

## The Calculus of Variations

The variational principles of mechanics are firmly rooted in the soil of that great century of Liberalism which starts with Descartes and ends with the French Revolution and which has witnessed the lives of Leibniz, Spinoza, Goethe, and Johann Sebastian Bach. It is the only period of cosmic thinking in the entire history of Europe since the time of the Greeks.<sup>1</sup>

The calculus of variations studies the extreme and critical points of functions. It has its roots in many areas, from geometry to optimization to mechanics, and it has grown so large that it is difficult to describe with any sort of completeness.

Perhaps the most basic problem in the calculus of variations is this: given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that is bounded from below, find a point  $\bar{x} \in \mathbb{R}^n$  (if one exists) such that

$$f(\bar{x}) = \inf_{x \in \mathbb{R}^n} f(x).$$

There are two main approaches to this problem. One is the ‘direct method,’ in which we take a sequence of points such that the sequence of values of  $f$  converges to the infimum of  $f$ , and then try to showing that the sequence, or a subsequence of it, converges to a minimizer. Typically, this requires some sort of compactness to show that there is a convergent subsequence of minimizers, and some sort of lower semi-continuity of the function to show that the limit is a minimizer.

The other approach is the ‘indirect method,’ in which we use the fact that any interior point where  $f$  is differentiable and attains a minimum is a critical, or stationary, point of  $f$ , meaning that the derivative of  $f$  is zero. We then examine the critical points of  $f$ , together with any boundary points and points where  $f$  is not differentiable, for a minimum.

Here, we will focus on the indirect method for functionals, that is, scalar-valued functions of functions. In particular, we will derive differential equations, called the Euler-Lagrange equations, that are satisfied by the critical points of certain functionals, and study some of the associated variational problems.

We will begin by explaining how the calculus of variations provides a formulation of one of the most basic systems in classical mechanics, a point particle moving in a conservative force field. See Arnold [6] for an extensive account of classical mechanics.

---

<sup>1</sup>Cornelius Lanczos, *The Variational Principles of Mechanics*.

### 1. Motion of a particle in a conservative force field

Consider a particle of constant mass  $m$  moving in  $n$ -space dimensions in a spatially-dependent force field  $\vec{F}(\vec{x})$ . The force field is said to be conservative if

$$\vec{F}(\vec{x}) = -\nabla V(\vec{x})$$

for a smooth potential function  $V : \mathbb{R}^n \rightarrow \mathbb{R}$ , where  $\nabla$  denotes the gradient with respect to  $\vec{x}$ . Equivalently, the force field is conservative if the work done by  $\vec{F}$  on the particle as it moves from  $\vec{x}_0$  to  $\vec{x}_1$ ,

$$\int_{\Gamma(\vec{x}_0, \vec{x}_1)} \vec{F} \cdot d\vec{x},$$

is independent of the path  $\Gamma(\vec{x}_0, \vec{x}_1)$  between the two endpoints.

Abusing notation, we denote the position of the particle at time  $a \leq t \leq b$  by  $\vec{x}(t)$ . We refer to a function  $\vec{x} : [a, b] \rightarrow \mathbb{R}^n$  as a particle trajectory. Then, according to Newton's second law, a trajectory satisfies

$$(3.1) \quad m\ddot{\vec{x}} = -\nabla V(\vec{x})$$

where a dot denotes the derivative with respect to  $t$ .

Taking the scalar product of (3.1) with respect to  $\dot{\vec{x}}$ , and rewriting the result, we find that

$$\frac{d}{dt} \left\{ \frac{1}{2} m |\dot{\vec{x}}|^2 + V(\vec{x}) \right\} = 0.$$

Thus, the total energy of the particle

$$E = T(\dot{\vec{x}}) + V(\vec{x}),$$

where  $V(\vec{x})$  is the potential energy and

$$T(\vec{v}) = \frac{1}{2} m |\vec{v}|^2$$

is the kinetic energy, is constant in time.

**Example 3.1.** The position  $x(t) : [a, b] \rightarrow \mathbb{R}$  of a one-dimensional oscillator moving in a potential  $V : \mathbb{R} \rightarrow \mathbb{R}$  satisfies the ODE

$$m\ddot{x} + V'(x) = 0$$

where the prime denotes a derivative with respect to  $x$ . The solutions lie on the curves in the  $(x, \dot{x})$ -phase plane given by

$$\frac{1}{2} m \dot{x}^2 + V(x) = E.$$

The equilibrium solutions are the critical points of the potential  $V$ . Local minima of  $V$  correspond to stable equilibria, while other critical points correspond to unstable equilibria. For example, the quadratic potential  $V(x) = \frac{1}{2} kx^2$  gives the linear simple harmonic oscillator,  $\ddot{x} + \omega^2 x = 0$ , with frequency  $\omega = \sqrt{k/m}$ . Its solution curves in the phase plane are ellipses, and the origin is a stable equilibrium.

**Example 3.2.** The position  $\vec{x} : [a, b] \rightarrow \mathbb{R}^3$  of a mass  $m$  moving in three space dimensions that is acted on by an inverse-square gravitational force of a fixed mass  $M$  at the origin satisfies

$$\ddot{\vec{x}} = -GM \frac{\vec{x}}{|\vec{x}|^3},$$

where  $G$  is the gravitational constant. The solutions are conic sections with the origin as a focus, as one can show by writing the equations in terms of polar coordinates in the plane of the particle motion, and integrating the resulting ODEs.

**Example 3.3.** Consider  $n$  particles of mass  $m_i$  and positions  $\vec{x}_i(t)$ , where  $i = 1, 2, \dots, n$ , that interact in three space dimensions through an inverse-square gravitational force. The equations of motion,

$$\ddot{\vec{x}}_i = -G \sum_{j=1}^n m_j \frac{\vec{x}_i - \vec{x}_j}{|\vec{x}_i - \vec{x}_j|^3} \quad \text{for } 1 \leq i \leq n,$$

are a system of  $3n$  nonlinear, second-order ODEs. The system is completely integrable for  $n = 2$ , when it can be reduced to the Kepler problem, but it is non-integrable for  $n \geq 3$ , and extremely difficult to analyze. One of the main results is KAM theory, named after Kolmogorov, Arnold and Moser, on the persistence of invariant tori for nonintegrable perturbations of integrable systems [6].

**Example 3.4.** The configuration of a particle may be described by a point in some other manifold than  $\mathbb{R}^n$ . For example, consider a pendulum of length  $\ell$  and mass  $m$  in a gravitational field with acceleration  $g$ . We may describe its configuration by an angle  $\theta \in \mathbb{T}$  where  $\mathbb{T} = \mathbb{R}/(2\pi\mathbb{Z})$  is the one-dimensional torus (or, equivalently, the circle  $\mathbb{S}^1$ ). The corresponding equation of motion is the pendulum equation

$$\ell \ddot{\theta} + g \sin \theta = 0.$$

### 1.1. The principle of stationary action

To give a variational formulation of (3.1), we define a function

$$L : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R},$$

called the *Lagrangian*, by

$$(3.2) \quad L(\vec{x}, \vec{v}) = T(\vec{v}) - V(\vec{x}).$$

Thus,  $L(\vec{x}, \vec{v})$  is the *difference* between the kinetic and potential energies of the particle, expressed as a function of position  $\vec{x}$  and velocity  $\vec{v}$ .

If  $\vec{x} : [a, b] \rightarrow \mathbb{R}^n$  is a trajectory, we define the *action* of  $\vec{x}(t)$  on  $[a, b]$  by

$$(3.3) \quad \mathcal{S}(\vec{x}) = \int_a^b L(\vec{x}(t), \dot{\vec{x}}(t)) dt.$$

Thus, the action  $\mathcal{S}$  is a real-valued function defined on a space of trajectories  $\{\vec{x} : [a, b] \rightarrow \mathbb{R}^n\}$ . A scalar-valued function of functions, such as the action, is often called a functional.

The *principle of stationary action* (also called *Hamilton's principle* or, somewhat incorrectly, the *principle of least action*) states that, for fixed initial and final positions  $\vec{x}(a)$  and  $\vec{x}(b)$ , the trajectory of the particle  $\vec{x}(t)$  is a stationary point of the action.

To explain what this means in more detail, suppose that  $\vec{h} : [a, b] \rightarrow \mathbb{R}^n$  is a trajectory with  $\vec{h}(a) = \vec{h}(b) = 0$ . The directional (or Gâteaux) derivative of  $\mathcal{S}$  at  $\vec{x}(t)$  in the direction  $\vec{h}(t)$  is defined by

$$(3.4) \quad d\mathcal{S}(\vec{x}) \vec{h} = \left. \frac{d}{d\varepsilon} \mathcal{S}(\vec{x} + \varepsilon \vec{h}) \right|_{\varepsilon=0}.$$

The (Fréchet) derivative of  $\mathcal{S}$  at  $\vec{x}(t)$  is the linear functional  $d\mathcal{S}(\vec{x})$  that maps  $\vec{h}(t)$  to the directional derivative of  $\mathcal{S}$  at  $\vec{x}(t)$  in the direction  $\vec{h}(t)$ .

**Remark 3.5.** Simple examples show that, even for functions  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , the existence of directional derivatives at a point does not guarantee the existence of a Fréchet derivative that provides a local linear approximation of  $f$ . In fact, it does not even guarantee the continuity of the function; for example, consider

$$f(x, y) = \frac{xy^2}{x^2 + y^4} \quad \text{if } (x, y) \neq (0, 0)$$

with  $f(0, 0) = 0$ . For sufficiently smooth functions, however, such as the action functional we consider here, the existence of directional derivatives does imply the existence of the derivative, and the Gâteaux and Fréchet derivatives agree, so we do not need to worry about the distinction.

A trajectory  $\vec{x}(t)$  is a *stationary point* of  $\mathcal{S}$  if it is a critical point, meaning that  $d\mathcal{S}(\vec{x}) = 0$ . Explicitly, this means that

$$\left. \frac{d}{d\varepsilon} \mathcal{S}(\vec{x} + \varepsilon \vec{h}) \right|_{\varepsilon=0} = 0$$

for every smooth function  $\vec{h} : [a, b] \rightarrow \mathbb{R}^n$  that vanishes at  $t = a, b$ . Thus, small variations in the trajectory of the order  $\varepsilon$  that keep its endpoints fixed, lead to variations in the action of the order  $\varepsilon^2$ .

**Remark 3.6.** Remarkably, the motion of any conservative, classical physical system can be described by a principle of stationary action. Examples include ideal fluid mechanics, elasticity, magnetohydrodynamics, electromagnetics, and general relativity. All that is required to specify the dynamics of a system is an appropriate configuration space to describe its state and a Lagrangian.

**Remark 3.7.** This meaning of the principle of stationary action is rather mysterious, but we will verify that it leads to Newton's second law. One way to interpret the principle is that it expresses a lack of distinction between different forms of energy (kinetic and potential): any variation of a stationary trajectory leads to an equal gain, or loss, of kinetic and potential energies. An alternative explanation, from quantum mechanics, is that the trajectories with stationary action are the ones with a minimal cancelation of quantum-mechanical amplitudes. Whether this makes the principle less, or more, mysterious is not so clear.

## 1.2. Equivalence with Newton's second law

To derive the differential equation satisfied by a stationary point  $\vec{x}(t)$  of the action  $\mathcal{S}$  defined in (3.2)–(3.3), we differentiate the equation

$$\mathcal{S}(\vec{x} + \varepsilon \vec{h}) = \int_a^b \left\{ \frac{1}{2} m \left| \dot{\vec{x}}(t) + \varepsilon \dot{\vec{h}}(t) \right|^2 - V(\vec{x}(t) + \varepsilon \vec{h}(t)) \right\} dt$$

with respect to  $\varepsilon$ , and set  $\varepsilon = 0$ , as in (3.4). This gives

$$d\mathcal{S}(\vec{x}) \vec{h} = \int_a^b \left\{ m \dot{\vec{x}} \cdot \dot{\vec{h}} - \nabla V(\vec{x}) \cdot \vec{h} \right\} dt.$$

Integrating the first term by parts, and using the fact that the boundary terms vanish because  $\vec{h}(a) = \vec{h}(b) = 0$ , we get

$$(3.5) \quad d\mathcal{S}(\vec{x}) \vec{h} = - \int_a^b \left\{ m\ddot{\vec{x}} + \nabla V(\vec{x}) \right\} \cdot \vec{h} dt.$$

If this integral vanishes for arbitrary  $\vec{h}(t)$ , it follows from the du Bois-Reymond lemma (1879) that the integrand vanishes. Thus,  $\vec{x}(t)$  satisfies

$$m\ddot{\vec{x}} + \nabla V(\vec{x}) = 0$$

for  $a \leq t \leq b$ . Hence, we recover Newton's second law (3.1).

### 1.3. The variational derivative

A convenient way to write the derivative of the action is in terms of the variational, or functional, derivative. The variational derivative of  $\mathcal{S}$  at  $\vec{x}(t)$  is the function

$$\frac{\delta\mathcal{S}}{\delta\vec{x}} : [a, b] \rightarrow \mathbb{R}^n$$

such that

$$d\mathcal{S}(\vec{x}) \vec{h} = \int_a^b \frac{\delta\mathcal{S}}{\delta\vec{x}(t)} \cdot \vec{h}(t) dt.$$

Here, we use the notation

$$\frac{\delta\mathcal{S}}{\delta\vec{x}(t)}$$

to denote the value of the variational derivative at  $t$ . Note that the variational derivative depends on the trajectory  $\vec{x}$  at which we evaluate  $d\mathcal{S}(\vec{x})$ , although the notation does not show this explicitly.

Thus, for the action functional (3.2)–(3.3), equation (3.5) implies that

$$\frac{\delta\mathcal{S}}{\delta\vec{x}} = - \left\{ m\ddot{\vec{x}} + \nabla V(\vec{x}) \right\}.$$

A trajectory  $\vec{x}(t)$  is a stationary point of  $\mathcal{S}$  if the variational derivative of  $\mathcal{S}$  vanishes at  $\vec{x}(t)$ .

The variational derivative of a functional is analogous to the gradient of a function. If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a scalar-valued function on  $n$ -dimensional Euclidean space, then the gradient  $\nabla f$  is defined by

$$\left. \frac{d}{d\varepsilon} f(\vec{x} + \varepsilon\vec{h}) \right|_{\varepsilon=0} = \nabla f(\vec{x}) \cdot \vec{h}$$

where ‘ $\cdot$ ’ denotes the Euclidean inner product. Thus, we use the inner product to identify the derivative at a point, which is a linear map belonging to the dual space of  $\mathbb{R}^n$ , with a corresponding gradient vector belonging to  $\mathbb{R}^n$ . For the variational derivative, we replace the Euclidean inner product of vectors by the  $L^2$ -inner product of functions,

$$\langle \vec{x}, \vec{y} \rangle = \int_a^b \vec{x}(t) \cdot \vec{y}(t) dt,$$

and define the variational derivative by

$$d\mathcal{S}(\vec{x}) \vec{h} = \left\langle \frac{\delta\mathcal{S}}{\delta\vec{x}}, \vec{h} \right\rangle.$$

**Remark 3.8.** Considering the scalar case  $x : [a, b] \rightarrow \mathbb{R}$  for simplicity, and taking  $h(t) = \delta_{t_0}(t)$ , where  $\delta_{t_0}(t) = \delta(t - t_0)$  is the delta function supported at  $t_0$ , we have formally that

$$\frac{\delta \mathcal{S}}{\delta x(t_0)} = \left. \frac{d}{d\varepsilon} \mathcal{S}(x + \varepsilon \delta_{t_0}) \right|_{\varepsilon=0}.$$

Thus, we may interpret the value of the functional derivative  $\delta \mathcal{S} / \delta x$  at  $t$  as describing the change in the values of the functional  $\mathcal{S}(x)$  due to changes in the function  $x$  at the point  $t$ .

#### 1.4. Examples from mechanics

Let us return to the examples considered in Section 1.

