINITION 1.28. *Gradient systems are differential equations that have the form*

$$\frac{dx}{dt} = -\nabla V(x), \tag{1.28}$$

with  $V$  (called the potential function) being a real-valued function.

In order to guarantee that the right-hand side of (1.28) is a continuously differentiable function of  $x$ , one requires that  $V$  is twice-continuously differentiable.

On solutions to (1.28) one has

$$\frac{d}{dt}V(x(t)) = \nabla V(x(t)) \cdot \frac{dx}{dt} = -|\nabla V(x)|^2.$$

A differential equation

$$\frac{dx}{dt} = f(x) = (f_1(x), \dots, f_d(x)) \tag{1.29}$$

is a gradient system if and only if there exists a scalar-valued function  $V(x)$  so that

$$-(f_1(x), \dots, f_d(x)) = \left( \frac{\partial V}{\partial x_1}(x), \dots, \frac{\partial V}{\partial x_d}(x) \right).$$

In dimension  $d = 1$ , one can always choose an antiderivative  $V$  of  $-f$  so that

$$\frac{dV}{dx}(x) = -f(x).$$

Equation (1.29) is always a gradient system in dimension one.

In dimension two, a system

$$\begin{cases} \frac{dx_1}{dt} = f_1(x_1, x_2), \\ \frac{dx_2}{dt} = f_2(x_1, x_2), \end{cases} \tag{1.30}$$

is a gradient system if and only if there is a potential  $V(x_1, x_2)$  so that

$$\frac{\partial V}{\partial x_1} = -f_1, \quad \frac{\partial V}{\partial x_2} = -f_2. \tag{1.31}$$

A necessary and sufficient condition for solvability of (1.31) is the equality of mixed partials,

$$\frac{\partial f_1}{\partial x_2} = \frac{\partial f_2}{\partial x_1}.$$

In the general case, the necessary and sufficient condition is again equality of mixed partials expressed as

$$\frac{\partial f_i}{\partial x_j} = \frac{\partial f_j}{\partial x_i} \quad \text{for all } 1 \leq i < j \leq d.$$

LEMMA 1.29. *The Hamiltonian system (1.23) is a **gradient system** if and only if the function  $H$  is harmonic.*



PROOF. Suppose that  $H$  is harmonic, *i.e.*,

$$\frac{\partial^2 H}{\partial p^2} + \frac{\partial^2 H}{\partial q^2} = 0.$$

Then the Jacobian of  $J^{-1}\nabla H$  given by

$$(J^{-1}\nabla H)' = \begin{pmatrix} -\frac{\partial^2 H}{\partial p \partial q} & -\frac{\partial^2 H}{\partial q^2} \\ \frac{\partial^2 H}{\partial p^2} & \frac{\partial^2 H}{\partial p \partial q} \end{pmatrix}$$

is symmetric. The integrability lemma shows that there exists  $V$  such that  $J^{-1}\nabla H = \nabla V$  and therefore, the Hamiltonian system is a gradient system.

Suppose that the Hamiltonian system is a gradient system. Then, there exists  $V$  such that

$$\frac{\partial V}{\partial p} = \frac{\partial H}{\partial q} \quad \text{and} \quad \frac{\partial V}{\partial q} = -\frac{\partial H}{\partial p}.$$

Therefore,

$$\Delta H := \frac{\partial^2 H}{\partial p^2} + \frac{\partial^2 H}{\partial q^2} = 0.$$

□

EXAMPLE 1.30. *The Hamiltonian system with  $H(p, q) = p^2 - q^2$  is a gradient system.*

**1.3.5. Hamilton-Jacobi equation.** The **Hamilton-Jacobi equation** is used to generate particular symplectic transformations that simplify Hamiltonian systems.

Let  $d = 1$  and let

$$H(p, q) = \frac{1}{2}p^2 + V(q).$$

Consider the Hamiltonian-Jacobi equation

$$\begin{cases} \frac{\partial u}{\partial t} + H\left(\frac{\partial u}{\partial q}, q\right) = 0, & q \in \mathbb{R}, t \geq 0, \\ u(0, q) = u_0(q), & q \in \mathbb{R}. \end{cases} \quad (1.32)$$

A smooth function  $u(P, q, t)$  satisfying (1.32) can be used to map the variables  $(p, q)$  to a set of variables  $(P, Q)$  that are constants over time. Let  $p = \frac{\partial u}{\partial q}$ , and define  $Q = \frac{\partial u}{\partial P}$ . Then,  $(p, q) \mapsto (P, Q)$  is symplectic. Moreover, in the new coordinates  $(P, Q)$ , the Hamiltonian system (1.23) reduces to

$$\begin{cases} \frac{dP}{dt} = 0, \\ \frac{dQ}{dt} = 0, \end{cases} \quad (1.33)$$

and becomes trivial to solve.

#### 1.4. Problems

PROBLEM 1.31 (**Exact equations**). *Consider the equation  $F(t, x) = 0$ , where  $F \in C^2(\mathbb{R}^2, \mathbb{R})$ . Suppose  $x(t)$  solves this equation.*

(i) *Show that  $x(t)$  satisfies*

$$g(t, x) \frac{dx}{dt} + f(t, x) = 0, \quad (1.34)$$

where

$$g(t, x) = \frac{\partial F(t, x)}{\partial x} \quad \text{and} \quad f(t, x) = \frac{\partial F(t, x)}{\partial t}.$$

(ii) *Show that we have*

$$\frac{\partial g(t, x)}{\partial t} = \frac{\partial f(t, x)}{\partial x},$$

and deduce that (1.34) is exact.

- (iii) *Conversely, show that if a first-order equation as (1.34) is exact, then there is a corresponding function  $F$  as above. Find an explicit formula for  $F$  in terms of  $f$  and  $g$ . Is  $F$  uniquely determined by  $f$  and  $g$ ?*
- (iv) *Show that*

$$(2tx + 3t + 5)\frac{dx}{dt} + 3t^2 + t + x^2 + 3x = 0$$

*is exact. Find  $F$  and find the solution.*

**PROBLEM 1.32 (Method of integrating factor).** *Consider*

$$g(t, x)\frac{dx}{dt} + f(t, x) = 0.$$

- (i) *Prove that  $\mu(t, x)$  is an integrating factor if*

$$\mu(t, x)g(t, x)\frac{dx}{dt} + \mu(t, x)f(t, x) = 0$$

*is exact.*

- (ii) *Consider*

$$t\frac{dx}{dt} + 3t - 2\frac{dx}{dt} = 0$$

*and look for an integrating factor  $\mu$  depending only on  $t$ . Solve the equation.*

**PROBLEM 1.33.** (i) *Prove that a smooth differential map  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is symplectic if and only if  $\det g' = 1$ .*

- (ii) *Find a counterexample to the statement in (i) in  $\mathbb{R}^{2d}$  for  $d > 1$ .*

**PROBLEM 1.34.** *Consider the system of linear equations*

$$\begin{cases} \frac{dX}{dt} = AX(t), \\ X(0) = X_0, \end{cases}$$

*where  $X, X_0$ , and  $A$  are  $d \times d$  real matrices.*

- (i) *Prove that if  $A$  is a skew-symmetric matrix then  $X^\top X$  is an invariant of the system.*
- (ii) *Prove that if  $X_0$  is orthogonal then the solution  $X(t)$  is orthogonal for all  $t \geq 0$ .*

**PROBLEM 1.35 (Transport theorem).** *Let  $\phi_t$  denote the flow of the system  $dx/dt = f(x)$ ,  $x \in \mathbb{R}^d$ , and let  $\Omega$  be a bounded domain in  $\mathbb{R}^d$ . Define*

$$V(t) = \int_{\phi_t(\Omega)} dx_1 \dots dx_d,$$

*and recall that the divergence of a vector field  $f = (f_1, \dots, f_d)^\top$  is*

$$\nabla \cdot f = \sum_{i=1}^d \frac{\partial f_i}{\partial x_i}.$$

- (i) *Use Liouville's theorem and the change of variables formula for multiple integrals to prove that*

$$\frac{dV}{dt} = \int_{\phi_t(\Omega)} (\nabla \cdot f) dx_1 \dots dx_d.$$

- (ii) *Prove that the flow of a vector field whose divergence is everywhere negative contracts volume.*
- (iii) *Suppose that  $g : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable. Prove the transport theorem*

$$\frac{d}{dt} \int_{\phi_t(\Omega)} g(t, x) dx_1 \dots dx_d = \int_{\phi_t(\Omega)} \left[ \frac{\partial g}{\partial t} + \nabla \cdot (gf) \right] dx_1 \dots dx_d.$$

## CHAPTER 2

# Existence, uniqueness, and regularity in the Lipschitz case

### 2.1. Banach fixed point theorem

**DEFINITION 2.1 (Contraction).** Let  $(X, d)$  be a metric space. A mapping  $F : X \rightarrow X$  is a **contraction** if there exists  $0 < \lambda < 1$  such that

$$d(F(x), F(y)) \leq \lambda d(x, y)$$

for all  $x, y \in X$ .

**THEOREM 2.2 (Banach fixed point theorem).** Let  $(X, d)$  be a **complete** metric space (i.e., every Cauchy sequence of elements of  $X$  is convergent) and let  $F : X \rightarrow X$  be a contraction. Then there exists a unique  $x \in X$  such that

$$F(x) = x.$$

### 2.2. Gronwall's lemma

**LEMMA 2.3 (Gronwall's lemma).** Let  $I = [0, T]$  and let  $\phi \in \mathcal{C}^0(I)$ . If there exist two constants  $\alpha, \beta \in \mathbb{R}$ ,  $\beta \geq 0$ , such that

$$\phi(t) \leq \alpha + \beta \int_0^t \phi(s) ds \quad \text{for all } t \in I, \tag{2.1}$$

then

$$\phi(t) \leq \alpha e^{\beta t} \quad \text{for all } t \in I.$$

**PROOF.** Let  $\varphi : I \rightarrow \mathbb{R}$  be the function

$$\varphi(t) := \alpha + \beta \int_0^t \phi(s) ds.$$

Since  $\phi \in \mathcal{C}^0$ , we conclude that  $\varphi \in \mathcal{C}^1$ , and

$$\frac{d\varphi}{dt} = \beta\phi(t) \quad \text{for all } t \in I.$$

By using (2.1), it follows that

$$\frac{d\varphi}{dt} \leq \beta\varphi.$$

Let  $\psi(t) := \exp(-\beta t)\varphi(t)$  for  $t \in I$ . Then

$$\begin{aligned} \frac{d\psi}{dt} &= -\beta e^{-\beta t}\varphi(t) + e^{-\beta t} \frac{d\varphi}{dt} \\ &= e^{-\beta t} \left( -\beta\varphi(t) + \frac{d\varphi}{dt} \right) \leq 0. \end{aligned}$$

Since  $\psi(0) = \varphi(0) = \alpha$ , we have  $\psi(t) \leq \alpha$  for  $t \in I$ , and hence

$$\varphi(t) \leq \alpha e^{\beta t},$$

which implies that  $\phi(t) \leq \varphi(t) \leq \alpha e^{\beta t}$  for all  $t \in I$ . □

### 2.3. Cauchy-Lipschitz theorem

Let  $I = [0, T]$ , let  $d$  be a positive integer, and let  $f : I \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ . Suppose that  $f \in \mathcal{C}^0(I \times \mathbb{R}^d)$ .

**DEFINITION 2.4** (Lipschitz condition). *If there exists a constant  $C_f \geq 0$  such that, for any  $x_1, x_2 \in \mathbb{R}^d$  and any  $t \in I$ , the following inequality holds:*

$$|f(t, x_1) - f(t, x_2)| \leq C_f |x_1 - x_2|, \quad (2.2)$$

*then we say that  $f$  satisfies a **Lipschitz condition** on  $I$ . The constant  $C_f$  is called the **Lipschitz constant** for  $f$ .*

**THEOREM 2.5** (Cauchy-Lipschitz theorem). *Consider the initial value problem*

$$\begin{cases} \frac{dx}{dt} = f(t, x), & t \in [0, T], \\ x(0) = x_0, & x_0 \in \mathbb{R}^d. \end{cases} \quad (2.3)$$

*If  $f \in \mathcal{C}^0(I \times \mathbb{R}^d)$  satisfies the Lipschitz condition (2.2) on  $[0, T]$ , then there exists a unique solution  $x \in \mathcal{C}^1(I)$  to (2.3) on  $[0, T]$ .*

**PROOF.** By (2.3), we have

$$x(t) = x_0 + \int_0^t f(s, x(s)) ds, \quad \forall t \in [0, T].$$

Define the functional  $F : \mathcal{C}^0([0, T]; \mathbb{R}^d) \rightarrow \mathcal{C}^0([0, T]; \mathbb{R}^d)$  by

$$F(y) := x_0 + \int_0^t f(s, y(s)) ds.$$

For  $y \in \mathcal{C}^0([0, T]; \mathbb{R}^d)$ , defined the norm of  $y$  by

$$\|y\| := \sup_{t \in [0, T]} \{|y(t)| e^{-C_f t}\}, \quad (2.4)$$

where  $C_f$  is the Lipschitz constant for  $f$ . It is easy to prove that (2.4) is equivalent to the usual norm  $\sup_{t \in [0, T]} |y(t)|$  and hence,  $\mathcal{C}^0([0, T]; \mathbb{R}^d)$  equipped with (2.4) is complete.

With (2.4), we compute

$$\begin{aligned} \|F[y_1] - F[y_2]\| &= \sup_{t \in [0, T]} |F[y_1](t) - F[y_2](t)| e^{-C_f t} \\ &\leq \sup_{t \in [0, T]} e^{-C_f t} \int_0^t |f(s, y_1(s)) - f(s, y_2(s))| ds \\ &\leq \sup_{t \in [0, T]} e^{-C_f t} C_f \int_0^t |y_1(s) - y_2(s)| ds \\ &\leq \sup_{t \in [0, T]} e^{-C_f t} C_f \int_0^t e^{C_f s} e^{-C_f s} |y_1(s) - y_2(s)| ds \\ &\leq \sup_{t \in [0, T]} \{e^{-C_f t} C_f \int_0^t e^{C_f s} ds\} \|y_1 - y_2\| \\ &\leq (1 - e^{-C_f T}) \|y_1 - y_2\|. \end{aligned}$$

By Banach fixed point theorem in a complete metric space (Theorem 2.2), there exists a unique  $y \in \mathcal{C}^0([0, T]; \mathbb{R}^d)$  such that  $F(y) = y$ . The **Picard iteration**

$$y^{(n+1)} = F[y^{(n)}]$$

is a Cauchy sequence and converges to the unique fixed point  $y$ . Therefore, there exists a unique solution to (2.3).  $\square$

REMARK 2.6. *Theorem 2.5 holds true if  $\mathbb{R}^d$  is replaced with a **Banach space** (a complete normed vector space). The proof is the same.*

If  $f$  is continuous, there is no guarantee that the initial value problem (2.3) possesses a unique solution.

EXAMPLE 2.7. *Consider*

$$\frac{dx}{dt} = x^{\frac{2}{3}}, \quad x(0) = 0. \quad (2.5)$$

*Then there are two solutions to (2.5) given by  $x_1(t) = \frac{t^3}{27}$  and  $x_2(t) = 0$ .*

THEOREM 2.8 (Cauchy-Peano existence theorem). *If  $f$  is continuous, then (2.3) admits a solution  $x(t)$  that is, at least, defined for small  $t$ .*

This theorem can be proved by using the **Arzela-Ascoli theorem**.

DEFINITION 2.9 (Equicontinuity). *A family of functions  $\mathcal{F}$  is said to be **equicontinuous** on  $[a, b]$  if for any given  $\epsilon > 0$ , there exists  $\delta > 0$  such that*

$$|f(t) - f(s)| < \epsilon$$

*whenever  $|t - s| < \delta$  for every function  $f \in \mathcal{F}$  and  $t, s \in [a, b]$ .*

DEFINITION 2.10 (Uniform boundedness). *A family of continuous functions  $\mathcal{F}$  on  $[a, b]$  is **uniformly bounded** if there exists a positive number  $M$  such that  $|f(t)| \leq M$  for every function  $f \in \mathcal{F}$  and  $t \in [a, b]$ .*

THEOREM 2.11 (Arzela-Ascoli). *Suppose that the sequence of functions  $\{f_n(t)\}_{n \in \mathbb{N}}$  on  $[a, b]$  is uniformly bounded and equicontinuous, then there exists a subsequence  $\{f_{n_k}(t)\}_{k \in \mathbb{N}}$  that is uniformly convergent on  $[a, b]$ .*

EXAMPLE 2.12. *Consider*

$$\frac{dx}{dt} = x^2, \quad x(0) = x_0 \neq 0.$$

*By separation of variables, we obtain*

$$\frac{dx}{x^2} = dt.$$

*Thus,*

$$-\frac{1}{x} = \int \frac{dx}{x^2} = t + C,$$

*and hence,*

$$x = -\frac{1}{t + C}.$$

*Since  $x(0) = x_0$ ,*

$$x(t) = \frac{x_0}{1 - x_0 t}.$$

*If  $x_0 > 0$ ,  $x(t)$  blows up when  $t \rightarrow \frac{1}{x_0}$  from below. If  $x_0 < 0$ , the singularity is in the past ( $t < 0$ ). The only solution defined for all positive and negative  $t$  is the constant solution  $x(t) = 0$ , corresponding to  $x_0 = 0$ .*

REMARK 2.13 (**Local existence and uniqueness theorem**). *If  $f(t, x)$  satisfies a Lipschitz condition in a bounded domain, then a unique solution exists in a limited region.*

THEOREM 2.14. *Let  $x_0 \in \mathbb{R}$ . Assume that  $f$  is continuous and satisfies the Lipschitz condition in the closed domain  $K := \{|x - x_0| \leq k\}$  and  $t \in [0, T]$ ,*

$$|f(t, x_1) - f(t, x_2)| \leq C_f |x_1 - x_2|, \quad \text{for all } x, y \in K, t \in [0, T],$$

then the equation

$$\begin{cases} \frac{dx}{dt} = f(t, x), & t \in [0, T], \\ x(0) = x_0, \end{cases}$$

has a unique solution in  $t \in [0, \min\{T, \frac{k}{M}\}]$ , where

$$M := \sup_{x \in K, t \in [0, T]} |f(t, x)|.$$

EXAMPLE 2.15. *The initial value problem*

$$\begin{cases} \frac{dx}{dt} = 1 + x^2, & t \in [0, 1], \\ x(0) = 0, \end{cases}$$

in the region  $\{(x, t) : |x| \leq 1, 0 \leq t \leq 1\}$  has a unique solution for  $0 \leq t \leq 1/2$ .

Now we turn to the continuity of the solution of (2.3).

**THEOREM 2.16 (Continuity with respect to the initial data).** *Suppose that  $f$  satisfies the Lipschitz condition (2.2). Let  $x_1(t)$  and  $x_2(t)$  be the solutions of (2.3) corresponding to the initial data  $x_1(0)$  and  $x_2(0)$ , respectively. Then we have*

$$|x_1(t) - x_2(t)| \leq e^{C_f t} |x_1(0) - x_2(0)| \quad \text{for all } t \in [0, T]. \quad (2.6)$$

PROOF. Since

$$\begin{aligned} \frac{d}{dt} |x_1(t) - x_2(t)|^2 &= 2(f(t, x_1(t)) - f(t, x_2(t)))(x_1(t) - x_2(t)) \\ &\leq 2C_f |x_1(t) - x_2(t)|^2, \quad t \in [0, T], \end{aligned}$$

we have

$$\frac{d}{dt} \left( |x_1(t) - x_2(t)|^2 e^{-2C_f t} \right) \leq 0. \quad (2.7)$$

Integrating (2.7) from 0 to  $t$  gives

$$|x_1(t) - x_2(t)|^2 e^{-2C_f t} \leq |x_1(0) - x_2(0)|^2,$$

or equivalently,

$$|x_1(t) - x_2(t)| \leq |x_1(0) - x_2(0)| e^{C_f t},$$

which yields the desired inequality.  $\square$

Next we discuss the differentiability of the solution of (2.3) with respect to the initial data.

Formally, taking the derivative of the solution  $x$  of (2.3) with respect to the initial data, we obtain that  $\partial x(t)/\partial x_0$  is the solution of the linear equation

$$\begin{cases} \frac{d}{dt} \frac{\partial x(t)}{\partial x_0} = \frac{\partial f}{\partial x}(t, x(t)) \frac{\partial x(t)}{\partial x_0}, \\ \frac{\partial x(t)}{\partial x_0} = 1. \end{cases} \quad (2.8)$$

**THEOREM 2.17.** *Suppose that  $f$  is of class  $C^1$ . Then  $x_0 \mapsto x(t)$  is differentiable and  $\partial x(t)/\partial x_0$  is the unique solution of the linear equation (2.8).*

PROOF. Let  $\Delta x(t, x_0, h) := x(t, x_0 + h) - x(t, x_0)$  be the difference quotient. By using the **mean-value theorem**, we have

$$\begin{aligned}\Delta x(t, x_0, h) &= h + \int_0^t (f(s, x(s, x_0 + h)) - f(s, x(s, x_0))) ds \\ &= h + \int_0^t (f(s, x(s, x_0) + \Delta x(s, x_0, h)) - f(s, x(s, x_0))) ds \\ &= h + \int_0^t \frac{\partial f}{\partial x}(s, x(s, x_0) + \tau \Delta x) \Delta x ds,\end{aligned}$$

where  $\tau = \tau(s, x_0, h) \in [0, 1]$ . Since there exists a positive constant  $M$  such that  $|\frac{\partial f}{\partial x}| \leq M$ , it holds that

$$|\Delta x| \leq |h| + M \int_0^t |\Delta x(s, x_0, h)| ds,$$

By Gronwall's lemma (Lemma 2.3),

$$|\Delta x(t, x_0, h)| \leq |h| e^{MT}.$$

Let  $v(t)$  be the unique solution of (2.8). We compute

$$\begin{aligned}\frac{\Delta x(t, x_0, h)}{h} - v(t) &= \int_0^t \left( \frac{f(s, x(s, x_0 + h)) - f(s, x(s, x_0))}{h} - \frac{\partial f}{\partial x}(s, x(s, x_0)) v(s) \right) ds \\ &= \int_0^t \frac{\Delta x(s, x_0, h)}{h} \left[ \frac{\partial f}{\partial x}(s, x(s, x_0) + \tau \Delta x(s, x_0, h)) - \frac{\partial f}{\partial x}(s, x(s, x_0)) \right] ds \\ &\quad + \int_0^t \frac{\partial f}{\partial x}(s, x(s, x_0)) \left( \frac{\Delta x(s, x_0, h)}{h} - v(s) \right) ds.\end{aligned}$$

By using the uniform continuity of  $\frac{\partial f}{\partial x}$ , we have that for any  $\epsilon > 0$  there exists  $h_0 > 0$  such that, for any  $|h| \leq h_0$ , the first term on the right-hand side is of order  $O(\epsilon)$ . Then, again by Gronwall's lemma, there exists a positive constant  $M'$  such that

$$\left| \frac{\Delta x(t, x_0, h)}{h} - v \right| \leq M' \epsilon e^{MT},$$

for  $|h|$  small enough, which proves that  $x_0 \mapsto x(t)$  is differentiable and its derivative is given by

$$\frac{\partial x}{\partial x_0} = v,$$

where  $v$  is the solution of (2.8). □

## 2.4. Stability

**THEOREM 2.18** (Strong continuity theorem). *Let*

$$\frac{dx}{dt} = f(t, x) \quad \text{and} \quad \frac{dy}{dt} = g(t, y)$$

*be two ODEs on  $[0, T]$ . If  $f$  satisfies the Lipschitz condition (2.2) on  $[0, T]$  and there exists  $\epsilon > 0$  such that, for any  $x \in \mathbb{R}^d$ ,  $t \in [0, T]$ ,*

$$|f(t, x) - g(t, x)| \leq \epsilon,$$

*then the following inequality holds:*

$$|x(t) - y(t)| \leq |x(0) - y(0)| e^{C_f t} + \frac{\epsilon}{C_f} (e^{C_f t} - 1), \quad t \in [0, T].$$

**REMARK 2.19.** *The function  $g$  may not satisfy a Lipschitz condition.*

PROOF. Since

$$\begin{aligned} \frac{d}{dt}|x(t) - y(t)|^2 &= 2(f(t, x(t)) - g(t, y(t)))(x(t) - y(t)) \\ &= 2(f(t, x(t)) - f(t, y(t)))(x(t) - y(t)) + 2(f(t, y(t)) - g(t, y(t)))(x(t) - y(t)), \end{aligned}$$

we have

$$\begin{aligned} \frac{d}{dt}|x(t) - y(t)|^2 &\leq \left| \frac{d}{dt}|x(t) - y(t)|^2 \right| \\ &\leq 2|f(t, x(t)) - f(t, y(t))| |x(t) - y(t)| + 2|f(t, y(t)) - g(t, y(t))| |x(t) - y(t)| \\ &\leq 2C_f |x(t) - y(t)|^2 + 2\epsilon |x(t) - y(t)| \\ &\leq 2C_f |x(t) - y(t)|^2 + 2\epsilon \sqrt{|x(t) - y(t)|^2}. \end{aligned}$$

If we denote by  $h(t) := |x(t) - y(t)|^2$ , then

$$\frac{dh}{dt} \leq 2C_f h + 2\epsilon \sqrt{h}.$$

Consider the following initial value problem:

$$\begin{cases} \frac{du}{dt} = 2C_f u + 2\epsilon \sqrt{u}, \\ u(0) = |x(0) - y(0)|^2. \end{cases} \quad (2.9)$$

Since  $C_f > 0$ ,  $u(0) \geq 0$ , it follows that  $\frac{du}{dt}$  is always non-negative when  $t \geq 0$ , and hence  $u$  is increasing.

Let  $z(t) := \sqrt{u(t)}$  and suppose that  $h(0) > 0$ . Then (2.9) is equivalent to

$$\begin{cases} \frac{dz}{dt} - C_f z = \epsilon, & t \in [0, T], \\ z(0) = \sqrt{u(0)}. \end{cases}$$

This gives the solution of (2.4):

$$\sqrt{u(t)} = z(t) = \sqrt{u(0)}e^{C_f t} + \frac{\epsilon}{C_f}(e^{C_f t} - 1).$$

Moreover,

$$\begin{aligned} \frac{d}{dt}(h(t) - u(t)) &\leq 2C_f(h(t) - u(t)) + 2\epsilon(\sqrt{h(t)} - \sqrt{u(t)}) \\ &= 2C_f(h(t) - u(t)) + 2\epsilon \frac{h(t) - u(t)}{\sqrt{h(t)} + \sqrt{u(t)}}. \end{aligned}$$

Suppose that there exists  $t_1$  such that  $h(t_1) > u(t_1)$ . Let  $t_0 := \sup\{t : 0 \leq t \leq t_1, h(t) \leq u(t)\}$ . By the continuity of  $h$  and  $u$ , we must have  $h(t_0) = u(t_0)$ . Since  $u(t_0) > 0$ , we obtain for  $t_0 \leq t \leq t_1$ , that

$$\begin{aligned} \frac{d}{dt}(h(t) - u(t)) &\leq 2C_f(h(t) - u(t)) + 2\epsilon \frac{h(t) - u(t)}{\sqrt{u(0)}} \\ &= (2C_f + \frac{2\epsilon}{\sqrt{u(0)}})(h(t) - u(t)). \end{aligned}$$

Hence,

$$\frac{d}{dt} \left( (h(t) - u(t)) \exp\left(-\left(2C_f + \frac{2\epsilon}{\sqrt{u(0)}}\right)t\right) \right) \leq 0.$$

Integrating from  $t_0$  to  $t$  gives  $h(t) \leq u(t)$  for  $t_0 \leq t \leq t_1$ , which is a contradiction to  $h(t_1) > u(t_1)$ .

Therefore, it follows that for all  $t \in [0, T]$ ,

$$\frac{d}{dt} \left( (h(t) - u(t)) \exp\left(-\left(2C_f + \frac{2\epsilon}{\sqrt{u(0)}}\right)t\right) \right) \leq 0.$$



By integrating now the last inequality from 0 to  $t$ , we obtain

$$(h(t) - u(t)) \exp\left(-\left(2C_f + \frac{2\epsilon}{\sqrt{u(0)}}\right)t\right) \leq h(0) - u(0).$$

Since  $u(0) = h(0)$ , we have  $h(t) \leq u(t)$  for  $t \in [0, T]$ , and hence

$$\begin{aligned} |x(t) - y(t)| &\leq \sqrt{u(t)} \\ &= \sqrt{u(0)}e^{C_f t} + \frac{\epsilon}{C_f}(e^{C_f t} - 1) \\ &= \sqrt{h(0)}e^{C_f t} + \frac{\epsilon}{C_f}(e^{C_f t} - 1). \end{aligned}$$

Therefore, the desired estimate

$$|x(t) - y(t)| \leq |x(0) - y(0)|e^{C_f t} + \frac{\epsilon}{C_f}(e^{C_f t} - 1)$$

holds.

If  $h(0) = 0$ , then, instead of (2.9), we consider the following equation:

$$\begin{cases} \frac{du_n}{dt} = 2C_f u_n + 2\epsilon\sqrt{u_n}, & t \in [0, T], \\ u_n(0) = \frac{1}{n}, \end{cases} \quad (2.10)$$

which, analogously to (2.9), has the explicit solution

$$u_n(t) = \left[ \frac{1}{\sqrt{n}}e^{C_f t} + \frac{\epsilon}{C_f}(e^{C_f t} - 1) \right]^2.$$

We only need to prove that for each  $n \in \mathbb{N}$ ,

$$h(t) \leq u_n(t) \quad (2.11)$$

holds for  $t \in [0, T]$ . Then by letting  $n \rightarrow +\infty$ ,  $u_n \rightarrow u$ , where  $u$  is the solution to (2.9), and hence  $h(t) \leq u(t)$ .

Inequality (2.11) can be proved by contradiction. Suppose that there exists  $t_1 > 0$  such that  $h(t_1) > u_n(t_1)$ . Let  $t_0$  be the largest  $t$  in the interval  $0 < t \leq t_1$  such that  $h(t_0) \leq u_n(t_0)$ . By the continuity of  $h(t)$  and  $u_n(t)$ , we assert that

$$h(t_0) = u_n(t_0) > 0,$$

and  $h(t) > u_n(t)$  on  $(t_0, t_0 + \epsilon)$ , a small right-neighborhood of  $t_0$ . But this is impossible according to the discussion in the case where  $h(0) > 0$  by replacing 0 by  $t_0$ . The proof of the theorem is now complete.  $\square$

## 2.5. Regularity

**THEOREM 2.20.** *If  $f \in \mathcal{C}^n$  for  $n \geq 0$ , then the solution  $x$  of (2.3) is of class  $\mathcal{C}^{n+1}$ .*

**PROOF.** The proof is by induction, the case  $n = 0$  being clear. If  $f \in \mathcal{C}^n$  then  $x$  is at least of class  $\mathcal{C}^n$ , by the inductive assumption. Then the function  $t \mapsto f(t, x(t)) = dx(t)/dt$  is also of class  $\mathcal{C}^n$ . The function  $x(t)$  is then of class  $\mathcal{C}^{n+1}$ .  $\square$

**REMARK 2.21.** *If  $f$  is a real analytic function, then it can be proved that  $x$  is also real analytic.*

### 2.6. Problems

PROBLEM 2.22 (**Generalized Gronwall's inequality**). Suppose  $\phi(t)$  satisfies

$$\phi(t) \leq \alpha(t) + \int_0^t \beta(s)\phi(s) ds \quad \text{for all } t \in [0, T],$$

with  $\alpha(t) \in \mathbb{R}$  and  $\beta(t) \geq 0$ .

(i) Prove that

$$\phi(t) \leq \alpha(t) + \int_0^t \alpha(s)\beta(s)e^{\int_s^t \beta(\tau) d\tau} ds.$$

(ii) Prove that, if in addition  $\alpha(s) \leq \alpha(t)$  for  $s \leq t$ , then

$$\phi(t) \leq \alpha(t)e^{\int_0^t \beta(s) ds}, \quad \text{for all } t \in [0, T].$$

PROBLEM 2.23. Let  $d = 1$  and let  $f(t, x)$  be a continuous function satisfying the Lipschitz condition (2.2). Let  $M := \sup_{x \in \mathbb{R}, t \in [0, T]} |f(t, x)|$ . Let  $x$  be the solution to (2.3) and let  $x^{(n)}$  be the  $n$ th term in its Picard's approximation. Prove that

$$|x(t) - x^{(n)}(t)| \leq \frac{MC_f^n}{(n+1)!} t^{n+1} \quad \text{for } t \in [0, T].$$

PROBLEM 2.24. State and prove a uniqueness theorem for the differential equation

$$\begin{cases} \frac{d^2x}{dt^2} = f(t, x, \frac{dx}{dt}), & t \in [0, T], \\ x(0) = x_0, \quad \frac{dx}{dt}(0) = x'_0, & x_0, x'_0 \in \mathbb{R}. \end{cases}$$

## CHAPTER 3

# Linear systems

### 3.1. Exponential of a matrix

Let  $\mathbb{M}_d(\mathbb{C})$  be the vector space of  $d \times d$  matrices with entries in  $\mathbb{C}$ . Let  $GL_d(\mathbb{C}) \subset \mathbb{M}_d(\mathbb{C})$  be the group of invertible matrices.

**DEFINITION 3.1** (Matrix norm). *The matrix norm of  $A \in \mathbb{M}_d(\mathbb{C})$  is*

$$\|A\| = \max_{|y|=1} |Ay|.$$

**LEMMA 3.2.** *The matrix norm has the following properties:*

- (i)  $|Ay| \leq \|A\||y|$  for all  $y \in \mathbb{C}^d$ ;
- (ii)  $\|A + B\| \leq \|A\| + \|B\|$  for all  $A, B \in \mathbb{M}_d(\mathbb{C})$ ;
- (iii)  $\|AB\| \leq \|A\| \|B\|$  for all  $A, B \in \mathbb{M}_d(\mathbb{C})$ .

**LEMMA 3.3** (Jordan-Chevalley decomposition). *Let  $A \in \mathbb{M}_d(\mathbb{C})$ . Then there exists  $C \in GL_d(\mathbb{C})$  such that  $A$  has a unique decomposition*

$$C^{-1}AC = D + N,$$

where  $D$  is diagonal,  $N$  is nilpotent (i.e.,  $N^d = 0$ ), and  $ND = DN$ .

We now define the **exponential of a matrix**.