**Example 3.9.** The action for the one-dimensional oscillator in Example 3.1 is

$$\mathcal{S}(x) = \int_a^b \left\{ \frac{1}{2} m \dot{x}^2 - V(x) \right\} dt,$$

and its variational derivative is

$$\frac{\delta \mathcal{S}}{\delta x} = -[m\ddot{x} + V'(x)].$$

**Example 3.10.** The potential energy  $V : \mathbb{R}^3 \setminus \{0\} \rightarrow \mathbb{R}$  for a central inverse-square force is given by

$$V(\vec{x}) = -\frac{GMm}{|\vec{x}|}.$$

The action of a trajectory  $\vec{x} : [a, b] \rightarrow \mathbb{R}^3$  is

$$\mathcal{S}(\vec{x}) = \int_a^b \left\{ \frac{1}{2} m |\dot{\vec{x}}|^2 + \frac{GMm}{|\vec{x}|} \right\} dt.$$

**Example 3.11.** The action for the  $n$ -body problem in Example 3.3 is

$$\mathcal{S}(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n) = \int_a^b \left\{ \frac{1}{2} \sum_{i=1}^n m_i |\dot{\vec{x}}_i|^2 + \frac{1}{2} \sum_{i,j=1}^n \frac{Gm_i m_j}{|\vec{x}_i - \vec{x}_j|} \right\} dt.$$

The equations of motion are obtained from the requirement that  $\mathcal{S}$  is stationary with respect to independent variations of  $\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$ .

**Example 3.12.** The configuration of a particle may be described by a point in some other manifold than  $\mathbb{R}^n$ . For example, consider a pendulum of length  $\ell$  and mass  $m$  in a gravitational field with acceleration  $g$ . We may describe its configuration by an angle  $\theta \in \mathbb{T}$ . The action is

$$\mathcal{S} = \int_a^b \left\{ \frac{1}{2} m \ell^2 \dot{\theta}^2 - mg\ell(1 - \cos \theta) \right\} dt,$$

and the corresponding equation of motion is the pendulum equation

$$\ell \ddot{\theta} + g \sin \theta = 0.$$

The following example connects mechanics and the calculus of variations with Riemannian geometry.

**Example 3.13.** Consider a particle moving freely on a Riemannian manifold  $M$  with metric  $g$ . If  $x = (x^1, x^2, \dots, x^n)$  are local coordinates on  $M$ , then the arclength  $ds$  on  $M$  is given by

$$ds^2 = g_{ij}(x) dx^i dx^j$$

where  $g_{ij}$  are the metric components. The metric is required to be symmetric, so  $g_{ij} = g_{ji}$ , and non-singular. We use the summation convention, meaning that repeated upper and lower indices are summed from 1 to  $n$ . A trajectory  $\gamma : [a, b] \rightarrow M$  has kinetic energy

$$T(\gamma, \dot{\gamma}) = \frac{1}{2} g_{ij}(\gamma) \dot{\gamma}^i \dot{\gamma}^j.$$

The corresponding action is

$$\mathcal{S} = \frac{1}{2} \int_a^b g_{ij}(\gamma) \dot{\gamma}^i \dot{\gamma}^j dt.$$

The principle of stationary action leads to the equation

$$g_{ij}(\gamma) \ddot{\gamma}^j + \Gamma_{jki}(\gamma) \dot{\gamma}^j \dot{\gamma}^k = 0 \quad i = 1, 2, \dots, n$$

where the connection coefficients, or Christoffel symbols,  $\Gamma_{jki}$  are defined by

$$\Gamma_{jki} = \frac{1}{2} \left( \frac{\partial g_{ij}}{\partial x^k} + \frac{\partial g_{ik}}{\partial x^j} - \frac{\partial g_{jk}}{\partial x^i} \right).$$

Since the metric is invertible, we may solve this equation for  $\ddot{\gamma}$  to get

$$(3.6) \quad \ddot{\gamma}^i + \Gamma_{jk}^i(\gamma) \dot{\gamma}^j \dot{\gamma}^k = 0 \quad i = 1, 2, \dots, n$$

where

$$\Gamma_{jk}^i = g^{ip} \Gamma_{jkp}$$

and  $g^{ij}$  denotes the components of the inverse matrix of  $g_{ij}$  such that

$$g^{ij} g_{jk} = \delta_k^i.$$

The solutions of the second-order system of ODEs (3.6) are the geodesics of the manifold.

## 2. The Euler-Lagrange equation

In the mechanical problems considered above, the Lagrangian is a quadratic function of the velocity. Here, we consider Lagrangians with a more general dependence on the derivative.

Let  $\mathcal{F}$  be a functional of scalar-valued functions  $u : [a, b] \rightarrow \mathbb{R}$  of the form

$$(3.7) \quad \mathcal{F}(u) = \int_a^b F(x, u(x), u'(x)) dx,$$

$$F : [a, b] \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R},$$

where  $F$  is a smooth function.

It is convenient to use the same notation for the variables

$$(x, u, u') \in [a, b] \times \mathbb{R} \times \mathbb{R}$$

on which  $F$  depends and the functions  $u(x), u'(x)$ . We denote the partial derivatives of  $F(x, u, u')$  by

$$F_x = \frac{\partial F}{\partial x} \Big|_{u, u'}, \quad F_u = \frac{\partial F}{\partial u} \Big|_{x, u'}, \quad F_{u'} = \frac{\partial F}{\partial u'} \Big|_{x, u}.$$

If  $h : [a, b] \rightarrow \mathbb{R}$  is a smooth function that vanishes at  $x = a, b$ , then

$$(3.8) \quad \begin{aligned} d\mathcal{F}(\vec{u})h &= \left. \frac{d}{d\varepsilon} \int_a^b F(x, u(x) + \varepsilon h(x), u'(x) + \varepsilon h'(x)) dx \right|_{\varepsilon=0} \\ &= \int_a^b \{F_u(x, u(x), u'(x))h(x) + F_{u'}(x, u(x), u'(x))h'(x)\} dx. \end{aligned}$$

It follows that a necessary condition for a  $C^1$ -function  $u(x)$  to be a stationary point of (3.7) in a space of functions with given values at the endpoints is that

$$(3.9) \quad \int_a^b \{F_u(x, u(x), u'(x))h(x) + F_{u'}(x, u(x), u'(x))h'(x)\} dx = 0$$

for all smooth functions  $h(x)$  that vanish at  $x = a, b$ .

If the function  $u$  in (3.8) is  $C^2$ , then we may integrate by parts to get

$$d\mathcal{F}(\vec{u})h = \int_a^b \left\{ F_u(x, u(x), u'(x)) - \frac{d}{dx} F_{u'}(x, u(x), u'(x)) \right\} h(x) dx.$$

It follows that the variational derivative of  $\mathcal{F}$  is given by

$$\frac{\delta\mathcal{F}}{\delta u} = -\frac{d}{dx} F_{u'}(x, u, u') + F_u(x, u, u').$$

Moreover, if a  $C^2$ -function  $u(x)$  is a stationary point of  $\mathcal{F}$ , then it must satisfy the ODE

$$(3.10) \quad -\frac{d}{dx} F_{u'}(x, u, u') + F_u(x, u, u') = 0.$$

Equation (3.10) is the *Euler-Lagrange equation* associated with the functional (3.7). It is a necessary, but not sufficient, condition that any smooth minimizer of (3.7) must satisfy. Equation (3.9) is the weak form of (3.10); it is satisfied by any  $C^1$ -minimizer (or, more generally, by any minimizer that belongs to a suitable Sobolev space  $W^{1,p}(a, b)$ ).

Note that  $d/dx$  in (3.10) is the total derivative with respect to  $x$ , meaning that the derivative is taken after the substitution of the functions  $u(x)$  and  $u'(x)$  into the arguments of  $F$ . Thus,

$$\frac{d}{dx} f(x, u, u') = f_x(x, u, u') + f_u(x, u, u')u' + f_{u'}(x, u, u')u''.$$

The coefficient of  $u''$  in (3.10) is equal to  $F_{u'u'}$ . The ODE is therefore of second order provided that

$$F_{u'u'}(x, u, u') \neq 0.$$

The derivation of the Euler-Lagrange equation extends straightforwardly to Lagrangians that depend on higher derivatives and to systems. For example, the Euler-Lagrange equation for the scalar functional

$$\mathcal{F}(u) = \int_a^b F(x, u(x), u'(x), u''(x)) dx,$$

where  $F : [a, b] \times \mathbb{R} \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ , is

$$\frac{d^2}{dx^2} F_{u''} - \frac{d}{dx} F_{u'} + F_u = 0.$$

This is a fourth-order ODE if  $F_{u''u''} \neq 0$ .

The Euler-Lagrange equation for a vector functional

$$\mathcal{F}(\vec{u}) = \int_a^b F(x, \vec{u}(x), \vec{u}'(x)) dx,$$

where  $F : [a, b] \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , is

$$-\frac{d}{dx} F_{u'_i} + F_{u_i} = 0 \quad \text{for } i = 1, 2, \dots, n.$$

This is an  $n \times n$  system of ODEs for  $\vec{u} = (u_1, u_2, \dots, u_n)$ . The system is second-order if the  $n \times n$  matrix with components  $F_{u'_i u'_j}$  is invertible.

The extension to functionals that involve more than one independent variable is less straightforward, and some examples will be considered below. In that case, the Euler-Lagrange equation is a PDE.

The question of whether a solution of the Euler-Lagrange equation is an extreme point of the functional is quite subtle even in the one-dimensional case. For example, the application of a second-derivative test, familiar from calculus for functions on finite-dimensional spaces, is not entirely straightforward. We will not discuss these questions here; see [11], for example, for more information.

### 3. Newton's problem of minimal resistance

If in a rare medium, consisting of equal particles freely disposed at equal distance from each other, a globe and a cylinder described on equal diameter move with equal velocities in the direction of the axis of the cylinder, the resistance of the globe will be half as great as that of the cylinder. . . I reckon that this proposition will not be without application in the building of ships.<sup>2</sup>

Many variational problems arise from optimization problems in which we seek to minimize (or maximize) some functional. We consider here a problem proposed and solved by Newton (1685) of finding the shape of a body with minimal resistance in a rarified gas. This was one of the first problems in the calculus of variations to be solved.

#### 3.1. Derivation of Newton's resistance functional

Following Newton, let us imagine that the gas is composed of uniformly distributed, non-interacting particles that reflect elastically off the body. We suppose that the particles have number-density  $n$ , mass  $m$ , and constant velocity  $v$  the downward  $z$ -direction, in a frame of reference moving with the body.

We assume that the body is cylindrically symmetric with a maximum radius of  $a$  and height  $h$ . We write the equation of the body surface in cylindrical polar coordinates as  $z = u(r)$ , where  $0 \leq r \leq a$  and

$$u(0) = h, \quad u(a) = 0.$$

Let  $\theta(r)$  denote the angle of the tangent line to the  $r$ -axis of this curve at the point  $(r, u(r))$ . Since the angle of reflection of a particle off the body is equal to the angle of incidence,  $\pi/2 - \theta$ , the reflected particle path makes an angle  $2\theta$  to the  $z$ -axis.

<sup>2</sup>I. Newton in *Principia Mathematica*, quoted from [11].

The change in momentum of the particle in the  $z$ -direction when it reflects off the body is therefore

$$mv(1 + \cos 2\theta).$$

For example, this is equal to  $2mv$  for normal incidence ( $\theta = 0$ ), and 0 for grazing incidence ( $\theta = \pi/2$ ).

The number of particles per unit time, per unit distance in the radial direction that hit the body is equal to

$$2\pi nvr.$$

Note that  $[2\pi nvr] = (1/L^3) \cdot (L/T) \cdot (L) = 1/(LT)$  as it should.

The rate at which the particles transfer momentum to the body per unit time, which is equal to force  $F$  exerted by the gas on the body, is given by

$$F = 2\pi nmv^2 \int_0^a r(1 + \cos 2\theta) dr.$$

Using the fact that  $\tan \theta = u'$  to eliminate  $\theta$ , we get that the resistance force on a profile  $z = u(r)$  is given by

$$F = 4\pi nma^2v^2 \mathcal{F}(u),$$

where the resistance functional  $\mathcal{F}$  is defined by

$$(3.11) \quad \mathcal{F}(u) = \frac{1}{a^2} \int_0^a \frac{r}{1 + [u'(r)]^2} dr.$$

Introducing dimensionless variables  $\tilde{r} = r/a$ ,  $\tilde{u} = u/a$  in (3.11), and dropping the tildes, we get the nondimensionalized resistance functional

$$(3.12) \quad \mathcal{F}(u) = \int_0^1 \frac{r}{1 + [u'(r)]^2} dr.$$

As we will see, this resistance functional does not provide the the most convincing physical results, although it has been used as a model for rarified flows and hypersonic flows. It is nevertheless remarkable that Newton was able to formulate and solve this problem long before a systematic development of the theory of fluid mechanics.

### 3.2. Resistances of some simple shapes

To see how the resistance functional  $\mathcal{F}$  in (3.11) behaves and formulate an appropriate optimization problem for it, let us consider some examples. Clearly, we have  $0 < \mathcal{F}(u) \leq 1/2$  for any  $u : [a, b] \rightarrow \mathbb{R}$ .

**Example 3.14.** For a vertical cylinder of radius  $a$ , we have  $u(r) = h$  for  $0 \leq r < a$  and  $u(a) = 0$ . The integrand in the functional (3.11) is small when  $u'$  is large, so we can approximate this discontinuous function by smooth functions whose resistance is arbitrarily close to the resistance of the cylinder. Setting  $u' = 0$  in (3.11), we get  $\mathcal{F} = 1/2$ . Thus, a blunt cylinder has the maximum possible resistance. The resistance is independent of the cylinder height  $h$ , since the gas particles graze the sides of the cylinder and exert no force upon it.

**Example 3.15.** For a sphere, with  $r^2 + z^2 = a^2$  and  $u(r) = \sqrt{a^2 - r^2}$ , we get  $\mathcal{F} = 1/4$ . As Newton observed, this is half the resistance of the cylinder. More generally, consider an ellipsoid of radius  $a$  and height  $h$ , with aspect ratio

$$(3.13) \quad M = \frac{h}{a},$$

and equation

$$\frac{r^2}{a^2} + \frac{z^2}{h^2} = 1, \quad u(r) = M\sqrt{a^2 - r^2}.$$

Using this expression for  $u$  in (3.11), and assuming that  $M \neq 1$ , we get the resistance

$$\mathcal{F}(u) = \frac{M^2 \log M^2 - (M^2 - 1)}{2(M^2 - 1)^2}.$$

The limit of this expression as  $M \rightarrow 0$  is equal to  $1/2$ , the resistance of the cylinder, and the limit as  $M \rightarrow 1$  is  $1/4$ , the resistance of the sphere. As  $M \rightarrow \infty$ , the resistance approaches zero. Thus, the resistance becomes arbitrarily small for a sufficiently tall, thin ellipsoid, and there is no profile that minimizes the resistance without a constraint on the aspect ratio.

**Example 3.16.** The equation of a circular cone with base  $a$  and height  $h$  is  $z = u(r)$  with  $u(r) = M(a - r)$ , where  $M$  is given by (3.13) as before. In this case  $u' = M$  is constant, and

$$\mathcal{F}(u) = \frac{1}{2(1 + M^2)}$$

As  $M \rightarrow 0$ , the resistance approaches  $1/2$ , and as  $M \rightarrow \infty$ , the resistance approaches 0.

**Example 3.17.** Suppose that  $u_n(r)$  consists of  $(n + 1/2)$  ‘tent’ functions of height  $h$  and base  $2b_n$  where

$$b_n = \frac{a}{2n + 1}.$$

Then, except at the ‘corners,’ we have  $|u'_n| = h/b_n$ , and therefore

$$\mathcal{F}(u_n) = \frac{1}{2[1 + (2n + 1)^2 M^2]}.$$

As before, we can approximate this piecewise smooth function by smooth functions with an arbitrarily small increase in the resistance. Thus,  $\mathcal{F}(u_n) \rightarrow 0$  as  $n \rightarrow \infty$ , even though  $0 \leq u_n(r) \leq h$  and the heights of the bodies are uniformly bounded. To eliminate this kind of oscillatory behavior, which would lead to multiple impacts of particles on the body contrary to what is assumed in the derivation of the resistance formula, we will impose the reasonable requirement that  $u'(r) \leq 0$  for  $0 \leq r \leq a$ .

### 3.3. The variational problem

We fix the aspect ratio  $M > 0$ , and seek to minimize  $\mathcal{F}$  over the space of functions

$$X_M = \{u \in W^{1,\infty}(0, 1) : [0, 1] \rightarrow \mathbb{R} \mid u(0) = M, \quad u(1) = 0, \quad u'(r) \leq 0\}.$$

Here,  $W^{1,\infty}(0, 1)$  denotes the Sobolev space of functions whose weak, or distributional, derivative is a bounded function  $u' \in L^\infty(0, 1)$ . Equivalently, this means that  $u$  is Lipschitz continuous with  $|u(x) - u(y)| \leq M|x - y|$ , where  $M = \|u\|_\infty$ . We could minimize  $\mathcal{F}$  over the larger space  $W^{1,1}(0, 1)$  of absolutely continuous functions with  $u' \in L^1(0, 1)$ , and get the same result. As we shall see, however, the smaller space  $C^1[0, 1]$  of continuously differentiable functions would not be adequate because the minimizer has a ‘corner’ and is not continuously differentiable.