**DEFINITION 3.4.** *For a matrix  $A \in \mathbb{M}_d(\mathbb{C})$ , we define*

$$e^A = \sum_{n \geq 0} \frac{A^n}{n!}.$$

We list some properties of the exponential of a matrix.

**LEMMA 3.5.** *The exponential of a matrix has the following properties:*

- (i) (exponential of the sum) Let  $A, B \in \mathbb{M}_d(\mathbb{C})$ . If  $AB = BA$ , then  $e^{A+B} = e^A e^B$ ;
- (ii) (conjugation and exponentiation) Let  $A, B \in \mathbb{M}_d(\mathbb{C})$  and  $C \in GL_d(\mathbb{C})$  be such that  $A = C^{-1}BC$ . Then we have

$$e^A = C^{-1}e^B C.$$

In fact,

$$e^A = \sum_{n \geq 0} \frac{A^n}{n!} = \sum_{n \geq 0} \frac{(C^{-1}BC)^n}{n!} = \sum_{n \geq 0} \frac{C^{-1}B^n C}{n!} = C^{-1}e^B C;$$

- (iii) (exponential of a diagonalizable matrix) If  $A$  is a diagonalizable matrix of the form

$$A = C^{-1} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{pmatrix} C,$$

where  $\lambda_1, \dots, \lambda_d \in \mathbb{C}$  and  $C \in GL_d(\mathbb{C})$ , then

$$e^A = C^{-1} \begin{pmatrix} e^{\lambda_1} & & 0 \\ & \ddots & \\ 0 & & e^{\lambda_d} \end{pmatrix} C;$$

(iv) (*exponential of a block matrix*) Let  $A_j \in \mathbb{M}_{h_j}(\mathbb{C})$  for  $j = 1, \dots, p$ . Let  $A$  be a block matrix of the form

$$A = \begin{pmatrix} A_1 & & 0 \\ & \ddots & \\ 0 & & A_p \end{pmatrix}.$$

Then

$$e^A = \begin{pmatrix} e^{A_1} & & 0 \\ & \ddots & \\ 0 & & e^{A_p} \end{pmatrix};$$

(v) (*derivative*) Let  $A \in \mathbb{M}_d(\mathbb{C})$ . We have

$$\frac{d}{dt} e^{tA} = A e^{tA} = e^{tA} A.$$

In addition, to the matrix exponential we will also need its inverse. That is, given a matrix  $A$  we want to find a matrix  $B$  such that

$$A = e^B.$$

In this case, we will call  $B = \log A$  a **matrix logarithm** of  $A$ . Note that  $B$  is not unique.

LEMMA 3.6. *A matrix  $A$  has a logarithm if and only if  $\det A \neq 0$ . Moreover, if  $A$  is real and all real eigenvalues are positive, then there is a real logarithm.*

### 3.2. Linear systems with constant coefficients

Let  $A \in \mathbb{M}_d(\mathbb{C})$  be independent of  $t$ . Let  $f \in \mathcal{C}^0([0, T])$ . Consider the following linear ODE with constant coefficients:

$$\begin{cases} \frac{dx}{dt} = Ax(t) + f(t), & t \in [0, T], \\ x(0) = x_0 \in \mathbb{R}^d. \end{cases} \quad (3.1)$$

Since

$$|A(x - y)| \leq \|A\| \|x - y\| \quad \text{for all } x, y \in \mathbb{C}^d,$$

by the Cauchy-Lipschitz theorem there exists a unique solution  $x$  to (3.1). If  $f = 0$ , then the system of equations (3.1) is an **autonomous** system.

If  $d = 1$  (i.e.,  $A = a \in \mathbb{C}$ ), then by the method of integrating factors,

$$x(t) = e^{at} x_0 + \int_0^t e^{a(t-s)} f(s) ds. \quad (3.2)$$

In the general case ( $d \geq 1$ ), if  $f = 0$ , then, from Lemma 3.5 (v), it follows that the solution  $x$  of (3.1) is  $x(t) = e^{tA} x_0$ .

For an arbitrary  $f$ , we have

$$\frac{d}{dt} (e^{-tA} x) = e^{-tA} f(t),$$

and hence the solution  $x(t)$  of (3.1) is given by

$$x(t) = e^{tA} x_0 + \int_0^t e^{(t-s)A} f(s) ds. \quad (3.3)$$

Observe that the solution of (3.1) has been reduced in (3.3) to matrix calculations and integration.

EXAMPLE 3.7. *An important class of linear system with constant coefficients are those that can be converted into diagonal form. Suppose that we are given a system  $dx/dt = Ax$  such that the eigenvalues  $\lambda_j$  of  $A$  are distinct. Then we can find an invertible matrix  $C$  such that  $C^{-1}AC$  is diagonal. If we choose a set of coordinates  $y = C^{-1}x$ , then in the new coordinates the equation becomes*

$$\frac{dy}{dt} = C^{-1}ACy = Dy, \quad y(0) = y_0. \quad (3.4)$$

By construction,  $D$  in (3.4) is diagonal and

$$y(t) = \begin{pmatrix} e^{\lambda_1 t} & & 0 \\ & \ddots & \\ 0 & & e^{\lambda_d t} \end{pmatrix} y_0.$$

EXAMPLE 3.8. Consider (3.1) with  $d = 2$ ,  $A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ . Then since  $A^2 = 0$  and hence  $e^{tA} = \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix}$ , the solution  $x(t)$  is given by

$$x(t) = \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix} x_0 + \int_0^t \begin{pmatrix} 1 & t-s \\ 0 & 1 \end{pmatrix} f(s) ds.$$

EXAMPLE 3.9. Consider (3.1) with  $d = 2$ ,  $A = \begin{pmatrix} 0 & \omega \\ -\omega & 0 \end{pmatrix}$  for some  $\omega \in \mathbb{R}$ ,  $\omega \neq 0$ . Then

$$e^{tA} = \begin{pmatrix} \cos \omega t & \sin \omega t \\ -\sin \omega t & \cos \omega t \end{pmatrix}.$$

This expression for  $e^{tA}$  can be verified by differentiation:

$$\frac{d}{dt} e^{tA} = \begin{pmatrix} -\omega \sin \omega t & \omega \cos \omega t \\ -\omega \cos \omega t & -\omega \sin \omega t \end{pmatrix} = \begin{pmatrix} 0 & \omega \\ -\omega & 0 \end{pmatrix} \begin{pmatrix} \cos \omega t & \sin \omega t \\ -\sin \omega t & \cos \omega t \end{pmatrix} = A e^{tA}.$$

The solution  $x(t)$  to (3.1) is then given by

$$x(t) = \begin{pmatrix} \cos \omega t & \sin \omega t \\ -\sin \omega t & \cos \omega t \end{pmatrix} x_0 + \int_0^t \begin{pmatrix} \cos \omega(t-s) & \sin \omega(t-s) \\ -\sin \omega(t-s) & \cos \omega(t-s) \end{pmatrix} f(s) ds.$$

### 3.3. Linear system with non-constant real coefficients

**3.3.1. The homogeneous case.** Let  $\mathbb{M}_d(\mathbb{R})$  be the vector space of  $d \times d$  matrices with entries in  $\mathbb{R}$ .

PROPOSITION 3.10. Let  $A : [0, T] \rightarrow \mathbb{M}_d(\mathbb{R})$  be continuous. The set  $S$  of solutions of  $dx/dt = A(t)x$  defined by

$$S = \left\{ x \in C^1([0, T]; \mathbb{R}^d) : x \text{ satisfies } \frac{dx}{dt} = A(t)x \right\} \quad (3.5)$$

is a linear subspace of  $C^1([0, T]; \mathbb{R}^d)$  of dimension  $d$ .

PROOF. If  $x, y \in S$ , then, for any  $\alpha, \beta \in \mathbb{R}$ ,  $\alpha x + \beta y \in C^1([0, T]; \mathbb{R}^d)$  is also a solution. Then  $S$  is a linear subspace of  $C^1([0, T]; \mathbb{R}^d)$ . We show that the dimension of  $S$  is  $d$ . Let the mapping  $F : S \rightarrow \mathbb{R}^d$  be defined by

$$F[x] = x(t_0) \quad (3.6)$$

for some  $t_0 \in [0, T]$ . Then  $F$  is linear:  $F[\alpha x + \beta y] = \alpha x(t_0) + \beta y(t_0) = \alpha F[x] + \beta F[y]$ .  $F$  is injective, i.e.,  $F[x] = 0$  implies that  $x = 0$ . In fact,  $x$  solves  $\frac{dx}{dt} = A(t)x(t)$  with the initial condition  $x(t_0) = 0$ . The solution to this problem is unique (by the Cauchy-Lipschitz theorem) and  $0$  is a solution. Then  $x = 0$ . Finally,  $F$  is surjective because for any  $x_0 \in \mathbb{R}^d$  the equation

$$\begin{cases} \frac{dx}{dt} = A(t)x(t), & t \in [0, T], \\ x(t_0) = x_0, \end{cases} \quad (3.7)$$

has a solution  $x \in C^1([0, T]; \mathbb{R}^d)$ . □

PROPOSITION 3.11. Let  $S$  be defined by (3.5) and let  $x_1, \dots, x_d \in S$ . The following statements are equivalent:

- (i)  $\{x_1, \dots, x_d\}$  is a basis of  $S$ ;
- (ii)  $\det[x_1(t), \dots, x_d(t)] \neq 0$  for all  $t \in [0, T]$ .
- (iii)  $\det[x_1(t_0), \dots, x_d(t_0)] \neq 0$  for some  $t_0 \in [0, T]$ .

Here,  $\det$  denotes the determinant of a matrix and  $[x_1, \dots, x_d]$  is the  $d \times d$  matrix with columns  $x_1, \dots, x_d \in \mathbb{R}^d$ .

PROOF. It is clear that (i) is equivalent to (ii). To see that (i) implies (iii), let  $\{x_1, \dots, x_d\}$  be a basis of  $S$ . Then  $\{F[x_1], \dots, F[x_d]\}$  forms a basis of  $\mathbb{R}^d$ , where the isomorphism  $F$  relative to  $t_0$  is defined by (3.6). Next let us check that (iii) implies (i). Let  $t_0$  be such that (iii) holds and let  $F : S \rightarrow \mathbb{R}^d$  be the isomorphism relative to  $t_0$  defined by (3.6). Then the inverse  $F^{-1} : \mathbb{R}^d \rightarrow S$  is also an isomorphism. It follows that  $x_1 = F^{-1}[x_1(t_0)], \dots, x_d = F^{-1}[x_d(t_0)]$  is a basis of  $S$ .  $\square$

DEFINITION 3.12 (**Fundamental matrix**). *If one of the three equivalent conditions of Proposition 3.11 holds, then the functions  $x_1, \dots, x_d$  are called a **fundamental system** of solutions of the differential equation  $\frac{dx}{dt} = A(t)x$ . The matrix  $X = [x_1, \dots, x_d]$  is then called a **fundamental matrix** of the equation.*

We now introduce the **Wronskian determinant**.

DEFINITION 3.13 (**Wronskian determinant**). *Let  $x_1, \dots, x_d \in S$ . The Wronskian determinant  $w \in C^1([0, T]; \mathbb{R})$  of  $x_1, \dots, x_d$  is defined by*

$$w(t) = \det[x_1(t), \dots, x_d(t)].$$

THEOREM 3.14. *Let  $x_1, \dots, x_d \in S$  and let  $w \in C^1([0, T]; \mathbb{R})$  be the Wronskian determinant of  $x_1, \dots, x_d$ . Then  $w$  solves the differential equation*

$$\frac{dw}{dt} = (\operatorname{tr} A(t))w \quad \text{for } t \in [0, T]. \quad (3.8)$$

Here,  $\operatorname{tr}$  denotes the trace of a matrix.

PROOF. If  $x_1, \dots, x_d$  are linearly dependent, then  $w = 0$  and (3.8) trivially holds. Suppose that  $x_1, \dots, x_d$  are linearly independent, i.e.,  $w(t) \neq 0$  for all  $t \in [0, T]$ .

Let  $X : [0, T] \rightarrow \mathbb{M}_d(\mathbb{R})$  be the fundamental matrix having as columns the solutions  $x_1, \dots, x_d$ , i.e.,

$$X(t) = (x_{ij}(t))_{i,j=1,\dots,d}, \quad t \in [0, T],$$

where  $x_j = (x_{1j}, \dots, x_{dj})^\top$  for  $j = 1, \dots, d$ .

Let  $z_j$  be the solution of

$$\begin{cases} \frac{dz_j}{dt} = A(t)z_j(t), \\ z_j(t_0) = e_j, \end{cases}$$

where  $\{e_j\}_{j=1,\dots,d}$  is the standard unit orthonormal basis in  $\mathbb{R}^d$ .

Then  $\{z_1, \dots, z_d\}$  is a basis of the space of solutions to  $dz/dt = Az$ . Moreover, there exists  $C \in GL_d(\mathbb{R}^d)$  such that

$$X(t) = Z(t)C, \quad t \in [0, T],$$

where  $Z = [z_1, \dots, z_d]$ . Since a fundamental matrix is uniquely determined by an initial condition,  $C = Z(t_0)^{-1}X(t_0)$ .

Let  $v(t) := \det Z(t)$ . Then  $v$  solves

$$\frac{dv}{dt}(t_0) = \operatorname{tr} A(t_0).$$

In fact, by the definition of the determinant of a matrix, we have

$$\frac{dv}{dt}(t) = \frac{d}{dt} \sum_{\sigma \in S_d} (-1)^{\operatorname{sgn} \sigma} \prod_{i=1}^d z_{i\sigma(i)}(t) = \sum_{\sigma \in S_d} (-1)^{\operatorname{sgn} \sigma} \sum_{j=1}^d \frac{d}{dt} z_{j\sigma(j)}(t) \prod_{i \neq j} z_{i\sigma(i)}(t),$$

where  $S_d$  is the set of all permutations of the  $d$  elements  $\{1, 2, \dots, d\}$  and  $\text{sgn } \sigma$  is the signature of the permutation  $\sigma$ . Note that

$$\prod_{i \neq j} z_{i\sigma(i)}(t_0) = 0 \quad \text{unless } \sigma = \text{identity},$$

and

$$\begin{aligned} \frac{dz_{jj}}{dt}(t_0) &= (A(t_0)z_j(t_0))_j \\ &= \sum_{h=1}^d a_{jh}(t_0)z_{hj}(t_0) = \sum_{h=1}^d a_{jh}(t_0)\delta_{hj}(t_0) \\ &= a_{jj}(t_0). \end{aligned}$$

Therefore,

$$\frac{dv}{dt}(t_0) = \sum_{j=1}^d a_{jj}(t_0) = \text{tr}A(t_0).$$

Now the general result follows from the differentiation of the following identity:

$$w = \det X = \det(ZC) = (\det C) \det Z = (\det C)v.$$

In fact, we have

$$\frac{dw}{dt}(t_0) = (\det C) \frac{dv}{dt}(t_0) = (\det C) \text{tr}A(t_0).$$

Therefore,

$$\frac{dw}{dt}(t_0) = \text{tr}A(t_0)w(t_0),$$

since  $v(t_0) = 1$ . □

**REMARK 3.15.** Let  $t_0 \in [0, T]$ . From (3.8), it follows that

$$w(t) = w(t_0)e^{\int_{t_0}^t \text{tr}A(s) ds} \quad \text{for } t \in [0, T]. \quad (3.9)$$

This is known as **Abel's identity** or **Liouville's formula**. Identity (3.9) shows that it suffices to check that the determinant of the fundamental matrix is nonzero for one  $t_0 \in [0, T]$ .

**3.3.2. The inhomogeneous case.** Consider the inhomogeneous linear differential equation of the form

$$\begin{cases} \frac{dx}{dt} = A(t)x + f(t), \end{cases} \quad (3.10)$$

where  $A(t) \in \mathcal{C}^0([0, T]; \mathbb{M}_d(\mathbb{R}))$  and  $f \in \mathcal{C}^0([0, T]; \mathbb{R}^d)$ .

Let  $X$  be a fundamental matrix for the homogeneous equation  $dx(t)/dt = A(t)x(t)$ , i.e.,

$$\frac{dX}{dt} = AX \quad \text{and} \quad \det X \neq 0 \quad \text{for all } t \in [0, T].$$

Then, any solution  $x$  to the homogeneous equation is of the form

$$x(t) = X(t)c, \quad t \in [0, T], \quad (3.11)$$

for some (column) vector  $c \in \mathbb{R}^d$ .

By using the method of integrating factors, we look for a solution to (3.10) of the form (3.11) with  $c \in \mathcal{C}^1([0, T]; \mathbb{R}^d)$ . In this case, we have

$$\frac{dx}{dt} = \frac{dX}{dt}c + X \frac{dc}{dt} = AXc + X \frac{dc}{dt} = Ax + X \frac{dc}{dt},$$

which implies  $X \frac{dc}{dt} = f(t)$ . Since  $X$  is invertible, we obtain

$$\frac{dc}{dt} = X^{-1}f(t).$$

Therefore, we find

$$c(t) = c_0 + \int_0^t X(s)^{-1} f(s) ds,$$

for some  $c_0 \in \mathbb{R}^d$ .

**THEOREM 3.16.** *Let  $X$  be a fundamental matrix for the homogeneous equation  $dx/dt = Ax$ . Then, for all  $c_0 \in \mathbb{R}^d$ , the function*

$$x(t) = X(t) \left( c_0 + \int_0^t X(s)^{-1} f(s) ds \right) \quad (3.12)$$

*is a solution to (3.10). Moreover, any solution to (3.10) is of the form (3.12) for some  $c_0 \in \mathbb{R}^d$ .*

**PROOF.** The first statement is already proved. To prove the second statement, let  $x_2$  be a solution to (3.10). Since

$$\frac{d}{dt}(x_2 - x(t)) = A(x_2 - x),$$

where  $x$  is given by (3.12), we get  $x_2 - x = Xc_1$  for some  $c_1 \in \mathbb{R}^d$  and the claim follows.  $\square$

Formula (3.12) is called **Duhamel's formula**.

### 3.4. Second order linear equations

Let  $d = 1$  and consider the following second order ODE:

$$\frac{d^2 x}{dt^2} = f\left(t, x, \frac{dx}{dt}\right),$$

for a given scalar function  $f$ . The above ODE is linear if  $f$  is linear in  $x$  and  $dx/dt$ , namely,

$$f\left(t, x, \frac{dx}{dt}\right) = g(t) - p(t) \frac{dx}{dt} - q(t)x,$$

where  $g, p, q$  are (scalar) functions of  $t$  but do not depend on  $x$ . Then the ODE becomes

$$\frac{d^2 x}{dt^2} + p(t) \frac{dx}{dt} + q(t)x = g(t). \quad (3.13)$$

The initial value problem consists of (3.13) together with a pair of initial conditions

$$x(t_0) = x_0, \quad \frac{dx}{dt}(t_0) = x'_0, \quad x_0, x'_0 \in \mathbb{R}. \quad (3.14)$$

The second order ODE (3.13) is called **homogeneous** if  $g = 0$  and **inhomogeneous** otherwise. If  $p(t)$  and  $q(t)$  are constant, then (3.13) is called linear ODE with constant coefficients.

Suppose that

$$p, q \in \mathcal{C}^0([0, T]). \quad (3.15)$$

If the condition (3.15) fails, then the points at which either  $p$  or  $q$  fail to be continuous are called **singular points**. The following are important examples:

$$\text{Bessel's equation: } p(t) = \frac{1}{t}, q(t) = 1 - \frac{\nu}{t^2}, \quad (\text{at } t = 0);$$

$$\text{Legendre's equation: } p(t) = \frac{2t}{1-t^2}, q(t) = \frac{n(n+1)}{1-t^2}, n \in \mathbb{N} \quad (\text{at } t = \pm 1).$$

**THEOREM 3.17.** *Suppose that  $p, q, g \in \mathcal{C}^0([0, T], \mathbb{R}^d)$ . Then there exists a unique solution  $x(t)$  on  $[0, T]$  to (3.13) with the initial conditions (3.14).*



**3.4.1. Structure of the general solution.** Here we discuss the structure of the general solution to the second order ODE (3.13).

First we consider the homogeneous case. We need the following results regarding the Wronskian determinant.

**DEFINITION 3.18.** *Two functions  $x_1$  and  $x_2$  on  $[0, T]$  are called **linearly independent** if neither of them is a multiple of the other. Otherwise, they are called **linearly dependent**.*

**PROPOSITION 3.19.** *Let  $w$  be the Wronskian determinant given by*

$$w(t) := x_1(t) \frac{dx_2}{dt}(t) - x_2(t) \frac{dx_1}{dt}(t) = \det \begin{pmatrix} x_1 & x_2 \\ \frac{dx_1}{dt} & \frac{dx_2}{dt} \end{pmatrix}.$$

*If  $w(t)$  is not zero at some  $t_0 \in [0, T]$ , then  $x_1$  and  $x_2$  are linearly independent.*

**PROOF.** Let us prove that if  $x_1$  and  $x_2$  are linearly dependent, then  $w(t) = 0$  for all  $t \in [0, T]$ . Suppose that  $x_1$  and  $x_2$  are linearly dependent. Then, with respect to  $(\alpha_1, \alpha_2)$ , the following system:

$$\begin{cases} \alpha_1 x_1 + \alpha_2 x_2 = 0, \\ \alpha_1 \frac{dx_1}{dt} + \alpha_2 \frac{dx_2}{dt} = 0, \end{cases} \quad \text{for all } t \in [0, T],$$

has a non-trivial solution. Therefore,

$$w(t) = \det \begin{pmatrix} x_1 & x_2 \\ \frac{dx_1}{dt} & \frac{dx_2}{dt} \end{pmatrix} = 0, \quad \text{for all } t \in [0, T].$$

This completes the proof.  $\square$

**PROPOSITION 3.20.** *If  $x_1$  and  $x_2$  solve (3.13) on  $[0, T]$  then  $w(t)$  is either identically zero or not equal to zero at any point of  $[0, T]$ .*

**PROOF.** We have

$$w'(t) = x_1 \frac{d^2 x_2}{dt^2} - x_2 \frac{d^2 x_1}{dt^2}.$$

We also have, by the assumption that  $x_1, x_2$  solve (3.13), that

$$\frac{d^2 x_i}{dt^2} = -p(t) \frac{dx_i}{dt} - q(t) x_i, \quad i = 1, 2.$$

So we get

$$\frac{dw}{dt} = -p(t) \left( x_1 \frac{dx_2}{dt} - \frac{dx_1}{dt} x_2 \right) = -p(t) w(t).$$

Therefore  $w(t) = w(t_0) e^{-\int_{t_0}^t p(s) ds}$ , which is either identically zero or never vanishes depending on  $w(t_0)$ .  $\square$

Now we discuss the structure of the general solution to the homogeneous system.

**THEOREM 3.21.** *Suppose that  $x_1$  and  $x_2$  solve the equation (3.13) with  $g = 0$ . Suppose also that  $x_1$  and  $x_2$  are linearly independent. Then the general solution is of the form  $c_1 x_1 + c_2 x_2$ , where  $c_1$  and  $c_2$  are constant coefficients.*

**PROOF.** Let  $\tilde{x}$  be an arbitrary solution with the initial condition  $\tilde{x}(t_0) = \tilde{x}_0, d\tilde{x}/dt(t_0) = \tilde{x}'_0$ . Consider the system of equations for  $(c_1, c_2)$

$$\begin{cases} c_1 x_1(t_0) + c_2 x_2(t_0) = \tilde{x}_0, \\ c_1 \frac{dx_1}{dt}(t_0) + c_2 \frac{dx_2}{dt}(t_0) = \tilde{x}'_0. \end{cases}$$

Since  $x_1 \frac{dx_2}{dt} - x_2 \frac{dx_1}{dt} \neq 0$  at  $t = t_0$ , there exists a unique nontrivial solution  $(c_1, c_2) = (\tilde{c}_1, \tilde{c}_2)$  to the above system. Then, by the existence and uniqueness theorem for the initial value problem of the second order ODE, we conclude that  $\tilde{c}_1 x_1 + \tilde{c}_2 x_2 = \tilde{x}$ .  $\square$

**3.4.2. Linear  $n$ -th order ODE with constant coefficients.** Here we discuss the approach to solving a linear  $n$ -th order ODE with constant coefficients. Consider

$$\frac{d^n x}{dt^n} + a_{n-1} \frac{d^{n-1} x}{dt^{n-1}} + \dots + a_1 \frac{dx}{dt} + a_0 x = 0, \quad (3.16)$$

where  $a_i \in \mathbb{R}$  for  $i = 0, \dots, n-1$ .

The general solution has the form

$$x(t) = c_1 x_1 + \dots + c_n x_n,$$

where  $\{x_i\}_{i=1}^n$  is the set of linearly independent solutions (a fundamental set of solutions) and  $c_i$  are constant coefficients.

Let  $w(t)$  be the Wronskian determinant of the set  $\{x_1, \dots, x_n\}$ , *i.e.*,

$$w(t) = \det \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ \frac{dx_1}{dt} & & & \\ \vdots & & & \vdots \\ \frac{d^{n-1} x_1}{dt^{n-1}} & \frac{d^{n-1} x_2}{dt^{n-1}} & \dots & \frac{d^{n-1} x_n}{dt^{n-1}} \end{bmatrix}.$$

If  $w(t_0) \neq 0$  for some  $t_0$ , then  $(x_1, \dots, x_n)$  forms a fundamental set of solution.

We solve the equation through the characteristic equation

$$\lambda^n + a_{n-1} \lambda^{n-1} + \dots + a_1 \lambda + a_0 = 0. \quad (3.17)$$

This equation is derived by guessing a solution  $x(t)$  has the form  $e^{\lambda t}$  with  $\lambda \in \mathbb{C}$ .

The characteristic equation (3.17) has  $n$  complex roots  $\hat{\lambda}_j$  counted with their multiplicities  $l_j$ . In other words, equation (3.17) can be rewritten in the form

$$\prod_{j=1}^m (\lambda - \hat{\lambda}_j)^{l_j} = 0$$

with  $\sum_{j=1}^m l_j = n$ . In fact, the general solution  $x(t)$  is a linear combination of  $t^k e^{\hat{\lambda}_j t}$  for  $0 \leq k < l_j$  and  $j = 1, \dots, m$ . In particular, if  $m = n$ , then  $x(t)$  is a linear combination of  $e^{\hat{\lambda}_j t}$ .

**THEOREM 3.22.** *Let  $\hat{\lambda}_j, 1 \leq j \leq m$ , be the zeros of the characteristic polynomial (3.17) associated with (3.16) and let  $l_j$  be the corresponding multiplicities. Then the functions*

$$x_{j,k}(t) = t^k e^{\hat{\lambda}_j t}, \quad 0 \leq k < l_j, \quad 1 \leq j \leq m, \quad (3.18)$$

*are  $n$  linearly independent solutions of (3.16). In particular, any other solution can be written as a linear combination of these solutions.*

**REMARK 3.23.** *Let  $y = (x, dx/dt, \dots, d^{n-1}x/dt^{n-1})^\top$ . Then (3.16) can be rewritten as*

$$\frac{dy}{dt} = Ay \quad \text{with } A := \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 1 \\ -a_0 & -a_1 & \dots & -a_{n-1} \end{pmatrix}.$$

*The characteristic polynomial of  $A$ ,  $P(\lambda) := \det(A - \lambda I)$ , is given by  $P(\lambda) = \prod_{j=1}^m (\lambda - \hat{\lambda}_j)^{l_j}$ . The algebraic multiplicity of the eigenvalue  $\hat{\lambda}_j$  of  $A$  is  $l_j$ . If  $A$  has a basis of eigenvectors, there*

will only be in  $y(t)$  terms of the form  $e^{\lambda_j t}$ . In general, let  $J$  be the Jordan bloc form of  $A$ . Then  $e^{tA} = C^{-1}e^{tJ}C$  for some invertible matrix  $C$ , where

$$e^{tJ} = \begin{pmatrix} e^{tJ_1} & & 0 \\ & \ddots & \\ 0 & & e^{tJ_k} \end{pmatrix},$$

and

$$e^{tJ_i} = e^{t(\lambda_i I + N_i)} = e^{t\lambda_i} \left( I + tN_i + \dots + \frac{t^{l_i-1}}{(l_i-1)!} N_i^{l_i-1} \right).$$

Therefore, as stated in Theorem 3.22, in general the solution will be the sum of terms of the form  $t^k e^{\lambda_j t}$ ,  $k < l_j$ .

**3.4.3. Reduction of order.** Here we discuss a method for finding a second solution to the homogeneous second order ODE when a first solution is known by reducing the order.

Suppose that  $x_1$  a solution of (3.13). Let

$$x(t) = v(t)x_1(t).$$

Then

$$\frac{dx}{dt}(t) = \frac{dv}{dt}x_1 + v\frac{dx_1}{dt}$$

and

$$\frac{d^2x}{dt^2}(t) = \frac{d^2v}{dt^2}x_1 + 2\frac{dv}{dt}\frac{dx_1}{dt} + v\frac{d^2x_1}{dt^2}.$$

So, we get

$$\frac{d^2v}{dt^2} + \left( p + 2\frac{(dx_1/dt)}{x_1} \right) \frac{dv}{dt} = 0. \quad (3.19)$$

By letting  $u = dv/dt$ , the equation above can be rewritten as a first order ODE

$$\frac{du}{dt} + \left( p + 2\frac{(dx_1/dt)}{x_1} \right) u = 0.$$

Therefore,

$$u(t) = ce^{-\int^t (p + 2\frac{(dx_1/dt)}{x_1}) ds} = \frac{c}{(x_1(t))^2} e^{-\int^t p(s) ds}. \quad (3.20)$$

Since  $v = \int^t u(s) ds$ , we get

$$x(t) = x_1(t) \int^t u(s) ds. \quad (3.21)$$

In conclusion, if one solution to (3.13) is known, then a second solution can be found and it is expressed by (3.21), where  $u$  is given by (3.20).

EXAMPLE 3.24. Consider the differential equation

$$\frac{d^2x}{dt^2} - 2t\frac{dx}{dt} - 2x = 0, \quad (3.22)$$

and observe that  $x_1(t) = e^{t^2}$  is a solution. Hence we can set  $x(t) = e^{t^2}v(t)$ . Then from (3.19), it follows that

$$\frac{d^2v}{dt^2} + 2t\frac{dv}{dt} = 0. \quad (3.23)$$

The solution of (3.23) is given by

$$\frac{dv}{dt} = e^{-t^2},$$

implying that

$$v(t) = \int_0^t e^{-s^2} ds = \frac{\sqrt{\pi}}{2} \operatorname{erf}(t),$$

where erf is the Gauss error function. Hence a second solution to (3.22) is given by

$$x_2(t) = e^{t^2} \operatorname{erf}(t).$$

### 3.5. Linearization and stability for autonomous systems

**3.5.1. Linear systems.** Let  $A \in \mathbb{M}_d(\mathbb{R})$  be independent of  $t$ . Consider the following linear system of ODEs:

$$\begin{cases} \frac{dx}{dt} = Ax(t), & t \in [0, +\infty[, \\ x(0) = x_0 \in \mathbb{R}^d. \end{cases} \quad (3.24)$$

By Lemma 3.3, there exists  $C \in GL_d(\mathbb{C})$  such that

$$C^{-1}AC = D + N,$$

where  $D$  is diagonal,  $N$  is nilpotent, and  $ND = DN$ . Let  $\lambda_j, j = 1, \dots, J$  be the (distinct) eigenvalues of  $A$ . Let  $l_j$  be the (algebraic) multiplicity of  $\lambda_j$  and denote by  $E_j = \ker(A - \lambda_j I)^{l_j}$  the characteristic subspace associated with  $\lambda_j$  (called also generalized eigenspace). We have  $\oplus E_j = \mathbb{C}^d$ . Moreover, each  $E_j$  is invariant under  $A$ .

The system (3.24) is said to be **stable** if there exists a positive constant  $C_0$  such that

$$|x(t)| \leq C_0 |x_0| \quad \text{for all } t \in [0, +\infty[. \quad (3.25)$$

**LEMMA 3.25.** *The system (3.24) is stable if and only if  $\Re \lambda_j < 0$  or  $\Re \lambda_j = 0$  and  $N|_{E_j} = 0$  for  $j = 1, \dots, J$ .*

**PROOF.** Let  $\tilde{x}(t) = C^{-1}x(t)$  and  $\tilde{x}_0 = C^{-1}x_0$ . By Lemma 3.5,

$$\tilde{x}(t) = e^{tD+tN}\tilde{x}_0, \quad t \in [0, +\infty[. \quad (3.26)$$

Since  $DN = ND$ , (3.26) yields

$$\tilde{x}(t) = \left( \sum_{i=0}^{d-1} \frac{(tN)^i}{i!} \right) e^{tD} \tilde{x}_0, \quad t \in [0, +\infty[. \quad (3.27)$$

If  $\tilde{x}_0 \in E_j$ , then

$$\tilde{x}(t) = e^{t\lambda_j} \left( \sum_{i=0}^{d-1} \frac{(tN)^i}{i!} \right) \tilde{x}_0, \quad t \in [0, +\infty[. \quad (3.28)$$

Therefore,  $x(t)$  satisfies (3.25) for some positive constant  $C_0$  if and only if  $\Re \lambda_j < 0$  or  $\Re \lambda_j = 0$  and  $N|_{E_j} = 0$ .  $\square$

**3.5.2. Nonlinear systems.** Consider the autonomous system

$$\begin{cases} \frac{dx}{dt} = f(x), \\ x(0) = x_0 \in \mathbb{R}^d, \end{cases} \quad (3.29)$$

where  $f$  is  $\mathcal{C}^1$ . Suppose that  $x^*$  is an equilibrium point for (3.29), i.e.,  $f(x^*) = 0$ .

**THEOREM 3.26 (Local stability).** *Suppose that all the eigenvalues  $\lambda$  of the Jacobian of  $f$  at  $x^*$ ,  $f'(x^*)$ , are with negative real parts. Then, there exists  $\delta > 0$  such that if  $|x_0 - x^*| \leq \delta$ , then  $|x(t) - x^*| \rightarrow 0$  as  $t \rightarrow +\infty$ .*

**PROOF.** Let  $A = f'(x^*)$  and consider the linearized system

$$\begin{cases} \frac{dy(t)}{dt} = Ay(t), & t \geq 0, \\ y(0) = x_0 - x^*, \end{cases} \quad (3.30)$$

which, in view of (3.3), has the explicit solution  $y(t) = e^{tA}(x_0 - x^*)$  for  $t \geq 0$ . Suppose that  $\Re \lambda < 0$  for any eigenvalue  $\lambda$  of  $f'(x^*)$ . From (3.27), it follows that there exists  $r > 0$  such that

$$|e^{tA}z| \leq C_0 e^{-rt}|z| \quad \text{for all } z \in \mathbb{R}^d,$$

where the constant  $C_0$  depends only on  $f$ .

Now, rewrite (3.29) as a small perturbation of the linearized system

$$\begin{cases} \frac{dx}{dt} = A(x - x^*) + g(x), \\ x(0) = x^*, \end{cases} \quad (3.31)$$

where

$$g(x) = |x - x^*|\epsilon(x), \quad \text{with } \epsilon \in \mathcal{C}^0 \quad \text{and } \epsilon(x^*) = 0. \quad (3.32)$$

Observe that there exists  $\delta_0 > 0$  such that for all  $\delta \in ]0, \delta_0[$ ,

$$\sup\{|g(x)| : |x - x^*| \leq \delta\} < \frac{r\delta}{C_0}. \quad (3.33)$$

To conclude it suffices to prove that if  $|x_0 - x^*| < \min(\delta, \delta/C_0)$ , then

$$|x(t) - x^*| \leq \delta \quad \text{for all } t \geq 0.$$

From (3.31), it follows that

$$x(t) - x^* = e^{tA}(x_0 - x^*) + \int_0^t e^{(t-s)A}g(x(s)) ds,$$

and hence, (3.32) yields

$$\begin{aligned} |x(t) - x^*| &\leq e^{-rt}C_0|x_0 - x^*| + \int_0^t e^{-r(t-s)}C_0|g(x(s))| ds \\ &\leq e^{-rt}C_0|x_0 - x^*| + (1 - e^{-rt})\frac{C_0}{r} \sup\{|g(x(s))| : 0 \leq s \leq t\}. \end{aligned}$$

Thus, for all  $t \geq 0$ ,

$$|x(t) - x^*| \leq \max\left(C_0|x_0 - x^*|, \frac{C_0}{r} \sup\{|g(x(s))| : 0 \leq s \leq t\}\right).$$

Introduce

$$T := \inf\{t > 0 : |x(t) - x^*| \geq \delta\}.$$

If we assume that  $T$  is finite, we would obtain

$$|x(t) - x^*| \leq \delta \quad \text{for all } t \in [0, T], \quad |x(T) - x^*| = \delta.$$

In view of (3.32), we arrive at a contradiction by using (3.33).  $\square$

**DEFINITION 3.27.** A function  $V \in \mathcal{C}^1(\mathbb{R}^d, \mathbb{R})$  is said to be a **Lyapunov function** for (3.29) if

$$V(x^*) < V(x) \quad \text{for any } x \neq x^*, \quad (3.34)$$

and

$$f(x) \cdot V'(x) \leq 0 \quad \text{for any } x \in \mathbb{R}^d. \quad (3.35)$$

**EXAMPLE 3.28.** (i) Consider the system

$$\begin{cases} \frac{dx_1}{dt} = x_2, \\ \frac{dx_2}{dt} = -2x_1 - x_2. \end{cases} \quad (3.36)$$

Then  $x^* = (0, 0)$  is an equilibrium point and

$$V(x) = x_1^2 + \frac{1}{2}x_2^2$$

is a Lyapunov function for (3.36).

- (ii) *For the gradient systems introduced in Subsection 1.3.4, there is a natural candidate for a Lyapunov function. Suppose that  $f(x) = -\nabla\Phi(x)$ . Suppose that the potential  $\Phi$  is smooth and there exists  $x^*$  such that  $\Phi(x^*) < \Phi(x)$  for any  $x \neq x^*$ . Then  $V = \Phi$  is a Lyapunov function for (3.29).*

**THEOREM 3.29.** *Suppose that there exists a Lyapunov function  $V$ . Then, for any  $\epsilon > 0$ , there exists  $\delta > 0$ , such that*

$$\sup_{t \geq 0} |x(t) - x^*| \leq \epsilon,$$

*provided that  $|x_0 - x^*| \leq \delta$ .*

**PROOF.** Condition (3.34) on  $V$  implies that for fixed  $\epsilon > 0$ , there exists  $\gamma > 0$  (sufficiently small) such that

$$\{x : |x - x^*| \leq 2\epsilon, V(x) \leq V(x^*) + \gamma\} \subset \{x : |x - x^*| \leq \epsilon\}.$$

Choose  $\delta$  ( $0 < \delta < \epsilon$ ) such that

$$\{x : |x - x^*| \leq \delta\} \subset \{x : |x - x^*| \leq 2\epsilon, V(x) \leq V(x^*) + \gamma\}.$$

By using the fundamental property of a Lyapunov function  $V$

$$\frac{d}{dt}V(x(t)) = f(x(t)) \cdot V'(x(t)) \leq 0, \quad t \geq 0, \quad (3.37)$$

we obtain that

$$V(x(t)) \leq V(x_0) \leq V(x^*) + \gamma \quad \text{if } |x_0 - x^*| \leq \delta.$$

In fact, we have

$$|x(s) - x^*| \leq 2\epsilon \quad \text{for all } s \geq 0,$$

since otherwise, there would exist  $t > 0$  such that  $|x(t) - x^*| = 2\epsilon$ . From  $V(x(t)) \leq V(x^*) + \gamma$  we would arrive at a contradiction.  $\square$

**THEOREM 3.30** (Global stability). *Suppose that there exists  $V \in \mathcal{C}^1(\mathbb{R}^d, \mathbb{R})$  satisfying (3.34) such that*

$$f(x) \cdot V'(x) < 0 \quad \text{for any } x \neq x^*, \quad (3.38)$$

*and the set  $\{x : V(x) \leq V(x_0)\}$  is bounded. Then the solution  $x(t)$  of (3.29) converges to  $x^*$  as  $t \rightarrow +\infty$ .*

**PROOF.** As in the proof of Theorem 3.29, we have  $V(x(t)) \leq V(x_0)$  and thus  $\{x(t) : t \geq 0\}$  is bounded. More precisely, (3.37) yields

$$\int_0^{+\infty} |f(x(t)) \cdot V'(x(t))| dt = \int_0^{+\infty} -f(x(t)) \cdot V'(x(t)) dt \leq V(x_0) - V^*,$$

where  $V^* := \lim_{t \rightarrow +\infty} V(x(t))$ . Note that  $V^* > -\infty$  since  $(x(t))_{t \geq 0}$  is bounded.

Therefore, we can choose  $(t_n)_{n \in \mathbb{N}}$  such that  $x(t_n) \rightarrow \tilde{x}$  and  $f(x(t_n)) \cdot V'(x(t_n)) \rightarrow 0$  as  $n \rightarrow +\infty$ . Hence,

$$f(\tilde{x}) \cdot V'(\tilde{x}) = 0,$$

which, by (3.38), gives  $\tilde{x} = x^*$ .  $\square$

**EXAMPLE 3.31.** *Consider the equation  $dx/dt = f(x)$  with the initial condition  $x(0) = x_0$ , where  $f(0) = 0$  and  $x^\top \cdot f(x) < 0$  if  $x \neq 0$ . Then  $x^* = 0$  is the unique equilibrium point. Let  $V(x) := |x|^2$ . We have  $V(x) > V(0)$  for  $x \neq 0$  and  $dV/dt = 2x^\top \cdot (dx/dt) = 2x^\top \cdot f(x) < 0$ . Moreover,  $\{x : V(x) \leq V(x_0)\}$  is bounded since  $V(x) \rightarrow +\infty$  if  $|x| \rightarrow +\infty$ . Therefore, it follows from Theorem 3.30 that  $\lim_{t \rightarrow +\infty} x(t) = 0$ .*

### 3.6. Periodic linear systems

In this section, we consider the equation

$$\frac{dx(t)}{dt} = A(t)x(t)$$

in the special case where the matrix  $A(t)$  is periodic,

$$A(t+T) = A(t), \quad T > 0.$$

This periodicity condition implies that  $x(t+T)$  is again a solution if  $x(t)$  is. A first naive guess would be that  $x(t+T) = x(t)$ . However, this is too much to hope for since it already fails with  $A(t)$  a constant matrix (see Example 3.36). Nevertheless, as it will be shown later,  $x(t)$  exhibits an exponential behavior if we move on by one period. If we factor out this exponential term, the remainder is periodic.

For  $t_0 \in \mathbb{R}$ , let the matrix  $Y(t, t_0)$  be the unique solution of

$$\begin{cases} \frac{dY(t, t_0)}{dt} = A(t)Y(t, t_0), & t > t_0, \\ Y(t_0, t_0) = I. \end{cases} \quad (3.39)$$

LEMMA 3.32. *Suppose that  $A(t)$  is periodic with period  $T$ . Then  $Y(t, t_0)$  satisfies*

$$Y(t+T, t_0+T) = Y(t, t_0).$$

PROOF. By

$$\frac{dY(t+T, t_0+T)}{dt} = A(t+T)Y(t+T, t_0+T) = A(t)Y(t+T, t_0+T)$$

and  $Y(t_0+T, t_0+T) = I$ , we see that  $Y(t+T, t_0+T)$  solves (3.39). Thus it is equal to  $Y(t, t_0)$  by uniqueness.  $\square$

Let  $Z(t_0) = Y(t_0+T, t_0)$ . By Lemma 3.32,  $Z$  is periodic,

$$Z(t_0+T) = Z(t_0).$$

LEMMA 3.33. *For all  $l \in \mathbb{N}$ , we have*

$$Y(t_0+lT, t_0) = Z(t_0)^l.$$

PROOF. We have

$$\begin{aligned} Y(t_0+lT, t_0) &= Y(t_0+lT, t_0+(l-1)T)Y(t_0+(l-1)T, t_0) \\ &= Z(t_0+(l-1)T)Y(t_0+(l-1)T, t_0) \\ &= Z(t_0)Y(t_0+(l-1)T, t_0) \\ &= Z(t_0)^l Y(t_0, t_0) = Z(t_0)^l. \end{aligned}$$

$\square$

From Liouville's formula (3.9), it follows that

$$\det Z(t_0) = e^{\int_{t_0}^{t_0+T} \text{tr}(A(s)) ds} = e^{\int_0^T \text{tr}(A(s)) ds},$$

which is independent of  $t_0$  and positive.

Therefore, by Lemma 3.6, we can find a matrix  $Q(t_0)$  such that

$$Z(t_0) = e^{TQ(t_0)}, \quad Q(t_0+T) = Q(t_0).$$

Note that  $Q(t_0)$  is not unique. Note also that  $Q(t_0)$  is complex even if  $A(t)$  is real unless all real eigenvalues of  $Z(t_0)$  are positive.

Now, writing

$$Y(t, t_0) = P(t, t_0)e^{(t-t_0)Q(t_0)},$$

a straightforward computation shows that

$$\begin{aligned} P(t+T, t_0) &= Y(t+T, t_0)M(t_0)^{-1}e^{-(t-t_0)Q(t_0)} \\ &= Y(t+T, t_0+T)e^{-(t-t_0)Q(t_0)} \\ &= Y(t, t_0)e^{-(t-t_0)Q(t_0)} = P(t, t_0). \end{aligned}$$

In summary, we have proven **Floquet's theorem**.

**THEOREM 3.34.** *Suppose that  $A(t)$  is periodic. Then  $Y(t, t_0)$  defined by (3.39) has the form*

$$Y(t, t_0) = P(t, t_0)e^{(t-t_0)Q(t_0)},$$

where  $P(\cdot, t_0)$  has the same period as  $A(\cdot)$  and  $P(t_0, t_0) = I$ .

As a consequence of Floquet's theorem we obtain the following result.

**COROLLARY 3.35.** *The transformation  $y(t) = P(t, t_0)^{-1}x(t)$  renders the system  $dx/dt = A(t)x$  into one with constant coefficients,*

$$\frac{dy(t)}{dt} = Z(t_0)y(t).$$

**EXAMPLE 3.36.** *Consider the one-dimensional case*

$$\frac{dx}{dt} = a(t)x, \quad a(t+T) = a(t).$$

Then

$$Y(t, t_0) = e^{\int_{t_0}^t a(s) ds}$$

and

$$Z(t_0) = e^{\int_{t_0}^{t_0+T} a(s) ds} = e^{T\langle a \rangle}, \quad \langle a \rangle = \frac{1}{T} \int_0^T a(s) ds.$$

Moreover,

$$P(t, t_0) = e^{\int_{t_0}^t (a(s) - \langle a \rangle) ds}, \quad Q(t_0) = \langle a \rangle.$$

The eigenvalues  $\rho_j$  of  $Z(t_0)$  are known as **Floquet multipliers** and the eigenvalues  $\gamma_j$  of  $Q(t_0)$  are known as **Floquet exponents**.  $\rho_j$  and  $\gamma_j$  are related via  $\rho_j = e^{T\gamma_j}$ . Since the periodic part  $P(t, t_0)$  is bounded, we obtain the following result as another consequence of Floquet's theorem.

**THEOREM 3.37.** *A periodic linear system is stable if all Floquet multipliers satisfy  $|\rho_j| \leq 1$  (respectively all Floquet exponents satisfy  $\Re\gamma_j \leq 0$ ) and for all Floquet multipliers with  $|\rho_j| = 1$  (respectively all Floquet exponents with  $\Re\gamma_j = 0$ ) the algebraic and geometric multiplicities are equal.*

**EXAMPLE 3.38.** *Consider **Hill's equation***

$$\frac{d^2x(t)}{dt^2} + q(t)x(t) = 0, \quad q(t+T) = q(t). \quad (3.40)$$

In this case, the associated system is with

$$A(t) = \begin{pmatrix} 0 & 1 \\ -q(t) & 0 \end{pmatrix}.$$

Let  $x_1$  and  $x_2$  be the solutions of (3.40) corresponding to the initial conditions

$$x_1(t_0, t_0) = 1, \frac{dx_1}{dt}(t_0, t_0) = 0 \quad \text{and} \quad x_2(t_0, t_0) = 0, \frac{dx_2}{dt}(t_0, t_0) = 1.$$

Then

$$Y(t, t_0) = \begin{pmatrix} x_1(t, t_0) & x_2(t, t_0) \\ \frac{dx_1}{dt}(t, t_0) & \frac{dx_2}{dt}(t, t_0) \end{pmatrix}.$$



Liouville's formula (3.9) shows that  $\det Y(t, t_0) = 1$  and hence the characteristic equation for

$$Z(t_0) = \begin{pmatrix} x_1(t_0 + T, t_0) & x_2(t_0 + T, t_0) \\ \frac{dx_1}{dt}(t_0 + T, t_0) & \frac{dx_2}{dt}(t_0 + T, t_0) \end{pmatrix}$$

is given by

$$\lambda^2 - 2\Delta\lambda + 1 = 0, \quad \Delta := \frac{\operatorname{tr}(Z(t_0))}{2}.$$

Therefore, by Theorem 3.37, Hill's equation (3.40) is stable if  $|\Delta| < 1$  and unstable if  $|\Delta| > 1$ .

### 3.7. Problems

PROBLEM 3.39 (**Laplace transform**). (i) Prove that if  $A \in \mathbb{M}_d(\mathbb{R})$ , then

$$e^{tA} - I = \int_0^t A e^{sA} ds.$$

(ii) Prove that if all eigenvalues of  $A$  have negative real parts, then

$$-A^{-1} = \int_0^{+\infty} e^{sA} ds.$$

(iii) Prove that if  $s \in \mathbb{R}$  is sufficiently large, then

$$(sA - I)^{-1} = \int_0^{+\infty} e^{s(A-tI)} ds,$$

that is, the Laplace transform of  $e^{tA}$  is  $(sI - A)^{-1}$ .

PROBLEM 3.40. Let  $A \in \mathbb{M}_d(\mathbb{R})$ .

(i) Apply the Jacobi formula

$$\frac{d}{dt} \det B(t) = (\det B(t)) \operatorname{tr}(B(t)^{-1} \frac{dB}{dt}(t)) \quad (3.41)$$

for  $B(t) = e^{tA}$  to prove that

$$\det e^A = e^{\operatorname{tr} A}.$$

(ii) Prove that a vector  $u$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\alpha$  if and only if  $u$  is an eigenvector of  $e^A$  corresponding to the eigenvalue  $e^\alpha$ .

(iii) Prove that if  $\det A(t) \neq 0$ , then

$$\frac{d}{dt} A^{-1}(t) = -A^{-1}(t) \frac{dA}{dt}(t) A^{-1}(t).$$

(iv) Prove that

$$\det(I + \epsilon A + o(\epsilon)) = 1 + \epsilon \operatorname{tr} A + o(\epsilon),$$

where  $o(\epsilon)$  (Landau symbol) collects terms which vanish faster than  $\epsilon$  as  $\epsilon \rightarrow 0$ .

PROBLEM 3.41 (**Reduction of order**). Use reduction of order to find the general solution of the following equations:

(i)

$$t \frac{d^2 x}{dt^2} - 2(t+1) \frac{dx}{dt} + (t+2)x = 0, \quad x_1(t) = e^t.$$

(ii)

$$t^2 \frac{d^2 x}{dt^2} - 3t \frac{dx}{dt} + 4x = 0, \quad x_1(t) = t^2.$$

PROBLEM 3.42. (i) Verify that the second-order equation

$$\frac{d^2x}{dt^2} + (1 - t^2)x = 0, \quad (3.42)$$

can be factorized as

$$\left(\frac{d}{dt} - t\right)\left(\frac{d}{dt} + t\right)x = 0. \quad (3.43)$$

(ii) By solving two first-order problems, find the solution of (3.42).

PROBLEM 3.43. Let  $A \in \mathbb{M}_d(\mathbb{R})$  be independent of  $t$ . Consider the linear system of ODEs (3.24).

(i) Assume that there exist two positive definite matrices  $P$  and  $Q$  such that

$$A^\top P + PA = -Q. \quad (3.44)$$

Prove that  $V(x) := x^\top Px$  is a Lyapunov function for (3.24).

(ii) Define

$$r := \min_{x \neq 0} \frac{x^\top Qx}{x^\top Px}.$$

Prove that  $V(x(t)) \leq e^{-rt}V(x_0)$ , where  $x(t)$  is the solution of (3.24).

(iii) Assume that every eigenvalue of  $A$  has a negative real part. Prove that given  $Q$ , the solution  $P$  to (3.44) can be written as

$$P = \int_0^{+\infty} e^{tA^\top} Q e^{tA} dt.$$

PROBLEM 3.44 (**Convergence of the gradient algorithm for finding a local minimum of a function**). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and assume that  $x^*$  is a local minimum, i.e.,  $f(x^*) < f(x)$  for  $x$  close enough but not equal to  $x^*$ . Assume that  $f$  is continuously differentiable and let  $x(t)$  be the solution to

$$\begin{cases} \frac{dx}{dt} = -\nabla f(x), & t \in [0, +\infty[, \\ x(0) = x_0 \in \mathbb{R}^d. \end{cases}$$

(i) Prove that if  $x_0$  is close to  $x^*$  then  $\lim_{t \rightarrow +\infty} x(t) = x^*$ .

(ii) Let  $f(x) = \frac{1}{2}x^\top Qx$ , where  $Q$  is symmetric, positive definite. Show directly that  $x(t)$  converges to zero ( $= x^*$ ).

PROBLEM 3.45. Consider

$$\frac{dx(t)}{dt} = a(t)Ax(t),$$

where  $a(t)$  is a scalar periodic function with period  $T$  and  $A$  is a constant  $2 \times 2$  matrix. Compute the Floquet exponent, and find  $P(t, t_0)$  and  $Q(t_0)$  in this case.

## Numerical solution of ordinary differential equations

### 4.1. Introduction

This chapter is concerned with the numerical solution of initial value problems for systems of ordinary differential equations. Since there is no hope of solving the vast majority of differential equations in explicit and analytic form, the design of suitable numerical schemes for accurately approximating solutions is essential. Explicit solutions, when they are known, can also be used as test cases for tracking the reliability and accuracy of a chosen numerical scheme. In this chapter, we survey the most basic numerical methods for solving initial value problems. It goes without saying that some equations are more difficult to accurately approximate than others, and a variety of more specialized techniques are employed when confronted with a recalcitrant system. However, all of the more advanced developments build on the basic schemes and ideas laid out in this chapter.

### 4.2. The general explicit one-step method

**4.2.1. Consistency, stability and convergence.** Consider the initial value problem

$$\begin{cases} \frac{dx}{dt} = f(t, x), & t \in [0, T], \\ x(0) = x_0, & x_0 \in \mathbb{R}, \end{cases} \quad (4.1)$$

where  $f \in C^0([0, T] \times \mathbb{R})$  satisfies the Lipschitz condition (2.2).

Starting at the initial time  $t = 0$ , we introduce successive discretization points

$$t_0 = 0 < t_1 < t_2 < \dots,$$

continuing on until we reach the final time  $T$ . To keep the analysis as simple as possible, we use a uniform **step size**, and so

$$\Delta t := t_{k+1} - t_k > 0, \quad (4.2)$$

does not depend on  $k$  and is assumed to be relatively small, with  $t_k = k\Delta t$ . We also suppose that  $K = T/(\Delta t)$  is an integer.

A general **explicit one-step method** may be written in the form:

$$x^{k+1} = x^k + \Delta t \Phi(t_k, x^k, \Delta t), \quad (4.3)$$

for some continuous function  $\Phi(t, x, h)$ . In (4.3), taking in succession  $k = 0, 1, \dots, K - 1$ , **one-step** at a time, the approximate values  $x^k$  of  $x$  at  $t_k$  can be easily obtained. Scheme (4.3) is called **explicit** since  $x^{k+1}$  is obtained from  $x^k$ .  $x^{k+1}$  appears only on the left-hand side of (4.3).

We define the **truncation error** of the numerical scheme (4.3) by

$$T_k(\Delta t) = \frac{x(t_{k+1}) - x(t_k)}{\Delta t} - \Phi(t_k, x(t_k), \Delta t). \quad (4.4)$$

As  $\Delta t \rightarrow 0, k \rightarrow +\infty, k\Delta t = t$ ,

$$T_k(\Delta t) \rightarrow \frac{dx}{dt} - \Phi(t, x, 0).$$

**DEFINITION 4.1 (Consistency).** *The numerical scheme (4.3) is **consistent** with (4.1) if*

$$\Phi(t, x, 0) = f(t, x) \quad \text{for all } t \in [0, T] \text{ and } x \in \mathbb{R}.$$

DEFINITION 4.2 (**Stability**). *The numerical scheme (4.3) for solving (4.1) is **stable** if  $\Phi$  is Lipschitz continuous in  $x$ , i.e., there exist positive constants  $C_\Phi$  and  $h_0$  such that*

$$|\Phi(t, x, h) - \Phi(t, y, h)| \leq C_\Phi |x - y|, \quad t \in [0, T], h \in [0, h_0], x, y \in \mathbb{R}. \quad (4.5)$$

Define the **global error** of the numerical scheme (4.3) by

$$e_k = x^k - x(t_k). \quad (4.6)$$

DEFINITION 4.3 (**Convergence**). *The numerical scheme (4.3) for solving (4.1) is **convergent** if*

$$|e_k| \rightarrow 0 \quad \text{as } \Delta t \rightarrow 0, k \rightarrow +\infty, k\Delta t = t \in [0, T].$$

THEOREM 4.4 (Dahlquist-Lax equivalence theorem). *The numerical scheme (4.3) is convergent if and only if it is consistent and stable.*

PROOF. From (4.1), it follows that

$$x(t_{k+1}) - x(t_k) = \int_{t_k}^{t_{k+1}} f(s, x(s)) ds,$$

which gives

$$x(t_{k+1}) - x(t_k) = (\Delta t)f(t_k, x(t_k)) + \int_{t_k}^{t_{k+1}} [f(s, x(s)) - f(t_k, x(t_k))] ds.$$

Therefore,

$$\left| x(t_{k+1}) - x(t_k) - (\Delta t)f(t_k, x(t_k)) \right| = \left| \int_{t_k}^{t_{k+1}} [f(s, x(s)) - f(t_k, x(t_k))] ds \right| \leq (\Delta t) \omega_1(\Delta t), \quad (4.7)$$

where

$$\omega_1(\Delta t) := \sup \{ |f(t, x(t)) - f(s, x(s))|, 0 \leq s, t \leq T, |t - s| \leq \Delta t \}. \quad (4.8)$$

Note that  $\omega_1(\Delta t) \rightarrow 0$  as  $\Delta t \rightarrow 0$ . Moreover, if  $f$  is Lipschitz in  $t$ , then  $\omega_1(\Delta t) = O(\Delta t)$ .

From (4.3) and

$$e_{k+1} - e_k = x^{k+1} - x^k - (x(t_{k+1}) - x(t_k)),$$

we obtain

$$e_{k+1} - e_k = \Delta t \Phi(t_k, x^k, \Delta t) - (x(t_{k+1}) - x(t_k)),$$

or equivalently,

$$e_{k+1} - e_k = \Delta t [\Phi(t_k, x^k, \Delta t) - f(t_k, x(t_k))] - [x(t_{k+1}) - x(t_k) - \Delta t f(t_k, x(t_k))].$$

Write

$$\begin{aligned} e_{k+1} - e_k &= \Delta t [\Phi(t_k, x^k, \Delta t) - \Phi(t_k, x(t_k), \Delta t) + \Phi(t_k, x(t_k), \Delta t) \\ &\quad - f(t_k, x(t_k))] - [x(t_{k+1}) - x(t_k) - \Delta t f(t_k, x(t_k))]. \end{aligned} \quad (4.9)$$

Let

$$\omega_2(\Delta t) := \sup \{ |\Phi(t, x, h) - f(t, x)|, t \in [0, T], x \in \mathbb{R}, 0 < h \leq (\Delta t) \}. \quad (4.10)$$

Since the numerical scheme is consistent,

$$\left| \Phi(t_k, x(t_k), \Delta t) - f(t_k, x(t_k)) \right| \leq \omega_2(\Delta t) \rightarrow 0 \quad \text{as } \Delta t \rightarrow 0. \quad (4.11)$$

On the other hand, from the stability condition (4.5), it follows that

$$\left| \Phi(t_k, x^k, \Delta t) - \Phi(t_k, x(t_k), \Delta t) \right| \leq C_\Phi |e_k|. \quad (4.12)$$

Combining (4.7), (4.9), (4.11), and (4.12) yields

$$|e_{k+1}| \leq (1 + C_\Phi \Delta t) |e_k| + \Delta t \omega_3(\Delta t), \quad 0 \leq k \leq K - 1, \quad (4.13)$$

where  $K = T/(\Delta t)$  and  $\omega_3(\Delta t) := \omega_1(\Delta t) + \omega_2(\Delta t) \rightarrow 0$  as  $\Delta t \rightarrow 0$ . By induction, we deduce from (4.13) that

$$|e_{k+1}| \leq (1 + C_\Phi \Delta t)^k |e_0| + (\Delta t) \omega_3(\Delta t) \sum_{l=0}^{k-1} (1 + C_\Phi \Delta t)^l, \quad 0 \leq k \leq K. \quad (4.14)$$

Estimate (4.14) together with

$$\sum_{l=0}^{k-1} (1 + C_\Phi \Delta t)^l = \frac{(1 + C_\Phi \Delta t)^k - 1}{C_\Phi \Delta t},$$

and

$$(1 + C_\Phi \Delta t)^K \leq (1 + C_\Phi \frac{T}{K})^K \leq e^{C_\Phi T},$$

yields

$$|e_k| \leq e^{C_\Phi T} |e_0| + \frac{e^{C_\Phi T} - 1}{C_\Phi} \omega_3(\Delta t). \quad (4.15)$$

Therefore, if  $e_0 = 0$ , then as  $\Delta t \rightarrow 0, k \rightarrow +\infty$  such that  $k\Delta t = t \in [0, T]$

$$\lim_{k \rightarrow +\infty} |e_k| = 0,$$

which shows that the scheme is in fact convergent.  $\square$

**DEFINITION 4.5.** *An explicit one-step method is said to be of **order**  $p$  if there exist positive constants  $h_0$  and  $C$  such that*

$$|T_k(\Delta t)| \leq C(\Delta t)^p, \quad 0 < \Delta t \leq h_0, k = 0, \dots, K-1,$$

where the truncation error  $T_k(\Delta t)$  is defined by (4.4).

If the explicit one-step method is stable, then the global error is bounded by the truncation error.

**PROPOSITION 4.6.** *Consider the explicit one-step scheme (4.3), where  $\Phi$  satisfies the stability condition (4.5). Suppose that  $e_0 = 0$ . Then*

$$|e_{k+1}| \leq \frac{(e^{C_\Phi T} - 1)}{C_\Phi} \max_{0 \leq l \leq k} |T_l(\Delta t)| \quad \text{for } k = 0, \dots, K-1, \quad (4.16)$$

where the truncation error  $T_l$  and the global error  $e_k$  are defined by (4.4) and (4.6), respectively.

**PROOF.** From (4.9), we have

$$e_{k+1} - e_k = -(\Delta t)T_k(\Delta t) + (\Delta t) \left[ \Phi(t_k, x^k, \Delta t) - \Phi(t_k, x(t_k), \Delta t) \right],$$

so we get

$$\begin{aligned} |e_{k+1}| &\leq (1 + C_\Phi(\Delta t))|e_k| + (\Delta t)|T_k(\Delta t)| \\ &\leq (1 + C_\Phi(\Delta t))|e_k| + (\Delta t) \max_{0 \leq l \leq k} |T_l(\Delta t)|. \end{aligned}$$

In exactly the same manner as in the proof of Theorem 4.4, we obtain estimate (4.16).  $\square$

**4.2.2. Explicit Euler's method.** Let  $\Phi(t, x, h) = f(t, x)$ . The numerical method (4.3) reduces to

$$x^{k+1} = x^k + (\Delta t)f(t, x^k). \quad (4.17)$$

The numerical method (4.17) is called the **explicit Euler scheme**.

**THEOREM 4.7.** *Consider the initial value problem (4.1). Suppose that  $f$  satisfies the Lipschitz condition (2.2) and  $f$  is Lipschitz with respect to  $t$ . Then the explicit Euler scheme (4.17) is convergent and the global error  $e_k$  is of order  $\Delta t$ . If  $f \in C^1$ , then (4.17) is of order one.*

**PROOF.** Since  $f$  satisfies the Lipschitz condition (2.2) then the numerical scheme with  $\Phi(t, x, h) = f(t, x)$  is stable. Moreover, it is consistent since  $\Phi(t, x, 0) = f(t, x)$  for all  $t \in [0, T]$  and  $x \in \mathbb{R}$ . Therefore, by Theorem 4.4, (4.17) is convergent. Furthermore, since  $f$  is Lipschitz in  $t$ ,  $\omega_1(\Delta t) = O(\Delta t)$ , where  $\omega_1$  is defined by (4.8). On the other hand,  $\omega_2(\Delta t) = 0$ , and hence  $\omega_3(\Delta t) = O(\Delta t)$ , where  $\omega_2$  is defined by (4.10) and  $\omega_3 = \omega_1 + \omega_2$ . Then, from (4.15), we have  $|e_k| = O(\Delta t)$  for  $1 \leq k \leq K$ . Now if  $f \in C^1$ , then from Theorem 2.20  $x \in C^2$ . By using the mean-value theorem, we have

$$\begin{aligned} T_k(\Delta t) &= \frac{1}{\Delta t} \left( x(t_{k+1}) - x(t_k) \right) - f(t_k, x(t_k)) \\ &= \frac{1}{\Delta t} \left( x(t_k) + (\Delta t) \frac{dx}{dt}(t_k) + \frac{(\Delta t)^2}{2} \frac{d^2x}{dt^2}(\tau) - x(t_k) \right) - f(t_k, x(t_k)) \\ &= \frac{\Delta t}{2} \frac{d^2x}{dt^2}(\tau), \end{aligned} \quad (4.18)$$

for some  $\tau \in [t_k, t_{k+1}]$ , which shows that (4.17) is of first order.  $\square$

**REMARK 4.8 (Round off error effects).** *Theorem 4.7 is true provided the arithmetic in calculating the numerical approximation is perfect, that is, when performing the operations required by (4.17) no errors occur. However computers always round off real numbers. In numerical methods rounding errors become important when the step size  $\Delta t$  is comparable with the precision of the computations. Thus, when running Euler's method (4.17), the best we can do is to compute the solution of the perturbed scheme:*

$$\tilde{x}^{k+1} = \tilde{x}^k + \Delta t f(t_k, \tilde{x}^k) + (\Delta t)\mu^k + \rho^k,$$

where  $\mu^k$  and  $\rho^k$  represent the errors in  $f$  and in the assembling, respectively. Assume that  $|\mu^k| \leq \mu$  and  $|\rho^k| \leq \rho$  for all  $k$  and  $f \in C^1$ . Defining  $\tilde{e}^k = x(t_k) - \tilde{x}^k$ , we have

$$|\tilde{e}^{k+1}| \leq (1 + C_f \Delta t) \tilde{e}^k + (\Delta t)\mu + \rho,$$

and hence

$$|\tilde{e}^k| \leq e^{C_f T} |\tilde{e}^0| + (\Delta t) e^{C_f T} \int_0^T \left| \frac{d^2x}{dt^2} \right|(s) ds + \mu(\Delta t) \frac{e^{C_f T}}{C_f} + \rho \frac{T}{\Delta t} e^{C_f T},$$

where  $C_f$  is the Lipschitz constant for  $f$ .

Introduce

$$\varphi(\Delta t) = \frac{\mu e^{C_f T}}{C_f} \Delta t + \frac{T \rho e^{C_f T}}{\Delta t}.$$

One can see that  $\varphi$  attains its minimum at  $\sqrt{\rho C_f T / \mu}$  and diverges for  $\Delta t \rightarrow 0$ . From a practical point of view, it is better to take time steps that are larger than  $\sqrt{\rho C_f T / \mu}$ .

**REMARK 4.9 (Control of the time step).** *In (4.17) the time step is uniform and is chosen such that the global error  $|e_k|$  is smaller than a given tolerance. In view of (4.18) this supposes a good knowledge of the exact solution. An alternative method consists in computing the numerical solution for an arbitrary  $\Delta t$  and then for  $2\Delta t$ . If the discrepancy between the two numerical solutions is smaller than the tolerance, we keep  $\Delta t$ . If not, we restart the calculations with a smaller step size, say  $\Delta t/2$ , until we reach the target.*

**4.2.3. High-order methods.** In general, the order of a numerical solution method governs both the accuracy of its approximations and the speed at which they converge to the true solution as the step size  $\Delta t \rightarrow 0$ . Although the explicit Euler method is simple and easy to implement, it is only a first order scheme as shown in Theorem 4.7, and therefore of limited use. So, the goal is to devise simple numerical methods that enjoy a higher order of accuracy. The higher its order, the more accurate the numerical scheme, and hence the larger the step size that can be used to produce the solution to a desired accuracy. However, this should be balanced with the fact that higher order methods inevitably require more computational effort at each step.