Also note that, as the examples above illustrate, it is necessary to impose a constraint, such as  $u' \leq 0$ , on the admissible functions, otherwise (as pointed out by Legendre in 1788) we could make the resistance as small as we wish by taking

profiles with rapid ‘zig-zags’ and large slopes, although the infimum  $\mathcal{F} = 0$  is not attained for any profile.

The functional (3.12) is a pathological one from the perspective of the general theory of the calculus of variations. First, it is not coercive, because

$$\frac{r}{1 + [u']^2} \rightarrow 0 \quad \text{as } |u'| \rightarrow \infty.$$

As a result, minimizing sequences need not be bounded, and, in the absence of constraints, minimizers can ‘escape’ to infinity. Second, it is not convex. A function  $\mathcal{F} : X \rightarrow \mathbb{R}$  on a real vector space  $X$  is convex if, for all  $u, v \in X$  and  $\lambda \in [0, 1]$ ,

$$\mathcal{F}(\lambda u + (1 - \lambda)v) \leq \lambda \mathcal{F}(u) + (1 - \lambda)\mathcal{F}(v).$$

In general, convex functions have good lower semicontinuity properties and convex minimization problems are typically well-behaved. The behavior of non-convex optimization problems can be much nastier.

#### 3.4. The Euler-Lagrange equation

The Euler-Lagrange equation for (3.12) is

$$\frac{d}{dr} \left\{ \frac{ru'}{[1 + (u')^2]^2} \right\} = 0.$$

Since the Lagrangian is independent of  $u$ , this has an immediate first integral,

$$(3.14) \quad ru' = -c [1 + (u')^2]^2$$

where  $c \geq 0$  is a constant of integration.

If  $c = 0$  in (3.14), then we get  $u' = 0$ , or  $u = \text{constant}$ . This solution corresponds to the cylinder with maximum resistance  $1/2$ . The maximum is not attained, however, within the class absolutely continuous functions  $u \in X_M$ , since for such functions if  $u'$  is zero almost everywhere with respect to Lebesgue measure, then  $u$  is constant, and it cannot satisfy both boundary conditions  $u(0) = M$ ,  $u(1) = 0$ .

If  $c > 0$  in (3.14), then it is convenient to parametrize the solution curve by  $p = u' < 0$ . From (3.14), the radial coordinate  $r$  is given in terms of  $p$  by

$$(3.15) \quad r = -\frac{c(1 + p^2)^2}{p}.$$

Using this equation to express  $dr$  in terms of  $dp$  in the integral

$$u = \int p dr,$$

and evaluating the result, we get

$$(3.16) \quad u = u_0 - c \left( -\log |p| + p^2 + \frac{3}{4}p^4 \right),$$

where  $u_0$  is a constant of integration.

From (3.15), we see that the minimum value of  $r(p)$  for  $p < 0$  is

$$r_0 = \frac{16\sqrt{3}c}{9}$$

at  $p = -1/\sqrt{3}$ . Thus, although this solution minimizes the resistance, we cannot use it over the whole interval  $0 \leq r \leq 1$ , only for  $r_0 \leq r \leq 1$ . In the remaining part of the interval, we use  $u = \text{constant}$ , and we obtain the lowest global resistance by placing the blunt part of the body around the nose  $r = 0$ , where it contributes least to the area and resistance.

While this plausibility argument seems reasonable, it is not entirely convincing, since the flat nose locally maximizes the resistance, and it is far from a proof. Nevertheless, with additional work, it is possible to prove that it does give the correct solution  $u \in X_M$  with minimal resistance.

This minimizing solution has the form

$$u(r) = \begin{cases} M & \text{for } 0 \leq r \leq r_0, \\ u_0 - c(-\log |p| + p^2 + \frac{3}{4}p^4) & \text{for } p_1 \leq p \leq -1/\sqrt{3}, \end{cases}$$

where  $r(p_1) = 1$ .

Imposing continuity of the solution at  $r = r_0$ ,  $p = 1/\sqrt{3}$  and the boundary condition  $u(1) = 0$ , with  $p = p_1$ , we get

$$\begin{aligned} M &= u_0 - c \left( \log \sqrt{3} + \frac{5}{12} \right), \\ p_1 &= -c(1 + p_1^2)^2, \\ u_0 &= c \left( -\log |p_1| + p_1^2 + \frac{3}{4}p_1^4 \right). \end{aligned}$$

Eliminating  $u_0$ , we may write the solution as

$$u(r) = M - c \left( p^2 + \frac{3}{4}p^4 - \log |\sqrt{3}p| - \frac{5}{12} \right)$$

for  $p_1 \leq p \leq -1/\sqrt{3}$ , where

$$M = c \left( p_1^2 + \frac{3}{4}p_1^4 - \log |\sqrt{3}p_1| - \frac{5}{12} \right), \quad p_1 = -c(1 + p_1^2)^2.$$

Thus,  $p_1$  is the solution of

$$\frac{p_1 \left( \log |\sqrt{3}p_1| - p_1^2 - \frac{3}{4}p_1^4 + \frac{5}{12} \right)}{(1 + p_1^2)^2} = M,$$

and  $r_0$  is given in terms of  $p_1$  by

$$r_0 = -\frac{16\sqrt{3}p_1}{9(1 + p_1^2)^2}.$$

Denoting by

$$C_0 = 2 \int_0^1 \frac{r}{1 + (u')^2} dr$$

the ratio of the minimal resistance to the maximal resistance of a cylinder, one gets the numerical values shown below [12]. Moreover, one can show that

$$r_0 \sim \frac{27}{16} \frac{1}{M^3}, \quad C_0 \sim \frac{27}{32} \frac{1}{M^2} \quad \text{as } M \rightarrow \infty.$$

Thus, as the aspect ratio increases, the radius of the blunt nose decreases and the total resistance of the body approaches zero (see Figure 1).

	$M = 1$	$M = 2$	$M = 3$	$M = 4$
$r_0$	0.35	0.12	0.048	0.023
$C_0$	0.37	0.16	0.0082	0.0049

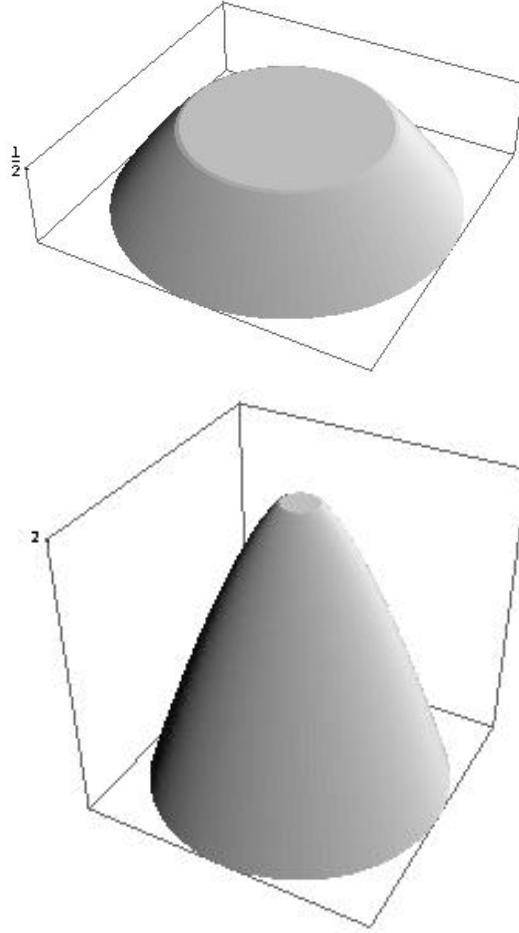


FIGURE 1. Minimal resistance surfaces for  $M = 0.5$  and  $M = 2.0$  (from Mark Peltier).

### 3.5. Non-radially symmetric solutions

The radially symmetric problem may be generalized to a two-dimensional, non-radially symmetric problem as follows. Suppose that  $\Omega \subset \mathbb{R}^2$  is a given domain (a bounded, open, connected set). Find a bounded, nonnegative convex function  $u : \Omega \rightarrow \mathbb{R}$  that minimizes

$$\mathcal{F}(u) = \int_{\Omega} \frac{1}{1 + |\nabla u|^2} dx dy.$$

In this case, the shape of the body is given by  $z = u(x, y)$ .

In the discussion above, we obtained the minimizer among radially symmetric bodies when  $\Omega$  is a disc  $D$ . It might seem natural to suppose that this radially symmetric solution minimizes the resistance among non-radially symmetric admissible functions  $u : D \rightarrow \mathbb{R}$ . It is interesting to note that this is not true. Brock, Ferroni, and Kawohl (1996) showed that there are non-radially symmetric convex functions on the disc that give a lower resistance than the radial solution found above.

#### 4. Constrained variational principles

It often occurs that we want to minimize a functional subject to a constraint. Constraints can take many forms. First, consider the minimization of a functional

$$\mathcal{F}(u) = \int_a^b F(x, u, u') dx,$$

over functions such that  $u(a) = 0$ ,  $u(b) = 0$ , subject to an integral constraint of the form

$$\mathcal{G} = \int_a^b G(x, u, u') dx.$$

Variational problems with integral constraints are called isoperimetric problems after the prototypical problem of finding the curve (a circle) that encloses the maximum area subject to the constraint that its length is fixed.<sup>3</sup>

We may solve this problem by introducing a Lagrange multiplier  $\lambda \in \mathbb{R}$  and seeking stationary points of the unconstrained functional

$$\mathcal{F}(u) - \lambda \mathcal{G}(u) = \int_a^b \{F(x, u, u') - \lambda G(x, u, u')\} dx.$$

The condition that this functional is stationary with respect to  $\lambda$  implies that  $\mathcal{G}(u) = 0$ , so a stationary point satisfies the constraint.

The Euler-Lagrange equation for stationarity of the functional with respect to variations in  $u$  is

$$-\frac{d}{dx} F_{u'}(x, u, u') + F_u(x, u, u') = \lambda \left[ -\frac{d}{dx} G_{u'}(x, u, u') + G_u(x, u, u') \right].$$

In principle, we solve this problem for  $u(x)$  and  $\lambda$  subject to the boundary conditions  $u(a) = 0$ ,  $u(b) = 0$  and the constraint  $\mathcal{G}(u) = 0$ .

##### 4.1. Eigenvalue problems

Consider the following Rayleigh quotient

$$\mathcal{Q}(u) = \frac{\int_a^b \{p(x)u'^2 + q(x)u^2\} dx}{\int_a^b u^2 dx}$$

where  $p(x)$ ,  $q(x)$  are given coefficient functions.

Since  $\mathcal{Q}(u)$  is homogeneous in  $u$ , the minimization of  $\mathcal{Q}(u)$  over nonzero functions  $u$  is equivalent to the minimization of the numerator subject to the constraint

<sup>3</sup>According to Virgil's Aeneid, Dido was given as much land as she could enclose with an ox hide to found the city of Carthage. She cut the hide into a thin strip, and used it to enclose a large circular hill.

that the denominator is equal to one; or, in other words, to the minimization of  $\mathcal{F}(u)$  subject to the constraint  $\mathcal{G}(u) = 0$  where

$$\mathcal{F}(u) = \frac{1}{2} \int_a^b \{p(x)u'^2 + q(x)u^2\} dx, \quad \mathcal{G}(u) = \frac{1}{2} \left\{ \int_a^b u^2 dx - 1 \right\}.$$

The corresponding Euler-Lagrange equation for the stationarity of  $\mathcal{F}(u) - \lambda\mathcal{G}(u)$  with respect to  $u$  is

$$- [p(x)u']' + q(x)u = \lambda u.$$

This is a Sturm-Liouville eigenvalue problem in which the Lagrange multiplier  $\lambda$  is an eigenvalue.

## 5. Elastic rods

As an example of the use of constrained variational principles, we will derive equations for the equilibria of an inextensible elastic rod and describe some applications.

Consider a thin, inextensible elastic rod that resists bending. Suppose that the cross-sections of the rod are isotropic and that we can ignore any twisting. We model the spatial configuration of the rod by a curve  $\vec{r}(s)$ ,

$$\vec{r}: [a, b] \rightarrow \mathbb{R}^3,$$

where it is convenient to parametrize the curve by arclength  $a \leq s \leq b$ .

We can model the twisting of a rod by introducing additional vectors that describe the orientation of its cross-section, leading to the Kirchoff and Cosserat theories [4], but we will not consider such generalizations here.

### 5.1. Kinematics

We introduce an orthonormal frame of vectors  $\{\vec{t}, \vec{n}, \vec{b}\}$  along the curve, consisting of the unit tangent, normal and binormal vectors, respectively. We have  $\vec{t} = \vec{r}'$  and  $\vec{b} = \vec{t} \times \vec{n}$ . According to the the Frenet-Serret formulas, these vectors satisfy

$$\vec{t}' = \kappa\vec{n}, \quad \vec{n}' = -\kappa\vec{t} + \tau\vec{b}, \quad \vec{b}' = -\tau\vec{n}$$

where  $\kappa(s)$  is the curvature and  $\tau(s)$  is the torsion of the curve.

These equations may also be written as

$$\begin{pmatrix} \vec{t} \\ \vec{n} \\ \vec{b} \end{pmatrix}' = \begin{pmatrix} 0 & \kappa & 0 \\ -\kappa & 0 & \tau \\ 0 & -\tau & 0 \end{pmatrix} \begin{pmatrix} \vec{t} \\ \vec{n} \\ \vec{b} \end{pmatrix}.$$

The skew-symmetric matrix on the right-hand side is the infinitesimal generator of the rotations of the orthonormal frame  $\{\vec{t}, \vec{n}, \vec{b}\}$  as it is transported along the curve.

### 5.2. A variational principle

We will derive equilibrium equations for the configuration of the rod from the condition that they minimize the energy.

We assume that the energy density of a rod configuration is proportional to the square of its curvature. This constitutive equation, and the model of a rod as an

‘elastic line,’ or *elastica*, was introduced and developed by James Bernoulli<sup>4</sup> (1694), Daniel Bernoulli (1728), and Euler (1727, 1732).

The curvature is given by  $\kappa^2 = \vec{t}' \cdot \vec{t}'$ , so the total energy of the rod is given by

$$(3.17) \quad \mathcal{E}(\vec{r}) = \int_a^b \frac{1}{2} J \vec{r}'' \cdot \vec{r}'' ds,$$

where the material function  $J : [a, b] \rightarrow \mathbb{R}$  gives the proportionality between the square of the curvature and the energy density due to bending.

Equations for the equilibrium configuration of the rod follow by minimizing the energy (3.17) subject to the constraint that  $s$  is arclength, meaning that

$$\vec{r}' \cdot \vec{r}' = 1.$$

This constraint is a pointwise constraint, rather than an integral, so we impose it by introducing a function  $\lambda : [a, b] \rightarrow \mathbb{R}$  as a Lagrange multiplier, and seeking stationary points of the functional

$$\mathcal{F}(\vec{r}, \lambda) = \int_a^b \frac{1}{2} \{J \vec{r}'' \cdot \vec{r}'' - \lambda (\vec{r}' \cdot \vec{r}' - 1)\} ds.$$

The Euler-Lagrange equation obtained by varying  $\vec{r}$  is

$$(3.18) \quad (J\vec{r}'')' + (\lambda\vec{r}')' = 0,$$

while we recover the constraint by varying  $\lambda$ . Integrating (3.18) once, and writing  $\vec{r}' = \vec{t}$ , we get

$$(3.19) \quad (J\vec{t}')' + \lambda\vec{t} = \vec{F}$$

where  $\vec{F}$  is a constant vector of integration. It corresponds to the contact force exerted by one part of the rod on another, which is constant in an inextensible rod which is not acted on by an external force.