4.2.3.1. *Taylor methods.* The explicit Euler scheme is based on a first order Taylor approximation to the solution. The Taylor expansion of the solution  $x(t)$  at the discretization points  $t_{k+1}$  has the form

$$x(t_{k+1}) = x(t_k + \Delta t) = x(t_k) + (\Delta t) \frac{dx}{dt}(t_k) + \frac{(\Delta t)^2}{2} \frac{d^2x}{dt^2}(t_k) + \frac{(\Delta t)^3}{6} \frac{d^3x}{dt^3}(t_k) + \dots \quad (4.19)$$

We can evaluate the first derivative term by using the differential equation

$$\frac{dx}{dt} = f(t, x). \quad (4.20)$$

The second derivative can be found by differentiating the equation with respect to  $t$ . Invoking the chain rule,

$$\frac{d^2x}{dt^2} = \frac{d}{dt} f(t, x) = \frac{\partial f}{\partial t}(t, x) + \frac{\partial f}{\partial x}(t, x) \frac{dx}{dt}. \quad (4.21)$$

Substituting (4.20) and (4.21) into (4.19) and truncating at order  $(\Delta t)^2$  leads to the **second order Taylor method**

$$x^{k+1} = x^k + (\Delta t) f(t_k, x^k) + \frac{(\Delta t)^2}{2} \left( \frac{\partial f}{\partial t}(t_k, x^k) + \frac{\partial f}{\partial x}(t_k, x^k) f(t_k, x^k) \right), \quad (4.22)$$

in which we have replaced the solution value  $x(t_k)$  by its computed approximation  $x^k$ . The resulting method is of second order.

PROPOSITION 4.10. *Suppose that  $f \in \mathcal{C}^2$ . Then (4.22) is of second order.*

PROOF. If  $f$  is of class  $\mathcal{C}^2$ , then by Theorem 2.20  $x \in \mathcal{C}^3$ . Therefore, by using the Taylor expansion (4.19), we obtain that the truncation error  $T_k$  is given by

$$T_k(\Delta t) = \frac{(\Delta t)^2}{6} \frac{d^3x}{dt^3}(\tau),$$

for some  $\tau \in [t_k, t_{k+1}]$  and so, (4.22) is of second order.  $\square$

Higher order Taylor methods are obtained by including further terms in the expansion (4.19). Whereas higher order Taylor methods are easy to motivate, they are rarely used in practice. There are two principal difficulties:

- (i) Owing to their dependence upon the partial derivatives of  $f$ ,  $f$  needs to be smooth;
- (ii) Efficient evaluation of the terms in the Taylor approximation and avoidance of round off errors are significant concerns.

4.2.3.2. *Integral equation method.* In order to design high-order numerical schemes that avoid the complications inherent in a direct Taylor expansion, we replace the differential equation by an equivalent **integral equation**. The solution  $x(t)$  of (4.1) coincides with the solution to the **integral equation**

$$x(t) = x_0 + \int_0^t f(s, x(s)) ds, \quad t \in [0, T]. \quad (4.23)$$

Starting at the discretization point  $t_k$  instead of 0, and integrating until time  $t = t_{k+1}$  gives an expression

$$x(t_{k+1}) = x(t_k) + \int_{t_k}^{t_{k+1}} f(s, x(s)) ds, \quad (4.24)$$

that implicitly computes the value of the solution at the subsequent discretization point. Comparing formula (4.24) with the explicit Euler method

$$x^{k+1} = x^k + (\Delta t)f(t_k, x^k),$$

where  $\Delta t$  is defined by (4.2) and assuming for the moment that  $x^k = x(t_k)$  is exact, we see that we are merely approximating the integral by

$$\int_{t_k}^{t_{k+1}} f(s, x(s)) ds \approx (\Delta t)f(t_k, x(t_k)), \quad (4.25)$$

which is the **left endpoint rule** for numerical integration—that approximates the integral of  $f(t, x(t))$  between  $t_k \leq t \leq t_{k+1}$  by the area of the rectangle whose height  $f(t_k, x(t_k))$  is prescribed by the left endpoint of the curve  $t \mapsto f(t, x(t))$ . Approximation (4.25) is not an especially accurate method of numerical integration. Better methods include the **Trapezoid rule**, which approximates the integral of the function  $f(t, x(t))$  between  $t_k \leq t \leq t_{k+1}$  by the area of the trapezoid obtained by connecting the points  $f(t_k, x(t_k))$  and  $f(t_{k+1}, x(t_{k+1}))$  of the curve  $t \mapsto f(t, x(t))$  by a straight line.

We recall the following basic numerical integration formulas for continuous functions.

(i) **Trapezoidal rule:**

$$\int_{t_k}^{t_{k+1}} g(s) ds \approx \frac{\Delta t}{2} \left( g(t_{k+1}) + g(t_k) \right); \quad (4.26)$$

(ii) **Simpson's rule:**

$$\int_{t_k}^{t_{k+1}} g(s) ds \approx \frac{\Delta t}{6} \left( g(t_{k+1}) + 4g\left(\frac{t_k + t_{k+1}}{2}\right) + g(t_k) \right); \quad (4.27)$$

(iii) The Trapezoidal rule is **exact** for polynomials of order one, while the Simpson's rule is exact for polynomials of second order.

Replacing (4.25) by the more accurate Trapezoidal approximation

$$\int_{t_k}^{t_{k+1}} f(s, x(s)) ds \approx \frac{(\Delta t)}{2} \left[ f(t_k, x(t_k)) + f(t_{k+1}, x(t_{k+1})) \right], \quad (4.28)$$

and substituting (4.28) into the integral equation (4.24) leads to the **Trapezoidal scheme**

$$x^{k+1} = x^k + \frac{(\Delta t)}{2} \left[ f(t_k, x^k) + f(t_{k+1}, x^{k+1}) \right]. \quad (4.29)$$

The Trapezoidal scheme is an **implicit numerical method**, since the updated value  $x^{k+1}$  appears on both sides of the equation, and hence is only defined implicitly. Only for very simple functions  $f(t, x)$  can one expect to solve (4.29) explicitly for  $x^{k+1}$  given  $t_k, x^k$ , and  $t_{k+1}$ .

PROPOSITION 4.11. *Suppose that  $f \in \mathcal{C}^2$  and*

$$\frac{(\Delta t)C_f}{2} < 1, \quad (4.30)$$

where  $C_f$  is the Lipschitz constant for  $f$  in  $x$  defined by (2.2). Then the Trapezoidal scheme (4.29) is convergent and is of second order.

PROOF. Let  $\Phi$  be defined **implicitly** by

$$\Phi(t, x, \Delta t) := \frac{1}{2} \left[ f(t, x) + f(t + \Delta t, x + (\Delta t)\Phi(t, x, \Delta t)) \right].$$

The scheme (4.29) is clearly consistent. In order to show that it converges, according to Theorem 4.4, we must establish the stability condition (4.5). We have

$$|\Phi(t, x, \Delta t) - \Phi(t, y, \Delta t)| \leq C_f |x - y| + \frac{\Delta t}{2} C_f |\Phi(t, x, \Delta t) - \Phi(t, y, \Delta t)|.$$



Hence

$$\left(1 - \frac{(\Delta t)C_f}{2}\right) |\Phi(t, x, \Delta t) - \Phi(t, y, \Delta t)| \leq C_f |x - y|,$$

and therefore, (4.5) holds with

$$C_\Phi = \frac{C_f}{1 - \frac{(\Delta t)C_f}{2}},$$

provided that  $\Delta t$  satisfies (4.30). Now we prove that (4.29) is of second order.

By the mean-value theorem,

$$\begin{aligned} T_k(\Delta t) &= \frac{x(t_{k+1}) - x(t_k)}{\Delta t} - \frac{1}{2} \left[ f(t_k, x(t_k)) + f(t_{k+1}, x(t_{k+1})) \right] \\ &= -\frac{1}{12} (\Delta t)^2 \frac{d^3 x}{dt^3}(\tau), \end{aligned} \quad (4.31)$$

for some  $\tau \in [t_k, t_{k+1}]$ , and therefore (4.29) is of second order, provided that  $f \in \mathcal{C}^2$  (and consequently  $x \in \mathcal{C}^3$ ).  $\square$

An alternative is to replace in (4.29)  $x^{k+1}$  by  $x^k + (\Delta t)f(t_k, x^k)$ . This yields the **improved Euler scheme**

$$x^{k+1} = x^k + \frac{(\Delta t)}{2} \left[ f(t_k, x^k) + f(t_{k+1}, \mathbf{x}^k + (\Delta t)\mathbf{f}(t_k, \mathbf{x}^k)) \right]. \quad (4.32)$$

**PROPOSITION 4.12.** *The numerical scheme (4.32) is convergent and is of second order.*

The improved Euler scheme (4.32) performs comparably to the Trapezoidal scheme (4.29), and significantly better than the Euler scheme (4.17). The improved Euler scheme (4.32) is the simplest of a large family of so-called **predictor-corrector algorithms**. In general, one begins by using a relatively crude method—in this case the explicit Euler method—to predict a first approximation  $\tilde{x}^{k+1}$  to the desired solution value  $x(t_{k+1})$ . One then employs a more sophisticated, typically implicit, method to correct the original prediction, by replacing the required update  $x^{k+1}$  on the right-hand side of the implicit scheme by a less accurate prediction  $\tilde{x}^{k+1}$ . The resulting explicit, corrected value  $x^{k+1}$  will be an improved approximation of the true solution, provided the method has been designed with due care.

We can design a range of numerical solution schemes by implementing alternative numerical approximations to the integral equation (4.24). For example, the **midpoint rule** approximates the integral of  $f(t, x(t))$  between  $t_k \leq t \leq t_{k+1}$  by the area of the rectangle whose height is the value of  $f$  at the midpoint  $t = t_k + (\Delta t)/2$

$$\int_{t_k}^{t_{k+1}} f(s, x(s)) ds \approx (\Delta t) f\left(t_k + \frac{\Delta t}{2}, x\left(t_k + \frac{\Delta t}{2}\right)\right). \quad (4.33)$$

The midpoint rule has the same order of accuracy as the trapezoid rule. Substituting (4.33) into (4.24) leads to the **midpoint scheme**

$$x^{k+1} = x^k + (\Delta t) f\left(t_k + \frac{\Delta t}{2}, x^k + \frac{\Delta t}{2} f(t_k, x^k)\right), \quad (4.34)$$

where we have approximated  $x(t_k + \frac{\Delta t}{2})$  by  $x^k + \frac{\Delta t}{2} f(t_k, x^k)$ .

A comparison between the terms in the Taylor expansion (4.19) of  $x(t_{k+1})$  and (4.34) reveals that the midpoint scheme is also of second order.

### 4.3. Example of linear systems

Consider the linear system of ODEs (3.24), where  $A \in \mathbb{M}_d(\mathbb{C})$  is independent of  $t$ .

A one-step numerical scheme for solving (3.24) is said to be **stable** if there exists a positive constant  $C_0$  such that

$$|x^{k+1}| \leq C_0 |x^0| \quad \text{for all } k \in \mathbb{N}. \quad (4.35)$$

Consider the following schemes for solving (3.1):

(i) Explicit Euler's scheme

$$x^{k+1} = x^k + (\Delta t)Ax^k; \quad (4.36)$$

(ii) Implicit Euler's scheme

$$x^{k+1} = x^k + (\Delta t)Ax^{k+1}; \quad (4.37)$$

(iii) Trapezoidal scheme:

$$x^{k+1} = x^k + \frac{(\Delta t)}{2} \left[ Ax^k + Ax^{k+1} \right], \quad (4.38)$$

where  $k \in \mathbb{N}$ , and  $x^0 = x_0$ .

**PROPOSITION 4.13.** *Suppose that  $\Re \lambda_j < 0$  for all  $j$ . The following results hold:*

- (i) *The explicit Euler scheme (4.36) is stable for  $\Delta t$  small enough;*
- (ii) *The implicit Euler scheme is unconditionally stable;*
- (iii) *The Trapezoidal scheme (4.38) is unconditionally stable.*

**PROOF.** Consider the explicit Euler scheme (4.36). By a change of basis, we have

$$\tilde{x}^k = (I + \Delta t(D + N))^k \tilde{x}^0,$$

where  $\tilde{x}^k = Cx^k$ . If  $\tilde{x}^0 \in E_j$ , then

$$\tilde{x}^k = \sum_{l=0}^{\min\{k,d\}} C_k^l (1 + \Delta t \lambda_j)^{k-l} (\Delta t)^l N^l \tilde{x}^0,$$

where  $C_k^l$  is the binomial coefficient.

If  $|1 + (\Delta t)\lambda_j| < 1$ , then  $\tilde{x}^k$  is bounded. If  $|1 + (\Delta t)\lambda_j| > 1$ , then one can find  $\tilde{x}^0$  such that  $|\tilde{x}^k| \rightarrow +\infty$  (exponentially) as  $k \rightarrow +\infty$ . If  $|1 + (\Delta t)\lambda_j| = 1$  and  $N \neq 0$ , then for all  $\tilde{x}^0$  such that  $N\tilde{x}^0 \neq 0$ ,  $N^2\tilde{x}^0 = 0$ , it can be seen that

$$\tilde{x}^k = (1 + (\Delta t)\lambda_j)^k \tilde{x}^0 + (1 + (\Delta t)\lambda_j)^{k-1} k \Delta t N \tilde{x}^0$$

goes to infinity as  $k \rightarrow +\infty$ .

The stability condition  $|1 + (\Delta t)\lambda_j| < 1$  is equivalent to

$$\Delta t < -2 \frac{\Re \lambda_j}{|\lambda_j|^2},$$

and therefore holds for  $\Delta t$  small enough.

For the implicit Euler scheme (4.36), we have

$$\tilde{x}^k = (I - \Delta t(D + N))^{-k} \tilde{x}^0.$$

Note that all the eigenvalues of the matrix  $(I - \Delta t(D + N))^{-1}$  are of modulus strictly smaller than 1. Therefore, the implicit Euler scheme (4.36) is unconditionally stable.

For the Trapezoidal scheme, we have

$$\tilde{x}^k = \left( I - \frac{(\Delta t)}{2}(D + N) \right)^{-k} \left( I + \frac{(\Delta t)}{2}(D + N) \right)^k \tilde{x}^0.$$

Therefore, the stability condition is

$$\left| 1 + \frac{(\Delta t)}{2} \lambda_j \right| < \left| 1 - \frac{(\Delta t)}{2} \lambda_j \right|,$$

which holds for all  $\Delta t > 0$  since  $\Re\lambda_j < 0$ .

□

Note that while the explicit and implicit Euler schemes are of order one, the Trapezoidal scheme is of order two.

**REMARK 4.14.** *If  $\Re\lambda_j = 0$  for some  $j$ , then the explicit Euler scheme may be unstable for any  $\Delta t > 0$ . Consider the second order linear equation*

$$\begin{cases} \frac{d^2x}{dt^2} + x = 0, & t \in [0, +\infty[, \\ x(0) = x_0, \frac{dx}{dt}(0) = x'_0, & x_0, x'_0 \in \mathbb{R}^d. \end{cases} \quad (4.39)$$

We first reduce (4.39) to the first order linear equation

$$\begin{cases} \frac{dX}{dt} = AX, & t \in [0, +\infty[, \\ X(0) = (x_0, x'_0)^\top \in \mathbb{R}^{2d}, \end{cases} \quad (4.40)$$

where  $X = (x, dx/dt)^\top$  and  $A = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$ . The eigenvalues of  $A$  are  $\pm i$ . Consequently, the explicit Euler scheme is unstable since  $|1 \pm i\Delta t| > 1$  for any  $\Delta t > 0$ . However, the implicit Euler scheme is stable since  $|1 \pm i\Delta t|^{-1} < 1$ .

#### 4.4. Runge-Kutta methods

The **Runge-Kutta methods** are by far the most popular and powerful general-purpose numerical methods for integrating ordinary differential equations.

The idea behind the Runge-Kutta methods is to evaluate  $f$  at carefully chosen values of its arguments,  $t$  and  $x$ , in order to create an approximation that is as accurate as a higher-order Taylor expansion of  $x(t + \Delta t)$  without evaluating derivatives of  $f$ . Runge-Kutta schemes are time-stepping schemes that can be derived by matching **multivariable Taylor series expansions** of  $f(t, x)$  with the Taylor series expansion of  $x(t + \Delta t)$ . To find the right values of  $t$  and  $x$  at which to evaluate  $f$ , we need to take a Taylor expansion of  $f$  evaluated at these (unknown) values, and then match the resulting numerical scheme to a Taylor series expansion of  $x(t + \Delta t)$  around  $t$ . Towards this, we state a generalization of Taylor's theorem to functions of two variables.

**THEOREM 4.15.** *Let  $f(t, x) \in \mathcal{C}^{n+1}([0, T] \times \mathbb{R})$ . Let  $(t_0, x_0) \in [0, T] \times \mathbb{R}$ . There exist  $t_0 \leq \tau \leq t$ ,  $x_0 \leq \xi \leq x$ , such that*

$$f(t, x) = P_n(t, x) + R_n(t, x),$$

where  $P_n(t, x)$  is the  $n$ th Taylor polynomial of  $f$  around  $(t_0, x_0)$ ,

$$\begin{aligned} P_n(t, x) &= f(t_0, x_0) + \left[ (t - t_0) \frac{\partial f}{\partial t}(t_0, x_0) + (x - x_0) \frac{\partial f}{\partial x}(t_0, x_0) \right] \\ &+ \left[ \frac{(t - t_0)^2}{2} \frac{\partial^2 f}{\partial t^2}(t_0, x_0) + (t - t_0)(x - x_0) \frac{\partial^2 f}{\partial t \partial x}(t_0, x_0) + \frac{(x - x_0)^2}{2} \frac{\partial^2 f}{\partial x^2}(t_0, x_0) \right] \\ &\dots + \left[ \frac{1}{n!} \sum_{j=0}^n C_j^n (t - t_0)^{n-j} (x - x_0)^j \frac{\partial^n f}{\partial t^{n-j} \partial x^j}(t_0, x_0) \right], \end{aligned}$$

and  $R_n(t, x)$  is the remainder term associated with  $P_n(t, x)$ ,

$$R_n(t, x) = \frac{1}{(n+1)!} \sum_{j=0}^{n+1} C_j^{n+1} (t - t_0)^{n+1-j} (x - x_0)^j \frac{\partial^{n+1} f}{\partial t^{n+1-j} \partial x^j}(\tau, \xi).$$

We now illustrate the proposed approach in order to obtain a second-order accurate method, that is, its local truncation error is  $O((\Delta t)^2)$ . This involves matching

$$x + \Delta t f(t, x) + \frac{(\Delta t)^2}{2} \left[ \frac{\partial f}{\partial t}(t, x) + \frac{\partial f}{\partial x}(t, x) f(t, x) \right] + \frac{(\Delta t)^3}{6} \frac{d^2}{dt^2} [f(\tau, x)]$$

to

$$x + (\Delta t) f(t + \alpha_1, x + \beta_1),$$

where  $\tau \in [t, t + \Delta t]$  and  $\alpha_1$  and  $\beta_1$  are to be found. After simplifying by removing terms that already match, we see that we only need to match

$$f(t, x) + \frac{(\Delta t)}{2} \left[ \frac{\partial f}{\partial t}(t, x) + \frac{\partial f}{\partial x}(t, x) f(t, x) \right] + \frac{(\Delta t)^2}{6} \frac{d^2}{dt^2} [f(t, x)]$$

with  $f(t + \alpha_1, x + \beta_1)$  at least up to terms of the order of  $O(\Delta t)$ , so that the local truncation error will be  $O((\Delta t)^2)$ . Applying the multivariable version of Taylor's theorem to  $f$ , we obtain

$$f(t + \alpha_1, x + \beta_1) = f(t, x) + \alpha_1 \frac{\partial f}{\partial t}(t, x) + \beta_1 \frac{\partial f}{\partial x}(t, x) + \frac{\alpha_1^2}{2} \frac{\partial^2 f}{\partial t^2}(\tau, \xi) + \alpha_1 \beta_1 \frac{\partial^2 f}{\partial t \partial x}(\tau, \xi) + \frac{\beta_1^2}{2} \frac{\partial^2 f}{\partial x^2}(\tau, \xi),$$

where  $t \leq \tau \leq t + \alpha_1$  and  $x \leq \xi \leq x + \beta_1$ . Hence comparing terms yields

$$\alpha_1 = \frac{\Delta t}{2} \quad \text{and} \quad \beta_1 = \frac{\Delta t}{2} f(t, x).$$

The resulting numerical scheme is therefore the **explicit midpoint method** (4.34), which is the simplest example of a Runge-Kutta method of second order. The **improved Euler method** (4.32) is also another often-used Runge-Kutta method.

The most general Runge-Kutta method takes the form

$$x^{k+1} = x^k + \Delta t \sum_{i=1}^m c_i f(t_{i,k}, x_{i,k}), \quad (4.41)$$

where  $m$  stands for the number of terms in the method. Each  $t_{i,k}$  denotes a point in  $[t_k, t_{k+1}]$ . The second argument  $x_{i,k} \approx x(t_{i,k})$  can be viewed as an approximation to the solution at the point  $t_{i,k}$ , and so is computed by a similar but simpler formula of the same type. To construct an  $n$ th order Runge-Kutta method, we need to take at least  $m \geq n$  terms in (4.41).

The best-known Runge-Kutta method is the **fourth-order Runge-Kutta method**, which uses four evaluations of  $f$  during each step. The method proceeds as follows:

$$\begin{cases} \kappa_1 := f(t_k, x^k), \\ \kappa_2 := f(t_k + \frac{\Delta t}{2}, x^k + \frac{\Delta t}{2} \kappa_1), \\ \kappa_3 := f(t_k + \frac{\Delta t}{2}, x^k + \frac{\Delta t}{2} \kappa_2), \\ \kappa_4 := f(t_{k+1}, x^k + \Delta t \kappa_3), \\ x^{k+1} = x^k + \frac{(\Delta t)}{6} (\kappa_1 + 2\kappa_2 + 2\kappa_3 + \kappa_4). \end{cases} \quad (4.42)$$

In (4.42), the values of  $f$  at the midpoint in time are given four times as much weight as values at the endpoints  $t_k$  and  $t_{k+1}$ , which is similar to Simpson's rule (4.27) from numerical integration.

**4.4.1. Construction of Runge-Kutta methods.** In this subsection we first construct Runge-Kutta methods by generalizing **collocation methods**. Then we discuss their consistency, stability, and order.

4.4.1.1. *Collocation methods.* Let  $\mathcal{P}_m$  denote the space of real polynomials of degree  $\leq m$ . Given a set of  $m$  **distinct** quadrature points  $c_1 < c_2 < \dots < c_m$  in  $\mathbb{R}$ , and corresponding data  $g_1, \dots, g_m$ , there exists a unique polynomial, called the **interpolating polynomial**,  $P(t) \in \mathcal{P}_{m-1}$  satisfying  $P(c_i) = g_i, i = 1, \dots, m$ .

Define the  $i$ th **Lagrange interpolating polynomial**  $l_i(t), i = 1, \dots, m$ , for the set of quadrature points  $\{c_j\}$  by

$$l_i(t) := \prod_{j \neq i, j=1}^m \frac{t - c_j}{c_i - c_j}.$$

The set of Lagrange interpolating polynomials form a basis of  $\mathcal{P}_{m-1}$  and the interpolating polynomial  $P$  corresponding to the data  $\{g_j\}$  is given by

$$P(t) := \sum_{i=1}^m g_i l_i(t). \quad (4.43)$$

Consider first a smooth function  $g$  on  $[0, 1]$ . We can approximate the integral of  $g$  on  $[0, 1]$  by exactly integrating the Lagrange interpolating polynomial of order  $m - 1$  based on  $m$  **quadrature points**  $0 \leq c_1 < c_2 < \dots < c_m \leq 1$ . The data are the values of  $g$  at the quadrature points  $g_i = g(c_i), i = 1, \dots, m$ .

Define the weights

$$b_i = \int_0^1 l_i(s) ds. \quad (4.44)$$

The **quadrature formula** is

$$\int_0^1 g(s) ds \approx \int_0^1 P(s) ds = \sum_{i=1}^m b_i g(c_i),$$

where  $P$  is defined by (4.43).

Now let  $f$  be a smooth function on  $[0, T]$  and let  $t_k = k\Delta t$  for  $k = 0, \dots, K = T/(\Delta t)$ , be the discretization points in  $[0, T]$ . The integral  $\int_{t_k}^{t_{k+1}} f(s) ds$  can be approximated by

$$\int_{t_k}^{t_{k+1}} f(s) ds = (\Delta t) \int_0^1 f(t_k + \Delta t\tau) d\tau \approx (\Delta t) \sum_{i=1}^m b_i f(t_k + (\Delta t)c_i). \quad (4.45)$$

Next let  $x$  be a polynomial of degree  $m$  satisfying

$$\begin{cases} x(0) = x_0, \\ \frac{dx}{dt}(c_i \Delta t) = F_i, \end{cases} \quad (4.46)$$

where  $F_i \in \mathbb{R}, i = 1, \dots, m$ .

From the Lagrange interpolation formula (4.43), it follows that for  $t$  in the first time-step interval  $[0, \Delta t]$ ,

$$\frac{dx}{dt}(t) = \sum_{i=1}^m F_i l_i\left(\frac{t}{\Delta t}\right). \quad (4.47)$$

Integrating (4.47) over the intervals  $[0, c_i \Delta t]$  gives

$$x(c_i \Delta t) = x_0 + (\Delta t) \sum_{j=1}^m F_j \int_0^{c_i} l_j(s) ds = x_0 + (\Delta t) \sum_{j=1}^m a_{ij} F_j, \quad i = 1, \dots, m, \quad (4.48)$$

where

$$a_{ij} := \int_0^{c_i} l_j(s) ds. \quad (4.49)$$

Integrating (4.47) over  $[0, \Delta t]$  yields

$$x(\Delta t) = x_0 + (\Delta t) \sum_{i=1}^m F_i \int_0^1 l_i(s) ds = x_0 + (\Delta t) \sum_{i=1}^m b_i F_i, \quad (4.50)$$

where  $b_i$  is defined by (4.44).

Writing  $dx/dt = f(x(t))$ , we obtain from (4.48) and (4.50) on the first time step interval  $[0, \Delta t]$

$$\begin{cases} F_i = f(x_0 + (\Delta t) \sum_{j=1}^m a_{ij} F_j), & i = 1, \dots, m, \\ x(\Delta t) = x_0 + (\Delta t) \sum_{i=1}^m b_i F_i. \end{cases} \quad (4.51)$$

Similarly, we have on  $[t_k, t_{k+1}]$

$$\begin{cases} F_{i,k} = f(x(t_k) + (\Delta t) \sum_{j=1}^m a_{ij} F_{j,k}), & i = 1, \dots, m, \\ x(t_{k+1}) = x(t_k) + (\Delta t) \sum_{i=1}^m b_i F_{i,k}. \end{cases} \quad (4.52)$$

In the **collocation method** (4.52), one first solves the coupled nonlinear system to obtain  $F_{i,k}$ ,  $i = 1, \dots, m$ , and then computes  $x(t_{k+1})$  from  $x(t_k)$ .

REMARK 4.16. *Since*

$$t^{l-1} = \sum_{i=1}^m c_i^{l-1} l_i(t), \quad t \in [0, 1], l = 1, \dots, m,$$

*we have*

$$\sum_{i=1}^m b_i c_i^{l-1} = \frac{1}{l}, \quad l = 1, \dots, m,$$

*and*

$$\sum_{j=1}^m a_{ij} c_j^{l-1} = \frac{c_i^l}{l}, \quad i, l = 1, \dots, m.$$

**4.4.2. Runge-Kutta methods as generalized collocation methods.** In (4.52), the coefficients  $b_i$  and  $a_{ij}$  are defined by certain integrals of the Lagrange interpolating polynomials associated with a chosen set of quadrature nodes  $c_i$ ,  $i = 1, \dots, m$ .

A natural generalization of collocation methods is obtained by allowing the coefficients  $c_i$ ,  $b_i$ , and  $a_{ij}$  to take arbitrary values, not necessary related to quadrature formulas. In fact, we no longer assume the  $c_i$  to be distinct. However, we should assume that

$$c_i = \sum_{j=1}^m a_{ij}, \quad i = 1, \dots, m. \quad (4.53)$$

The result is the class of Runge-Kutta methods for solving (4.1), which can be written as

$$\begin{cases} F_{i,k} = f(t_{i,k}, x^k + (\Delta t) \sum_{j=1}^m a_{ij} F_{j,k}), \\ x^{k+1} = x^k + (\Delta t) \sum_{i=1}^m b_i F_{i,k}, \end{cases} \quad (4.54)$$

where  $t_{i,k} = t_k + c_i \Delta t$ , or equivalently,

$$\begin{cases} x_{i,k} = x^k + (\Delta t) \sum_{j=1}^m a_{ij} f(t_{j,k}, x_{j,k}), \\ x^{k+1} = x^k + (\Delta t) \sum_{i=1}^m b_i f(t_{i,k}, x_{i,k}). \end{cases} \quad (4.55)$$

Let

$$\kappa_j := f(t + c_j \Delta t, x_j), \quad (4.56)$$

and define the function  $\Phi$  by

$$\begin{cases} x_i = x + (\Delta t) \sum_{j=1}^m a_{ij} \kappa_j, \\ \Phi(t, x, \Delta t) = \sum_{i=1}^m b_i f(t + c_i \Delta t, x_i). \end{cases} \quad (4.57)$$

One can see that the scheme (4.55) is a one step method. Moreover, if  $a_{ij} = 0$  for  $j \geq i$ , then (4.55) is explicit.

It is also easy to see that with definition (4.55), explicit Euler's method and Trapezoidal scheme are Runge-Kutta methods. For example, explicit Euler's method (4.17) can be put into the form (4.55) with  $m = 1, b_1 = 1, a_{11} = 0$ . The Trapezoidal scheme (4.29) has  $m = 2, b_1 = b_2 = 1/2, a_{11} = a_{12} = 0, a_{21} = a_{22} = 1/2$ . Finally, for the fourth-order Runge-Kutta method (4.42), we have  $m = 4, c_1 = 0, c_2 = c_3 = 1/2, c_4 = 1, b_1 = 1/6, b_2 = b_3 = 1/3, b_4 = 1/6, a_{21} = a_{32} = 1/2, a_{43} = 1$ , and all the other  $a_{ij}$  entries are zero.

#### 4.4.3. Consistency, stability, convergence, and order of Runge-Kutta methods.

From (4.57), the Runge-Kutta scheme is consistent if and only if

$$\sum_{j=1}^m b_j = 1. \quad (4.58)$$

Let  $|A|$  be the matrix defined by  $(|a_{ij}|)_{i,j=1}^m$ . Let the **spectral radius**  $\rho(|A|)$  of the matrix  $|A|$  be defined by

$$\rho(|A|) := \max\{|\lambda_j|, \lambda_j \text{ is an eigenvalue of } |A|\}. \quad (4.59)$$

The following stability result holds.

**THEOREM 4.17.** *Let  $C_f$  be the Lipschitz constant for  $f$ . Suppose that*

$$(\Delta t) C_f \rho(|A|) < 1. \quad (4.60)$$

*Then the Runge-Kutta method (4.55) for solving (4.1) is stable.*

**PROOF.** Let  $\Phi$  be defined by (4.57). We have

$$\Phi(t, x, \Delta t) - \Phi(t, y, \Delta t) = \sum_{i=1}^m b_i \left[ f(t + c_i \Delta t, x_i) - f(t + c_i \Delta t, y_i) \right], \quad (4.61)$$

where

$$x_i = x + (\Delta t) \sum_{j=1}^m a_{ij} f(t + c_j \Delta t, x_j), \quad (4.62)$$

and

$$y_i = y + (\Delta t) \sum_{j=1}^m a_{ij} f(t + c_j \Delta t, y_j). \quad (4.63)$$

Subtracting (4.63) from (4.62) yields

$$x_i - y_i = x - y + (\Delta t) \sum_{j=1}^m a_{ij} \left[ f(t + c_j \Delta t, x_j) - f(t + c_j \Delta t, y_j) \right]. \quad (4.64)$$

Therefore, for  $i = 1, \dots, m$ ,

$$|x_i - y_i| \leq |x - y| + (\Delta t) C_f \sum_{j=1}^m |a_{ij}| |x_j - y_j|, \quad (4.65)$$

where  $C_f$  is the Lipschitz constant for  $f$ . Let the vectors  $X$  and  $Y$  be defined by

$$X = \begin{bmatrix} |x_1 - y_1| \\ \vdots \\ |x_m - y_m| \end{bmatrix} \quad \text{and} \quad Y = \begin{bmatrix} |x - y| \\ \vdots \\ |x - y| \end{bmatrix}.$$

From (4.65), it follows that

$$X \leq Y + (\Delta t) C_f |A| X, \quad (4.66)$$

and therefore,

$$X \leq (I - (\Delta t) C_f |A|)^{-1} Y, \quad (4.67)$$

provided that condition (4.60) holds. Finally, combining (4.61) and (4.67) yields the stability of the Runge-Kutta scheme (4.55).  $\square$

By the Dahlquist-Lax equivalence theorem (Theorem 4.4), it follows that the Runge-Kutta scheme (4.55) is convergent provided that (4.58) and (4.60) hold.

In order to establish the order of the Runge-Kutta scheme (4.55), we compute the order as  $\Delta t \rightarrow 0$  of the truncation error

$$T_k(\Delta t) = \frac{x(t_{k+1}) - x(t_k)}{\Delta t} - \Phi(t_k, x(t_k), \Delta t),$$

where  $\Phi$  is defined by (4.57). We write

$$T_k(\Delta t) = \frac{x(t_{k+1}) - x(t_k)}{\Delta t} - \sum_{i=1}^m b_i f(t_k + c_i \Delta t, x(t_k)) + \Delta t \sum_{j=1}^m a_{ij} \kappa_j.$$

Suppose that  $f$  is smooth enough. We have

$$f(t_k + c_i \Delta t, x(t_k) + \Delta t \sum_{j=1}^m a_{ij} \kappa_j) = f(t_k, x(t_k)) + \Delta t \left[ c_i \frac{\partial f}{\partial t}(t_k, x(t_k)) + \left( \sum_{j=1}^m a_{ij} \kappa_j \right) \frac{\partial f}{\partial x}(t_k, x(t_k)) \right] + O((\Delta t)^2).$$

Suppose that (4.53) holds. Then, from

$$\sum_{j=1}^m a_{ij} \kappa_j = \left( \sum_{j=1}^m a_{ij} \right) f(t_k, x(t_k)) + O(\Delta t) = c_i f(t_k, x(t_k)) + O(\Delta t),$$

it follows that

$$f(t_k + c_i \Delta t, x(t_k) + \Delta t \sum_{j=1}^m a_{ij} \kappa_j) = f(t_k, x(t_k)) + \Delta t c_i \left[ \frac{\partial f}{\partial t}(t_k, x(t_k)) + \frac{\partial f}{\partial x}(t_k, x(t_k)) f(t_k, x(t_k)) \right] + O((\Delta t)^2).$$

Therefore, we obtain the following theorem.

**THEOREM 4.18.** *Assume that  $f$  is smooth enough. Then the Runge-Kutta scheme (4.55) for solving (4.1) is of order 2 provided that the conditions (4.58) and*

$$\sum_{i=1}^m b_i c_i = \frac{1}{2} \quad (4.68)$$

*hold.*

One can prove by higher-order Taylor expansions that the following results hold.



**THEOREM 4.19.** *Assume that  $f$  is smooth enough. Then the Runge-Kutta scheme (4.55) for solving (4.1) is of order 3 provided that the conditions (4.58), (4.68), and*

$$\sum_{i=1}^m b_i c_i^2 = \frac{1}{3}, \quad \sum_{i=1}^m \sum_{j=1}^m b_i a_{ij} c_j = \frac{1}{6} \quad (4.69)$$

*hold. It is of order 4 provided that (4.58), (4.68), (4.69), and*

$$\sum_{i=1}^m b_i c_i^3 = \frac{1}{4}, \quad \sum_{i=1}^m \sum_{j=1}^m b_i c_i a_{ij} c_j = \frac{1}{8}, \quad \sum_{i=1}^m \sum_{j=1}^m b_i a_{ij} c_j^2 = \frac{1}{12}, \quad \sum_{i=1}^m \sum_{j=1}^m \sum_{l=1}^m b_i a_{ij} a_{jl} c_l = \frac{1}{24} \quad (4.70)$$

*hold.*

The Runge-Kutta scheme (4.42) satisfies the four conditions (4.58), (4.68), (4.69), and (4.70). Hence, (4.42) is of order 4.

#### 4.5. Multi-step methods

While Runge-Kutta methods present an improvement over Euler's methods in terms of accuracy, this is achieved by investing additional computational effort. For example, the fourth-order method (4.42) involves four function evaluations per step. For comparison, by considering three consecutive points  $t_{k-1}, t_k, t_{k+1}$ , integrating the differential equation between  $t_{k-1}$  and  $t_{k+1}$ , and applying **Simpson's rule** to approximate the resulting integral yields

$$\begin{aligned} x(t_{k+1}) &= x(t_{k-1}) + \int_{t_{k-1}}^{t_{k+1}} f(s, x(s)) ds \\ &\approx x(t_{k-1}) + \frac{(\Delta t)}{3} \left[ f(t_{k-1}, x(t_{k-1})) + 4f(t_k, x(t_k)) + f(t_{k+1}, x(t_{k+1})) \right], \end{aligned}$$

which leads to the method

$$x^{k+1} = x^{k-1} + \frac{(\Delta t)}{3} \left[ f(t_{k-1}, x^{k-1}) + 4f(t_k, x^k) + f(t_{k+1}, x^{k+1}) \right]. \quad (4.71)$$

In contrast with the one-step methods considered in the previous sections where only a single value of  $x^k$  was required to compute the next approximation  $x^{k+1}$ , in (4.71) we need two preceding values,  $x^k$  and  $x^{k-1}$  in order to calculate  $x^{k+1}$ , and therefore (4.71) is a **two-step method**.

A general  $n$ -step method is of the form

$$\sum_{j=0}^n \alpha_j x^{k+j} = (\Delta t) \sum_{j=0}^n \beta_j f(t_{k+j}, x^{k+j}), \quad (4.72)$$

where the coefficients  $\alpha_j$  and  $\beta_j$  are real constants and  $\alpha_n \neq 0$ .

If  $\beta_n = 0$ , then  $x^{k+n}$  is obtained explicitly from previous values of  $x^j$  and  $f(t_j, x^j)$ , and the  $n$ -step method is **explicit**. Otherwise, the  $n$ -step method is **implicit**.

In multi-step methods we need a starting procedure which provides approximations to the exact solution at the points  $t_1, \dots, t_{n-1}$ . One possibility for obtaining these missing starting values is the use of any one-step method, e.g., a Runge-Kutta method.

The following are classical examples of multi-step methods:

**EXAMPLE 4.20.** (i) *The two-step **Adams-Bashforth method***

$$x^{k+2} = x^{k+1} + \frac{(\Delta t)}{2} \left[ 3f(t_{k+1}, x^{k+1}) - f(t_k, x^k) \right] \quad (4.73)$$

*is an example of an **explicit two-step method**;*

(ii) *The three-step Adams-Bashforth method*

$$x^{k+3} = x^{k+2} + \frac{(\Delta t)}{12} \left[ 23f(t_{k+2}, x^{k+2}) - 16f(t_{k+1}, x^{k+1}) + 5f(t_k, x^k) \right] \quad (4.74)$$

*is an example of an explicit three-step method;*

(iii) *The four-step Adams-Bashforth method*

$$x^{k+4} = x^{k+3} + \frac{(\Delta t)}{24} \left[ 55f(t_{k+3}, x^{k+3}) - 59f(t_{k+2}, x^{k+2}) + 37f(t_{k+1}, x^{k+1}) - 9f(t_k, x^k) \right] \quad (4.75)$$

*is an example of an explicit four-step method;*

(iv) *The two-step Adams-Moulton method*

$$x^{k+2} = x^{k+1} + \frac{(\Delta t)}{12} \left[ 5f(t_{k+2}, x^{k+2}) + 8f(t_{k+1}, x^{k+1}) - f(t_k, x^k) \right] \quad (4.76)$$

*is an example of an implicit two-step method;*

(v) *The three-step Adams-Moulton method*

$$x^{k+3} = x^{k+2} + \frac{(\Delta t)}{24} \left[ 9f(t_{k+3}, x^{k+3}) + 19f(t_{k+2}, x^{k+2}) + 5f(t_{k+1}, x^{k+1}) - 9f(t_k, x^k) \right] \quad (4.77)$$

*is an example of an implicit three-step method.*

The construction of general classes of linear multi-step methods is discussed in the next subsection.

**4.5.1. Construction of linear multi-step methods.** Suppose that  $x^k, k \in \mathbb{N}$ , is a sequence of real numbers. We introduce the **shift operator**  $E$ , the **forward difference operator**  $\Delta_+$  and the **backward difference operator**  $\Delta_-$  by

$$E : x^k \mapsto x^{k+1}, \quad \Delta_+ : x^k \mapsto x^{k+1} - x^k, \quad \Delta_- : x^k \mapsto x^k - x^{k-1}.$$

Since  $\Delta_+ = E - I$  and  $\Delta_- = I - E^{-1}$ , it follows that, for any  $n \in \mathbb{N}$ ,

$$(E - I)^n = \sum_{j=0}^n (-1)^j C_j^n E^{n-j},$$

and

$$(I - E^{-1})^n = \sum_{j=0}^n (-1)^j C_j^n E^{-j}.$$

Therefore,

$$\Delta_+^n x^k = \sum_{j=0}^n (-1)^j C_j^n x^{k+n-j}$$

and

$$\Delta_-^n x^k = \sum_{j=0}^n (-1)^j C_j^n x^{k-j}.$$

Now let  $y(t) \in C^\infty(\mathbb{R})$  and let  $t_k = k\Delta t, \Delta t > 0$ . By applying the Taylor series we find that, for any  $s \in \mathbb{N}$ ,

$$E^s y(t_k) = y(t_k + s\Delta t) = \left( \sum_{l=0}^{+\infty} \frac{1}{l!} (s\Delta t \frac{\partial}{\partial t})^l y \right) (t_k) = (e^{s(\Delta t) \frac{\partial}{\partial t}} y) (t_k),$$

and hence

$$E^s = e^{s(\Delta t) \frac{\partial}{\partial t}}.$$

Thus, formally,

$$(\Delta t) \frac{\partial}{\partial t} = \ln E = -\ln(I - \Delta_-) = \Delta_- + \frac{1}{2}\Delta_-^2 + \frac{1}{3}\Delta_-^3 + \dots \quad (4.78)$$

Therefore, if  $x(t)$  is the solution of (4.1), then by using (4.78) we find that

$$(\Delta t)f(t_k, x(t_k)) = \left( \Delta_- + \frac{1}{2}\Delta_-^2 + \frac{1}{3}\Delta_-^3 + \dots \right) x(t_k). \quad (4.79)$$

The successive truncation of the infinite series on the right-hand side of (4.79) yields

$$\begin{aligned} x^k - x^{k-1} &= (\Delta t)f(t_k, x^k), \\ \frac{3}{2}x^k - 2x^{k-1} + \frac{1}{2}x^{k-2} &= (\Delta t)f(t_k, x^k), \\ \frac{11}{6}x^k - 3x^{k-1} + \frac{3}{2}x^{k-2} - \frac{1}{3}x^{k-3} &= (\Delta t)f(t_k, x^k), \end{aligned} \quad (4.80)$$

and so on. This gives rise to a class of implicit multi-step methods called **backward differentiation formulas**.

Similarly,

$$E^{-1}((\Delta t)\frac{\partial}{\partial t}) = (\Delta t)\frac{\partial}{\partial t}E^{-1} = -(I - \Delta_-)\ln(I - \Delta_-),$$

and hence,

$$((\Delta t)\frac{\partial}{\partial t}) = -E(I - \Delta_-)\ln(I - \Delta_-) = -(I - \Delta_-)\ln(I - \Delta_-)E. \quad (4.81)$$

Therefore, if  $x(t)$  is the solution of (4.1), then we find that

$$(\Delta t)f(t_k, x(t_k)) = \left( \Delta_- - \frac{1}{2}\Delta_-^2 - \frac{1}{6}\Delta_-^3 + \dots \right) x(t_{k+1}). \quad (4.82)$$

The successive truncation of the infinite series on the right-hand side of (4.82) yields the following explicit numerical schemes:

$$\begin{aligned} x^{k+1} - x^k &= (\Delta t)f(t_k, x^k), \\ \frac{1}{2}x^{k+1} - \frac{1}{2}x^{k-1} &= (\Delta t)f(t_k, x^k), \\ \frac{1}{3}x^{k+1} + \frac{1}{2}x^k - x^{k-1} + \frac{1}{6}x^{k-2} &= (\Delta t)f(t_k, x^k), \\ &\vdots \end{aligned} \quad (4.83)$$

The first of these numerical scheme is the explicit Euler method, while the second is the explicit mid-point method.

In order to construct further classes of multi-step methods, we define, for  $y \in \mathcal{C}^\infty$ ,

$$D^{-1}y(t_k) = y(t_0) + \int_{t_0}^{t_k} y(s) ds,$$

and observe that

$$(E - I)D^{-1}y(t_k) = \int_{t_k}^{t_{k+1}} y(s) ds.$$

Now, from

$$(E - I)D^{-1} = \Delta_+ D^{-1} = E\Delta_- D^{-1} = (\Delta t)E\Delta_- ((\Delta t)D)^{-1},$$

it follows that

$$(E - I)D^{-1} = -(\Delta t)E\Delta_- (\ln(I - \Delta_-))^{-1}. \quad (4.84)$$

Furthermore,

$$(E - I)D^{-1} = E\Delta_- D^{-1} = \Delta_- ED^{-1} = \Delta_- (DE^{-1})^{-1} = (\Delta t)\Delta_- ((\Delta t)DE^{-1})^{-1}.$$

Thus,

$$(E - I)D^{-1} = -(\Delta t)\Delta_- \left( (I - \Delta_-)\ln(I - \Delta_-) \right)^{-1}. \quad (4.85)$$

By using (4.84) and (4.85), we deduce from

$$x(t_{k+1}) - x(t_k) = \int_{t_k}^{t_{k+1}} f(s, x(s)) ds = (E - I)D^{-1}f(t_k, x(t_k)),$$

that

$$x(t_{k+1}) - x(t_k) = \begin{cases} -(\Delta t)\Delta_-((I - \Delta_-)\ln(I - \Delta_-))^{-1}f(t_k, x(t_k)) \\ -(\Delta t)E\Delta_-(\ln(I - \Delta_-))^{-1}f(t_k, x(t_k)), \end{cases} \quad (4.86)$$

where  $x(t)$  is the solution of (4.1).

On expanding  $\ln(I - \Delta_-)$  into a Taylor series on the right-hand side of (4.86) we find that

$$x(t_{k+1}) - x(t_k) = (\Delta t) \left[ I + \frac{1}{2}\Delta_- + \frac{5}{12}\Delta_-^2 + \frac{3}{8}\Delta_-^3 + \dots \right] f(t_k, x(t_k)), \quad (4.87)$$

and

$$x(t_{k+1}) - x(t_k) = (\Delta t) \left[ I - \frac{1}{2}\Delta_- - \frac{1}{12}\Delta_-^2 - \frac{1}{24}\Delta_-^3 + \dots \right] f(t_{k+1}, x(t_{k+1})). \quad (4.88)$$

The successive truncation of (4.87) yields the family (4.75) of (explicit) Adams-Bashforth methods, while similar successive truncation of (4.88) gives rise to the family (4.77) of (implicit) Adams-Moulton methods.

**4.5.2. Consistency, stability, and convergence.** In this subsection, we introduce the concepts of consistency, stability, and convergence for analyzing linear multi-step methods.

**DEFINITION 4.21 (Consistency).** *The  $n$ -step method (4.72) is **consistent** with (4.1) if the **truncation error** defined by*

$$T_k(\Delta t) = \frac{\sum_{j=0}^n [\alpha_j x(t_{k+j}) - (\Delta t)\beta_j \frac{dx}{dt}(t_{k+j})]}{(\Delta t)}$$

is such that for any  $\epsilon > 0$  there exists  $h_0$  for which

$$|T_k(\Delta t)| \leq \epsilon \quad \text{for } 0 < \Delta t \leq h_0 \quad (4.89)$$

and any  $(n+1)$  points  $((t_j, x(t_j)), \dots, (t_{j+n}, x(t_{j+n})))$  on any solution  $x(t)$ .

**THEOREM 4.22.** *The  $n$ -step method (4.72) is consistent if and only if the two conditions*

$$\sum_{j=0}^n \alpha_j = 0 \quad \text{and} \quad \sum_{j=0}^n j\alpha_j = \sum_{j=0}^n \beta_j, \quad (4.90)$$

hold. Furthermore, it is of order  $p$  if and only if

$$\frac{1}{l} \sum_{j=0}^n j^l \alpha_j = \sum_{j=0}^n j^{l-1} \beta_j, \quad \text{for all } l = 1, \dots, p, \quad (4.91)$$

and

$$\frac{1}{p+1} \sum_{j=0}^n j^{p+1} \alpha_j \neq \sum_{j=0}^n j^p \beta_j. \quad (4.92)$$

**PROOF.** Assume that  $f \in C^\infty$ . Using the Taylor expansions for both  $x$  and  $dx/dt$ ,

$$x(t_{k+j}) = \sum_{l=0}^{+\infty} \frac{1}{l!} (j\Delta t)^l x^{(l)}(t_k), \quad \frac{dx}{dt}(t_{k+j}) = \sum_{l=0}^{+\infty} \frac{1}{l!} (j\Delta t)^l x^{(l+1)}(t_k),$$

we obtain

$$\sum_{j=0}^n [\alpha_j x(t_{k+j}) - (\Delta t)\beta_j \frac{dx}{dt}(t_{k+j})] = \sum_{j=0}^n \left[ \alpha_j \sum_{l=0}^{+\infty} \frac{1}{l!} (j\Delta t)^l x^{(l)}(t_k) - (\Delta t)\beta_j \sum_{l=0}^{+\infty} \frac{1}{l!} (j\Delta t)^l x^{(l+1)}(t_k) \right]$$

$$= \left( \sum_{j=0}^n \alpha_j \right) x(t_k) + \left( \sum_{j=0}^n [j\alpha_j - \beta_j] \right) \Delta t \frac{dx}{dt}(t_k) + \sum_{l=2}^{+\infty} \left( \sum_{j=0}^n \left[ \frac{j^l}{l!} \alpha_j - \frac{j^{l-1}}{(l-1)!} \beta_j \right] \right) (\Delta t)^l x^{(l)}(t_k),$$

which yields the result.  $\square$

In view of Theorem 4.22, one can easily check that (4.71) is of order 4, (4.73) is of order 2, (4.74) is of order 3, (4.75) is of order 4, (4.76) is of order 3, and (4.77) is of order 4.

**DEFINITION 4.23 (Stability).** *The  $n$ -step method (4.72) is **stable** if there exists a constant  $C$  such that, for any two sequences  $(x^k)$  and  $(\tilde{x}^k)$  which have been generated by the same formulas but different initial data  $x^0, x^1, \dots, x^{n-1}$  and  $\tilde{x}^0, \tilde{x}^1, \dots, \tilde{x}^{n-1}$ , respectively, we have*

$$|x^k - \tilde{x}^k| \leq C \max\{|x^0 - \tilde{x}^0|, |x^1 - \tilde{x}^1|, \dots, |x^{n-1} - \tilde{x}^{n-1}|\} \quad (4.93)$$

as  $\Delta t \rightarrow 0$  for all  $k \geq n$ .

**THEOREM 4.24 (Convergence).** *Suppose that the  $n$ -step method (4.72) is consistent with (4.1). The stability condition (4.93) is necessary and sufficient for the convergence. Moreover, if  $x \in C^{p+1}$  and the truncation error is  $O((\Delta t)^p)$ , then the global error  $e_k = x(t_k) - x^k$  is also  $O((\Delta t)^p)$ .*

**PROOF.** One way to prove Theorem 4.24 is to rewrite (4.72) as a one-step method in a higher dimensional space. For this, let  $\phi(t_k, x^k, \dots, x^{k+n-1}, \Delta t)$  be defined implicitly by

$$\phi = \sum_{j=0}^{n-1} \beta'_j f(t_{k+j}, x^{k+j}) + \beta'_n f(t_{k+n}, (\Delta t)\phi - \sum_{j=0}^{n-1} \alpha'_j x^{k+j}),$$

where  $\alpha'_j = \alpha_j/\alpha_n$  and  $\beta'_j = \beta_j/\alpha_n$ . Then, (4.72) can be written as

$$x^{k+n} = - \sum_{j=0}^{n-1} \alpha'_j x^{k+j} + (\Delta t)\phi.$$

Introduce the  $n$ -dimensional vectors

$$X^k = (x^{k+n-1}, \dots, x^k)^\top, \quad \Phi(t_k, X^k, \Delta t) = \phi(t_k, x^k, \dots, x^{k+n-1}, \Delta t)(1, 0, \dots, 0)^\top,$$

and the  $n \times n$  matrix

$$A = \begin{pmatrix} -\alpha'_{n-1} & -\alpha'_{n-2} & \cdots & \cdot & -\alpha'_0 \\ 1 & 0 & \cdots & \cdot & 0 \\ & 1 & \cdots & \vdots & 0 \\ & & \ddots & \vdots & \vdots \\ & & & 1 & 0 \end{pmatrix}.$$

The  $n$ -step method (4.72) can be rewritten as

$$X^{k+1} = AX^k + \Delta t \Phi(t_k, X^k, \Delta t),$$

and the concepts of consistency and stability can be expressed in this new notation. In fact, let  $x(t)$  be the exact solution and denote by  $X(t_k) = (x(t_{k+n-1}), \dots, x(t_k))^\top$ . The consistency condition (4.89) implies that

$$|X(t_{k+1}) - AX(t_k) - \Delta t \Phi(t_k, X(t_k), \Delta t)| \rightarrow 0 \text{ as } \Delta t \rightarrow 0.$$

Moreover, if (4.72) is of order  $p$  then

$$|X(t_{k+1}) - AX(t_k) - \Delta t \Phi(t_k, X(t_k), \Delta t)| = O((\Delta t)^p)$$

as  $\Delta t \rightarrow 0$ . Furthermore, the stability condition (4.93) implies that there exists a matrix norm such that  $\|A\| \leq 1$ . The rest of the proof is similar to the proof of Theorem 4.4.  $\square$

#### 4.6. Stiff equations and systems

Let  $\epsilon > 0$  be a small parameter. Consider the initial value problem

$$\begin{cases} \frac{dx(t)}{dt} = -\frac{1}{\epsilon}x(t), & t \in [0, T], \\ x(0) = 1, \end{cases} \quad (4.94)$$

which has an exponential solution  $x(t) = e^{-t/\epsilon}$ . The explicit Euler method with step size  $\Delta t$  relies on the iterative scheme

$$x^{k+1} = \left(1 - \frac{\Delta t}{\epsilon}\right)x^k, \quad x^0 = 1, \quad (4.95)$$

with solution

$$x^k = \left(1 - \frac{\Delta t}{\epsilon}\right)^k.$$

Since  $\epsilon > 0$  the exact solution is exponentially decaying and positive. But now, if  $1 - \frac{\Delta t}{\epsilon} < -1$ , then the iterates (4.95) grow exponentially fast in magnitude, with alternating signs. In this case, the numerical solution is nowhere close to the true solution. If  $-1 < 1 - \frac{\Delta t}{\epsilon} < 0$ , then the numerical solution decays in magnitude, but continues to alternate between positive and negative values. Thus, to correctly model the qualitative features of the solution and obtain a numerically accurate solution, we need to choose the step size  $\Delta t$  so as to ensure that  $1 - \frac{\Delta t}{\epsilon} > 0$ , and hence  $\Delta t < \epsilon$ .

Equation (4.94) is the simplest example of what is known as a **stiff differential equation**. In general, an equation or system is stiff if it has one or more very rapidly decaying solutions. In the case of the autonomous constant coefficient linear system (3.24), stiffness occurs whenever the coefficient matrix  $A$  has an eigenvalue  $\lambda_{j_0}$  with large negative real part:  $\Re \lambda_{j_0} \ll 0$ , resulting in a very rapidly decaying eigensolution. It only takes one such eigensolution to render the equation stiff, and ruin the numerical computation of even well behaved solutions. Even though the component of the actual solution corresponding to  $\lambda_{j_0}$  is almost irrelevant, as it becomes almost instantaneously tiny, its presence continues to render the numerical solution to the system very difficult. Stiff equations require more sophisticated numerical schemes to integrate.

Most of the numerical methods derived above also suffer from instability due to stiffness of (4.94) for sufficiently small positive  $\epsilon$ . Interestingly, stability of (4.94) suffices to characterize acceptable step sizes  $\Delta t$ , depending on the size of  $-1/\epsilon$ , which, in the case of linear systems, is the eigenvalue. Applying the Trapezoidal scheme (4.29) to (4.94) leads to

$$x^{k+1} = x^k - \frac{\Delta t}{2\epsilon}(x^k + x^{k+1}), \quad x^0 = 1, \quad (4.96)$$

which we solve for

$$x^{k+1} = \frac{1 - \frac{\Delta t}{2\epsilon}}{1 + \frac{\Delta t}{2\epsilon}}x^k, \quad x^0 = 1. \quad (4.97)$$

Thus, the behavior of the numerical solution is entirely determined by the size of the coefficient

$$\mu := \frac{1 - \frac{\Delta t}{2\epsilon}}{1 + \frac{\Delta t}{2\epsilon}}.$$

Since  $|\mu| < 1$  for all  $\epsilon > 0$ , the Trapezoidal scheme (4.96) is not affected by stiffness.

In the system of equations (1.5), the parameter satisfies  $0 < a \ll 1$ . This makes (1.5) a stiff system of ODEs.

#### 4.7. Perturbation theories for differential equations

**4.7.1. Regular perturbation theory.** Let  $\epsilon > 0$  be a small parameter and consider the differential equation

$$\begin{cases} \frac{dx}{dt} = f(t, x, \epsilon), & t \in [0, T], \\ x(0) = x_0, & x_0 \in \mathbb{R}. \end{cases} \quad (4.98)$$

If we suppose that  $f \in \mathcal{C}^1$ , then (4.98) is a **regular perturbation problem**. The solution  $x(t, \epsilon)$  is in  $\mathcal{C}^1$  and has the following Taylor expansion:

$$x(t, \epsilon) = x^{(0)}(t) + \epsilon x^{(1)}(t) + o(\epsilon) \quad (4.99)$$

with respect to  $\epsilon$  in a neighborhood of 0.

Clearly, the unperturbed term  $x^{(0)}$  is given as the solution of the unperturbed equation

$$\begin{cases} \frac{dx^{(0)}}{dt} = f_0(t, x^{(0)}), & t \in [0, T], \\ x^{(0)}(0) = x_0, & x_0 \in \mathbb{R}, \end{cases} \quad (4.100)$$

where  $f_0(t, x) := f(t, x, 0)$ . Moreover, the first-order correction term  $x^{(1)}$ , which is the derivative of  $x(t, \epsilon)$  with respect to  $\epsilon$  at 0,

$$x^{(1)}(t) = \frac{\partial x}{\partial \epsilon}(t, 0),$$

solves the equation

$$\begin{cases} \frac{dx^{(1)}}{dt} = \frac{\partial f}{\partial x}(t, x^{(0)}, 0)x^{(1)} + \frac{\partial f}{\partial \epsilon}(t, x^{(0)}, 0), & t \in [0, T], \\ x^{(1)}(0) = 0. \end{cases} \quad (4.101)$$

The initial condition  $x^{(1)}(0) = 0$  follows from the fact that the initial condition  $x_0$  does not depend on  $\epsilon$ .

The numerical methods described in Section 4.4 can be used to efficiently compute the unperturbed solution  $x^{(0)}$  and the first-order correction  $x^{(1)}$ .

**REMARK 4.25.** *Consider the equation*

$$\begin{cases} \frac{dx}{dt} = -\epsilon x + 1, & t \in [0, +\infty[, \\ x(0) = 0. \end{cases} \quad (4.102)$$

*The solution can be easily found*

$$x(t, \epsilon) = \frac{e^{-\epsilon t} - 1}{\epsilon}. \quad (4.103)$$

*If we apply the perturbation theory to (4.102), then by solving (4.100) and (4.101) with*

$$f(t, x, \epsilon) = -\epsilon x + 1,$$

*we find*

$$x^{(0)}(t) = -t \quad \text{and} \quad x^{(1)}(t) = \frac{t^2}{2},$$

*which gives*

$$x(t, \epsilon) = -t + \epsilon \frac{t^2}{2} + o(\epsilon). \quad (4.104)$$

*The approximation (4.104) of course coincides with the Taylor expansion of the exact solution given by (4.103). However, note that the approximation is valid only for fixed  $t = O(1)$  and diverges to  $+\infty$  as  $t$  increases while the exact solution converges to  $-1/\epsilon$ . The limits  $\epsilon \rightarrow 0$  and  $t \rightarrow +\infty$  do not commute. Expansion (4.104) is not uniformly valid in time.*

**4.7.2. Singular perturbation theory.** In this subsection we consider a system of ordinary differential equations (together with appropriate boundary conditions) in which the highest derivative is multiplied by a small, positive parameter  $\epsilon$ . In what follows we give the general (nonlinear) form of the system:

$$\begin{cases} \epsilon \frac{d^2 x}{dt^2} = f(t, x, \frac{dx}{dt}), & t \in [0, T], \\ x(0) = x_0, & x(T) = x_1. \end{cases} \quad (4.105)$$

The problem above is called a **singular perturbation problem**, and is characterized by the fact that its order reduces when the problem parameter  $\epsilon$  equals zero. In such a situation, the problem becomes singular since, in general, not all of the original boundary conditions can be satisfied by the reduced problem. Singular perturbed problems form a particular class of **stiff problems**.

Consider the following linear, scalar and of second-order ODE which is subject to Dirichlet boundary conditions:

$$\begin{cases} \epsilon \frac{d^2 x}{dt^2} + 2 \frac{dx}{dt} + x = 0, & t \in [0, 1], \\ x(0) = 0, & x(1) = 1. \end{cases} \quad (4.106)$$

Let

$$\alpha(\epsilon) := \frac{1 - \sqrt{1 - \epsilon}}{\epsilon} \quad \text{and} \quad \beta(\epsilon) := 1 + \sqrt{1 - \epsilon}.$$

The solution of equation (4.106) is given by

$$x(t, \epsilon) = \frac{e^{-\alpha t} - e^{-\beta t/\epsilon}}{e^{-\alpha} - e^{-\beta/\epsilon}}, \quad t \in [0, 1]. \quad (4.107)$$

The solution  $x(t, \epsilon)$  involves two terms which vary on widely different length-scales. Let us consider the behavior of  $x(t, \epsilon)$  as  $\epsilon \rightarrow 0^+$ . The asymptotic behavior is nonuniform, and there are two cases, which lead to matching **outer** and **inner** solutions.

- (i) **Outer limit:**  $t > 0$  fixed and  $\epsilon \rightarrow 0^+$ . Then  $x(t, \epsilon) \rightarrow x^{(0)}(t)$ , where

$$x^{(0)}(t) := e^{(1-t)/2}. \quad (4.108)$$

This leading-order **outer solution** satisfies the boundary condition at  $t = 1$  but not the boundary condition at  $t = 0$ . Indeed,  $x^{(0)}(0) = e^{1/2}$ .

- (ii) **Inner limit:**  $t/\epsilon = \tau$  fixed and  $\epsilon \rightarrow 0^+$ . Then  $x(\epsilon\tau, \epsilon) \rightarrow X^{(0)}(\tau) := e^{1/2}(1 - e^{-2\tau})$ . This leading-order **inner solution** satisfies the boundary condition at  $t = 0$  but not the one at  $t = 1$ , which corresponds to  $\tau = 1/\epsilon$ . Indeed,  $\lim_{\tau \rightarrow +\infty} X^{(0)}(\tau) = e^{1/2}$ .

- (iii) **Matching:** Both the inner and outer expansions are valid in the region  $\epsilon \ll t \ll 1$ , corresponding to  $t \rightarrow 0$  and  $\tau \rightarrow +\infty$  as  $\epsilon \rightarrow 0^+$ . They satisfy the **matching condition**

$$\lim_{t \rightarrow 0^+} x^{(0)}(t) = \lim_{\tau \rightarrow +\infty} X^{(0)}(\tau). \quad (4.109)$$

Let us now construct an asymptotic solution of (4.106) without relying on the fact that we can solve it exactly.

We begin with the outer solution. We look for a straightforward expansion

$$x(t, \epsilon) = x^{(0)}(t) + \epsilon x^{(1)}(t) + O(\epsilon^2). \quad (4.110)$$

We use this expansion in (4.106) and equate the coefficients of the leading-order terms to zero. Guided by our analysis of the exact solution, we only impose the boundary condition at  $t = 1$ . We will see later that matching is impossible if, instead, we attempt to impose the boundary condition at  $t = 0$ . We obtain that

$$\begin{cases} 2 \frac{dx^{(0)}}{dt} + x^{(0)} = 0, & t \in [0, 1], \\ x^{(0)}(1) = 1. \end{cases} \quad (4.111)$$

The solution of (4.111) is given by (4.108), in agreement with the expansion of the exact solution  $x(t, \epsilon)$ .



Next we consider the inner solution. We suppose that there is a **boundary layer** at  $t = 0$  of width  $\delta(\epsilon)$ , and introduce a **stretched variable**  $\tau = t/\delta$ . We look for an inner solution  $X(\tau, \epsilon) = x(t, \epsilon)$ . Since

$$\frac{d}{dt} = \frac{1}{\delta} \frac{d}{d\tau},$$

we find from (4.106) that  $X$  satisfies

$$\frac{\epsilon}{\delta^2} \frac{d^2 X}{d\tau^2} + \frac{2}{\delta} \frac{dX}{d\tau} + X = 0.$$

There are two possible dominant balances in this equation:

- (i)  $\delta = 1$ , leading to the outer solution;
- (ii)  $\delta = \epsilon$ , leading to the inner solution.

Thus we conclude that the boundary layer thickness is of the order of  $\epsilon$ , and the appropriate inner variable is  $\tau = t/\epsilon$ . The equation for  $X$  is then

$$\begin{cases} \frac{d^2 X}{d\tau^2} + 2 \frac{dX}{d\tau} + \epsilon X = 0, \\ X(0, \epsilon) = 0. \end{cases}$$

We impose only the boundary condition at  $\tau = 0$ , since we do not expect the inner expansion to be valid outside the boundary layer where  $t = O(\epsilon)$ .

We seek an inner expansion

$$X(\tau, \epsilon) = X^{(0)}(\tau) + \epsilon X^{(1)}(\tau) + O(\epsilon^2)$$

and find that

$$\begin{cases} \frac{d^2 X^{(0)}}{d\tau^2} + 2 \frac{dX^{(0)}}{d\tau} = 0, \\ X^{(0)}(0) = 0. \end{cases} \quad (4.112)$$

The general solution of (4.112) is

$$X^{(0)}(\tau) = c(1 - e^{-2\tau}), \quad (4.113)$$

where  $c$  is an arbitrary constant of integration.

We can determine the unknown constant  $c$  in (4.113) by requiring that the inner solution (4.113) matches with the outer solution (4.108). Here the matching condition is simply

$$\lim_{t \rightarrow 0^+} x^{(0)}(t) = \lim_{\tau \rightarrow +\infty} X^{(0)}(\tau),$$

which implies that  $c = e^{1/2}$ .

In summary, the asymptotic solution as  $\epsilon \rightarrow 0^+$  is given by

$$x(t, \epsilon) = \begin{cases} e^{1/2}(1 - e^{-2\tau}) & \text{as } \epsilon \rightarrow 0^+ \text{ with } t/\epsilon \text{ fixed,} \\ e^{(1-t)/2} & \text{as } \epsilon \rightarrow 0^+ \text{ with } t \text{ fixed.} \end{cases}$$

### 4.7.3. WKB approximations.

4.7.3.1. *Schrödinger equation.* Consider the **Schrödinger equation**

$$\begin{cases} i\epsilon \frac{\partial \Psi}{\partial t}(t, x) + \epsilon^2 \frac{\partial^2 \Psi}{\partial x^2}(t, x) - V(x)\Psi(t, x) = 0, & x \in \mathbb{R}, t \geq 0, \\ \Psi(0, x) = \Psi_0(x), & x \in \mathbb{R}, \end{cases} \quad (4.114)$$

where  $\epsilon \ll 1$  and  $V(x) > 0$ .