We obtain an expression for the Lagrange multiplier  $\lambda$  by imposing the constraint that  $\vec{t}$  is a unit vector on solutions of (3.19). Taking the inner product of (3.19) with  $\vec{t}$ , and rewriting the result, we get

$$(J\vec{t} \cdot \vec{t}')' - J\vec{t}' \cdot \vec{t} + \lambda\vec{t} \cdot \vec{t} = \vec{F} \cdot \vec{t}.$$

Using  $\vec{t} \cdot \vec{t} = 1$  and  $\vec{t}' \cdot \vec{t} = 0$  in this equation, we get

$$\lambda = \vec{F} \cdot \vec{t} + J\vec{t}' \cdot \vec{t}'.$$

Thus, (3.19) may be written as

$$(3.20) \quad (J\vec{t}')' + J\kappa^2\vec{t} = \vec{F} - (\vec{F} \cdot \vec{t})\vec{t}, \quad \kappa^2 = \vec{t}' \cdot \vec{t}'.$$

Equation (3.20) is a second order ODE for the tangent vector  $\vec{t}(s)$ . We supplement it with suitable boundary conditions at the ends of rod. For example, if the ends are fully clamped, we specify the directions  $\vec{t}(a)$ ,  $\vec{t}(b)$  of the rod at each endpoint. Given a solution for  $\vec{t}$ , we may then recover the position of the rod by integrating the equation  $\vec{r}' = \vec{t}$ . Note that, in this case, we cannot expect to also

<sup>4</sup>There were a lot of Bernoulli's. The main ones were the older brother James (1654-1705), the younger brother Johann (1667-1748), and Johann's son Daniel (1700-1782). James and Johann has a prolonged feud over the priority of their mathematical results, and, after James died, Johann became jealous of his son Daniel's work, in particular on Bernoulli's law in hydrodynamics.

specify the position of both endpoints. In general, the issue of what boundary conditions to use in rod theories is somewhat subtle (see [4] for further discussion).

Taking the cross product of (3.20) with  $\vec{t}$ , and using the fact that  $\vec{t} \times \vec{t}' = \kappa \vec{b}$ , we get

$$\vec{m}' = \vec{t} \times \vec{F}, \quad \text{where } \vec{m} = J\kappa \vec{b}.$$

This equation expresses a balance of moments in the rod due to the constant contact force  $\vec{F}$  and a contact couple  $\vec{m}$ . The couple is proportional to the curvature, as proposed by Bernoulli and Euler, corresponding to the constitutive assumption that the energy density is a quadratic function of the curvature. Thus, we obtain the same equations from the Euler-Lagrange equations of the variational principle as we would by balancing the forces and moments acting on the rod.

### 5.3. Dimensional considerations

From (3.17), the material function  $J$  has the dimension of energy  $\cdot$  length. It is often written as  $J = EI$  where  $E$  is Young's modulus for the elastic material making up the rod, and  $I$  is the moment of inertia of a cross-section.

Young's modulus gives the ratio of tensile stress to tensile strain in an elastic solid. Strain, which measures a deformed length to an undeformed length, is dimensionless, so  $E$  has the dimension of stress, force/area, meaning that

$$[E] = \frac{M}{LT^2}.$$

For example, the Young's modulus of steel is approximately 200 kN/mm<sup>2</sup>.

The moment of inertia, in this context, is a second area moment of the rod cross-section, and has the dimension of  $L^4$ . The term 'moment of inertia' is also used to describe the relationship between angular velocity and angular momentum for a rotating rigid body; the moment of inertia here corresponds to this notion with mass replaced by area.

Explicitly, we define the components of a second-order, area-moment tensor of a region  $\Omega \subset \mathbb{R}^2$  in the plane, with Cartesian coordinates  $x_i$ ,  $i = 1, 2$ , by

$$I_{ij} = \int_{\Omega} x_i x_j dA.$$

In general, this symmetric, positive-definite tensor has two positive real eigenvalues, corresponding to the moments of inertia about the principal axes defined by the corresponding eigenvectors. If these eigenvalues coincide, then we get the isotropic case with  $I_{ij} = I\delta_{ij}$  where  $I$  is the moment of inertia. For example, if  $\Omega$  is a disc of radius  $a$ , then  $I = \pi a^4/4$ .

Thus,

$$[EI] = \frac{ML^2}{T^2} \cdot L,$$

consistent with the dimension of  $J$ . In general,  $J$  may depend upon  $s$ , for example because the cross-sectional area of the rod, and therefore moment of inertia, varies along its length.

### 5.4. The persistence length of DNA

An interesting application of rod theories is to the modeling of polymers whose molecular chains resist bending, such as DNA. A statistical mechanics of flexible polymers may be derived by supposing that the polymer chain undergoes a random

walk due to thermal fluctuations. Such polymers typically coil up because there are more coiled configurations than straight ones, so coiling is entropically favored.

If a polymer has elastic rigidity, then the increase in entropy that favors its coiling is opposed by the bending energy required to coil. As a result, the tangent vector of the polymer chain is highly correlated over distances short enough that significant bending energies are required to change its direction, while it is decorrelated over much longer distances. A typical lengthscale over which the tangent vector is correlated is called the *persistence length* of the polymer.

According to statistical mechanics, the probability that a system at absolute temperature  $T$  has a specific configuration with energy  $E$  is proportional to

$$(3.21) \quad e^{-E/kT}$$

where  $k$  is Boltzmann's constant. Boltzmann's constant has the approximate value  $k = 1.38 \times 10^{-23} \text{ JK}^{-1}$ . The quantity  $kT$  is an order of magnitude for the random thermal energy of a single microscopic degree of freedom at temperature  $T$ .

The bending energy of an elastic rod is set by the coefficient  $J$  in (3.17), with dimension energy  $\cdot$  length. Thus, the quantity

$$A = \frac{J}{kT}$$

is a lengthscale over which thermal and bending energies are comparable, and it provides a measure of the persistence length. For DNA, a typical value of this length at standard conditions is  $A \approx 50 \text{ nm}$ , or about 150 base pairs of the double helix.

The statistical mechanics of an elastica, or 'worm-like chain,' may be described, formally at least, in terms of path integrals (integrals over an infinite-dimensional space of functions). The expected value  $\mathbf{E}[\mathcal{F}(\vec{r})]$  of some functional  $\mathcal{F}(\vec{r})$  of the elastica configuration is given by

$$\mathbf{E}[\mathcal{F}(\vec{r})] = \frac{1}{Z} \int \mathcal{F}(\vec{r}) e^{-\mathcal{E}(\vec{r})/kT} D\vec{r},$$

where the right-hand side is a path integral over a path space of configurations  $\vec{r}(s)$  using the Boltzmann factor (3.21) and the elastica energy (3.17). The factor  $Z$  is inserted to normalize the Boltzmann distribution to a probability distribution.

These path integrals are difficult to evaluate in general, but in some cases the energy functional may be approximated by a quadratic functional, and the resulting (infinite-dimensional) Gaussian integrals may be evaluated exactly. This leads to results which are in reasonable agreement with the experimentally observed properties of DNA [36]. One can also include other effects in the model, such as the twisting energy of DNA.

## 6. Buckling and bifurcation theory

Let us consider planar deformations of an elastic rod of length  $L$ . In this case, we may write

$$\vec{t} = (\cos \theta, \sin \theta)$$

in (3.20), where  $\theta(s)$  is the angle of the rod to the  $x$ -axis. We assume that the rod is uniform, so that  $J = EI$  is constant, and that the force  $\vec{F} = (F, 0)$  in the rod is directed along the  $x$ -axis, with  $F > 0$ , corresponding to a compression.

With these assumptions, equation (3.20) reduces to a scalar ODE

$$EI\theta'' + F \sin \theta = 0.$$

This ODE is the Euler-Lagrange equation of the functional

$$\mathcal{E}(\theta) = \int_0^L \left\{ \frac{1}{2}EI(\theta')^2 - F(1 - \cos \theta) \right\} ds$$

The first term is the bending energy of the rod, and the second term is the work done by the force in shortening the length of the rod in the  $x$ -direction.

This equation is identical in form to the pendulum equation. Here, however, the independent variable is arclength, rather than time, and we will impose boundary conditions, not initial conditions, on the solutions.

Let us suppose that the ends of the rod at  $s = 0$ ,  $s = L$  are horizontally clamped, so that  $\theta(0) = 0$ ,  $\theta(L) = 0$ . Introducing a dimensionless arclength variable  $\tilde{s} = s/L$ , and dropping the tildes, we may write this BVP as

$$(3.22) \quad \theta'' + \lambda \sin \theta = 0,$$

$$(3.23) \quad \theta(0) = 0, \quad \theta(1) = 0,$$

where the dimensionless force parameter  $\lambda > 0$  is defined by

$$\lambda = \frac{FL^2}{EI}.$$

This problem was studied by Euler, and is one of the original problems in the bifurcation theory of equilibria.

The problem (3.22)–(3.23) has the trivial solution  $\theta = 0$  for any value of  $\lambda$ , corresponding to the unbuckled state of the rod. This is the unique solution when  $\lambda$  is sufficiently small, but other non-trivial solutions bifurcate off the trivial solution as  $\lambda$  increases. This phenomenon corresponds to the buckling of the rod under an increased load.

The problem can be solved explicitly in terms of elliptic functions, as we will show below. First, however, we will obtain solutions by perturbing off the trivial solution. This method is applicable to more complicated problems which cannot be solved exactly.

### 6.1. The bifurcation equation

To study the bifurcation of non-zero solutions off the zero solution, we first linearize (3.22)–(3.23) about  $\theta = 0$ . This gives

$$(3.24) \quad \begin{aligned} \theta'' + \lambda_0 \theta &= 0, \\ \theta(0) &= 0, \quad \theta(1) = 0. \end{aligned}$$

We denote the eigenvalue parameter in the linearized problem by  $\lambda_0$ .

Equation (3.24) has a unique solution  $\theta = 0$  except when  $\lambda_0 = \lambda_0^{(n)}$ , where the eigenvalues  $\lambda_0^{(n)}$  are given by

$$\lambda_0^{(n)} = n^2 \pi^2 \quad \text{for } n \in \mathbb{N}.$$

The corresponding solutions are then  $\theta(s) = A\theta^{(n)}(s)$ , where

$$\theta^{(n)}(s) = \sin(n\pi s).$$

The implicit function theorem implies that if  $\bar{\lambda}$  is not an eigenvalue of the linearized problem, then the zero solution is the unique solution of the nonlinear problem for  $(\theta, \lambda)$  in a small enough neighborhood of  $(0, \bar{\lambda})$ .

On the other hand, non-trivial solutions can bifurcate off the zero solution at eigenvalues of the linearized problem. We will compute these solutions by expanding the nonlinear problem about an eigenvalue. As we discuss below, this formal computation can be made rigorous by use of a Lyapunov-Schmidt reduction.

Fix  $n \in \mathbb{N}$ , and let

$$\lambda_0 = n^2\pi^2$$

be the  $n^{\text{th}}$  eigenvalue. We drop the superscript  $n$  to simplify the notation.

We introduce a small parameter  $\varepsilon$ , and consider values of the eigenvalue parameter  $\lambda$  close to  $\lambda_0$ . We suppose that  $\lambda(\varepsilon)$  has the expansion

$$(3.25) \quad \lambda(\varepsilon) = \lambda_0 + \varepsilon^2\lambda_2 + \dots \quad \text{as } \varepsilon \rightarrow 0,$$

where we write  $\varepsilon^2$  instead of  $\varepsilon$  to simplify the subsequent equations.

We look for small-amplitude solutions  $\theta(s; \varepsilon)$  of (3.22)–(3.23) with an expansion of the form

$$(3.26) \quad \theta(s; \varepsilon) = \varepsilon\theta_1(s) + \varepsilon^3\theta_3(s) + \dots \text{ as } \varepsilon \rightarrow 0.$$

Using (3.25) and (3.26) in (3.22)–(3.23), Taylor expanding the result with respect to  $\varepsilon$ , and equating coefficients of  $\varepsilon$  and  $\varepsilon^3$  to zero, we find that

$$(3.27) \quad \begin{aligned} \theta_1'' + \lambda_0\theta_1 &= 0, \\ \theta_1(0) = 0, \quad \theta_1(1) &= 0, \end{aligned}$$

$$(3.28) \quad \begin{aligned} \theta_3'' + \lambda_0\theta_3 + \lambda_2\theta_1 - \frac{1}{6}\lambda_0\theta_1^3 &= 0, \\ \theta_3(0) = 0, \quad \theta_3(1) &= 0, \end{aligned}$$

The solution of (3.27) is

$$(3.29) \quad \theta_1(s) = A \sin(n\pi s),$$

where  $A$  is an arbitrary constant of integration.

Equation (3.28) then becomes

$$\begin{aligned} \theta_3'' + \lambda_0\theta_3 + \lambda_2 A \sin(n\pi s) - \frac{1}{6}\lambda_0 A^3 \sin^3(n\pi s) &= 0, \\ \theta_3(0) = 0, \quad \theta_3(1) &= 0, \end{aligned}$$

In general, this equation is not solvable for  $\theta_3$ . To derive the solvability condition, we multiply the ODE by the eigenfunction  $\sin(n\pi s)$  and integrate the result over  $0 \leq s \leq 1$ .

Integration by parts, or Green's formula, gives

$$\begin{aligned} \int_0^1 \sin(n\pi s) \{\theta_3'' + \lambda_0\theta_3\} ds - \int_0^1 \{\sin(n\pi s)'' + \lambda_0 \sin(n\pi s)\} \theta_3 ds \\ = [\sin(n\pi s) \theta_3' - \sin(n\pi s)' \theta_3]_0^1. \end{aligned}$$

It follows that

$$\int_0^1 \sin(n\pi s) \{\theta_3'' + \lambda_0\theta_3\} ds = 0,$$

and hence that

$$\lambda_2 A \int_0^1 \sin^2(n\pi s) ds = \frac{1}{6} \lambda_0 A^3 \int_0^1 \sin^4(n\pi s) ds.$$

Using the integrals

$$\int_0^1 \sin^2(n\pi s) ds = \frac{1}{2}, \quad \int_0^1 \sin^4(n\pi s) ds = \frac{3}{8},$$

we get

$$(3.30) \quad \lambda_2 A = \frac{1}{8} \lambda_0 A^3.$$

This is the bifurcation equation for the problem.

To rewrite (3.30) in terms of the original variables, let  $\alpha$  denote the maximum value of a solution  $\theta(s)$ . Then, from (3.26) and (3.29), we have

$$\alpha = \varepsilon A + O(\varepsilon^3).$$

Using (3.30) in (3.25), we get the bifurcation equation

$$(3.31) \quad (\lambda - \lambda_0) \alpha = \frac{1}{8} \lambda_0 \alpha^3 + O(\alpha^5) \quad \text{as } \alpha \rightarrow 0.$$

Thus, in addition to the trivial solution  $\alpha = 0$ , we have solutions with

$$(3.32) \quad \alpha^2 = \frac{8(\lambda - \lambda_0)}{\lambda_0} + O(\alpha^4)$$

branching from each of the linearized eigenvalues  $\lambda_0$  for  $\lambda > \lambda_0$ . This type of bifurcation is called a pitchfork bifurcation. It is supercritical because the new solutions appear for values of  $\lambda$  larger than the bifurcation value.

Thus, the original infinite-dimensional bifurcation problem (3.22)–(3.23) reduces to a one-dimensional bifurcation equation of the form  $F(\alpha, \lambda) = 0$  in a neighborhood of the bifurcation point  $(\theta, \lambda) = (0, \lambda_0)$ . The bifurcation equation has the Taylor expansion (3.31) as  $\alpha \rightarrow 0$  and  $\lambda \rightarrow \lambda_0$ .

## 6.2. Energy minimizers

For values of  $\lambda > \pi^2$ , solutions of the nonlinear BVP (3.22)–(3.23) are not unique. This poses the question of which solutions should be used. One criterion is that solutions of an equilibrium problem should be dynamically stable. We cannot address this question directly here, since we have not derived a set of time-dependent evolution equations. We can, however, use energy considerations.

The potential energy for (3.22) is

$$(3.33) \quad \mathcal{E}(\theta) = \int_0^1 \left\{ \frac{1}{2} (\theta')^2 - \lambda (1 - \cos \theta) \right\} ds.$$

We claim that the zero solution is a global minimizer of (3.33) when  $\lambda \leq \pi^2$ , with  $\mathcal{E}(0) = 0$ , but it is not a minimizer when  $\lambda > \pi$ . As a result, the zero solution loses stability as  $\lambda$  passes through the first eigenvalue  $\pi^2$ , after which the rod will buckle.

To show that  $\theta = 0$  is not a minimizer for  $\lambda > \pi^2$ , we compute the energy in the direction of the eigenvector of the first eigenvalue:

$$\begin{aligned}\mathcal{E}(\alpha \sin \pi s) &= \int_0^1 \left\{ \frac{1}{2} \alpha^2 \pi^2 \cos^2 \pi s - \lambda [1 - \cos(\alpha \sin \pi s)] \right\} ds \\ &= \int_0^1 \left\{ \frac{1}{2} \alpha^2 \pi^2 \cos^2 \pi s - \frac{1}{2} \alpha^2 \lambda \sin^2 \pi s \right\} ds + O(\alpha^4) \\ &= \frac{1}{4} \alpha^2 (\pi^2 - \lambda) + O(\alpha^4).\end{aligned}$$

It follows that we can have  $\mathcal{E}(\theta) < \mathcal{E}(0)$  when  $\lambda > \pi^2$ .