Write

$$\Psi(t, x) = e^{i \frac{S(t, x)}{\epsilon}}.$$

It follows that

$$-\frac{\partial S}{\partial t} - \left(\frac{\partial S}{\partial x}\right)^2 + i\epsilon \frac{\partial^2 S}{\partial x^2} - V(x)S = 0.$$

Hence, the leading order term in the asymptotic expansion with respect to  $\epsilon$

$$S(t, x) = S^{(0)}(t, x) + \epsilon S^{(1)}(t, x) + \dots$$

satisfies the Hamilton-Jacobi type equation

$$\frac{\partial S^{(0)}}{\partial t}(t, x) + \left(\frac{\partial S^{(0)}}{\partial x}\right)^2(t, x) + V(x)S^{(0)}(t, x) = 0.$$

**4.7.4. Wave equation.** Consider the **Helmholtz equation**

$$\left\{ \begin{array}{l} \epsilon^2 \frac{d^2 \Psi}{dx^2}(x) + V(x)\Psi(x) = 0, \quad x \in \mathbb{R}, \end{array} \right. \quad (4.115)$$

where  $\epsilon \ll 1$  and  $V(x) > 0$ .

Using the ansatz

$$\Psi(x) = a(x, \epsilon) e^{\frac{S(x)}{\epsilon}} = (a^{(0)}(x) + \epsilon a^{(1)}(x) + \dots) e^{\frac{S(x)}{\epsilon}},$$

it follows that

$$\left(-\left|\frac{dS}{dx}\right|^2 + V\right) + 2i\epsilon \frac{dS}{dx} \frac{da}{dx} + i\epsilon a \frac{d^2 S}{dx^2} + \epsilon^2 \frac{d^2 a}{dx^2} = 0.$$

Therefore, the phase  $S$  is solution to the **eikonal equation**

$$\left|\frac{dS}{dx}\right|^2(x) = V(x), \quad (4.116)$$

and the leading order term  $a^{(0)}$  in the asymptotic expansion of the amplitude  $a(x, \epsilon)$  with respect to  $\epsilon$  satisfies the **transport equation**

$$2 \frac{dS}{dx} \frac{da^{(0)}}{dx} + a^{(0)} \frac{d^2 S}{dx^2} = 0. \quad (4.117)$$

#### 4.8. Problems

PROBLEM 4.26. (i) *Prove Proposition 4.12.*

(ii) *Prove estimate (4.31).*

PROBLEM 4.27. *Prove that the so-called **explicit Milne's four-step method**:*

$$x^{k+1} = x^{k-3} + \frac{(\Delta t)}{3} \left[ 8f(t_k, x^k) - 4f(t_{k-1}, x^{k-1}) + 8f(t_{k-2}, x^{k-2}) \right]$$

*is of order 4.*

# Geometrical numerical integration methods for differential equations

## 5.1. Introduction

Geometric integration is the numerical integration of a differential equation, while preserving one or more of its geometric properties exactly, *i.e.*, to within round-off error. Many of these geometric properties are of crucial importance in physical applications: preservation of energy, momentum, volume, symmetries, time-reversal symmetry, dissipation, and symplectic structure being examples. The aim of this chapter is to present geometric numerical integration methods for ordinary differential equations. We concentrate mainly on Hamiltonian systems and on methods that preserve their symplectic structure, invariants, symmetries, or volume.

## 5.2. Structure preserving methods for Hamiltonian systems

The numerical methods discussed in Chapter 4 are designed for general differential equations, and a distinction was drawn only between stiff and nonstiff problems. As shown in Chapter 1, Hamiltonian systems are an important class of differential equations with a geometric structure (their flow has the geometric property of being symplectic), whose preservation in the numerical discretization leads to substantially better methods, especially when integrating over long times. In general, most geometric properties are not preserved by the standard numerical methods presented in Chapter 4.

Some of the reasons we are motivated to preserve structure are

- (i) it may yield methods that are faster, simpler, more stable, and/or more accurate for some types of ODEs;
- (ii) it may yield more robust and quantitatively better results than standard methods for the long-time integration of Hamiltonian systems.

The standard problem in numerical ODEs discussed in the previous chapter is to compute the solution to an initial value problem at a fixed time, to within a given global error, as efficiently as possible. The class of method, its order and local error, and choice of time steps are all tailored to this end. In contrast, a typical application of a geometric numerical method is to fix a (sometimes moderately large) time step and compute solutions with perhaps many different initial conditions over very long time intervals.

**5.2.1. Symplectic methods.** Consider the Hamiltonian system

$$\begin{cases} \frac{dp}{dt} = -\frac{\partial H}{\partial q}(p, q), \\ \frac{dq}{dt} = \frac{\partial H}{\partial p}(p, q), \\ p(0) = p_0, q(0) = q_0, \end{cases} \quad (5.1)$$

where  $p_0, q_0 \in \mathbb{R}^d$ , and the Hamiltonian function  $H : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth function.

Let  $x = (p, q)^\top$ . The Hamiltonian system of equations (5.1) can be rewritten as a first-order differential equation

$$\begin{cases} \frac{dx}{dt} = f(x), \\ x(0) = x_0 \in \mathbb{R}^{2d}, \end{cases} \quad (5.2)$$

where  $x_0 = (p_0, q_0)^\top$  and

$$\begin{aligned} f &: \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d} \\ x &\mapsto J^{-1} \nabla H(x). \end{aligned}$$

DEFINITION 5.1. Let  $J$  be defined by (1.26). A numerical one-step method  $(p^{k+1}, q^{k+1}) = \Phi_{\Delta t}(p^k, q^k)$  for solving (5.1) is called **symplectic** if the **numerical flow**  $\Phi_{\Delta t}$  is a symplectic map:

$$\Phi'_{\Delta t}(p, q)^\top J \Phi'_{\Delta t}(p, q) = J, \quad (5.3)$$

for all  $(p, q)$  and all step sizes  $\Delta t$ .

### 5.2.2. Symplectic Euler methods.

THEOREM 5.2. The implicit Euler method for solving (5.1)

$$\begin{cases} p^{k+1} = p^k - \Delta t \frac{\partial H}{\partial q}(p^{k+1}, q^k), \\ q^{k+1} = q^k + \Delta t \frac{\partial H}{\partial p}(p^{k+1}, q^k), \end{cases} \quad (5.4)$$

is symplectic. Moreover, if the Hamiltonian function  $H(p, q) = T(p) + V(q)$  is **separable**, then (5.4) is explicit.

PROOF. Let  $\Phi_{\Delta t}$  be the numerical flow associated with (5.4). We have

$$\Phi'_{\Delta t}(p^k, q^k) = \frac{\partial(p^{k+1}, q^{k+1})}{\partial(p^k, q^k)}.$$

From

$$\begin{pmatrix} I + \Delta t \frac{\partial^2 H}{\partial p \partial q} & 0 \\ -\Delta t \frac{\partial^2 H}{\partial p^2} & I \end{pmatrix} \Phi'_{\Delta t}(p^k, q^k) = \begin{pmatrix} I & -\Delta t \frac{\partial^2 H}{\partial q^2} \\ 0 & I + \Delta t \frac{\partial^2 H}{\partial p \partial q} \end{pmatrix}, \quad (5.5)$$

where the matrices  $\frac{\partial^2 H}{\partial p^2}$ ,  $\frac{\partial^2 H}{\partial q^2}$ , and  $\frac{\partial^2 H}{\partial p \partial q}$  are evaluated at  $(p^{k+1}, q^k)$ , one can easily verify by computing  $\Phi'_{\Delta t}(p^k, q^k)$  from (5.5) that the symplecticity condition (5.3) holds.  $\square$

A variant of (5.4) is

$$\begin{cases} p^{k+1} = p^k - \Delta t \frac{\partial H}{\partial q}(p^k, q^{k+1}), \\ q^{k+1} = q^k + \Delta t \frac{\partial H}{\partial p}(p^k, q^{k+1}). \end{cases} \quad (5.6)$$

Analogously to (5.4), the Euler method (5.6) is symplectic and turns out to be explicit for separable Hamiltonian functions.

### 5.2.3. Composition of symplectic methods.

**THEOREM 5.3.** *The composition of two symplectic one-step methods for solving (5.1) is also symplectic.*

**PROOF.** Let  $\Phi_{\Delta t}^{(1)}$  and  $\Phi_{\Delta t}^{(2)}$  be the numerical flows associated with two symplectic one-step methods for solving (5.1). Let  $\Phi_{\Delta t} := \Phi_{\Delta t}^{(2)} \circ \Phi_{\Delta t}^{(1)}$ . We have

$$\begin{aligned} (\Phi'_{\Delta t}(x))^\top J \Phi'_{\Delta t}(x) &= ((\Phi_{\Delta t}^{(2)})'(x^*)(\Phi_{\Delta t}^{(1)})'(x))^\top J (\Phi_{\Delta t}^{(2)})'(x^*)(\Phi_{\Delta t}^{(1)})'(x) \\ &= ((\Phi_{\Delta t}^{(1)})'(x))^\top ((\Phi_{\Delta t}^{(2)})'(x^*))^\top J (\Phi_{\Delta t}^{(2)})'(x^*)(\Phi_{\Delta t}^{(1)})'(x) \\ &= ((\Phi_{\Delta t}^{(1)})'(x))^\top J (\Phi_{\Delta t}^{(1)})'(x) = J, \end{aligned}$$

where  $x^* = \Phi_{\Delta t}^{(1)}(x)$ . That is, the composition of symplectic one-step methods is again a symplectic one-step method.  $\square$

**5.2.4. The adjoint method.** The flow  $\phi_t$  of an autonomous differential equation  $dx/dt = f(x)$  satisfies  $\phi_{-t}^{-1} = \phi_t$ . This property is in general not satisfied by the one-step map  $\Phi_{\Delta t}$  of a numerical method.

**DEFINITION 5.4.** *The adjoint method  $\Phi_{\Delta t}^*$  of a method  $\Phi_{\Delta t}$  is the inverse map of the original method with reversed time step  $-\Delta t$ , i.e.,*

$$\Phi_{\Delta t}^* := \Phi_{-\Delta t}^{-1}.$$

*In other terms,  $\Phi_{\Delta t}^*$  is defined by replacing, in the method associated with  $\Phi_{\Delta t}$ ,  $\Delta t$  by  $-\Delta t$  and exchanging the superscripts  $k$  and  $k+1$ .*

The adjoint method satisfies the usual properties.

**PROPOSITION 5.5.** *We have*

- (i)  $(\Phi_{\Delta t}^*)^* = \Phi_{\Delta t}$ ;
- (ii)  $(\Phi_{\Delta t}^{(2)} \circ \Phi_{\Delta t}^{(1)})^* = (\Phi_{\Delta t}^{(1)})^* \circ (\Phi_{\Delta t}^{(2)})^*$  for any two one-step methods  $\Phi_{\Delta t}^{(1)}$  and  $\Phi_{\Delta t}^{(2)}$ ;
- (iii)  $(\Phi_{\Delta t/2} \circ \Phi_{\Delta t/2}^*)^* = \Phi_{\Delta t/2} \circ \Phi_{\Delta t/2}^*$ .

**5.2.5. Leapfrog method.** Define the **leapfrog method** (Verlet method and Strömer-Verlet method are also often-used names) for solving the Hamiltonian system (5.1) by

$$\left\{ \begin{array}{l} p^{k+\frac{1}{2}} = p^k - \frac{\Delta t}{2} \frac{\partial H}{\partial q}(p^{k+\frac{1}{2}}, q^k), \\ q^{k+1} = q^k + \frac{\Delta t}{2} \left( \frac{\partial H}{\partial p}(p^{k+\frac{1}{2}}, q^k) + \frac{\partial H}{\partial p}(p^{k+\frac{1}{2}}, q^{k+1}) \right), \\ p^{k+1} = p^{k+\frac{1}{2}} - \frac{\Delta t}{2} \frac{\partial H}{\partial q}(p^{k+\frac{1}{2}}, q^{k+1}). \end{array} \right. \quad (5.7)$$

**THEOREM 5.6.** *The leapfrog method (5.7) for solving the Hamiltonian system (5.1) is symplectic.*

**PROOF.** The leapfrog method (5.7) can be interpreted as the composition of the symplectic Euler method

$$\left\{ \begin{array}{l} p^{k+\frac{1}{2}} = p^k - \frac{\Delta t}{2} \frac{\partial H}{\partial q}(p^{k+\frac{1}{2}}, q^k), \\ q^{k+\frac{1}{2}} = q^k + \frac{\Delta t}{2} \frac{\partial H}{\partial p}(p^{k+\frac{1}{2}}, q^k), \end{array} \right. \quad (5.8)$$

and its adjoint

$$\begin{cases} q^{k+1} &= q^{k+\frac{1}{2}} + \frac{\Delta t}{2} \frac{\partial H}{\partial p}(p^{k+\frac{1}{2}}, q^{k+1}), \\ p^{k+1} &= p^{k+\frac{1}{2}} - \frac{\Delta t}{2} \frac{\partial H}{\partial q}(p^{k+\frac{1}{2}}, q^{k+1}). \end{cases} \quad (5.9)$$

In other terms, if  $\Psi_{\Delta t}$  denotes the numerical flow associated with the leapfrog method and  $\Phi_{\Delta t}$  the one associated with the symplectic Euler method (5.4), then

$$\Psi_{\Delta t} = \Phi_{\Delta t/2}^* \circ \Phi_{\Delta t/2}. \quad (5.10)$$

The methods (5.8) and (5.9) are symplectic. Hence their composition (5.7) is also symplectic.  $\square$

### 5.2.6. Preserving time-reversal symmetry and invariants.

5.2.6.1. *Preserving time-reversal symmetry.* The leapfrog method (5.7) is symmetric with respect to changing the direction of time: replacing  $\Delta t$  by  $-\Delta t$  and exchanging the superscripts  $k$  and  $k+1$  results in the same method. In terms of the numerical one-step map  $\Phi_{\Delta t} : (p^k, q^k) \mapsto (p^{k+1}, q^{k+1})$ , the symmetry property is stated as follows.

DEFINITION 5.7. *The numerical one-step map  $\Phi_{\Delta t}$  is said to be **symmetric** if*

$$\Phi_{\Delta t} = \Phi_{\Delta t}^* (= \Phi_{-\Delta t}^{-1}). \quad (5.11)$$

Relation (5.11) does not hold for the symplectic Euler methods (5.8) and (5.9), where the time reflection transforms (5.8) to (5.9) and vice versa.

The time-symmetry of the leapfrog method (5.7), which follows from (5.10) and item (iii) in Proposition 5.5, implies an important geometric property of the numerical map, namely **reversibility**.

Assume that

$$H(-p, q) = H(p, q). \quad (5.12)$$

Then the system (5.1) has the property that inverting the direction of the initial  $p_0$  does not change the solution trajectory. The flow  $\phi_t$  associated with (5.1) satisfies

$$\phi_t(p_0, q_0) = (p, q) \Rightarrow \phi_t(-p, q) = (-p_0, q_0). \quad (5.13)$$

Relation (5.13) shows that  $\phi_t$  is **reversible** with respect to the reflection  $(p, q) \mapsto (-p, q)$ .

DEFINITION 5.8. *The numerical one-step map  $\Phi_{\Delta t}$  is said to be **reversible** if*

$$\Phi_{\Delta t}(p, q) = (\hat{p}, \hat{q}) \Rightarrow \Phi_{\Delta t}(-\hat{p}, \hat{q}) = (-p, q), \quad (5.14)$$

for all  $p, q$  and all  $\Delta t$ .

Since

$$\Phi_{\Delta t}(p, q) = (\hat{p}, \hat{q}) \Rightarrow \Phi_{-\Delta t}(-p, q) = (-\hat{p}, \hat{q}), \quad (5.15)$$

the symmetry (5.11) of the leapfrog method (5.7) is therefore equivalent to the reversibility (5.13).

THEOREM 5.9. *The leapfrog method (5.7) applied to (5.1) with  $H$  satisfying (5.12) is both symmetric and reversible, i.e., its one-step map satisfies (5.11) and (5.14).*

REMARK 5.10. *Consider a one-step method  $\Phi_{\Delta t}$  of order one. Then, formally,*

$$\Phi_{\Delta t}(x_0) = \varphi_{\Delta t}(x_0) + C(x_0)\Delta t + O((\Delta t)^2),$$

and

$$\Phi_{\Delta t}^*(x_0) = \varphi_{\Delta t}(x_0) - C(x_0)\Delta t + O((\Delta t)^2),$$

with  $\varphi_{\Delta t}$  being the exact flow. Therefore, if  $\Phi_{\Delta t}$  is symmetric, then it should be of order two since  $C(x_0)$  has to be zero.

REMARK 5.11. *From Remark 5.10, it follows that the composition with the adjoint method turns every consistent one-step method of order one into a second-order symmetric method*

$$\Psi_{\Delta t} = \Phi_{\Delta t/2} \circ \Phi_{\Delta t/2}^*.$$

## 5.2.6.2. Preserving invariants.

DEFINITION 5.12. A numerical one-step method  $\Phi_{\Delta t}$  for solving (5.2) is said to **preserve the invariant**  $F$  if  $F(\Phi_{\Delta t}(p, q)) = \text{Constant}$  for all  $p, q$  and all  $\Delta t$ . If  $F = H$ , then we say that the scheme preserves **energy**.

THEOREM 5.13. The leapfrog method (5.7) applied to (5.1) preserves linear invariants and quadratic invariants of the form

$$F(p, q) = p^\top (Bq + b). \quad (5.16)$$

PROOF. Let the linear invariant be  $F(p, q) = b^\top q + c^\top p$ , so that

$$b^\top \frac{\partial H}{\partial p}(p, q) - c^\top \frac{\partial H}{\partial q}(p, q) = 0,$$

for all  $p, q$ . Multiplying the formulas for  $\Phi_{\Delta t}(p, q)$  in (5.7) by  $(c, b)^\top$  thus yields the desired result on linear invariants.

Next we turn to the conservation by the leapfrog method of quadratic invariants of the form (5.16). In order to prove that (5.7) applied to (5.1) preserves quadratic invariants of the form  $F(p, q) = p^\top (Bq + b)$ , we write (5.7) as the composition of the two symplectic Euler methods (5.8) and (5.9). For the first half-step, we obtain

$$(p^{k+\frac{1}{2}})^\top (Bq^{k+\frac{1}{2}} + b) = (p^k)^\top (Bq^k + b).$$

For the second half-step, we obtain in the same way

$$(p^{k+1})^\top (Bq^{k+1} + b) = (p^{k+\frac{1}{2}})^\top (Bq^{k+\frac{1}{2}} + b),$$

and the result follows.  $\square$

The energy is generally not preserved by the leapfrog method (5.7). Consider  $H(p, q) = \frac{1}{2}(p^2 + q^2)$ . Applying (5.7) gives

$$\begin{pmatrix} p^{k+1} \\ q^{k+1} \end{pmatrix} = \begin{bmatrix} 1 - \frac{(\Delta t)^2}{2} & -\Delta t(1 - \frac{(\Delta t)^2}{4}) \\ \Delta t & 1 - \frac{(\Delta t)^2}{2} \end{bmatrix} \begin{pmatrix} p^k \\ q^k \end{pmatrix}. \quad (5.17)$$

Since the **propagation matrix** in (5.17) is not orthogonal,  $H(p, q)$  is not preserved along numerical solutions.

Consider the Hamiltonian

$$H(p, q) := \frac{1}{2}p^\top M^{-1}p + V(q), \quad (5.18)$$

where  $M$  is a symmetric positive definite matrix and the potential  $V$  is a smooth function.

In the particular case of the Hamiltonian (5.18), the leapfrog method (5.7) reduces to the explicit method

$$\begin{cases} p^{k+\frac{1}{2}} = p^k - \frac{\Delta t}{2} \nabla V(q^k), \\ q^{k+1} = q^k + \Delta t M^{-1} p^{k+\frac{1}{2}}, \\ p^{k+1} = p^{k+\frac{1}{2}} - \frac{\Delta t}{2} \nabla V(q^{k+1}). \end{cases} \quad (5.19)$$

Note that the Hamiltonian (5.18) is invariant under  $p \mapsto -p$  and the corresponding Hamiltonian system (5.1) is invariant under the transformation

$$\begin{bmatrix} p \\ t \end{bmatrix} \mapsto \begin{bmatrix} -p \\ -t \end{bmatrix}. \quad (5.20)$$

The **time-reversal symmetry** of (5.1) is preserved by the leapfrog method (5.19).

5.2.6.3. *Preserving volume.* Recall that, due to equality of mixed partial derivatives, (5.2) is divergence-free, *i.e.*,

$$\nabla \cdot f := \sum_{i=1}^{2d} \frac{\partial f_i}{\partial x_i} = 0.$$

A remarkable feature of divergence-free vector fields is that the associated flows are volume preserving.

Given a map  $\phi : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$  and a domain  $\Omega$ , by change of variables

$$\text{vol}(\phi(\Omega)) = \int_{\Omega} |\det \phi'(y)| dy,$$

where  $\phi'$  is the Jacobian of  $\phi$ . It follows that  $\phi$  preserves volume provided that

$$|\det \phi'(y)| = 1 \quad \text{for } y \in \Omega. \quad (5.21)$$

Let  $\phi_t$  be the flow associated with  $dx/dt = f(x)$ , where  $\nabla \cdot f = 0$ . Then  $\phi_t$  satisfies

$$\frac{d\phi_t(y)}{dt} = f(\phi_t(y)),$$

and therefore, its Jacobian  $\phi'_t$  satisfies

$$\frac{d\phi'_t(y)}{dt} = f'(\phi_t(y))\phi'_t(y).$$

Assuming  $\phi'_t$  is invertible yields

$$\text{tr} \left[ \frac{d\phi'_t(y)}{dt} \phi'_t(y)^{-1} \right] = \text{tr} f'(\phi_t(y)).$$

Combining  $\text{tr} f' = \nabla \cdot f = 0$  and Jacobi's formula (3.41) for the derivative of a determinant gives

$$\text{tr} \left[ \frac{d\phi'_t(y)}{dt} \phi'_t(y)^{-1} \right] = \frac{1}{\det \phi'_t(y)} \frac{d}{dt} \det \phi'_t(y) = 0.$$

Hence,

$$\det \phi'_t(y) = \det \phi'_{t=0}(y) = 1.$$

The following result holds.

**THEOREM 5.14 (Liouville's theorem).** *The flow  $\phi_t$  associated with the system*

$$\begin{cases} \frac{dx}{dt} = f(x), \\ x(0) = x_0 \in \mathbb{R}^{2d}, \end{cases} \quad (5.22)$$

where the  $C^1$  vector field  $f$  is **divergence-free**, is a **volume preserving map** (for all  $t$ ).

Note that if the system (5.22) is Hamiltonian, then Theorem 5.14 can be immediately obtained from the **symplecticity** of the associated flow. In fact, from

$$(\phi'_t)^\top J \phi'_t = J,$$

it follows that  $|\det \phi'_t|^2 = 1$  since  $\det J = 1$ . Moreover, using the facts that  $\det \phi'_{t=0} = 1$  and the continuity of the determinant, we obtain that  $\det \phi'_t = 1$  for all  $t$ .

**REMARK 5.15.** *Since*

$$\nabla \cdot J^{-1} \nabla H(x) = - \sum_{j=1}^d \frac{\partial^2 H}{\partial x_j \partial x_{d+j}} + \sum_{j=1}^d \frac{\partial^2 H}{\partial x_{d+j} \partial x_j} = 0$$

for any smooth function  $H$ , Hamiltonian systems are divergence free equations. If  $d = 1$ , all divergence-free systems are Hamiltonians since  $\nabla \cdot f = 0$  implies that  $f = \nabla \times H$  for some function  $H \in C^2$  (at least locally) and

$$J^{-1} \nabla = \nabla \times .$$



For  $d > 1$ , the previous identity is no longer true. Consequently, divergence-free systems are not necessary Hamiltonians.

**DEFINITION 5.16.** A numerical one-step method for solving (5.22) is said to be **volume preserving** if  $|\det \Phi'_{\Delta t}(p, q)| = 1$  for all  $p, q$ .

Note that if (5.22) is a Hamiltonian system, then any symplectic numerical method preserves the volume. However, no standard methods can be volume-preserving for all divergence-free vector fields.

**EXAMPLE 5.17.** Consider the divergence-free problem

$$\begin{cases} \frac{dx}{dt} = Ax, \\ x(0) = x_0 \in \mathbb{R}^{2d}, \end{cases} \quad (5.23)$$

where  $A \in \mathbb{M}_{2d}(\mathbb{R})$  and  $\text{tr} A = 0$ . The Explicit and implicit Euler's schemes for solving (5.23)

$$\begin{aligned} x^{k+1} &= x^k + \Delta t A x^k, \\ x^{k+1} &= x^k + \Delta t A x^{k+1}, \end{aligned}$$

are volume-preserving if and only if

$$|\det(I + \Delta t A)| = 1,$$

and

$$|\det(I - \Delta t A)| = 1,$$

respectively.

**5.2.7. Composition methods.** Now using the fact that (5.2) is divergence-free, we have (when  $f_{2d}$  is assumed for simplicity to depend only on  $x_{2d}$ ),

$$\begin{aligned} f_{2d}(x) &= f_{2d}(\bar{x}) + \int_{\bar{x}}^{x_{2d}} \frac{\partial f_{2d}}{\partial x_{2d}} dx_{2d} \\ &= f_{2d}(\bar{x}) - \int_{\bar{x}}^{x_{2d}} \left( \sum_{i=1}^{2d-1} \frac{\partial f_i(x)}{\partial x_i} \right) dx_{2d}, \end{aligned} \quad (5.24)$$

where  $\bar{x}$  is an arbitrary point which can be chosen conveniently (e.g., if possible such that  $f_{2d}(\bar{x}) = 0$ ).

Substituting (5.24) into (5.2) yields

$$\begin{aligned} \frac{dx_1}{dt} &= f_1(x), \\ &\vdots \\ \frac{dx_{2d-1}}{dt} &= f_{2d-1}(x), \\ \frac{dx_{2d}}{dt} &= f_{2d}(\bar{x}) - \sum_{i=1}^{2d-1} \int_{\bar{x}}^{x_{2d}} \frac{\partial f_i(x)}{\partial x_i} dx_{2d}. \end{aligned} \quad (5.25)$$

We now split this as the sum of  $2d - 1$  vector fields

$$\begin{aligned} \frac{dx_i}{dt} &= 0, \quad i \neq j, 2d - 1, \\ \frac{dx_j}{dt} &= f_j(x), \\ \frac{dx_{2d}}{dt} &= f_{2d}(\bar{x}) \delta_{j, 2d-1} - \int_{\bar{x}}^{x_{2d}} \frac{\partial f_j(x)}{\partial x_j} dx_{2d}, \end{aligned} \quad (5.26)$$

for  $j = 1, \dots, 2d - 1$ . Here  $\delta$  is the Kronecker delta function.

Note that each of the  $2d - 1$  vector fields is divergence-free. Moreover, we have split (5.25) into the  $2d - 1$  problems (5.26). Each of these problems has a simpler structure than (5.2). In fact, each of them corresponds to a two-dimensional Hamiltonian system

$$\begin{aligned}\frac{dx_j}{dt} &= -\frac{\partial H_j}{\partial x_{2d}}, \\ \frac{dx_{2d}}{dt} &= \frac{\partial H_j}{\partial x_j},\end{aligned}\tag{5.27}$$

with Hamiltonian

$$H_j(x) := f_{2d}(\bar{x})\delta_{j,2d-1}x_j - \int_{\bar{x}}^{x_{2d}} f_j(x) dx_{2d},\tag{5.28}$$

treating  $x_i$  for  $i \neq j, 2d$  as fixed parameters.

Each of the two-dimensional problems (5.27) can either be solved exactly (if possible), or approximated with a symplectic integrator  $\Phi_{\Delta t}^{(j)}$ . A volume-preserving integrator for  $f$  is then given by

$$\Phi_{\Delta t} = \Phi_{\Delta t}^{(1)} \circ \Phi_{\Delta t}^{(2)} \circ \dots \circ \Phi_{\Delta t}^{(2d-1)}.\tag{5.29}$$

**5.2.8. Splitting methods.** Consider a Hamiltonian system

$$\frac{dx}{dt} = f(x) = J^{-1}\nabla H(x), \quad H(x) = H_1(x) + H_2(x),\tag{5.30}$$

and suppose the flows

$$\frac{dx}{dt} = f_1(x) = J^{-1}\nabla H_1(x) \quad \text{and} \quad \frac{dx}{dt} = f_2(x) = J^{-1}\nabla H_2(x),\tag{5.31}$$

can be exactly integrated.

Let  $\phi_t^{(1)}$  and  $\phi_t^{(2)}$  be the exact flows associated with the equations in (5.31) and let  $\phi$  be the flow associated with (5.30).

Since the exact solution of a Hamiltonian system defines a symplectic map, we have

$$((\phi_t^{(1)})')^\top J(\phi_t^{(1)})' = J \quad \text{and} \quad ((\phi_t^{(2)})')^\top J(\phi_t^{(2)})' = J.$$

Next consider the numerical method defined by composing these two exact flows:

$$\Phi_{\Delta t}(x) := \phi_{\Delta t}^{(2)} \circ \phi_{\Delta t}^{(1)}(x).$$

This map is also symplectic, since

$$\begin{aligned}(\Phi_{\Delta t}'(x))^\top J\Phi_{\Delta t}'(x) &= ((\phi_{\Delta t}^{(2)})'(x^*)(\phi_{\Delta t}^{(1)})'(x))^\top J(\phi_{\Delta t}^{(2)})'(x^*)(\phi_{\Delta t}^{(1)})'(x) \\ &= ((\phi_{\Delta t}^{(1)})'(x))^\top ((\phi_{\Delta t}^{(2)})'(x^*))^\top J(\phi_{\Delta t}^{(2)})'(x^*)(\phi_{\Delta t}^{(1)})'(x) \\ &= ((\phi_{\Delta t}^{(1)})'(x))^\top J(\phi_{\Delta t}^{(1)})'(x) = J,\end{aligned}$$

where  $x^* = \phi_{\Delta t}^{(1)}(x)$ . That is, as shown in Theorem 5.3, the composition of symplectic maps is again a symplectic map.

If, from a given initial value  $x_0$ , we first solve the first system to obtain a value  $x_{\frac{1}{2}}$ , and from this value integrate the second system to obtain  $x_1$ , we get two numerical integrators where one is the adjoint of the other:

$$\Phi_{\Delta t} = \phi_{\Delta t}^{(2)} \circ \phi_{\Delta t}^{(1)} \quad \text{and} \quad \Phi_{\Delta t}^* = \phi_{\Delta t}^{(1)} \circ \phi_{\Delta t}^{(2)}.$$

By Taylor expansion, we find that

$$\phi_{\Delta t}^{(2)} \circ \phi_{\Delta t}^{(1)}(x_0) = \phi_{\Delta t}(x_0) + O((\Delta t)^2),$$

so that  $\Phi_{\Delta t}$  (and analogously  $\Phi_{\Delta t}^*$ ) gives approximation of order one to the solution of (5.30).

Another idea is to use a symmetric version and put

$$\Phi_{\Delta t} = \phi_{\Delta t/2}^{(1)} \circ \phi_{\Delta t}^{(2)} \circ \phi_{\Delta t/2}^{(1)}.\tag{5.32}$$

By breaking up in (5.32)

$$\phi_{\Delta t}^{(2)} = \phi_{\Delta t/2}^{(2)} \circ \phi_{\Delta t/2}^{(2)}$$

an using Taylor expansion, we see that (5.32) is symmetric and of order two.

EXAMPLE 5.18. Consider the separable Hamiltonian  $H(p, q) = U(p) + V(q)$ . Based on splitting the Hamiltonian  $H$  into  $U$  and  $V$ , we interpret the symplectic Euler methods and the leapfrog method for solving (5.2) as splitting methods.

To do so, we consider (5.30) as the sum of two Hamiltonians, the first one depending only on  $p$ , the second one only on  $q$ . The corresponding Hamiltonian systems

$$\left\{ \begin{array}{l} \frac{dp}{dt} = 0, \\ \frac{dq}{dt} = \frac{\partial U}{\partial p}(p), \\ p(0) = p_0, q(0) = q_0, \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} \frac{dp}{dt} = -\frac{\partial V}{\partial q}(q), \\ \frac{dq}{dt} = 0, \\ p(0) = p_0, q(0) = q_0, \end{array} \right.$$

can be solved explicitly

$$\left\{ \begin{array}{l} p(t) = p_0, \\ q(t) = q_0 + t \frac{\partial U}{\partial p}(p_0), \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} p(t) = p_0 - t \frac{\partial V}{\partial q}(q_0), \\ q(t) = q_0. \end{array} \right.$$

Denoting the flows of these two systems by  $\phi_t^U$  and  $\phi_t^V$ , we see that the symplectic Euler method

$$\left\{ \begin{array}{l} p^{k+1} = p^k - \Delta t \frac{\partial V}{\partial q}(q^k), \\ q^{k+1} = q^k + \Delta t \frac{\partial U}{\partial p}(p^{k+1}), \end{array} \right.$$

is just the composition

$$\phi_{\Delta t}^U \circ \phi_{\Delta t}^V, \tag{5.33}$$

and its adjoint is

$$\phi_{\Delta t}^V \circ \phi_{\Delta t}^U. \tag{5.34}$$

The leapfrog method

$$\left\{ \begin{array}{l} p^{k+\frac{1}{2}} = p^k - \frac{\Delta t}{2} \frac{\partial V}{\partial q}(q^k), \\ q^{k+1} = q^k + \Delta t \frac{\partial U}{\partial p}(p^{k+\frac{1}{2}}), \\ p^{k+1} = p^{k+\frac{1}{2}} - \frac{\Delta t}{2} \frac{\partial V}{\partial q}(q^{k+1}), \end{array} \right.$$

is

$$\phi_{\Delta t/2}^V \circ \phi_{\Delta t}^U \circ \phi_{\Delta t/2}^V. \tag{5.35}$$

Decompositions (5.33), (5.34), and (5.35) give second proofs of Theorems 5.2 and 5.6 in the case of a separable Hamiltonian. They also show that the symplectic Euler methods are of order one while the leapfrog method is order two.

### 5.3. Runge-Kutta methods

Now we turn to Runge-Kutta methods

$$\left\{ \begin{array}{l} x_{i,k} = x^k + (\Delta t) \sum_{j=1}^m a_{ij} f(x_{j,k}), \\ x^{k+1} = x^k + (\Delta t) \sum_{i=1}^m b_i f(x_{i,k}), \end{array} \right. \tag{5.36}$$

for solving (5.2).

**THEOREM 5.19.** (i) *All the Runge-Kutta methods (5.36) preserve linear invariants;*  
(ii) *The Runge-Kutta method (5.36) whose coefficients satisfy the condition*

$$b_i a_{ij} + b_j a_{ji} - b_i b_j = 0, \quad i, j = 1, \dots, m, \quad (5.37)$$

*preserves all quadratic invariants.*

**PROOF.** Define  $\Phi_{\Delta t}$  by  $x^{k+1} = \Phi_{\Delta t}(x^k)$ . Let  $F(x) = d^\top x$ , where  $d \in \mathbb{R}^{2d}$ . We compute

$$F(\Phi_{\Delta t}(x^k)) = d^\top \left( x^k + \Delta t \sum_{i=1}^m b_i f(x_{i,k}) \right) = d^\top x^k,$$

since  $d^\top x$  is assumed to be an invariant of (5.2) and hence  $d^\top f(x_{i,k}) = 0$ .

Next, let  $F(x) = x^\top C x$ , where  $C$  is a symmetric  $2d \times 2d$  matrix. Assume that  $F$  is an invariant of (5.2). We have

$$x^\top C f(x) = 0 \quad \text{for all } x. \quad (5.38)$$

On the other hand, we have

$$\begin{aligned} F(\Phi_{\Delta t}(x^k)) &= \left( x^k + \Delta t \sum_{j=1}^m b_j f(x_{j,k}) \right)^\top C \left( x^k + \Delta t \sum_{i=1}^m b_i f(x_{i,k}) \right) \\ &= (x^k)^\top C x^k + (\Delta t) \sum_{i=1}^m (x^k)^\top C b_i f(x_{i,k}) + (\Delta t) \sum_{j=1}^m b_j f(x_{j,k})^\top C x^k \\ &\quad + (\Delta t)^2 \sum_{i,j=1}^m b_i b_j f(x_{j,k})^\top C f(x_{i,k}). \end{aligned}$$

From (5.38), we obtain

$$(x_{i,k})^\top C f(x_{i,k}) = 0,$$

and hence, by writing

$$x^k = x^k + \Delta t \sum_{j=1}^m a_{ij} f(x_{j,k}) - \Delta t \sum_{j=1}^m a_{ij} f(x_{j,k}) = x_{i,k} - \Delta t \sum_{j=1}^m a_{ij} f(x_{j,k}),$$

we get

$$\begin{aligned} F(\Phi_{\Delta t}(x^k)) &= (x^k)^\top C x^k - (\Delta t)^2 \sum_{i,j=1}^m b_i a_{ij} f(x_{j,k})^\top C f(x_{i,k}) - (\Delta t)^2 \sum_{i,j=1}^m b_j a_{ji} f(x_{j,k})^\top C f(x_{i,k}) \\ &\quad + (\Delta t)^2 \sum_{i,j=1}^m b_i b_j f(x_{j,k})^\top C f(x_{i,k}) \\ &= (x^k)^\top C x^k - (\Delta t)^2 \left( \sum_{i,j=1}^m (b_i a_{ij} + b_j a_{ji} - b_i b_j) f(x_{j,k})^\top C f(x_{i,k}) \right). \end{aligned}$$

Therefore, the Runge-Kutta method (5.36) preserves the quadratic invariant  $F$  provided that (5.37) holds.  $\square$

Lemma 1.19 shows that  $H$  is an invariant of (5.2). If  $H$  is quadratic, then Theorem 5.19 says that the energy is preserved by the Runge-Kutta method (5.36) provided that condition (5.37) holds.

The following characterization of **symplectic Runge-Kutta methods** for solving (5.2) holds.

**THEOREM 5.20.** *The Runge-Kutta method (5.36) for solving (5.2) whose coefficients satisfy condition (5.37) is symplectic.*

PROOF. Theorem 1.25 shows that the flow  $\phi_t$  is a symplectic transformation (if  $H$  is smooth enough). Let  $\Psi(t) := \frac{\partial \phi_t(x_0)}{\partial x_0} = \phi'_t$ , where  $x_0$  is the initial condition. We have

$$\begin{cases} \frac{d\Psi}{dt} = f'(x)\Psi, \\ \Psi(0) = I. \end{cases} \quad (5.39)$$

Apply a Runge-Kutta method satisfying (5.37) to (5.2) and (5.39) to obtain the approximations  $x^{k+1}$  and  $\Psi^{k+1}$  from  $x^k$  and  $\Psi^k$ . Since  $\Psi^\top J \Psi$  is a quadratic invariant of the **augmented system** (5.36) and (5.39), we obtain

$$(\Psi^k)^\top J \Psi^k = J \quad \text{for all } k.$$

Suppose for a moment that

$$\Psi^{k+1} = \frac{\partial x^{k+1}}{\partial x^k}. \quad (5.40)$$

We obtain

$$\left(\frac{\partial x^{k+1}}{\partial x^k}\right)^\top J \frac{\partial x^{k+1}}{\partial x^k} = J,$$

which means that the Runge-Kutta method for solving (5.2) whose coefficients satisfy condition (5.37) is symplectic.

In order to complete the proof, we prove (5.40). We want to show that the result of first applying  $\Phi_{\Delta t}$  and then differentiating with respect to  $x^k$  is the same as applying the same Runge-Kutta method to (5.39).

In fact, on the one hand, by differentiating (5.36) with respect to  $x^k$  we obtain

$$\begin{cases} \frac{\partial x_{i,k}}{\partial x^k} = I + (\Delta t) \sum_{j=1}^m a_{ij} f'(x_{j,k}) \frac{\partial x_{j,k}}{\partial x^k}, \\ \frac{\partial x^{k+1}}{\partial x^k} = I + (\Delta t) \sum_{i=1}^m b_i f'(x_{i,k}) \frac{\partial x_{i,k}}{\partial x^k}. \end{cases} \quad (5.41)$$

Multiplying the first equation in (5.41) by  $f'(x_{i,k}) \frac{\partial x_{i,k}}{\partial x^k}$

$$f'(x_{i,k}) \frac{\partial x_{i,k}}{\partial x^k} = f'(x_{i,k}) \left( I + (\Delta t) \sum_{j=1}^m a_{ij} f'(x_{j,k}) \frac{\partial x_{j,k}}{\partial x^k} \right), \quad (5.42)$$

$$\frac{\partial x^{k+1}}{\partial x^k} = I + (\Delta t) \sum_{i=1}^m b_i f'(x_{i,k}) \frac{\partial x_{i,k}}{\partial x^k}. \quad (5.43)$$

On the other hand, applying the same Runge-Kutta method to (5.39) yields

$$\Psi_{i,k} = f'(x^k + \Delta t \sum_{j=1}^m a_{ij} x_{j,k}) \left( I + (\Delta t) \sum_{j=1}^m a_{ij} \Psi_{j,k} \right), \quad (5.44)$$

$$\Psi^{k+1} = I + (\Delta t) \sum_{i=1}^m b_i \Psi_{i,k}. \quad (5.45)$$

We conclude the proof by observing that (5.44) is the same system as (5.42) but in the unknowns  $\Psi_{i,k}$ ,  $i = 1, \dots, m$ . It is easily seen that this system has a unique solution for sufficiently small  $\Delta t$ , so it must be

$$\Psi_{i,k} = f'(x_{i,k}) \frac{\partial x_{i,k}}{\partial x^k} \quad \text{for } i = 1, \dots, m,$$

which, in view of (5.43) and (5.45), yields (5.40).  $\square$

For arbitrary Hamiltonians, the only known symplectic one-step numerical methods are the symplectic Runge-Kutta methods of the form (4.55) that satisfy the symplectic condition (5.37).

EXAMPLE 5.21. *The midpoint scheme for solving (5.2)*

$$x^{k+1} = x^k + \Delta t f\left(\frac{x^k + x^{k+1}}{2}\right), \quad (5.46)$$

*is symplectic and preserves linear and quadratic invariants. Moreover, it is time-reversible.*

#### 5.4. Long-time behaviour of numerical solutions

In (5.17) we have seen that the energy is not exactly preserved by the leapfrog method (5.7). In that example, it is however, approximately preserved. As shown in the following theorem, the symplecticity of a one-step numerical method yields an approximate conservation of energy over very long times for general Hamiltonian systems.

THEOREM 5.22. *For an analytic Hamiltonian  $H$  and a symplectic one-step numerical method  $\Phi_{\Delta t}$  of order  $n$ , if the numerical trajectory remains in a compact subset, then there exist  $h > 0$  and  $\Delta t^* > 0$  such that, for  $\Delta t \leq \Delta t^*$ ,*

$$H(p^k, q^k) = H(p^0, q^0) + O((\Delta t)^n), \quad (5.47)$$

*for exponentially long times  $k\Delta t \leq e^{\frac{h}{\Delta t}}$ . Here,  $(p^{k+1}, q^{k+1}) = \Phi_{\Delta t}(p^k, q^k)$ .*

Theorem (5.22) is based on symplecticity. It can be proved via backward error analysis. The idea is to deduce the long-time behavior estimate (5.47) from properties of the solution of the equation corresponding to an approximation  $H_{\Delta t}$  of the Hamiltonian  $H$ .

#### 5.5. Problems

PROBLEM 5.23. *Consider the flow  $\phi_t$  of (5.2). Given a one-step numerical scheme  $x^{k+1} = \Phi_{\Delta t}(x^k)$ , its adjoint*

$$x^{k+1} = \Phi_{\Delta t}^*(x^k)$$

*is the method defined by*

$$x^k = \Phi_{-\Delta t}(x^{k+1}),$$

*or equivalently,*

$$x^{k+1} = \Phi_{-\Delta t}^{-1}(x^k).$$

- (i) *Prove that  $\phi_t \circ \phi_s = \phi_{t+s}$  and hence,  $\phi_t \circ \phi_{-t} = I$ , for  $t, s \in \mathbb{R}$ .*
- (ii) *Prove that  $\Phi_{\Delta t}$  is symmetric if and only if  $\Phi_{\Delta t} = \Phi_{\Delta t}^*$ .*
- (iii) *Prove that  $(\Phi_{\Delta t}^*)^* = \Phi_{\Delta t}$ .*
- (iv) *Prove that for any one-step methods  $\Phi_{\Delta t}$  and  $\Psi_{\Delta t}$ ,*

$$(\Phi_{\Delta t} \circ \Psi_{\Delta t})^* = \Psi_{\Delta t}^* \circ \Phi_{\Delta t}^*.$$

- (v) *Prove that for any one-step method  $\Phi_{\Delta t}$ ,*

$$x^{k+1} = \Phi_{\Delta t/2} \circ \Phi_{\Delta t/2}^*(x^k)$$

*is a symmetric method.*

PROBLEM 5.24. *Consider the Runge-Kutta method that is consistent, i.e.,  $\sum_{i=1}^m b_i = 1$ , and with coefficients such that  $\sum_{j=1}^m a_{ij} = c_i$ , for  $1 \leq i \leq m$ .*

- (i) *Prove that the adjoint of the Runge-Kutta method is again a Runge-Kutta method, with coefficients given by*

$$a_{ij}^* = b_{m+1-j} - a_{m+1-i, m+1-j}, \quad b_i^* = b_{m+1-i} \quad \text{for } 1 \leq i, j \leq m.$$

- (ii) *Deduce that if the method is symmetric, then  $a_{ij} = b_j - a_{m+1-i, m+1-j}$  for all  $i, j = 1, \dots, m$ .*
- (iii) *Prove that, if the Runge-Kutta method is explicit, then it can not be symmetric.*

PROBLEM 5.25. Consider the average vector field method

$$x^{k+1} = x^k + \Delta t \int_0^1 f(\theta x^{k+1} + (1 - \theta)x^k) d\theta, \quad (5.48)$$

where the vector field  $f$  is Lipschitz continuous.

- (i) Prove that (5.48) is well-defined for a stepsize  $\Delta t$  small enough.
- (ii) Prove that (5.48) preserves exactly the energy of any Hamiltonian system.
- (iii) Suppose that the Hamiltonian function is a polynomial. Prove that there exists a quadrature formula  $(b_i, c_i)_{i=1, \dots, m}$ , with nodes  $c_i$  and weights  $b_i$ , such that

$$\int_0^1 f(\theta x^{k+1} + (1 - \theta)x^k) d\theta = \sum_{i=1}^m b_i f(x^k + c_i(x^{k+1} - x^k)),$$

where  $f(x) = J^{-1} \nabla H(x)$ .

- (iv) Construct a Runge-Kutta method that exactly preserves a given polynomial Hamiltonian  $H$ .





## Finite difference methods

### 6.1. Introduction

Finite difference methods are basic numerical solution methods for partial differential equations. They are obtained by replacing the derivatives in the equation by the appropriate numerical differentiation formulas. However, there is no guarantee that the resulting numerical scheme will accurately approximate the true solution. Further analysis is required. In this chapter, we establish some of the most basic finite difference schemes for the heat and the wave equations.

### 6.2. Numerical algorithms for the heat equation

#### 6.2.1. Finite difference approximations. Consider the heat equation

$$\begin{cases} \frac{\partial u}{\partial t} - \gamma \frac{\partial^2 u}{\partial x^2} = 0, & x \in [0, 1], t \geq 0, \\ u(0, t) = u(1, t) = 0, & t \geq 0, \\ u(x, 0) = u_0(x), & x \in [0, 1], \end{cases} \quad (6.1)$$

where  $\gamma > 0$  is the thermal conductivity.

In order to design a numerical approximation to the solution  $u$  of (6.1), we begin by introducing a rectangular mesh consisting of points  $(t_k, x_j)$  with

$$0 = t_0 < t_1 < t_2 < \dots \quad \text{and} \quad 0 = x_0 < x_1 < \dots < x_{N+1} = 1.$$

For simplicity, we maintain a uniform mesh spacing in both directions, with

$$\Delta t = t_{k+1} - t_k, \quad \Delta x = x_{j+1} - x_j = \frac{1}{N+1},$$

representing, respectively, the time step size and the spatial mesh size. We shall use the notation

$$u_j^k \approx u(x_j, t_k) \quad \text{where} \quad x_j = j\Delta x, \quad t_k = k\Delta t,$$

to denote the numerical approximation of  $u$  at the mesh point  $(x_j, t_k)$ .

The Dirichlet boundary conditions  $u(0, t) = u(1, t) = 0$ ,  $t \geq 0$ , yield

$$u_0^k = u_{N+1}^k = 0 \quad \text{for all } k > 0. \quad (6.2)$$

As a first attempt at designing a numerical method, we shall employ the simplest finite difference approximations to the derivatives. The time derivative can be approximated by

$$\frac{\partial u}{\partial t}(x_j, t_k) \approx \frac{u(x_j, t_{k+1}) - u(x_j, t_k)}{\Delta t} + O(\Delta t) \approx \frac{u_j^{k+1} - u_j^k}{\Delta t} + O(\Delta t). \quad (6.3)$$

Similarly, the second order space derivative is approximated by **centered differences**

$$\begin{aligned} \frac{\partial^2 u}{\partial x^2}(x_j, t_k) &\approx \frac{u(x_{j-1}, t_k) - 2u(x_j, t_k) + u(x_{j+1}, t_k))}{(\Delta x)^2} + O((\Delta x)^2) \\ &\approx \frac{u_{j-1}^k - 2u_j^k + u_{j+1}^k}{(\Delta x)^2} + O((\Delta x)^2). \end{aligned} \quad (6.4)$$

Replacing the derivatives in the heat equation (6.1) by their finite difference approximations (6.3) and (6.4), we end up with the **explicit scheme**

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} + \gamma \frac{-u_{j-1}^k + 2u_j^k - u_{j+1}^k}{(\Delta x)^2} = 0 \quad (6.5)$$

for  $k \geq 0$  and  $j \in \{1, \dots, N\}$ .

Let

$$\mu := \frac{\gamma \Delta t}{(\Delta x)^2}, \quad (6.6)$$

and let

$$u^{(k)} := (u_1^k, u_2^k, \dots, u_N^k)^\top \approx (u(x_1, t_k), u(x_2, t_k), \dots, u(x_N, t_k))^\top, \quad (6.7)$$

be the vector whose entries are the numerical approximations to the solution values at time  $t_k$  at the interior nodes.

The scheme (6.5) can be written in the matrix form

$$u^{(k+1)} = Au^{(k)}, \quad (6.8)$$

where

$$A := \begin{pmatrix} 1-2\mu & \mu & & & \\ \mu & 1-2\mu & \mu & & \\ & \mu & 1-2\mu & \mu & \\ & & \ddots & \ddots & \ddots \\ & & & \mu & 1-2\mu & \mu \\ & & & & \mu & 1-2\mu \end{pmatrix}. \quad (6.9)$$

The matrix  $A$  is symmetric and tridiagonal:  $A = \text{diag}(\mu, 1-2\mu, \mu) = I_N + \mu \text{diag}(1, -2, 1)$ . Here,  $I_N$  is the  $N \times N$  identity matrix.

**LEMMA 6.1.** *Let  $M := \text{diag}(b, a, b)$  be a  $N \times N$  tridiagonal symmetric matrix. The eigenvalues of  $M$  are*

$$\lambda_n = a + 2b \cos \theta_n, \quad n = 1, \dots, N, \quad (6.10)$$

and the corresponding eigenvectors are

$$v_n = \sqrt{2} (\sin \theta_n, \sin(2\theta_n), \dots, \sin(N\theta_n))^\top, \quad (6.11)$$

where

$$\theta_n = \frac{n\pi}{N+1}.$$

Moreover,  $\{v_n\}_{n=1}^N$  form an orthonormal basis of  $\mathbb{R}^N$  with respect to the (scaled) inner product  $\frac{1}{N} \sum_{i=1}^N u_i w_i$  for  $u = (u_1, \dots, u_N)^\top$  and  $w = (w_1, \dots, w_N)^\top$  in  $\mathbb{R}^N$ .

Applying Lemma 6.1 to  $A$  defined by (6.9) shows that the eigenvectors  $v_n$  of  $A$  are independent of  $\mu$ .

**REMARK 6.2.** *By using the following approximation of the time derivative instead of (6.3):*

$$\frac{\partial u}{\partial t}(x_j, t_k) \approx \frac{u(x_j, t_k) - u(x_j, t_{k-1})}{\Delta t} + O(\Delta t) \approx \frac{u_j^k - u_j^{k-1}}{\Delta t} + O(\Delta t) \quad (6.12)$$

for  $k \geq 1$ , we obtain the **implicit scheme**

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} + \gamma \frac{-u_{j-1}^{k+1} + 2u_j^{k+1} - u_{j+1}^{k+1}}{(\Delta x)^2} = 0 \quad (6.13)$$

for  $k \geq 0$  and  $j \in \{1, \dots, N\}$ .

With the same notation as in (6.7), the implicit scheme (6.13) can be written in the matrix form

$$Bu^{(k+1)} = u^{(k)}, \quad (6.14)$$

where

$$B := \begin{pmatrix} 1+2\mu & -\mu & & & \\ -\mu & 1+2\mu & -\mu & & \\ & -\mu & 1+2\mu & -\mu & \\ & & \ddots & \ddots & \ddots \\ & & & -\mu & 1+2\mu & -\mu \\ & & & & -\mu & 1+2\mu \end{pmatrix} = I_N - \mu \text{diag}(1, -2, 1). \quad (6.15)$$

The matrix  $B$  is symmetric and tridiagonal. Moreover, since it is diagonal dominant, it is positive definite and hence, invertible.

REMARK 6.3. A convex combination of the explicit and implicit schemes (6.5) and (6.13) yields the  $\theta$ -scheme, for  $0 \leq \theta \leq 1$ ,

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} + \theta \gamma \frac{-u_{j-1}^{k+1} + 2u_j^{k+1} - u_{j+1}^{k+1}}{(\Delta x)^2} + (1 - \theta) \gamma \frac{-u_{j-1}^k + 2u_j^k - u_{j+1}^k}{(\Delta x)^2} = 0 \quad (6.16)$$

for  $k \geq 0$  and  $j \in \{1, \dots, N\}$ . If  $\theta \neq 0$ , then the scheme is implicit. For  $\theta = 1/2$ , we obtain the **Crank-Nicolson** scheme.

REMARK 6.4. If we consider the heat equation with the periodic boundary conditions

$$u(0, t) = u(1, t) \quad \text{and} \quad \frac{\partial u}{\partial x}(0, t) = \frac{\partial u}{\partial x}(1, t) \quad \text{for } t \geq 0, \quad (6.17)$$

then (6.2) should be replaced with

$$u_0^k = u_{N+1}^k \quad \text{for all } k > 0. \quad (6.18)$$

If the Neumann boundary conditions,

$$\frac{\partial u}{\partial x}(0, t) = \frac{\partial u}{\partial x}(1, t) = 0 \quad \text{for } t \geq 0, \quad (6.19)$$

are imposed then one may approximate those conditions by

$$\frac{u_1^k - u_0^k}{\Delta x} = 0 \quad \text{and} \quad \frac{u_{N+1}^k - u_N^k}{\Delta x} = 0, \quad (6.20)$$

and eliminate  $u_0^k$  and  $u_{N+1}^k$  to calculate only  $(u_j^k)_{1 \leq j \leq N}$ . Note that (6.20) is a first-order approximation. The second-order approximations

$$\frac{u_1^k - u_{-1}^k}{2\Delta x} = 0 \quad \text{and} \quad \frac{u_{N+2}^k - u_N^k}{2\Delta x} = 0, \quad (6.21)$$

can be used through the introduction of the two fictitious points  $x_{-1}$  and  $x_{N+2}$ .

REMARK 6.5. Both the explicit and implicit schemes (6.3) and (6.12) are one-step methods. Higher step methods can be designed by employing appropriate finite difference approximations to the derivatives. Examples of two-step finite difference methods are

(i) The **Richardson scheme**:

$$\frac{u_j^{k+1} - u_j^{k-1}}{2\Delta t} + \gamma \frac{-u_{j-1}^k + 2u_j^k - u_{j+1}^k}{(\Delta x)^2} = 0; \quad (6.22)$$

(ii) The **DuFort-Frankel scheme**:

$$\frac{u_j^{k+1} - u_j^{k-1}}{2\Delta t} + \gamma \frac{-u_{j-1}^k + u_j^{k+1} + u_j^{k-1} - u_{j+1}^k}{(\Delta x)^2} = 0; \quad (6.23)$$

(iii) The **Gear scheme**:

$$\frac{3u_j^{k+1} - 4u_j^k + u_j^{k-1}}{2\Delta t} + \gamma \frac{-u_{j-1}^{k+1} + 2u_j^{k+1} - u_{j+1}^{k+1}}{(\Delta x)^2} = 0. \quad (6.24)$$

**6.2.2. Consistency, stability, and convergence.** A general **finite difference method** is defined by

$$F_{\Delta t, \Delta x}(\{u_{j+n}^{k+m}\}_{m^- \leq m \leq m^+, n^- \leq n \leq n^+}) = 0, \quad (6.25)$$

where the integers  $m^\pm, n^\pm$  define the width of the stencil of the scheme. Here,  $F_{\Delta t, \Delta x}$  is such that for any  $u$  not satisfying the heat equation,  $F_{\Delta t, \Delta x}(\{u(x_{j+n}, t_{k+m})\}_{m^- \leq m \leq m^+, n^- \leq n \leq n^+})$  does not converge to 0 as  $\Delta x, \Delta t \rightarrow 0$ .

**DEFINITION 6.6 (Consistency and order).** *The finite difference scheme (6.25) is consistent with the equation  $F(u) = 0$  if, for any smooth solution  $u(x, t)$ , the **truncation error** defined by*

$$F_{\Delta t, \Delta x}(\{u(x_{j+n}, t_{k+m})\}_{m^- \leq m \leq m^+, n^- \leq n \leq n^+}) \quad (6.26)$$

*goes to zero as  $\Delta t$  and  $\Delta x$  go to zero independently. Moreover, the scheme is said to be of order  $p$  in time and order  $q$  in space if the truncation error is of the order of  $O((\Delta t)^p + (\Delta x)^q)$  as  $\Delta t$  and  $\Delta x$  go to zero.*

**THEOREM 6.7.** *The explicit scheme (6.5) is consistent with the heat equation (6.1), of order one in time and two in space. Moreover, if*

$$\frac{\gamma \Delta t}{(\Delta x)^2} = \frac{1}{6}, \quad (6.27)$$

*then it is of order two in time and four in space.*

**PROOF.** Let  $v(x, t) \in C^6$ . By the Taylor expansion of  $v$  evaluated at  $(x, t)$ ,

$$\begin{aligned} \frac{v(x, t + \Delta t) - v(x, t)}{\Delta t} + \gamma \frac{-v(x - \Delta x, t) + 2v(x, t) - v(x + \Delta x, t)}{(\Delta x)^2} &= \left( \frac{\partial v}{\partial t} - \gamma \frac{\partial^2 v}{\partial x^2} \right)(x, t) \\ &+ \frac{\Delta t}{2} \frac{\partial^2 v}{\partial t^2}(x, t) - \frac{\gamma (\Delta x)^2}{12} \frac{\partial^4 v}{\partial x^4}(x, t) + O((\Delta t)^2 + (\Delta x)^4). \end{aligned} \quad (6.28)$$

If  $v$  is a solution to (6.1), then it follows from (6.28) that the truncation error goes to zero as  $\Delta t, \Delta x \rightarrow 0$  and hence, the explicit scheme is consistent. Moreover, it is of order 1 in time and 2 in space. If we suppose that (6.27) holds, then the terms in  $\Delta t$  and  $(\Delta x)^2$  cancel out since

$$\frac{\partial^2 v}{\partial t^2} = \gamma \frac{\partial^3 v}{\partial t \partial x^2} = \gamma^2 \frac{\partial^4 v}{\partial x^4}.$$

Thus, the explicit scheme is of order 2 in time and 4 in space.  $\square$

Analogously to Theorem 6.7, the following results can be proved.

**THEOREM 6.8.** (i) *The truncation error for the  $\theta$ -scheme (6.16) is of the order of  $O(\Delta t + (\Delta x)^2)$  for any  $0 \leq \theta \neq \frac{1}{2} \leq 1$  and is of order of  $O((\Delta t)^2 + (\Delta x)^2)$  for  $\theta = \frac{1}{2}$ , i.e., for the Crank-Nicolson scheme.*

(ii) *The truncation error for the Richardson scheme (6.22) is of order of  $O((\Delta t)^2 + (\Delta x)^2)$ .*

(iii) *The truncation error for the DuFort-Frankel scheme (6.23) is of the order of  $O(\frac{\Delta t}{\Delta x} + (\Delta x)^2)$  and hence, (6.23) is not consistent.*

(iv) *The truncation error for the Gear scheme (6.24) is of order of  $O((\Delta t)^2 + (\Delta x)^2)$ .*

**DEFINITION 6.9 (Stability).** *A finite difference scheme is stable with respect to the norm  $\|\cdot\|_r$  defined by*

$$\|u^{(k)}\|_r := \left( \sum_{j=1}^N \Delta x |u_j^k|^r \right)^{\frac{1}{r}}, \quad 1 \leq r \leq +\infty, \quad (6.29)$$

*where  $u^{(k)}$  is given by (6.7), if there exists a positive constant  $C$  independent of  $\Delta t$  and  $\Delta x$  such that*

$$\|u^{(k)}\|_r \leq C \|u^{(0)}\|_r \quad \text{for all } k \geq 0. \quad (6.30)$$

Note that

$$\|u^{(k)}\|_\infty := \sup_{1 \leq j \leq N} |u_j^k|.$$

**DEFINITION 6.10 (Linear scheme).** A finite difference scheme defined by (6.25) is said to be linear if (6.25) is linear with respect to its arguments  $u_{j+n}^{k+m}$ .

If a finite difference scheme is linear, then it can be written in the form

$$u^{(k+1)} = Au^{(k)}, \quad (6.31)$$

where  $A$  is the **iteration matrix**. From (6.31), it follows that

$$u^{(k+1)} = A^{k+1}u^{(0)},$$

and therefore, the stability of (6.31) is equivalent to

$$\|A^k u^{(0)}\|_r \leq C \|u^{(0)}\|_r, \quad \text{for all } k \geq 0 \text{ and } u^{(0)} \in \mathbb{R}^N. \quad (6.32)$$

Introduce the matrix norm

$$\|M\|_r = \sup_{u \in \mathbb{R}^N, u \neq 0} \frac{\|Mu\|_r}{\|u\|_r}.$$

The stability of (6.31) with respect to  $\|\cdot\|_r$  is equivalent to

$$\|A^k\|_r \leq C, \quad \text{for all } k \geq 0.$$

**REMARK 6.11.** Note that since we require (6.30) to hold uniformly in  $\Delta x$  as  $\Delta x$  together with the fact that  $N = O(1/\Delta x)$ , the norms  $\|\cdot\|_r$  defined by (6.29) are not equivalent.

**REMARK 6.12.** The  $\|\cdot\|_2$  is associated with the weighted scalar product

$$(u, v)_2 = (\Delta x) \sum_{i=1}^N u_i v_i, \quad (6.33)$$

where  $u_i$  and  $v_i$  are the components of the vectors  $u$  and  $v$ .

**REMARK 6.13.** Consider for instance the explicit scheme (6.5). Then (6.31) holds with  $A$  being defined by (6.9). Let  $\tilde{u}_j^k = u(x_j, t_k)$  and  $\tilde{u}^{(k)} = (\tilde{u}_1^k, \dots, \tilde{u}_N^k)^\top$ . Then the truncation error introduced in (6.26) is given by

$$\epsilon^{(k)} := \frac{\tilde{u}^{(k+1)} - \tilde{u}^{(k)}}{\Delta t} + \frac{(I_N - A)}{\Delta t} \tilde{u}^{(k)} = \frac{\tilde{u}^{(k+1)} - A\tilde{u}^{(k)}}{\Delta t}.$$

Therefore,

$$\tilde{u}^{(k+1)} = A\tilde{u}^{(k)} + (\Delta t)\epsilon^{(k)}.$$

**6.2.2.1. Stability in the  $L^\infty$  norm.** Recall that the **implicit scheme** given by (6.13) is well defined since  $u^{(k+1)}$  can be obtained from  $u^{(k)}$  by inverting the definite positive matrix  $B$  given by (6.15).

The following results hold.

**THEOREM 6.14.** (i) The explicit scheme (6.5) is stable with respect to the  $L^\infty$  norm if and only if the following **Courant-Friedrichs-Lewy (CFL) condition** holds:

$$2\gamma\Delta t \leq (\Delta x)^2. \quad (6.34)$$

(ii) The implicit scheme (6.13) is unconditionally stable with respect to the  $L^\infty$  norm.

Before proving Theorem 6.14, we first introduce the discrete maximum principle.

**DEFINITION 6.15.** We say that a finite difference scheme satisfies the discrete maximum principle if for all  $k \geq 0, 1 \leq j \leq N$ ,

$$\min(0, \min_{0 \leq j \leq N+1} u_j^0) \leq u_j^k \leq \max(0, \max_{0 \leq j \leq N+1} u_j^0) \quad (6.35)$$

for any initial data  $u^{(0)}$ .

Condition (6.35) prevents unbounded oscillations of the numerical solution. It is clearly a **sufficient condition** for the stability with respect to the  $L^\infty$  norm.

Now, under the CFL condition (6.34), the explicit scheme satisfies the discrete maximum principle. This can be easily verified by induction. In fact, we can rewrite the explicit scheme as follows:

$$u_j^{k+1} = \frac{\gamma\Delta t}{(\Delta x)^2}u_{j-1}^k + (1 - 2\frac{\gamma\Delta t}{(\Delta x)^2})u_j^k + \frac{\gamma\Delta t}{(\Delta x)^2}u_{j+1}^k, \quad (6.36)$$

which shows that if the CFL condition holds, then  $u_j^{k+1}$  is a convex combination of  $u_{j-1}^k, u_j^k, u_{j+1}^k$  since all the coefficients in (6.36) are positive and their sum is one. So if  $m \leq u_j^0 \leq M$  for all  $j$ , then  $m \leq u_j^k \leq M$  for all  $j$  and all  $k \geq 0$ . Moreover, assume that the CFL condition does not hold. Then by taking  $u_j^0 = (-1)^j$ , we find that

$$u_j^k = (-1)^j (1 - 4\frac{\gamma\Delta t}{(\Delta x)^2})^k.$$

Hence, from  $1 - 4\frac{\gamma\Delta t}{(\Delta x)^2} < -1$ , it follows that  $|u_j^k| \rightarrow +\infty$  as  $k \rightarrow +\infty$ .

To prove item (ii) in Theorem 6.14, we rewrite the implicit scheme as follows:

$$(1 + 2\mu)u_j^{k+1} = u_j^k + \mu u_{j-1}^{k+1} + \mu u_{j+1}^{k+1},$$

which shows that

$$(1 + 2\mu)|u_j^{k+1}| \leq \|u^{(k)}\|_\infty + 2\mu\|u^{(k+1)}\|_\infty,$$

and hence,

$$\|u^{(k+1)}\|_\infty \leq \|u^{(k)}\|_\infty.$$

The following stability results with respect to the  $L^\infty$  norm hold.

**THEOREM 6.16.** (i) *The Crank-Nicolson scheme is stable with respect to the  $L^\infty$  norm if  $\frac{\gamma\Delta t}{(\Delta x)^2} \leq 1$ .*

(ii) *The DuFort-Frankel scheme (6.23) is stable with respect to the  $L^\infty$  norm if  $\frac{2\Delta t}{(\Delta x)^2} \leq 1$ .*

In order to prove the stability of the Crank-Nicolson scheme with respect to the  $L^\infty$  norm under the CFL condition  $\frac{\gamma\Delta t}{(\Delta x)^2} \leq 1$ , we rewrite it as follows:

$$(I_N - \frac{\mu}{2}\text{diag}(1, -2, 1))u^{(k+1)} = (I_N + \frac{\mu}{2}\text{diag}(1, -2, 1))u^{(k)}.$$

By the unconditional stability of the implicit scheme, we have

$$\|u^{(k+1)}\|_\infty \leq \|(I_N - \frac{\mu}{2}\text{diag}(1, -2, 1))u^{(k+1)}\|_\infty.$$

On the other hand, under the CFL condition  $\frac{\gamma\Delta t}{(\Delta x)^2} \leq 1$ , we have from item (i) in Theorem 6.14,

$$\|(I_N + \frac{\mu}{2}\text{diag}(1, -2, 1))u^{(k)}\|_\infty \leq C\|u^{(k)}\|_\infty.$$

Combining the above two estimates yields the desired result.