For the converse, we use the Poincaré (or Wirtinger) inequality, which states that

$$\int_0^1 \theta^2 ds \leq \frac{1}{\pi^2} \int_0^1 \theta'^2 ds$$

for all smooth functions such that  $\theta(0) = 0$ ,  $\theta(1) = 0$ . (The inequality also holds for all  $\theta \in H_0^1(0, 1)$ .) The best constant,  $1/\pi^2$ , in this inequality is the reciprocal of the lowest eigenvalue of (3.24), and it may be obtained by minimization of the corresponding Rayleigh quotient.

Using the inequality

$$1 - \cos \theta \leq \frac{1}{2} \theta^2$$

in (3.33), followed by the Poincaré inequality, we see that

$$\mathcal{E}(\theta) \geq \int_0^1 \left\{ \frac{1}{2} (\theta')^2 - \frac{1}{2} \theta^2 \right\} ds \geq \frac{1}{2} \left( 1 - \frac{\lambda}{\pi^2} \right) \int_0^1 (\theta')^2 ds.$$

It follows that  $\mathcal{E}(\theta) \geq 0$  if  $\lambda < \pi^2$ , and  $\theta = 0$  is the unique global minimizer of  $\mathcal{E}$  among functions that vanish at the endpoints.

As the parameter  $\lambda$  passes through each eigenvalue  $\lambda_0^{(n)}$ , the energy function develops another direction (tangent to the corresponding eigenvector) in which it decreases as  $\theta$  moves away from the critical point 0. These results are connected to conjugate points and Morse theory (see [37]).

The branches that bifurcate from  $\lambda_0^{(n)}$  for  $n \geq 2$  are of less interest than the first branch, because for  $\lambda > \lambda_0^{(1)}$  we expect the solution to lie on one of the stable branches that bifurcates from  $\lambda_0^{(1)}$  rather than on the trivial branch. We are then interested in secondary bifurcations of solutions from the stable branch rather than further bifurcations from the unstable trivial branch.

### 6.3. Solution by elliptic functions

Let us return to the solution of (3.22)–(3.23) in terms of elliptic functions.

The pendulum equation (3.22) has the first integral

$$\frac{1}{2} (\theta')^2 + \lambda (1 - \cos \theta) = 2\lambda k^2$$

where  $k$  is a constant of integration; equivalently

$$(\theta')^2 = 4\lambda \left( k^2 - \sin^2 \frac{\theta}{2} \right).$$

Thus, if  $\alpha$  is the maximum value of  $\theta$ , we have

$$(3.34) \quad k = \sin \frac{\alpha}{2}.$$

Solving for  $\theta'$ , separating variables and integrating, we get

$$\int \frac{d\theta}{\sqrt{k^2 - \sin^2(\theta/2)}} = 2\sqrt{\lambda}s.$$

Here, the sign of the square root is chosen appropriately, and we neglect the constant of integration, which can be removed by a translation of  $s$ . Making the substitution  $ku = \sin(\theta/2)$  in the integral, we get

$$(3.35) \quad \int \frac{du}{\sqrt{(1-u^2)(1-k^2u^2)}} = \sqrt{\lambda}s.$$

Trigonometric functions arise as inverse functions of integrals of the form

$$\int \frac{du}{\sqrt{p(u)}}$$

where  $p(u)$  is a quadratic polynomial. In an analogous way, elliptic functions arise as inverse functions of integrals of the same form where  $p(u)$  is a nondegenerate cubic or quartic polynomial. The Jacobi elliptic function  $u \mapsto \text{sn}(u, k)$ , with modulus  $k$ , has the inverse function

$$\text{sn}^{-1}(u, k) = \int_0^u \frac{dt}{\sqrt{(1-t^2)(1-k^2t^2)}}.$$

Rewriting  $u$  in terms of  $\theta$ , it follows from (3.35) that  $u = \text{sn}(\sqrt{\lambda}s, k)$ , so solutions  $\theta(s)$  of (3.22) with  $\theta(0) = 0$  are given by

$$\sin\left(\frac{\theta}{2}\right) = k \text{sn}\left(\sqrt{\lambda}s, k\right).$$

The arclength  $\ell$  of this solution from the endpoint  $\theta = 0$  to the maximum deflection angle  $\theta = \alpha$  is given by

$$\ell = \int_0^\ell ds = \int_0^\alpha \frac{d\theta}{\theta'}.$$

Using the substitution  $ku = \sin(\theta/2)$ , we get

$$\ell = \frac{1}{\sqrt{\lambda}}K(k)$$

where  $K(k)$  is the complete elliptic integral of the first kind, defined by

$$(3.36) \quad K(k) = \int_0^1 \frac{du}{\sqrt{(1-u^2)(1-k^2u^2)}}.$$

This solution satisfies the boundary condition  $\theta(1) = 0$  if  $\ell = 1/(2n)$  for some integer  $n = 1, 2, 3, \dots$ , meaning that

$$(3.37) \quad \lambda = 4n^2K^2(k).$$

This is the exact bifurcation equation for the  $n^{\text{th}}$  branch that bifurcates off the trivial solution.

A Taylor expansion of this equation agrees with the result from perturbation theory. From (3.36), we have, as  $k \rightarrow 0$ ,

$$K(k) = \int_0^1 \frac{du}{\sqrt{1-u^2}} + \frac{1}{2}k^2 \int_0^1 \frac{u^2 du}{\sqrt{1-u^2}} + \dots = \frac{\pi}{2} \left( 1 + \frac{1}{4}k^2 + \dots \right).$$

Also, from (3.34), we have  $k = \alpha/2 + \dots$ . It follows that (3.37) has the expansion

$$\lambda = n^2 \pi^2 \left( 1 + \frac{1}{4}k^2 + \dots \right)^2 = n^2 \pi^2 \left( 1 + \frac{1}{8}\alpha^2 + \dots \right),$$

in agreement with (3.32).

There are also solutions with nonzero winding number, meaning that  $\theta(0) = 0$  and  $\theta(1) = 2\pi N$  for some nonzero integer  $N$ . These cannot be reached from the zero solution along a continuous branch, since the winding number is a discrete topological invariant.

#### 6.4. Lyapounov-Schmidt reduction

The Lyapounov-Schmidt method provides a general approach to the rigorous derivation of local equilibrium bifurcation equations, based on an application of the implicit function theorem. We will outline the method and then explain how it applies to the buckling problem considered above. The main idea is to project the equation into two parts, one which can be solved uniquely and the other which gives the bifurcation equation.

Suppose that  $X, Y, \Lambda$  are Banach spaces, and  $F : X \times \Lambda \rightarrow Y$  is a smooth map (at least  $C^1$ ; see [14], for example, for more about derivatives of maps on Banach spaces). We are interested in solving the equation

$$(3.38) \quad F(x, \lambda) = 0$$

for  $x$  as  $\lambda$  varies in the parameter space  $\Lambda$ .

We denote the partial derivatives of  $F$  at  $(x, \lambda) \in X \times \Lambda$  by

$$F_x(x, \lambda) : X \rightarrow Y, \quad F_\lambda(x, \lambda) : \Lambda \rightarrow Y.$$

These are bounded linear maps such that

$$F_x(x, \lambda)h = \left. \frac{d}{d\varepsilon} F(x + \varepsilon h, \lambda) \right|_{\varepsilon=0}, \quad F_\lambda(x, \lambda)\eta = \left. \frac{d}{d\varepsilon} F(x, \lambda + \varepsilon\eta) \right|_{\varepsilon=0}.$$

Suppose that  $(x_0, \lambda_0)$  is a solution of (3.38), and denote by

$$L = F_x(x_0, \lambda_0) : X \rightarrow Y$$

the derivative of  $F(x, \lambda)$  with respect to  $x$  at  $(x_0, \lambda_0)$ .

The implicit function theorem states that if the bounded linear map  $L$  has an inverse  $L^{-1} : Y \rightarrow X$ , then (3.38) has a unique solution  $x = f(\lambda)$  in some neighborhood of  $(x_0, \lambda_0)$ . Moreover, the solution is at least as smooth as  $F$ , meaning that if  $F$  is  $C^k$  in a neighborhood of  $(x_0, \lambda_0)$ , then  $f$  is  $C^k$  in a neighborhood of  $\lambda_0$ . Thus, roughly speaking, the nonlinear problem is locally uniquely solvable if the linearized problem is uniquely solvable.

It follows that a necessary condition for new solutions of (3.38) to bifurcate off a solution branch  $x = f(\lambda)$  at  $(x_0, \lambda_0)$ , where  $x_0 = f(\lambda_0)$ , is that  $F_x(x_0, \lambda_0)$  is not invertible.

Consider such a point, and suppose that the non-invertible map  $L : X \rightarrow Y$  is a Fredholm operator. This means that: (a) the null-space of  $L$ ,

$$N = \{h \in X : Lh = 0\},$$

has finite dimension, and we can write  $X = M \oplus N$  where  $M, N$  are closed subspaces of  $X$ ; (b) the range of  $L$ ,

$$R = \{k \in Y : k = Lh \text{ for some } h \in X\},$$

has finite codimension, and  $Y = R \oplus S$  where  $R, S$  are closed subspaces of  $Y$ .

The condition that the range  $R$  of  $L$  is a closed subspace is satisfied automatically for maps on finite-dimensional spaces, but it is a significant assumption for maps on infinite-dimensional spaces. The condition that  $R$  has finite codimension simply means that any complementary space, such as  $S$ , has finite dimension (in which case the dimension does not depend on the choice of  $S$ ).

We write  $x \in X$  as  $x = m + n$  where  $m \in M$  and  $n \in N$ , and let

$$Q : Y \rightarrow Y$$

denote the projection onto  $R$  along  $S$ . That is, if  $y = r + s$  is the unique decomposition of  $y \in Y$  into a sum of  $r \in R$  and  $s \in S$ , then  $Qy = r$ . Since  $R$  is closed, the linear map  $Q$  is bounded.

Equation (3.38) is equivalent to the pair of equations obtained by projecting it onto the range of  $L$  and the complementary space:

$$(3.39) \quad QF(m + n, \lambda) = 0,$$

$$(3.40) \quad (I - Q)F(m + n, \lambda) = 0.$$

We write (3.39) as

$$G(m, \nu) = 0,$$

where  $\nu = (n, \lambda) \in \Gamma$ , with  $\Gamma = N \oplus \Lambda$ , and  $G : M \times \Gamma \rightarrow R$  is defined by

$$G(m, \nu) = QF(m + n, \lambda).$$

Let  $x_0 = m_0 + n_0$  and  $\nu_0 = (n_0, \lambda_0)$ , so  $(m_0, \nu_0) \in M \times \Gamma$  corresponds to  $(x_0, \lambda_0) \in X \times \Lambda$ . It follows from our definitions that the derivative of  $G$

$$G_m(m_0, \nu_0) : M \rightarrow R$$

is an invertible linear map between Banach spaces. The implicit function theorem then implies that that (3.39) has a unique local solution for  $m$  of the form

$$m = g(n, \lambda)$$

where  $g : N \times \Lambda \rightarrow M$ .

Using this expression for  $m$  in (3.40), we find that  $(n, \lambda)$  satisfies an equation of the form

$$(3.41) \quad \Phi(n, \lambda) = 0$$

where  $\Phi : N \oplus \Lambda \rightarrow S$  is defined locally by

$$\Phi(n, \lambda) = (I - Q)F(g(n, \lambda) + n, \lambda).$$

Equation (3.41) is the bifurcation equation for (3.38). It describes all solutions of (3.38) in a neighborhood of a point  $(x_0, \lambda_0)$  where the derivative  $F_x(x_0, \lambda_0)$  is singular.

This result is sometimes expressed in the following way. The  $m$ -component the solution  $x = m + n$  is 'slaved' to the  $n$ -component; thus, if we can solve the

bifurcation equation for  $n$  in terms of  $\lambda$ , then  $m$  is determined by  $n$ . This allows us to reduce a larger bifurcation problem for  $x \in X$  to a smaller bifurcation problem for  $n \in N$ .

If the null-space of  $L$  has dimension  $p$  and the range has codimension  $q$ , then (3.41) is equivalent to a system of  $p$  equations for  $q$  unknowns, depending on a parameter  $\lambda \in \Lambda$ . The integer  $p - q$  is called the Fredholm index of  $L$ . In the commonly occurring case when the Fredholm index of  $L$  is zero, the bifurcation equation is a  $p \times p$  system of equations. Thus, we can reduce bifurcation problems on infinite-dimensional spaces to ones on finite-dimensional spaces; the number of unknowns is equal to the dimension of the null space of  $L$  at the bifurcation point.

Next, we show how this method applies to the buckling problem. We write (3.22)–(3.23) as an equation  $F(\theta, \lambda) = 0$ , where  $F : X \times \mathbb{R} \rightarrow Y$  is given by

$$F(\theta, \lambda) = \theta'' + \lambda \sin \theta$$

and

$$X = \{\theta \in H^2(0, 1) : \theta(0) = 0, \theta(1) = 0\}, \quad Y = L^2(0, 1).$$

Here,  $H^2(0, 1)$  denotes the Sobolev space of functions whose weak derivatives of order less than or equal to 2 are square-integrable on  $(0, 1)$ . Functions in  $H^2(0, 1)$  are continuously differentiable on  $[0, 1]$ , so the boundary conditions make sense pointwise. Other function spaces, such as spaces of Hölder continuous functions, could be used equally well.

Consider bifurcations off the trivial solution  $\theta = 0$ . The derivative

$$L = F_\theta(0, \lambda_0)$$

is given by

$$Lh = h'' + \lambda_0 h.$$

This is singular on  $X$  if  $\lambda_0 = n^2\pi^2$  for some  $n \in \mathbb{N}$ , so these are the only possible bifurcation points.

In this case, the null-space  $N$  of  $L$  is one-dimensional:

$$N = \{\alpha \sin(n\pi s) : \alpha \in \mathbb{R}\}.$$

We take as a closed complementary space

$$M = \left\{ \varphi \in X : \int_0^1 \varphi(s) \sin(n\pi s) ds = 0 \right\}$$

The range  $R$  of  $L$  consists of the  $L^2$ -functions that are orthogonal to  $\sin(n\pi s)$ , meaning that

$$R = \left\{ \rho \in L^2(0, 1) : \int_0^1 \rho(s) \sin(n\pi s) ds = 0 \right\}.$$

As a complementary space, we take

$$S = \{\alpha \sin(n\pi s) : \alpha \in \mathbb{R}\}.$$

The projection  $Q : L^2(0, 1) \rightarrow L^2(0, 1)$  onto  $R$  is then given by

$$(Q\rho)(s) = \rho(s) - \left[ 2 \int_0^1 \rho(t) \sin(n\pi t) dt \right] \sin(n\pi s).$$

We write

$$\theta(s) = \varphi(s) + \alpha \sin(n\pi s)$$

where  $\alpha$  is an arbitrary constant and  $\varphi \in M$ , so that

$$(3.42) \quad \int_0^1 \varphi(s) \sin(n\pi s) ds = 0.$$

In this case, equation (3.39) becomes

$$(3.43) \quad \begin{aligned} & \varphi'' + \lambda \sin[\varphi + \alpha \sin(n\pi s)] \\ & - 2\lambda \left\{ \int_0^1 \sin[\varphi(t) + \alpha \sin(n\pi t)] \sin(n\pi t) dt \right\} \sin(n\pi s) = 0, \end{aligned}$$

subject to the boundary conditions  $\varphi(0) = 0$ ,  $\varphi(1) = 0$ , and the projection condition (3.42).

Equation (3.43) has the form  $G(\varphi, \alpha, \lambda) = 0$ , where  $G : M \times \mathbb{R} \times \mathbb{R} \rightarrow R$ . The derivative  $G_\varphi(0, 0, \lambda_0) : M \rightarrow R$  is given by

$$G_\varphi(0, 0, \lambda_0) h(s) = h''(s) + \lambda_0 \left[ h(s) - \left( 2 \int_0^1 \sin(n\pi t) h(t) dt \right) \sin(n\pi s) \right].$$

It is one-to-one and onto, and has a bounded inverse. Therefore we can solve (3.43) locally for  $\varphi(s) = g(s; \alpha, \lambda)$ . Equation (3.40) then gives the bifurcation equation

$$(3.44) \quad 2\lambda \int_0^1 \sin[g(s; \alpha, \lambda) + \alpha \sin(n\pi s)] \sin(n\pi s) ds - \alpha \lambda_0 = 0.$$

A Taylor expansion of (3.43)–(3.44) in  $(\alpha, \lambda)$  about  $(0, \lambda_0)$  gives the same results as before.

Finally, we remark that these results, which are based on linearization and Taylor expansion, are local. There are also topological methods in bifurcation theory, introduced by Krasnoselski (1956) and Rabinowitz (1971), that use degree theory and provide global, but less explicit, results.

## 7. Laplace's equation

One of the most important variational principles for a PDE is Dirichlet's principle for the Laplace equation. We will show how Dirichlet's principle leads to the Laplace equation and describe how it arises in the potential theory of electrostatic fields.

### 7.1. Dirichlet principle

Let  $\Omega \subset \mathbb{R}^n$  be a domain and  $u : \bar{\Omega} \rightarrow \mathbb{R}$  a function. We assume that the domain and the function are sufficiently smooth.

The Dirichlet integral of  $u$  over  $\Omega$  is defined by

$$(3.45) \quad \mathcal{F}(u) = \int_{\Omega} \frac{1}{2} |\nabla u|^2 dx.$$

Let us derive the Euler-Lagrange equation that must be satisfied by a minimizer of  $\mathcal{F}$ . To be specific, we consider a minimizer of  $\mathcal{F}$  in a space of functions that satisfy Dirichlet conditions

$$u = f \quad \text{on } \partial\Omega$$

where  $f$  is a given function defined on the boundary  $\partial\Omega$  of  $\Omega$ .

If  $h : \bar{\Omega} \rightarrow \mathbb{R}$  is a function such that  $h = 0$  on  $\partial\Omega$ , then

$$d\mathcal{F}(u)h = \frac{d}{d\varepsilon} \int_{\Omega} \frac{1}{2} |\nabla u + \varepsilon \nabla h|^2 dx \Big|_{\varepsilon=0} = \int_{\Omega} \nabla u \cdot \nabla h dx.$$

Thus, any minimizer of the Dirichlet integral must satisfy

$$(3.46) \quad \int_{\Omega} \nabla u \cdot \nabla h \, dx = 0$$

for all smooth functions  $h$  that vanish on the boundary.

Using the identity

$$\nabla \cdot (h \nabla u) = h \Delta u + \nabla u \cdot \nabla h$$

and the divergence theorem, we get

$$\int_{\Omega} \nabla u \cdot \nabla h \, dx = - \int_{\Omega} (\Delta u) h \, dx + \int_{\partial \Omega} h \frac{\partial u}{\partial n} \, dS.$$

Since  $h = 0$  on  $\partial \Omega$ , the integral over the boundary is zero, and we get

$$d\mathcal{F}(u) h = - \int_{\Omega} (\Delta u) h \, dx$$

Thus, the variational derivative of  $\mathcal{F}$ , defined by

$$d\mathcal{F}(u) h = \int_{\Omega} \frac{\delta \mathcal{F}}{\delta u} h \, dx,$$

is given by

$$\frac{\delta \mathcal{F}}{\delta u} = -\Delta u.$$

Therefore, a smooth minimizer  $u$  of  $\mathcal{F}$  satisfies Laplace's equation

$$(3.47) \quad \Delta u = 0.$$

This is the classical form of Laplace's equation, while (3.46) is the weak form.

Similarly, a minimizer of the functional

$$\mathcal{F}(u) = \int_{\Omega} \left\{ \frac{1}{2} |\nabla u|^2 - f u \right\} dx,$$

where  $f : \bar{\Omega} \rightarrow \mathbb{R}$  is a given function, satisfies Poisson's equation

$$-\Delta u = f.$$

We will study the Laplace and Poisson equations in more detail later on.

## 7.2. The direct method

One of the simplest ways to prove the existence of solutions of the Laplace equation (3.47), subject, for example, to Dirichlet boundary conditions to show directly the existence of minimizers of the Dirichlet integral (3.45). We will not give any details here but we will make a few comments (see [13] for more information).

It was taken more-or-less taken for granted by Dirichlet, Gauss, and Riemann that since the Dirichlet functional (3.45) is a quadratic functional of  $u$ , which is bounded from below by zero, it attains its minimum for some function  $u$ , as would be the cases for such functions on  $\mathbb{R}^n$ . Weierstrass pointed out that this argument requires a nontrivial proof for functionals defined on infinite-dimensional spaces, because the Heine-Borel theorem that a bounded set is (strongly) precompact is not true in that case.

Let us give a few simple one-dimensional examples which illustrate the difficulties that can arise.

**Example 3.18.** Consider the functional (Weierstrass, 1895)

$$(3.48) \quad \mathcal{F}(u) = \frac{1}{2} \int_{-1}^1 x^2 [u'(x)]^2 dx$$

defined on functions

$$u : [-1, 1] \rightarrow \mathbb{R} \text{ such that } u(-1) = -1, u(1) = 1.$$

This functional is quadratic and bounded from below by zero. Furthermore, its infimum over smooth functions that satisfy the boundary conditions is equal to zero. To show this, for instance, let

$$u^\varepsilon(x) = \frac{\tan^{-1}(x/\varepsilon)}{\tan^{-1}(1/\varepsilon)} \quad \text{for } \varepsilon > 0.$$

A straightforward computation gives

$$\mathcal{F}(u^\varepsilon) = \frac{\varepsilon}{\tan^{-1}(1/\varepsilon)} \rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0^+.$$

The Euler-Lagrange equation for (3.48) is

$$- [x^2 u']' = 0.$$

Solutions  $u^+$ ,  $u^-$  that satisfy the boundary conditions  $u^+(1) = 1$ ,  $u^-(-1) = -1$  have the form

$$u^+(x) = 1 + c^+ \left(1 - \frac{1}{x}\right), \quad u^-(x) = -1 + c^- \left(1 + \frac{1}{x}\right)$$

for some constants  $c^\pm$ . However, we cannot satisfy both boundary conditions for any choice the constants. Thus, there is no smooth, or even absolutely continuous, function that minimizes  $\mathcal{F}$ . Note that

$$\mathcal{F}(u) = \int_{-1}^1 F(x, u') dx, \quad F(x, p) = \frac{1}{2} x^2 p^2.$$

The integrand  $F(x, p)$  is a strictly convex function of  $p$  for each  $x$ , with

$$F_{pp}(x, p) = x^2 > 0,$$

except when  $x = 0$ . This loss of strict convexity at  $x = 0$  is what leads to the singular behavior of the solutions of the Euler-Lagrange equations and the lack of a minimizer.

**Example 3.19.** Consider the functional

$$\mathcal{F}(u) = \int_0^1 x^{2/3} [u']^2 dx$$

defined on functions  $u : [0, 1] \rightarrow \mathbb{R}$  with  $u(0) = 0$ ,  $u(1) = 1$ . The infimum is equal to zero. This infimum is attained for the function  $u(x) = x^{1/3}$ , which is not differentiable at  $x = 0$ . Thus, we cannot find a minimizer if we restrict the functional to  $C^1$ -functions; but we can find a minimizer on the larger class of absolutely continuous functions with weak derivative in  $L^1(0, 1)$ . The minimizer is Hölder continuous with exponent  $1/3$ .

**Example 3.20.** Consider the non-convex functional

$$\mathcal{F}(u) = \int_0^1 \left(1 - [u']^2\right)^2 dx$$

defined on functions  $u : [0, 1] \rightarrow \mathbb{R}$  with  $u(0) = 0$ ,  $u(1) = 0$ . The infimum is equal to zero. This infimum is not attained at any  $C^1$ -function, but it is attained at any ‘zig-zag’ Lipschitz continuous function that vanishes at the endpoints and whose derivative is equal to  $\pm 1$  almost everywhere. If we change the functional to

$$\mathcal{F}(u) = \int_0^1 \left\{ u^2 + \left(1 - [u']^2\right)^2 \right\} dx$$

then the infimum is still zero (as can be seen by taking a sequence of functions  $u_n$  with  $n$  ‘zig-zags’ and small  $L^\infty$ -norm). This infimum, however, is not attained by any absolutely continuous function, since we cannot simultaneously make  $|u'| = 1$  and  $u = 0$ . The difficulty here is associated with a lack of weak lower semicontinuity of the non-convex functional  $\mathcal{F}$ ; for example, for the ‘zig-zag’ functions, we have  $u_n \rightharpoonup 0$  in  $W^{1,1}(0, 1)$ , but  $\mathcal{F}(0) > \liminf_{n \rightarrow \infty} \mathcal{F}(u_n)$ .

These difficulties were resolved for the Dirichlet functional by Hilbert (1900) and Lebesgue (1907), and Hilbert included several problems in the calculus of variations among his list of 23 problems at the 1900 ICM in Paris.

The Dirichlet functional is defined provided that  $\nabla u$  is square-integrable. Thus, it is natural to look for minimizers of (3.45) in the Sobolev space  $H^1(\Omega)$  of Lebesgue measurable, square-integrable functions  $u : \Omega \rightarrow \mathbb{R}$  such that  $u \in L^2(\Omega)$ , meaning that  $\int_\Omega u^2(x) dx < \infty$ , with square-integrable weak derivatives  $\partial_{x^i} u \in L^2(\Omega)$ . If  $g : \partial\Omega \rightarrow \mathbb{R}$  is a given boundary value that is attained by some function in  $H^1(\Omega)$ , then one can prove that there is a unique minimizer of (3.45) in the space

$$X = \{u \in H^1(\Omega) : \text{such that } u = g \text{ on } \partial\Omega\}.$$

The definition of the boundary values, or trace, of Sobolev functions requires a more careful discussion, but we will not go into the details here.

A further central issue in the calculus of variations is the regularity of minimizers. It is possible to prove that the minimizer of the Dirichlet functional is, in fact, a smooth function with continuous derivative of all orders inside  $\Omega$ . In particular, it follows that it is a classical solution of Laplace’s equation. Furthermore, if the boundary data and the domain are smooth, then the solution is also smooth on  $\overline{\Omega}$ .

### 7.3. Electrostatics

As an example of a physical problem leading to potential theory, consider a static electric field in a dielectric medium. (A dielectric medium is simply an insulator that does not conduct electricity, such as glass, air, or a vacuum.) We suppose that the dielectric has a charge-density  $\rho(\vec{x})$ , and that there is no magnetic field.

The electrostatic properties of the dielectric are characterized by two vector fields, the electric field  $\vec{E}(\vec{x})$  and the electric displacement  $\vec{D}(\vec{x})$ . According to Maxwell’s equations, these satisfy [28]

$$(3.49) \quad \text{curl } \vec{E} = 0,$$

$$(3.50) \quad \text{div } \vec{D} = \rho.$$

The integral form of these balance laws is

$$(3.51) \quad \int_{\Gamma} \vec{E} \cdot d\vec{x} = 0,$$

$$(3.52) \quad \int_{\partial\Omega} \vec{D} \cdot \vec{n} d\vec{x} = \int_{\Omega} \rho d\vec{x},$$

for any closed curve  $\Gamma$  and any bounded volume  $\Omega$ .

Equation (3.51) states that the circulation of  $\vec{E}$  around the closed curve  $\Gamma$  is equal to zero, since by Stokes' theorem it is equal to the flux of  $\text{curl } \vec{E}$  through a surface bounded by  $\Gamma$ . Equation (3.52) states that the flux of  $\vec{D}$  through a closed surface  $\partial\Omega$  is equal to the total charge in the enclosed volume  $\Omega$ .

On a simply connected domain, equation(3.49) implies that

$$\vec{E} = -\nabla\Phi.$$

for a suitable potential  $\Phi(\vec{x})$ .

The electric displacement is related to the electric field by a constitutive relation, which describes the response of the dielectric medium to an applied electric field. We will assume that it has the simplest linear, isotropic form

$$\vec{E} = \epsilon\vec{D}$$

where  $\epsilon$  is a constant, called the dielectric constant, or electric permittivity, of the medium. In a linear, anisotropic medium,  $\epsilon$  becomes a tensor; for large electric fields, it may be necessary to use a nonlinear constitutive relation.

It follows from these equations and (3.50) that  $\Phi$  satisfies Poisson's equation

$$(3.53) \quad -\epsilon\Delta\Phi = \rho.$$

This equation is supplemented by boundary conditions; for example, we require that  $\Phi$  is constant on a conducting boundary, and the normal derivative of  $\Phi$  is zero on an insulating boundary.

The energy of the electrostatic field in some region  $\Omega \subset \mathbb{R}^3$  is

$$\mathcal{E} = \int_{\Omega} \left\{ \frac{1}{2} \vec{E} \cdot \vec{D} - \rho\Phi \right\} d\vec{x} = \int_{\Omega} \left\{ \frac{1}{2} \epsilon |\nabla\Phi|^2 - \rho\Phi \right\} d\vec{x}.$$

The term proportional to  $\vec{E} \cdot \vec{D}$  is the energy of the field, while the term proportional to  $\rho\Phi$  is the work required to bring the charge distribution to the potential  $\Phi$ .

The potential  $\Phi$  minimizes this functional, and the condition that  $\mathcal{E}(\Phi)$  is stationary with respect to variations in  $\Phi$  leads to (3.53).

## 8. The Euler-Lagrange equation

A similar derivation of the Euler-Lagrange equation as the condition satisfied by a smooth stationary point applies to more general functionals. For example, consider the functional

$$\mathcal{F}(u) = \int_{\Omega} F(x, u, \nabla u) dx$$

where  $\Omega \subset \mathbb{R}^n$  and  $u : \bar{\Omega} \rightarrow \mathbb{R}^m$ . Then, writing

$$x = (x^1, x^2, \dots, x^n), \quad u = (u^1, u^2, \dots, u^m),$$

denoting a derivative with respect to  $x^j$  by  $\partial_j$ , and using the summation convention, we have

$$d\mathcal{F}(u)h = \int_{\Omega} \{F_{u^i}(x, u, \nabla u) h^i + F_{\partial_j u^i}(x, u, \nabla u) \partial_j h^i\} dx.$$

Thus,  $u$  is a stationary point of  $\mathcal{F}$  if

$$(3.54) \quad \int_{\Omega} \{F_{u^i}(x, u, \nabla u) h^i + F_{\partial_j u^i}(x, u, \nabla u) \partial_j h^i\} dx = 0$$

for all smooth test functions  $h : \bar{\Omega} \rightarrow \mathbb{R}$  that vanish on the boundary.

Using the divergence theorem, we find that

$$d\mathcal{F}(u)h = \int_{\Omega} \{F_{u^i}(x, u, \nabla u) - \partial_j [F_{\partial_j u^i}(x, u, \nabla u)]\} h^i dx.$$

Thus,

$$\frac{\delta \mathcal{F}}{\delta h^i} = -\partial_j [F_{\partial_j u^i}(x, u, \nabla u)] + F_{u^i}(x, u, \nabla u),$$

and a smooth stationary point  $u$  satisfies

$$-\partial_j [F_{\partial_j u^i}(x, u, \nabla u)] + F_{u^i}(x, u, \nabla u) = 0 \quad \text{for } i = 1, 2, \dots, n.$$

The weak form of this equation is (3.54).

### 8.1. The minimal surface equation

Suppose that a surface over a domain  $\Omega \subset \mathbb{R}^n$  is the graph of a smooth function  $z = u(x)$ , where  $u : \bar{\Omega} \rightarrow \mathbb{R}$ . The area  $\mathcal{A}$  of the surface is

$$\mathcal{A}(u) = \int_{\Omega} \sqrt{1 + |\nabla u|^2} dx.$$

The problem of finding a surface of minimal area that spans a given curve  $z = g(x)$  over the boundary, where  $g : \partial\Omega \rightarrow \mathbb{R}$ , is called Plateau's problem. Any smooth minimizer of the area functional  $\mathcal{A}(u)$  must satisfy the Euler-Lagrange equation, called the minimal surface problem,

$$\nabla \cdot \left[ \frac{\nabla u}{\sqrt{1 + |\nabla u|^2}} \right] = 0.$$

As a physical example, a film of soap has energy per unit area equal to its surface tension. Thus, a soap film on a wire frame is a minimal surface.

A full analysis of this problem is not easy. The PDE is elliptic, but it is nonlinear and it is not uniformly elliptic, and it has motivated a large amount of work on quasilinear elliptic PDEs. See [13] for more information.