**6.2.2.2. Stability in the  $L^2$  norm.** In order to investigate the stability of a finite difference scheme for solving the heat equation with respect to the  $L^2$  norm, we consider (6.1) with the **periodic boundary conditions**

$$u(x+1, t) = u(x, t) \quad \text{for all } x \in [0, 1], \quad t \geq 0.$$

For any  $u^{(k)} = (u_j^k)_{j=0, \dots, N}$ , we associate a piecewise constant function  $u^{(k)}(x)$ , periodic with period 1, defined on  $[0, 1]$  by

$$u^{(k)}(x) := u_j^k \quad \text{for } x_{j-\frac{1}{2}} < x < x_{j+\frac{1}{2}},$$

where

$$x_{j+\frac{1}{2}} = \left(j + \frac{1}{2}\right)\Delta x, \quad j = 0, \dots, N, \quad x_{-\frac{1}{2}} = 0, \quad x_{N+1+\frac{1}{2}} = 1.$$

The Fourier series of  $u^{(k)}$  reads

$$u^{(k)}(x) = \sum_{n \in \mathbb{Z}} \hat{u}_n^{(k)} e^{2\pi i n x},$$

where

$$\hat{u}_n^{(k)} := \int_0^1 u^{(k)}(x) e^{-2\pi i n x} dx.$$

Moreover, by **Plancherel's formula**, we have

$$\int_0^1 |u^{(k)}(x)|^2 dx = \sum_{n \in \mathbb{Z}} |\hat{u}_n^{(k)}|^2. \quad (6.37)$$

Furthermore, an important property of Fourier series of periodic functions is that

$$v^{(k)}(x) = u^{(k)}(x + \Delta x) \Rightarrow \hat{v}_n^{(k)} = \hat{u}_n^{(k)} e^{2\pi i n \Delta x}.$$

With this notation, one can rewrite the explicit scheme (6.5) in the form

$$\frac{u^{(k+1)}(x) - u^{(k)}(x)}{\Delta t} + \gamma \frac{-u^{(k)}(x - \Delta x) + 2u^{(k)}(x) - u^{(k)}(x + \Delta x)}{(\Delta x)^2} = 0. \quad (6.38)$$

Applying the Fourier transform yields

$$\hat{u}_n^{(k+1)} = \left(1 - \frac{\gamma \Delta t}{(\Delta x)^2} (-e^{-2\pi i n \Delta x} + 2 - e^{2\pi i n \Delta x})\right) \hat{u}_n^{(k)},$$

or equivalently,

$$\hat{u}_n^{(k+1)} = \alpha(n) \hat{u}_n^{(k)} = \alpha(n)^{k+1} \hat{u}_n^{(0)} \quad \text{with } \alpha(n) := 1 - \frac{4\gamma \Delta t}{(\Delta x)^2} (\sin(\pi n \Delta x))^2. \quad (6.39)$$

Therefore,  $\hat{u}_n^{(k)}$  is bounded as  $k \rightarrow +\infty$  if and only if the **amplification factor**  $\alpha(n)$  satisfies

$$|\alpha(n)| \leq 1 \quad \text{for all } n \in \mathbb{Z}. \quad (6.40)$$

Assume that (6.40) holds, *i.e.*,  $2\gamma \Delta t / (\Delta x)^2 \leq 1$ . Then from (6.37), it follows that

$$\|u^{(k)}\|_2^2 = \int_0^1 |u^{(k)}(x)|^2 dx = \sum_{n \in \mathbb{Z}} |\hat{u}_n^{(k)}|^2 \leq \sum_{n \in \mathbb{Z}} |\hat{u}_n^{(0)}|^2 = \|u^{(0)}\|_2^2,$$

and therefore the scheme is stable with respect to the  $L^2$  norm.

Similarly, the implicit scheme (6.13) can be rewritten in the form

$$\frac{u^{(k+1)}(x) - u^{(k)}(x)}{\Delta t} + \gamma \frac{-u^{(k+1)}(x - \Delta x) + 2u^{(k+1)}(x) - u^{(k+1)}(x + \Delta x)}{(\Delta x)^2} = 0. \quad (6.41)$$

Again, by applying the Fourier transform, it follows that

$$\hat{u}_n^{(k+1)} = \beta(n) \hat{u}_n^{(k)} = \beta(n)^{k+1} \hat{u}_n^{(0)},$$

where

$$\beta(n) := \left(1 + \frac{4\gamma \Delta t}{(\Delta x)^2} (\sin(\pi n \Delta x))^2\right)^{-1}.$$

Since the amplification factor  $\beta(n)$  satisfies  $0 \leq \beta(n) \leq 1$ , for all  $\Delta t > 0$  and  $\Delta x > 0$ , we obtain

$$\|u^{(k)}\|_2^2 \leq \|u^{(0)}\|_2^2$$

for all  $k \geq 0$ .

- THEOREM 6.17.** (i) *The explicit scheme (6.5) is stable with respect to the  $L^2$  norm if and only if the CFL condition (6.34) holds.*  
(ii) *The implicit scheme (6.13) is unconditionally stable with respect to the  $L^2$  norm.*

Note that the stability results for (6.5) and (6.13) with respect to the  $L^2$  norm are the same as those with respect to  $L^\infty$  norm. This is however not in general true for other finite difference schemes.

The following stability results for the  $\theta$ -scheme with respect to the  $L^2$  norm hold.

**THEOREM 6.18.** *The  $\theta$ -scheme (6.16) is unconditionally stable with respect to the  $L^2$  norm if  $\frac{1}{2} \leq \theta \leq 1$  and provided the CFL condition  $2(1 - 2\theta)\gamma\Delta t \leq (\Delta x)^2$  if  $0 \leq \theta < \frac{1}{2}$ .*

The method described here is called the **von Neumann stability analysis**.

Based on Lemma 6.1, there is a more direct (but equivalent) way for verifying the stability with respect to the  $L^2$  norm for the explicit, implicit, and Crank-Nicolson schemes. Such a technique extends to the heat equation with either Dirichlet or Neumann boundary conditions. For more general schemes, one uses the von Neumann analysis of stability.

To fix ideas, consider first the explicit scheme (6.5) for solving the heat equation (6.1) with the Dirichlet boundary conditions.

We expand  $w \in \mathbb{R}^N$  in the orthonormal basis of eigenvectors  $\{v_n\}_{n=1}^N$  of  $A$  (with respect to the weighted scalar product (6.33)) given by (6.11):

$$w = \sum_{n=1}^N \hat{w}_n v_n \quad \text{with} \quad \hat{w}_n = (w, v_n)_2 = (\Delta x) \sum_{i=1}^N w_i (v_n)_i,$$

where  $w = (w_1, \dots, w_N)^\top$  and  $v_n = ((v_n)_1, \dots, (v_n)_N)^\top$ .

The **discrete Parseval identity** is

$$\|w\|_2^2 = (\Delta x) \sum_{i=1}^N (w_i)^2 = \sum_{n=1}^N |\hat{w}_n|^2. \quad (6.42)$$

Since  $u^{(k+1)} = Au^{(k)}$  with  $A = \text{diag}(\mu, 1 - 2\mu, \mu)$ , the stability in the  $L^2$  norm is related to the spectral radius  $\rho(A)$ . That is

$$\|u^{(k+1)}\|_2 \leq \|A\|_2 \|u^{(k)}\|_2,$$

and since  $A$  is symmetric with respect to  $(\cdot)_2$ ,

$$\|A\|_2 = \rho(A) = \max_{1 \leq l \leq N} |\lambda_l(A)| = \max_{1 \leq l \leq N} |1 - 2\mu + 2\mu \cos \theta_l|.$$

The uniform stability with respect to  $N$  implies that  $\mu \leq \frac{1}{2}$ .

If we consider the implicit scheme (6.13), then since

$$\rho(B^{-1}) = \max_{1 \leq l \leq N} \frac{1}{|1 + 2\mu - 2\mu \cos \theta_l|} \leq 1,$$

for any  $\mu > 0$ , the implicit is unconditionally stable with respect to the  $L^2$ -norm.

Finally, we can easily check that the Crank-Nicolson scheme can be rewritten as

$$u^{(k+1)} = \tilde{B}^{-1} \tilde{A} u^{(k)},$$

where

$$\tilde{B} = \frac{1}{2} \text{diag}(-\mu, 2 + 2\mu, -\mu) \quad \text{and} \quad \tilde{A} = \frac{1}{2} \text{diag}(\mu, 2 - 2\mu, \mu).$$

Since  $\tilde{A}$  and  $\tilde{B}$  have the same eigenvectors (by Lemma 6.1), we have

$$\rho(\tilde{B}^{-1} \tilde{A}) = \max_{1 \leq l \leq N} \frac{\lambda_l(\tilde{A})}{\lambda_l(\tilde{B})} = \max_{1 \leq l \leq N} \frac{|1 - \mu + \mu \cos \theta_l|}{|1 + \mu - \mu \cos \theta_l|}.$$

Consequently,  $\rho(\tilde{B}^{-1} \tilde{A}) \leq 1$  for all  $\mu > 0$  and therefore, the Crank-Nicolson scheme is unconditionally stable with respect to the  $L^2$  norm.



### 6.2.3. Convergence.

**THEOREM 6.19 (Lax theorem).** *Let  $u$  be a smooth solution of the heat equation (6.1). Suppose that the finite difference scheme for computing the numerical solution  $u_j^k$  is linear, consistent, and stable with respect to the norm  $\|\cdot\|_r$ . Let  $e_j^k := u_j^k - u(x_j, t_k)$  and  $e^{(k)} = (e_1^k, e_2^k, \dots, e_N^k)^\top$ . Assume that  $u_j^0 = u_0(x_j)$ . Then,*

$$\lim_{\Delta t, \Delta x \rightarrow 0} \left( \sup_{t_k \leq T} \|e^{(k)}\|_r \right) = 0 \quad \text{for all } T > 0.$$

Moreover, if the scheme is of order  $p$  in time and  $q$  in space, then there exists a constant  $C_T > 0$  such that

$$\sup_{t_k \leq T} \|e^{(k)}\|_r \leq C_T ((\Delta t)^p + (\Delta x)^q).$$

**PROOF.** Let  $u^{(k+1)} = Au^{(k)}$ , where  $A$  is the iteration matrix, and let  $\tilde{u}_j^k = u(x_j, t_k)$ . Since the scheme is consistent, there exists  $\epsilon^{(k)}$  such that

$$\tilde{u}^{(k+1)} = A\tilde{u}^{(k)} + (\Delta t)\epsilon^{(k)} \quad \text{and} \quad \lim_{\Delta t, \Delta x \rightarrow 0} \|\epsilon^{(k)}\|_r = 0, \quad (6.43)$$

uniformly in  $k$ . If the scheme is of order  $p$  in time and  $q$  in space, then

$$\|\epsilon^{(k)}\|_r \leq C((\Delta t)^p + (\Delta x)^q);$$

see Remark 6.13.

By subtracting (6.43) from (6.31), we obtain

$$e^{(k+1)} = Ae^{(k)} - \Delta t\epsilon^{(k)}, \quad (6.44)$$

and therefore, by induction,

$$e^{(k)} = A^k e^{(0)} - \Delta t \sum_{l=1}^k A^{k-l} \epsilon^{(l-1)}. \quad (6.45)$$

The stability of the scheme yields

$$\|A^k\|_r \leq C'$$

for some positive constant  $C'$ . Therefore, since  $e^{(0)} = 0$ , (6.45) yields

$$\|e^{(k)}\|_r \leq (\Delta t)kCC'((\Delta t)^p + (\Delta x)^q) \leq TCC'((\Delta t)^p + (\Delta x)^q). \quad (6.46)$$

The proof is then complete.  $\square$

**6.2.4. Multi-step schemes.** Assume that  $u^{(k+1)}$  depends linearly on  $u^{(k)}$  and  $u^{(k-1)}$ , as for example in (6.22), (6.23), and (6.24). Then, we set

$$U^{(k)} = \begin{pmatrix} u^{(k)} \\ u^{(k-1)} \end{pmatrix}.$$

There exist then two  $N \times N$  matrices  $A_1$  and  $A_2$  such that

$$U^{(k+1)} = AU^{(k)} = \begin{pmatrix} A_1 & A_2 \\ I_N & 0 \end{pmatrix} U^{(k)},$$

where  $A$  is a  $2N \times 2N$  matrix and  $I_N$  is the  $N \times N$  identity matrix. As before, we obtain that  $U^{(k)} = A^k U^{(1)}$  and the stability of the scheme is equivalent to

$$\|A^k\|_r \leq C \quad \text{for all } k \geq 0.$$

For  $r = 2$  and  $A$  normal, the  $L^2$  stability condition reduces to the **von Neumann stability condition**

$$\rho(A) \leq 1 \quad (6.47)$$

with  $\rho(A)$  being the spectral radius. In general, we have  $\|A\|_2 \geq \rho(A)$  and therefore, the von Neumann stability condition is only a **necessary condition**.

LEMMA 6.20. *The Richardson scheme (6.22) is unstable with respect to the  $L^2$  norm.*

PROOF. With the same notation as in Subsection 6.2.2.2, the Richardson scheme (6.22) reads

$$\frac{u^{(k+1)}(x) - u^{(k-1)}(x)}{2\Delta t} + \gamma \frac{-u^{(k)}(x - \Delta x) + 2u^{(k)}(x) - u^{(k)}(x)}{(\Delta x)^2} = 0. \quad (6.48)$$

Then, applying the Fourier transform yields

$$\hat{u}_n^{(k+1)} + \frac{8\gamma\Delta t}{(\Delta x)^2} (\sin(\pi n\Delta x))^2 \hat{u}_n^{(k)} - \hat{u}_n^{(k-1)} = 0, \quad (6.49)$$

or in other words,

$$\hat{U}_n^{(k+1)} = \begin{pmatrix} \hat{u}_n^{(k+1)} \\ \hat{u}_n^{(k)} \end{pmatrix} = \begin{pmatrix} -\frac{8\gamma\Delta t}{(\Delta x)^2} (\sin(\pi n\Delta x))^2 & 1 \\ 1 & 0 \end{pmatrix} \hat{U}_n^{(k)} = A(n)\hat{U}_n^{(k)}. \quad (6.50)$$

Consequently,

$$\hat{U}_n^{(k+1)} = A(n)^k \hat{U}_n^{(1)}.$$

In (6.50),  $A(n)$  is a  $2 \times 2$  (amplification) matrix, while for a one-step method, it is a scalar; see (6.39).

For  $n \in \mathbb{Z}$ , the vector  $\hat{U}_n^{(k)}$  is bounded iff the amplification matrix  $A(n)$  satisfies

$$\|A(n)^k\|_2 \leq C \quad \text{for all } k \geq 1 \quad (6.51)$$

for some constant  $C$  independent of  $k$  and  $n$ . Since  $A(n)$  is real symmetric,  $\|A(n)\|_2 = \rho(A(n))$  and  $\|A(n)^k\|_2 = \|A(n)\|_2^k$ . Here,  $\rho(M)$  is the spectral radius of  $M$ . Therefore, (6.51) is satisfied iff  $\rho(A(n)) \leq 1$ . The eigenvalues of  $A(n)$  are roots of the second order polynomial

$$\lambda^2 + \frac{8\gamma\Delta t}{(\Delta x)^2} (\sin(\pi n\Delta x))^2 \lambda - 1 = 0,$$

which admits two distinct real roots with product equals to  $-1$ . Therefore,  $A(n)$  has an eigenvalue with modulus strictly larger than 1. Consequently, the Richardson scheme is unstable with respect to the  $L^2$  norm.  $\square$

For the DuFort-Frankel and Gear schemes, the following convergence results hold.

THEOREM 6.21. *We have*

- (i) *The DuFort-Frankel (6.23) is stable and hence convergent with respect to the  $L^2$  norm, provided that  $\Delta t/(\Delta x)^2$  stays bounded as  $\Delta t$  and  $\Delta x$  go to 0.*
- (ii) *The Gear scheme (6.24) is unconditionally stable and hence convergent with respect to the  $L^2$  norm.*

### 6.3. Numerical algorithms for the wave equation

We first consider the one-way wave equation given by

$$\begin{cases} \frac{\partial u}{\partial t} = c \frac{\partial u}{\partial x}, \\ u(x, 0) = u_0(x), \end{cases} \quad (6.52)$$

where  $c > 0$  is the wave speed. The solution of (6.52) is given by  $u(x, t) = u_0(x + ct)$ . Note that if a smooth function  $u$  satisfies the first equation in (6.52), then

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}. \quad (6.53)$$

There are three finite difference approximations of the solution:

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} = \begin{cases} c \frac{u_{j+1}^k - u_j^k}{\Delta x} & \text{upwind scheme,} \\ c \frac{u_j^k - u_{j-1}^k}{\Delta x} & \text{downwind scheme,} \\ c \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} & \text{centered scheme.} \end{cases}$$

Using the Taylor expansions of a smooth solution  $u$  to (6.52),

$$\begin{aligned} \frac{u(x, t + \Delta t) - u(x, t)}{\Delta t} &= \frac{\partial u}{\partial t}(x, t) + \frac{\Delta t}{2} \frac{\partial^2 u}{\partial t^2}(x, t) + O((\Delta t)^2), \\ \frac{u(x + \Delta x, t) - u(x, t)}{\Delta x} &= \frac{\partial u}{\partial x}(x, t) + \frac{\Delta x}{2} \frac{\partial^2 u}{\partial x^2}(x, t) + O((\Delta x)^2), \end{aligned}$$

and

$$\frac{u(x + \Delta x, t) - u(x - \Delta x, t)}{\Delta x} = \frac{\partial u}{\partial x}(x, t) + O((\Delta x)^2),$$

we obtain that the truncation error in the upwind scheme is  $O(\Delta t + \Delta x)$ . Analogously, the truncation error in the downwind scheme is  $O(\Delta t + \Delta x)$ , while the one in the centered scheme is  $O(\Delta t + (\Delta x)^2)$ . Note that if

$$c = \frac{\Delta x}{\Delta t},$$

then the truncation error in the upwind scheme is  $O((\Delta t)^2 + (\Delta x)^2)$ . This directly follows from (6.53).

Now, regarding the stability of these schemes, one can easily see that the upwind scheme is stable with respect to the  $L^2$  norm provided that the following CFL condition holds:

$$\frac{c\Delta t}{\Delta x} \leq 1,$$

while both the downwind and the centered schemes are unstable. In fact, with the notation of Subsection 6.2.2.2, we have for the centered scheme

$$\hat{u}_n^{(k+1)} = \left(1 + i \frac{c\Delta t}{\Delta x} \sin(2\pi n\Delta x)\right) \hat{u}_n^{(k)}.$$

One can write the following implicit version of the centered scheme which is consistent, of order one in time and two in space and is unconditionally stable with respect to the  $L^2$  norm:

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} - c \frac{u_{j+1}^{k+1} - u_{j-1}^{k+1}}{2\Delta x} = 0. \quad (6.54)$$

If we want to stay within the class of explicit centered schemes, we can use the **Lax-Friedrichs scheme**

$$\frac{2u_j^{k+1} - u_{j+1}^k - u_{j-1}^k}{2\Delta t} - c \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} = 0, \quad (6.55)$$

which is consistent if  $\Delta t/\Delta x$  is constant as  $\Delta t, \Delta x \rightarrow 0$ , stable in  $L^2$  under the CFL condition

$$c\Delta t \leq \Delta x, \quad (6.56)$$

and of order 1 in time and space. It is worth emphasizing that this scheme is not consistent in the sense of Definition 6.6, but is only conditionally consistent. In fact, the truncation error is given by

$$-\frac{(\Delta x)^2}{2\Delta t} \left(1 - \frac{(c\Delta t)^2}{(\Delta x)^2}\right) \frac{\partial^2 u}{\partial x^2}(x_j, t_k) + O((\Delta x)^2 + \frac{(\Delta x)^4}{\Delta t}). \quad (6.57)$$

To check its  $L^2$  stability properties under the CFL condition (6.56), we use Fourier analysis to obtain

$$\hat{u}_n^{(k+1)} = \left(\cos(2\pi n\Delta x) + i \frac{c\Delta t}{\Delta x} \sin(2\pi n\Delta x)\right) \hat{u}_n^{(k)}.$$

A centered, explicit scheme of higher order than the Lax-Friedrichs scheme is the **Lax-Wendroff scheme**

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} - c \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} - \frac{c^2 \Delta t}{2} \frac{u_{j-1}^k - 2u_j^k + u_{j+1}^k}{(\Delta x)^2} = 0, \quad (6.58)$$

which is consistent, stable in  $L^2$  under the CFL condition (6.56), and is of order 2 in time and space.

A general way to fix the stability issue for the centered scheme is to replace the centered scheme with

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} = c \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} + \theta \frac{u_{j+1}^k - 2u_j^k + u_{j-1}^k}{(\Delta x)^2}, \quad (6.59)$$

where  $\theta > 0$ , or equivalently, with

$$\frac{u_j^{k+1} - (\frac{\lambda}{2}u_{j+1}^k + (1-\lambda)u_j^k + \frac{\lambda}{2}u_{j-1}^k)}{\Delta t} = c \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x}.$$

Here,  $\lambda$  is defined by

$$\lambda = \frac{2\Delta t}{(\Delta x)^2} \theta.$$

For  $\theta = (\Delta x)^2/(2\Delta t)$  (i.e.,  $\lambda = 1$ ), (6.59) reduces to the Lax-Friedrichs scheme (6.55) while for  $\theta = c^2 \Delta t/2$ , (6.59) reduces to the Lax-Wendroff scheme (6.58). Moreover, the scheme (6.59) solves (approximately, up to order two in time and space) the equation

$$\frac{\partial u}{\partial t} = c \frac{\partial u}{\partial x} + \left(\theta - \frac{c^2 \Delta t}{2}\right) \frac{\partial^2 u}{\partial x^2}.$$

Here, we have used the fact that

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}.$$

Next, consider the wave equation (with periodic boundary conditions)

$$\begin{cases} \frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}, & 0 < x < 1, \quad t \geq 0, \\ u(x+1, t) = u(x, t), & 0 < x < 1, \quad t \geq 0, \\ u(x, 0) = u_0(x), & 0 < x < 1, \\ \frac{\partial u}{\partial t}(x, 0) = u_1(x), & 0 < x < 1. \end{cases} \quad (6.60)$$

To insure that the solution stays bounded in  $t$ , we suppose that

$$\int_0^1 u_1(x) dx = 0. \quad (6.61)$$

In fact, if  $u_0 = 0$  and  $u_1$  is equal to some constant  $C$ , then  $u(t, x) = Ct$ . To eliminate this effect we impose the normalization condition (6.61).

Similar to the numerical schemes for the heat equation, we can use differentiation formulas to arrive at a numerical scheme for the wave equation (6.60). Since both time and space derivatives are of second order, we use centered differences to approximate them. Analogously to (6.4), we have

$$\begin{aligned} \frac{\partial^2 u}{\partial t^2}(x_j, t_k) &\approx \frac{u(x_j, t_{k-1}) - 2u(x_j, t_k) + u(x_j, t_{k+1}))}{(\Delta t)^2} + O((\Delta t)^2) \\ &\approx \frac{u_j^{k-1} - 2u_j^k + u_j^{k+1}}{(\Delta t)^2} + O((\Delta t)^2). \end{aligned} \quad (6.62)$$

Then up to an error of order  $O((\Delta x)^2 + (\Delta t)^2)$  the solution to the wave equation (6.60) can be approximated by the following explicit finite difference scheme:

$$\frac{u_j^{k+1} - 2u_j^k + u_j^{k-1}}{(\Delta t)^2} = c^2 \frac{u_{j+1}^k - 2u_j^k + u_{j-1}^k}{(\Delta x)^2}. \quad (6.63)$$

One can prove that (6.63) is stable in the  $L^2$  norm provided that  $c(\Delta t)/(\Delta x) \leq 1$ .

Another standard finite difference scheme for solving (6.60) is the  **$\theta$ -centered scheme**

$$\left\{ \begin{array}{l} \frac{u_j^{k+1} - 2u_j^k + u_j^{k-1}}{(\Delta t)^2} + \theta c^2 \frac{-u_{j-1}^{k+1} + 2u_j^{k+1} - u_{j+1}^{k+1}}{(\Delta x)^2} \\ + (1 - 2\theta)c^2 \frac{-u_{j-1}^k + 2u_j^k - u_{j+1}^k}{(\Delta x)^2} + \theta c^2 \frac{-u_{j-1}^{k-1} + 2u_j^{k-1} - u_{j+1}^{k-1}}{(\Delta x)^2} = 0, \end{array} \right. \quad (6.64)$$

where  $0 \leq \theta \leq 1/2$ .

If  $\theta = 0$ , then the scheme is nothing else than the explicit scheme (6.63), while it is implicit if  $\theta \neq 0$ .

The initial conditions can be expressed by

$$u_j^0 = u_0(x_j) \quad \text{and} \quad \frac{u_j^1 - u_j^0}{\Delta t} = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u_1(x) dx,$$

which shows that (6.61) is satisfied by the numerical solution.

**THEOREM 6.22.** *If  $1/4 \leq \theta \leq 1/2$ , then the  $\theta$ -centered scheme (6.64) is unconditionally stable with respect to the  $L^2$  norm. If  $0 \leq \theta < 1/4$ , (6.64) is stable provided that the CFL condition*

$$\frac{c\Delta t}{\Delta x} < \sqrt{\frac{1}{1 - 4\theta}} \quad (6.65)$$

*holds and is unstable if  $c\Delta t/\Delta x > 1/\sqrt{1 - 4\theta}$ .*

**PROOF.** By using Fourier analysis, we obtain

$$\hat{u}_n^{(k+1)} - 2\hat{u}_n^{(k)} + \alpha(n)(\theta\hat{u}_n^{(k+1)} + (1 - 2\theta)\hat{u}_n^{(k)} + \theta\hat{u}_n^{(k-1)}) + \hat{u}_n^{(k-1)} = 0,$$

where

$$\alpha(n) = 4c^2 \left(\frac{\Delta t}{\Delta x}\right)^2 \sin^2(\pi n \Delta x).$$

Therefore,

$$\hat{U}_n^{(k+1)} = \begin{pmatrix} \hat{u}_n^{(k+1)} \\ \hat{u}_n^{(k)} \end{pmatrix} = \begin{pmatrix} \frac{2 - (1 - 2\theta)\alpha(n)}{1 + \theta\alpha(n)} & -1 \\ 1 & 0 \end{pmatrix} \hat{U}_n^{(k)} = A(n)\hat{U}_n^{(k)}.$$

The eigenvalues of  $A(n)$  are the roots of

$$\lambda^2 - \frac{2 - (1 - 2\theta)\alpha(n)}{1 + \theta\alpha(n)}\lambda + 1 = 0. \quad (6.66)$$

The discriminant of this second order equation is

$$\Delta = -\frac{\alpha(n)(4 - (1 - 4\theta)\alpha(n))}{(1 + \theta\alpha(n))^2}.$$

The study of the stability properties of (6.64) is quite delicate since the amplification matrix  $A(n)$  is not normal (i.e., it does not commute with its adjoint  $\overline{A}^\top$ ). Recall that for a non normal matrix, its  $L^2$  norm does not in general coincide with its spectral radius  $\rho(A)$ . Then, let us only check here the necessary condition  $\rho(A(n)) \leq 1$ . If  $c\Delta t/\Delta x > 1/\sqrt{1 - 4\theta}$ , choosing  $n$  such that  $\sin^2(\pi n \Delta x) \approx 1$  yields  $\Delta > 0$  and thus, there are two distinct real solutions to (6.66) with product 1. Hence,  $\rho(A(n)) > 1$  and the scheme is unstable. If  $c\Delta t/\Delta x < 1/\sqrt{1 - 4\theta}$ , then  $\Delta \leq 0$  for all

$n$  and the two roots are complex with modulus 1. Therefore,  $\rho(A(n)) = 1$  and the von Neumann stability condition (6.47) is satisfied.  $\square$

An important property of the wave equation is the conservation of energy.

LEMMA 6.23. *Suppose that  $u$  satisfies the the wave equation on  $(0, 1) \times (0, \infty)$  together with the boundary conditions*

$$u(0, t) = u(1, t) = 0,$$

and the initial conditions

$$u(x, 0) = g(x), \quad \frac{\partial u}{\partial t}(x, 0) = h(x).$$

Then, the energy

$$E(t) := \int_0^1 \left(\frac{\partial u}{\partial t}\right)^2 dx + \int_0^1 \left(\frac{\partial u}{\partial x}\right)^2 dx \quad (6.67)$$

is constant over time, i.e.,  $E(t) = E(0)$  for all  $t \geq 0$ .

PROOF. By multiplying the wave equation by  $\partial u / \partial t$  and integrating in  $x$  over  $(0, 1)$ , we obtain that  $dE(t)/dt = 0$ .  $\square$

In view of Lemma 6.23, the energy  $E(t)$  given by (6.67) is conserved. It is then desirable that a discrete version of the energy is conserved at the discrete level. For the  $\theta$ -scheme designed to solve the wave equation with periodic boundary conditions, we introduce the discrete energy

$$E^{k+1} = \Delta x \left[ \sum_{j=0}^N \left( \frac{u_j^{k+1} - u_j^k}{\Delta t} \right)^2 + a_{\Delta x} (u^{(k+1)}, u^{(k)}) + \theta a_{\Delta x} (u^{(k+1)} - u^{(k)}, u^{(k+1)} - u^{(k)}) \right]$$

with  $u^{(k)} = (u_0^k, \dots, u_N^k)^\top$  and

$$a_{\Delta x}(u, v) = c^2 \sum_{j=0}^N \left( \frac{u_{j+1} - u_j}{\Delta x} \right) \left( \frac{v_{j+1} - v_j}{\Delta x} \right)$$

with

$$u = (u_0, \dots, u_N)^\top \text{ and } v = (v_0, \dots, v_N)^\top \text{ and } u_{N+1} = u_0, v_{N+1} = v_0.$$

$E^{k+1}$  approximates  $E(t_{k+1})$  up to  $O(\Delta x + \Delta t)$ . We can show that  $E^k = E^0$  for all  $k \geq 0$  and therefore, the  $\theta$ -scheme preserves the conservation of energy property. The proof is based on the following **discrete integration by parts formula**:

$$\sum_{j=0}^N (-u_{j+1} + 2u_j - u_{j-1})v_j = \sum_{j=0}^N (u_{j+1} - u_j)(v_{j+1} - v_j) \quad \text{with } u_{-1} = u_N. \quad (6.68)$$

Another way to derive finite difference schemes for the wave equation is to rewrite (6.60) as a system of first order equations (by choosing  $v = \partial u / \partial t$  and  $w = \partial u / \partial x$ )

$$\begin{cases} \frac{\partial}{\partial t} \begin{pmatrix} v \\ w \end{pmatrix} = \begin{pmatrix} 0 & c \\ c & 0 \end{pmatrix} \frac{\partial}{\partial x} \begin{pmatrix} v \\ w \end{pmatrix}, & 0 < x < 1, \quad t \geq 0, \\ v(x+1, t) = v(x, t), \quad w(x+1, t) = w(x, t), & 0 < x < 1, \quad t \geq 0, \\ w(x, 0) = \frac{\partial u_0}{\partial x}(x), & 0 < x < 1, \\ v(x, 0) = u_1(x), & 0 < x < 1. \end{cases} \quad (6.69)$$

Hence, we can use the algorithms developed for the one-way wave equation in order to solve (6.60). For instance, the following scheme for solving (6.60) is of Lax-Friedrichs type:

$$\frac{1}{2\Delta t} \begin{pmatrix} 2v_j^{k+1} - v_{j+1}^k - v_{j-1}^k \\ 2w_j^{k+1} - w_{j+1}^k - w_{j-1}^k \end{pmatrix} - \frac{c}{2\Delta x} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} v_{j+1}^k - v_{j-1}^k \\ w_{j+1}^k - w_{j-1}^k \end{pmatrix} = 0, \quad (6.70)$$

while

$$\frac{1}{\Delta t} \begin{pmatrix} v_j^{k+1} - v_j^k \\ w_j^{k+1} - w_j^k \end{pmatrix} - \frac{c}{2\Delta x} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} v_{j+1}^k - v_{j-1}^k \\ w_{j+1}^k - w_{j-1}^k \end{pmatrix} - \frac{c^2 \Delta t}{2(\Delta x)^2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}^2 \begin{pmatrix} v_{j-1}^k - 2v_j^k + v_{j+1}^k \\ w_{j-1}^k - 2w_j^k + w_{j+1}^k \end{pmatrix} = 0 \quad (6.71)$$

is of Lax-Wendroff type.

PROBLEM 6.24. Consider the **advection equation**

$$\begin{cases} \frac{\partial u}{\partial t} = -v \frac{\partial u}{\partial x}, & 0 < x < 1, \quad t \geq 0, \\ u(t, x+1) = u(t, x), & 0 < x < 1, \quad t \geq 0, \\ u(0, x) = u_0(x), & 0 < x < 1, \end{cases} \quad (6.72)$$

where  $v > 0$ .

(i) Prove that the **centered explicit scheme**

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} + v \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} = 0$$

is unconditionally unstable in  $L^2$ .

(ii) Prove that the **Lax-Friedrichs scheme**

$$\frac{2u_j^{k+1} - u_{j+1}^k - u_{j-1}^k}{2\Delta t} + v \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} = 0$$

is consistent if  $\Delta t/\Delta x$  is constant as  $\Delta t, \Delta x \rightarrow 0$ , stable in  $L^2$  under the CFL condition

$$v\Delta t \leq \Delta x, \quad (6.73)$$

and of order 1 in time and space.

(iii) Prove that the **Lax-Wendroff scheme**

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} + v \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} - \frac{v^2 \Delta t}{2} \frac{u_{j-1}^k - 2u_j^k + u_{j+1}^k}{(\Delta x)^2} = 0$$

is consistent, stable in  $L^2$  under the CFL condition (6.73), and is of order 2 in time and space.

(iv) Prove that the **leapfrog scheme**

$$\frac{u_j^{k+1} - u_j^{k-1}}{2\Delta t} + v \frac{u_{j+1}^k - u_{j-1}^k}{2\Delta x} = 0$$

is consistent and is stable in  $L^2$  under the CFL condition

$$v\Delta t \leq M\Delta x, \quad (6.74)$$

with  $M < 1$ .





## Bibliography

- [1] E. Hairer, C. Lubich, and G. Wanner, *Geometric numerical integration. Structure-preserving algorithms for ordinary differential equations*. Second edition. Springer Series in Computational Mathematics, 31. Springer-Verlag, Berlin, 2006.
- [2] E. Süli and D.F. Mayers, *An introduction to numerical analysis*. Cambridge University Press, Cambridge, 2003.
- [3] G. Teschl, *Ordinary differential equations and dynamical systems*. Graduate Studies in Mathematics, 140. American Mathematical Society, Providence, RI, 2012.