### 8.2. Nonlinear elasticity

Consider an equilibrium deformation of an elastic body. We label material points by their location  $\vec{x} \in \mathcal{B}$  in a suitable reference configuration  $\mathcal{B} \subset \mathbb{R}^n$ . A deformation is described by an invertible function  $\vec{\varphi} : \mathcal{B} \rightarrow \mathbb{R}^n$ , where  $\vec{\varphi}(\vec{x})$  is the location of the material point  $\vec{x}$  in the deformed configuration of the body.

The deformation gradient

$$\mathbf{F} = \nabla \vec{\varphi}, \quad F_{ij} = \frac{\partial \varphi_i}{\partial x_j}$$

gives a linearized approximation of the deformation at each point, and therefore describes the local strain and rotation of the deformation.

An elastic material is said to be hyperelastic if the work required to deform it, per unit volume in the reference configuration, is given by a scalar-valued strain energy function. Assuming, for simplicity, that the body is homogeneous so the work does not depend explicitly on  $\vec{x}$ , the strain energy is a real-valued function  $W(\mathbf{F})$  of the deformation gradient. In the absence of external forces, the total energy of a deformation  $\vec{\varphi}$  is given by

$$\mathcal{W}(\vec{\varphi}) = \int_{\mathcal{B}} W(\nabla\vec{\varphi}(\vec{x})) \, d\vec{x}.$$

Equilibrium deformations are minimizers of the total energy, subject to suitable boundary conditions. Therefore, smooth minimizers satisfy the Euler-Lagrange equations

$$\nabla \cdot \mathbf{S}(\nabla\vec{\varphi}(\vec{x})) = 0, \quad \frac{\partial S_{ij}}{\partial x_j} = 0 \quad i = 1, \dots, n,$$

where we use the summation convention in the component form of the equations, and  $\mathbf{S}(\mathbf{F})$  is the Piola-Kirchoff stress tensor, given by

$$\mathbf{S} = \nabla_{\mathbf{F}} W, \quad S_{ij} = \frac{\partial W}{\partial F_{ij}}.$$

This is an  $n \times n$  nonlinear system of second-order PDEs for  $\vec{\varphi}(\vec{x})$ .

There are restrictions on how the strain energy  $W$  depends on  $\mathbf{F}$ . The principle of material frame indifference [25] implies that, for any material,

$$W(\mathbf{R}\mathbf{F}) = W(\mathbf{F})$$

for all orthogonal transformations  $\mathbf{R}$ . If the elastic material is isotropic, then

$$W(\mathbf{F}) = \tilde{W}(\mathbf{B})$$

depends only on the left Cauchy-Green strain tensor  $\mathbf{B} = \mathbf{F}\mathbf{F}^\top$ , and, in fact, only on the principle invariants of  $\mathbf{B}$ .

The constitutive restriction imply that  $W$  is not a convex function of  $\mathbf{F}$ . This creates a difficulty in the proof of the existence of minimizers by the use of direct methods, because one cannot use convexity to show that  $\mathcal{W}$  is weakly lower semicontinuous.

This difficult was overcome by Ball (1977). He observed that one can write

$$W(\mathbf{F}) = \hat{W}(\mathbf{F}, \text{cof } \mathbf{F}, \det \mathbf{F})$$

where  $\hat{W}$  is a convex function of  $\mathbf{F}$  and the cofactors  $\text{cof } \mathbf{F}$  of  $\mathbf{F}$ , including its determinant. A function with this property is said to be *polyconvex*.

According to the theory of compensated compactness, given suitable bounds on the derivatives of  $\vec{\varphi}$ , the cofactors of  $\mathbf{F}$  are weakly continuous, which is a very unusual property for nonlinear functions. Using this fact, combined with the observation that the strain energy is polyconvex, Ball was able to prove the existence of minimizers for nonlinear hyperelasticity.

### 9. The wave equation

Consider the motion of a medium whose displacement may be described by a scalar function  $u(x, t)$ , where  $x \in \mathbb{R}^n$  and  $t \in \mathbb{R}$ . For example, this function might represent the transverse displacement of a membrane  $z = u(x, y, t)$ .

Suppose that the kinetic energy  $\mathcal{T}$  and potential energy  $\mathcal{V}$  of the medium are given by

$$\mathcal{T}(u_t) = \frac{1}{2} \int \rho_0 u_t^2 dx, \quad \mathcal{V}(u) = \frac{1}{2} \int k |\nabla u|^2 dx,$$

where  $\rho_0(\vec{x})$  is a mass-density and  $k(\vec{x})$  is a stiffness, both assumed positive. The Lagrangian  $\mathcal{L} = \mathcal{T} - \mathcal{V}$  is

$$(3.55) \quad \mathcal{L}(u, u_t) = \int \frac{1}{2} \left\{ \rho_0 u_t^2 - k |\nabla u|^2 \right\} dx,$$

and the action — the time integral of the Lagrangian — is

$$\mathcal{S}(u) = \int \int \frac{1}{2} \left\{ \rho_0 u_t^2 - k |\nabla u|^2 \right\} dx dt.$$

Note that the kinetic and potential energies and the Lagrangian are functionals of the spatial field and velocity  $u(\cdot, t)$ ,  $u_t(\cdot, t)$  at each fixed time, whereas the action is a functional of the space-time field  $u(x, t)$ , obtained by integrating the Lagrangian with respect to time.

The Euler-Lagrange equation satisfied by a stationary point of this action is

$$(3.56) \quad \rho_0 u_{tt} - \nabla \cdot (k \nabla u) = 0.$$

If  $\rho_0, k$  are constants, then

$$(3.57) \quad u_{tt} - c_0^2 \Delta u = 0,$$

where  $c_0^2 = k/\rho_0$ . This is the linear wave equation with wave-speed  $c_0$ .

Unlike the energy for Laplace's equation, the action functional for the wave equation is not positive definite. We therefore cannot expect a solution of the wave equation to be a minimizer of the action, in general, only a critical point. As a result, direct methods are harder to implement for the wave equation (and other hyperbolic PDEs) than they are for Laplace's equation (and other elliptic PDEs), although there are 'mountain-pass' lemmas that can be used to establish the existence of stationary points. Moreover, in general, stationary points of functionals do not have the increased regularity that minimizers of convex functionals typically possess.

### 10. Hamiltonian mechanics

Let us return to the motion of a particle in a conservative force field considered in Section 1. We will give an alternative, Hamiltonian, formulation of its equations of motion.

Given a Lagrangian

$$L(\vec{x}, \vec{v}) = \frac{1}{2} m |\vec{v}|^2 - V(\vec{x}),$$

we define the momentum  $\vec{p}$  by

$$(3.58) \quad \vec{p} = \frac{\partial L}{\partial \vec{v}},$$

meaning that  $\vec{p} = m\vec{v}$ . Here, we use the notation

$$\frac{\partial L}{\partial \vec{v}} = \left( \frac{\partial L}{\partial v_1}, \frac{\partial L}{\partial v_2}, \dots, \frac{\partial L}{\partial v_n} \right)$$

to denote the derivative with respect to  $\vec{v}$ , keeping  $\vec{x}$  fixed, with a similar notation for the derivative  $\partial/\partial \vec{p}$  with respect to  $\vec{p}$ , keeping  $\vec{x}$  fixed. The derivative  $\partial/\partial \vec{x}$  is taken keeping  $\vec{v}$  or  $\vec{p}$  fixed, as appropriate.

We then define the Hamiltonian function  $H$  by

$$(3.59) \quad H(\vec{x}, \vec{p}) = \vec{p} \cdot \vec{v} - L(\vec{x}, \vec{v}),$$

where we express  $\vec{v} = \vec{p}/m$  on the right hand side in terms of  $\vec{p}$ . This gives

$$H(\vec{x}, \vec{p}) = \frac{1}{2m} \vec{p} \cdot \vec{p} + V(\vec{x}).$$

Thus, we transform  $L$  as a function of  $\vec{v}$  into  $H$  as a function of  $\vec{p}$ . The variable  $\vec{x}$  plays the role of a parameter in this transformation. The function  $H(\vec{x}, \vec{p})$ , given by (3.58)–(3.59) is the Legendre transform of  $L(\vec{x}, \vec{v})$  with respect to  $\vec{v}$ ; conversely,  $L$  is the Legendre transform of  $H$  with respect to  $\vec{p}$ .

Note that the the Hamiltonian is the total energy of the particle,

$$H(\vec{x}, \vec{p}) = T(\vec{p}) + V(\vec{x}),$$

where  $T$  is the kinetic energy expressed as a function of the momentum

$$T(\vec{p}) = \frac{1}{2m} |\vec{p}|^2.$$

The Lagrangian equation of motion (3.1) may then be written as a first order system for  $(\vec{x}, \vec{p})$ :

$$\dot{\vec{x}} = \frac{1}{m} \vec{p}, \quad \dot{\vec{p}} = -\frac{\partial V}{\partial \vec{x}}.$$

This system has the canonical Hamiltonian form

$$(3.60) \quad \dot{\vec{x}} = \frac{\partial H}{\partial \vec{p}}, \quad \dot{\vec{p}} = -\frac{\partial H}{\partial \vec{x}}.$$

Equation (3.60) is a  $2n$ -dimensional system of first-order equations. We refer to the space  $\mathbb{R}^{2n}$ , with coordinates  $(\vec{x}, \vec{p})$ , as the phase space of the system.

### 10.1. The Legendre transform

The above transformation, from the Lagrangian as a function of velocity to the Hamiltonian as a function of momentum, is an example of a Legendre transform. In that case, the functions involved were quadratic.

More generally, if  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we define the Legendre transform  $f^*(x^*)$  of the function  $f(x)$  as follows. Let

$$x^* = \frac{\partial f}{\partial x}(x),$$

and suppose we can invert this equation to get  $x = x(x^*)$ . This is the case, for example, if  $f$  is a smooth, convex function. We then define

$$f^*(x^*) = x^* \cdot x(x^*) - f(x(x^*)).$$

Note that, by the chain rule and the definition of  $x^*$ , we have

$$\frac{\partial f^*}{\partial x^*} = x + x^* \cdot \frac{\partial x}{\partial x^*} - \frac{\partial f}{\partial x} \frac{\partial x}{\partial x^*} = x + x^* \cdot \frac{\partial x}{\partial x^*} - x^* \cdot \frac{\partial x}{\partial x^*} = x,$$

and, from the definition of  $f^*$ ,

$$f(x) = x \cdot x^*(x) - f^*(x^*(x)).$$

Thus, if  $f$  is convex, the Legendre transform of  $f^*(x^*)$  is the original function  $f(x)$  (see [43] for more on convex analysis).

Consider a Lagrangian  $F(x, u, u')$ , where  $u : [a, b] \rightarrow \mathbb{R}^n$  and  $F : [a, b] \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ . Taking the Legendre transform of  $F$  with respect to the  $u'$ -variable, we get

$$p = \frac{\partial F}{\partial u'}, \quad H(x, u, p) = p \cdot u' - F(x, u, u').$$

It follows that

$$\begin{aligned} \frac{\partial H}{\partial u} &= p \cdot \frac{\partial u'}{\partial u} - \frac{\partial F}{\partial u} - \frac{\partial F}{\partial u'} \frac{\partial u'}{\partial u} = -\frac{\partial F}{\partial u}, \\ \frac{\partial H}{\partial p} &= u' + p \cdot \frac{\partial u'}{\partial p} - \frac{\partial F}{\partial u'} \frac{\partial u'}{\partial p} = u'. \end{aligned}$$

Hence, the Euler-Lagrange equation

$$-\frac{d}{dx} \left( \frac{\partial F}{\partial u'} \right) + \frac{\partial F}{\partial u} = 0$$

may be written as a Hamiltonian system

$$u' = \frac{\partial H}{\partial p}, \quad p' = -\frac{\partial H}{\partial u}.$$

In general, the Hamiltonian in these equations may depend on the independent variable  $x$  (or  $t$  in the mechanical problem above) as well as the dependent variables. For simplicity, we will consider below Hamiltonians that do not depend explicitly on the independent variable.

## 10.2. Canonical coordinates

It is important to recognize that there are two ingredients in the Hamiltonian system (3.60). One is obvious: the Hamiltonian function  $H(\vec{x}, \vec{p})$  itself. The other, less obvious, ingredient is a Hamiltonian structure that allows us to map the differential of a Hamiltonian function

$$dH = \frac{\partial H}{\partial \vec{x}} d\vec{x} + \frac{\partial H}{\partial \vec{p}} d\vec{p}$$

to the Hamiltonian vector field

$$X_H = \frac{\partial H}{\partial \vec{p}} \frac{\partial}{\partial \vec{x}} - \frac{\partial H}{\partial \vec{x}} \frac{\partial}{\partial \vec{p}}$$

that appears in (3.60).

We will not describe the symplectic geometry of Hamiltonian systems in any detail here (see [6] for more information, including an introduction to differential forms) but we will make a few comments to explain the role of canonical coordinates  $(\vec{x}, \vec{p})$  in the formulation of Hamiltonian equations.

The Hamiltonian structure of (3.60) is defined by a symplectic two-form

$$(3.61) \quad \omega = d\vec{x} \wedge d\vec{p}$$

on the phase space  $\mathbb{R}^{2n}$ . More generally, one can consider symplectic manifolds, which are manifolds, necessarily even-dimensional, equipped with a closed, nondegenerate two-form  $\omega$ .

The two-form (3.61) can be integrated over a two-dimensional submanifold  $S$  of  $\mathbb{R}^{2n}$  to give an ‘area’

$$\int_S \omega = \sum_{i=1}^n \int_S dx^i \wedge dp_i.$$

Roughly speaking, this integral is the sum of the oriented areas of the projections of  $S$ , counted according to multiplicity, onto the  $(x^i, p_i)$ -coordinate planes. Thus, the phase space of a Hamiltonian system has a notion of oriented area, defined by the *skew-symmetric* two-form  $\omega$ . In a somewhat analogous way, Euclidean space (or a Riemannian manifold) has a notion of length and angle, which is defined by a *symmetric* two-form, the metric  $g$ . The geometry of symplectic manifolds  $(M, \omega)$  is, however, completely different from the more familiar geometry of Riemannian manifolds  $(M, g)$ .

According to Darboux’s theorem, if  $\omega$  is a closed nondegenerate two-form, then there are local coordinates  $(\vec{x}, \vec{p})$  in which it is given by (3.61). Such coordinates are called canonical coordinates, and Hamilton’s equations take the canonical form (3.60) for every Hamiltonian  $H$  in any canonical system of coordinates. The canonical form of  $\omega$  and Hamilton’s equations, however, are not preserved under arbitrary transformations of the dependent variables.

A significant part of the theory of Hamiltonian systems, such as Hamilton-Jacobi theory, is concerned with finding canonical transformations that simplify Hamilton’s equations. For example, if, for a given Hamiltonian  $H(\vec{x}, \vec{p})$ , we can find a canonical change of coordinates such that

$$(\vec{x}, \vec{p}) \mapsto (\vec{x}', \vec{p}'), \quad H(\vec{x}, \vec{p}) \mapsto H(\vec{p}'),$$

meaning that the transformed Hamiltonian is independent of the position variable  $\vec{x}'$ , then we can solve the corresponding Hamiltonian equations explicitly. It is typically not possible to do this, but the completely integrable Hamiltonian systems for which it is possible form an important and interesting class of solvable equations. We will not discuss these ideas further here (see [24] for more information).

## 11. Poisson brackets

It can be inconvenient to use conjugate variables, and in some problems it may be difficult to identify which variables form conjugate pairs. The Poisson bracket provides a way to write Hamiltonian systems, as well as odd-order generalizations of the even-order canonical systems, which does not require the use of canonical variables. The Poisson bracket formulation is also particularly convenient for the description of Hamiltonian PDEs.

First we describe the Poisson-bracket formulation of the canonical equations. Let  $\vec{u} = (\vec{x}, \vec{p})^\top \in \mathbb{R}^{2n}$ . Then we may write (3.60) as

$$\dot{\vec{u}} = \mathbf{J} \frac{\partial H}{\partial \vec{u}}$$

where  $\mathbf{J} : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$  is the constant skew-symmetric linear map with matrix

$$(3.62) \quad \mathbf{J} = \begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix}.$$

If  $F, G : \mathbb{R}^{2n} \rightarrow \mathbb{R}$  are smooth functions, then we define their Poisson bracket  $\{F, G\}$ , which is also a function  $\{F, G\} : \mathbb{R}^{2n} \rightarrow \mathbb{R}$ , by

$$\{F, G\} = \frac{\partial F}{\partial \vec{u}} \cdot \mathbf{J} \frac{\partial G}{\partial \vec{u}}.$$

In terms of derivatives with respect to  $(\vec{x}, \vec{p})$ , the bracket is given by

$$(3.63) \quad \{F, G\} = \frac{\partial F}{\partial \vec{x}} \cdot \frac{\partial G}{\partial \vec{p}} - \frac{\partial F}{\partial \vec{p}} \cdot \frac{\partial G}{\partial \vec{x}} = \sum_{i=1}^n \left( \frac{\partial F}{\partial x^i} \frac{\partial G}{\partial p_i} - \frac{\partial F}{\partial p_i} \frac{\partial G}{\partial x^i} \right)$$

Hamilton's equations may be written as

$$\dot{\vec{u}} = \{\vec{u}, H\},$$

or, in component form,

$$u^i = \{u^i, H\} \quad 1 \leq i \leq 2n.$$

Moreover, if  $F(\vec{u})$  is any function, then

$$\dot{F} = \frac{\partial F}{\partial \vec{x}} \dot{\vec{x}} + \frac{\partial F}{\partial \vec{p}} \dot{\vec{p}} = \frac{\partial F}{\partial \vec{x}} \frac{\partial H}{\partial \vec{p}} - \frac{\partial F}{\partial \vec{p}} \frac{\partial H}{\partial \vec{x}}.$$

It follows that

$$\dot{F} = \{F, H\}.$$

Thus, a function  $F(\vec{x}, \vec{p})$  that does not depend explicitly on time  $t$  is a conserved quantity for Hamilton's equations if its Poisson bracket with the Hamiltonian vanishes; for example, the Poisson bracket of the Hamiltonian with itself vanishes, so the Hamiltonian is conserved.

The Poisson bracket in (3.63) has the properties that for any functions  $F, G, H$  and constants  $a, b$

$$(3.64) \quad \{F, G\} = -\{G, F\},$$

$$(3.65) \quad \{aF + bG, H\} = a\{F, H\} + b\{G, H\},$$

$$(3.66) \quad \{FG, H\} = F\{G, H\} + \{F, H\}G,$$

$$(3.67) \quad \{F, \{G, H\}\} + \{G, \{H, F\}\} + \{H, \{F, G\}\} = 0.$$

That is, it is skew-symmetric (3.64), bilinear (3.65), a derivation (3.66), and satisfies the Jacobi identity (3.67).

Any bracket with these properties that maps a pair of smooth functions  $F, G$  to a smooth function  $\{F, G\}$  defines a Poisson structure. The bracket corresponding to the matrix  $\mathbf{J}$  in (3.62) is the canonical bracket, but there are many other brackets. In particular, the skew-symmetric linear operator  $\mathbf{J}$  can depend on  $u$ , provided that the associated bracket satisfies the Jacobi identity.

## 12. Rigid body rotations

Consider a rigid body, such as a satellite, rotating about its center of mass in three space dimensions.

We label the material points of the body by their position  $\vec{a} \in \mathcal{B}$  in a given reference configuration  $\mathcal{B} \subset \mathbb{R}^3$ , and denote the mass-density of the body by

$$\rho : \mathcal{B} \rightarrow [0, \infty).$$

We use coordinates such that the center of mass of the body is at the origin, so that

$$\int_{\mathcal{B}} \rho(\vec{a}) \vec{a} d\vec{a} = 0.$$

Here,  $d\vec{a}$  denotes integration with respect to volume in the reference configuration.

The possible configurations of the body are rotations of the reference configuration, so the configuration space of the body may be identified with the rotation group. This is the special orthogonal group  $SO(3)$  of linear transformations  $R$  on  $\mathbb{R}^3$  such that

$$R^\top R = I, \quad \det R = 1.$$

The first condition is the orthogonality condition,  $R^\top = R^{-1}$ , which ensures that  $R$  preserves the Euclidean inner product of vectors, and therefore lengths and angles. The second condition restricts  $R$  to the ‘special’ transformations with determinant one. It rules out the orientation-reversing orthogonal transformations with  $\det R = -1$ , which are obtained by composing a reflection and a rotation.

First, we will define the angular velocity and angular momentum of the body in a spatial reference frame. Then we will ‘pull back’ these vectors to the body reference frame, in which the equations of motion simplify.

### 12.1. Spatial description

Consider the motion of the body in an inertial frame of reference whose origin is at the center of mass of the body. The position vector  $\vec{x}$  of a point  $\vec{a} \in \mathcal{B}$  at time  $t$  is given by

$$(3.68) \quad \vec{x}(\vec{a}, t) = R(t)\vec{a}$$

where  $R(t) \in SO(3)$  is a rotation. Thus, the motion of the rigid body is described by a curve of rotations  $R(t)$  in the configuration space  $SO(3)$ .

Differentiating (3.68) with respect to  $t$ , and using (3.68) in the result, we find that the velocity

$$\vec{v}(\vec{a}, t) = \dot{\vec{x}}(\vec{a}, t)$$

of the point  $\vec{a}$  is given by

$$(3.69) \quad \vec{v} = w\vec{x},$$

where

$$(3.70) \quad w = \dot{R}R^\top.$$

Differentiation of the equation  $RR^\top = I$  with respect to  $t$  implies that

$$\dot{R}R^\top + R\dot{R}^\top = 0.$$

Thus,  $w$  in (3.70) is skew-symmetric, meaning that  $w^\top = -w$ .

If  $W : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  is a skew-symmetric linear map on three-dimensional Euclidean space, then there is a unique vector  $\vec{\Omega} \in \mathbb{R}^3$  such that

$$(3.71) \quad W\vec{x} = \vec{\Omega} \times \vec{x}.$$

We denote this correspondence by  $\vec{\Omega} = \hat{W}$ . With respect to a right-handed orthonormal basis, the matrix of  $W$  and the components of  $\vec{\Omega}$  are related by

$$\begin{pmatrix} 0 & -\Omega_3 & \Omega_2 \\ \Omega_3 & 0 & -\Omega_1 \\ -\Omega_2 & \Omega_1 & 0 \end{pmatrix} \longleftrightarrow \begin{pmatrix} \Omega_1 \\ \Omega_2 \\ \Omega_3 \end{pmatrix}$$

We let  $\vec{\omega}(t) = \hat{w}(t)$  denote the vector associated with  $w$  in (3.70). Then, from (3.69), the velocity of the body in the spatial frame is

$$(3.72) \quad \vec{v} = \vec{\omega} \times \vec{x}.$$

Thus, the vector  $\vec{\omega}(t)$  is the angular velocity of the body at time  $t$ .

The angular momentum  $\pi$ , or moment of momentum, of the body is defined by

$$(3.73) \quad \pi(t) = \int_{\mathcal{B}} \rho(\vec{a}) [\vec{x}(\vec{a}, t) \times \vec{v}(\vec{a}, t)] d\vec{a}$$

Equivalently, making the change of variables  $\vec{a} \mapsto \vec{x}$  in (3.68) in the integral, whose Jacobian is equal to one, we get

$$\pi = \int_{\mathcal{B}_t} \rho(R^\top(t)\vec{a}) [\vec{x} \times (\vec{\omega}(t) \times \vec{x})] d\vec{x},$$

where  $\mathcal{B}_t = \vec{x}(\mathcal{B}, t)$  denotes the region occupied by the body at time  $t$ .

Conservation of angular momentum implies that, in the absence of external forces and couples,

$$(3.74) \quad \dot{\vec{\pi}} = 0.$$

This equation is not so convenient to solve for the motion of the body, because the angular momentum  $\vec{\pi}$  depends in a somewhat complicated way on the angular velocity  $\vec{\omega}$  and the rotation matrix  $R$ . We will rewrite it with respect to quantities defined with respect to the body frame, which leads to a system of ODEs for the angular momentum, or angular velocity, in the body frame.

## 12.2. Body description

The spatial coordinate  $\vec{x}$  is related to the body coordinate  $\vec{a}$  by  $\vec{x} = R\vec{a}$ . Similarly, we define a body frame velocity  $\vec{V}(\vec{a}, t)$ , angular velocity  $\vec{\Omega}(t)$ , and angular momentum  $\vec{\Pi}(t)$  in terms of the corresponding spatial vectors by

$$(3.75) \quad \vec{v} = R\vec{V}, \quad \vec{\omega} = R\vec{\Omega}, \quad \vec{\pi} = R\vec{\Pi}.$$

Thus, we rotate the spatial vectors back to the body frame.

First, from (3.69), we find that if  $\vec{v} = R\vec{V}$ , then

$$(3.76) \quad \vec{V} = W\vec{a}$$

where  $W$  is the skew-symmetric map

$$(3.77) \quad W = R^\top \dot{R}.$$

Therefore, denoting by  $\vec{\Omega} = \hat{W}$  the vector associated with  $W$ , we have

$$(3.78) \quad \vec{V} = \vec{\Omega} \times \vec{a}.$$

Since  $w = RWR^\top$ , it follows that  $\vec{\omega} = R\vec{\Omega}$ , as in (3.75).

Next, since rotations preserve cross-products, we have

$$\vec{x} \times \vec{v} = R(\vec{a} \times \vec{V}).$$

Using this equation, followed by (3.78), in (3.73), we find that  $\vec{\pi} = R\vec{\Pi}$  where

$$\vec{\Pi}(t) = \int_{\mathcal{B}} \rho(\vec{a}) [\vec{a} \times (\vec{\Omega}(t) \times \vec{a})] d\vec{a}.$$

This equation is a linear relation between the angular velocity and angular momentum. We write it as

$$(3.79) \quad \vec{\Pi} = \mathbf{I}\vec{\Omega}.$$

where  $\mathbf{I} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  is a constant linear map depending only on the mass distribution of the body. It is called the *inertia tensor*.

An explicit expression for the inertia tensor is

$$(3.80) \quad \mathbf{I} = \int_{\mathcal{B}} \rho(\vec{a}) [(\vec{a} \cdot \vec{a}) I - \vec{a} \otimes \vec{a}] d\vec{a},$$

where  $I$  denotes the identity transformation, or, in components,

$$I_{ij} = \int_{\mathcal{B}} \rho(\vec{a}) [(a_k a_k) \delta_{ij} - a_i a_j] d\vec{a}.$$

The inertia tensor is symmetric and positive definite. In the limiting case of a rod, or ‘rotator,’ idealized as a straight line with a mass density per unit length, the eigenvalue of  $\mathbf{I}$  corresponding to rotations about the axis of the rod is zero, and  $\mathbf{I}$  is singular. We will not consider that case here, and assume that  $\mathbf{I}$  is nonsingular.

The quantities in (3.79) have dimensions

$$[\vec{\Omega}] = \frac{1}{T}, \quad [\vec{\Pi}] = \frac{ML^2}{T}, \quad [\mathbf{I}] = ML^2,$$

so the equation is dimensionally consistent.

Using the equation  $\vec{\pi} = R\vec{\Pi}$  in the spatial equation of conservation of angular momentum (3.74), using (3.77) to write  $\dot{R}$  in terms of  $W$ , and using the fact that  $\vec{\Omega} = \hat{W}$ , we get the body form of conservation of angular momentum

$$\dot{\vec{\Pi}} + \Omega \times \vec{\Pi} = 0.$$

Together with (3.79), this equation provides a  $3 \times 3$  system of ODEs for either the body angular velocity  $\vec{\Omega}(t)$

$$\mathbf{I}\dot{\vec{\Omega}} + \vec{\Omega} \times (\mathbf{I}\vec{\Omega}) = 0,$$

or the body angular momentum  $\vec{\Pi}(t)$

$$(3.81) \quad \dot{\vec{\Pi}} + (\mathbf{I}^{-1}\vec{\Pi}) \times \vec{\Pi} = 0.$$

Once we have solved these equations for  $\vec{\Omega}(t)$ , and therefore  $W(t)$ , we may reconstruct the rotation  $R(t)$  by solving the matrix equation

$$\dot{R} = RW.$$

### 12.3. The kinetic energy

The kinetic energy  $T$  of the body is given by

$$T = \frac{1}{2} \int_{\mathcal{B}} \rho(\vec{a}) |\vec{v}(\vec{a}, t)|^2 d\vec{a}.$$

Since  $R$  is orthogonal and  $\vec{v} = R\vec{V}$ , we have  $|\vec{v}|^2 = |\vec{V}|^2$ . Therefore, using (3.78), the kinetic energy of the body is given in terms of the body angular velocity by

$$(3.82) \quad T = \frac{1}{2} \int_{\mathcal{B}} \rho(\vec{a}) |\Omega(t) \times \vec{a}|^2 d\vec{a}.$$

From (3.80), this expression may be written as

$$T = \frac{1}{2} \vec{\Omega} \cdot \mathbf{I} \vec{\Omega}.$$

Thus, the body angular momentum  $\vec{\Pi}$  is given by

$$\vec{\Pi} = \frac{\partial T}{\partial \vec{\Omega}}.$$

Note that this equation is dimensionally consistent, since  $[T] = ML^2/T^2$ . Expressed in terms of  $\vec{\Pi}$ , the kinetic energy is

$$T = \frac{1}{2} \vec{\Pi} \cdot (\mathbf{I}^{-1} \vec{\Pi}).$$

As we will show below, the kinetic energy  $T$  is conserved for solutions of (3.81).

#### 12.4. The rigid body Poisson bracket

The equations (3.81) are a  $3 \times 3$  system, so they cannot be canonical Hamiltonian equations, which are always even in number. We can, however, write them in Poisson form by use of a suitable noncanonical Poisson bracket.

We define a Poisson bracket of functions  $F, G : \mathbb{R}^3 \rightarrow \mathbb{R}$  by

$$(3.83) \quad \{F, G\} = -\vec{\Pi} \cdot \left( \frac{\partial F}{\partial \vec{\Pi}} \times \frac{\partial G}{\partial \vec{\Pi}} \right),$$

This bracket is a skew-symmetric, bilinear derivation. It also satisfies the Jacobi identity (3.67), as may be checked by a direct computation. The minus sign is not required in order for (3.83) to define a bracket, but it is included to agree with the usual sign convention, which is related to a difference between right and left invariance in the Lie group  $SO(3)$  underlying this problem.

For each  $\vec{\Pi} \in \mathbb{R}^3$ , we define a linear map  $\mathbf{J}(\vec{\Pi}) : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  by

$$\mathbf{J}(\vec{\Pi}) \vec{x} = \vec{\Pi} \times \vec{x}.$$

Then, using the cyclic symmetry of the scalar triple product, we may write the Poisson bracket as

$$\{F, G\} = \frac{\partial F}{\partial \vec{\Pi}} \cdot \mathbf{J}(\vec{\Pi}) \left[ \frac{\partial G}{\partial \vec{\Pi}} \right],$$

Equation (3.81) is then

$$\dot{\vec{\Pi}} = \mathbf{J}(\vec{\Pi}) \left[ \frac{\partial T}{\partial \vec{\Pi}} \right]$$

or, in Poisson bracket form,

$$\dot{\vec{\Pi}} = \left\{ \vec{\Pi}, T \right\}, \quad \dot{\Pi}_i = \{ \Pi_i, T \}.$$

Next let us derive the conserved quantities for this equation. Any function  $F$  such that  $\{F, T\} = 0$  is conserved. In particular, since the Poisson bracket is skew-symmetric, the kinetic energy  $T$  itself is conserved.

Let  $L : \mathbb{R}^3 \rightarrow \mathbb{R}$  denote the total angular momentum function

$$L(\vec{\Pi}) = \vec{\Pi} \cdot \vec{\Pi}.$$

Then, from (3.83),

$$\{F, L\} = -2\vec{\Pi} \cdot \left( \frac{\partial F}{\partial \vec{\Pi}} \times \vec{\Pi} \right) = -2 \frac{\partial F}{\partial \vec{\Pi}} \cdot (\vec{\Pi} \times \vec{\Pi}) = 0.$$

Thus,  $\{F, L\} = 0$  for *any* function  $F$ . Such a function  $L$  is called a *Casimir* (or distinguished) function of the Poisson bracket; it is a conserved quantity for any Hamiltonian with that Poisson bracket. In particular, it follows that  $L$  is a conserved quantity for (3.81)

The conservation of  $L$  is also easy to derive directly from (3.81). Taking the inner product of the equation with  $\pi$ , we get

$$\frac{d}{dt} \left( \frac{1}{2} \vec{\pi} \cdot \vec{\pi} \right) = 0,$$

and, since  $R$  is orthogonal,  $\vec{\pi} \cdot \vec{\pi} = \vec{\Pi} \cdot \vec{\Pi}$ .

Thus, the trajectories of (3.81) lie on the intersection of the invariant spheres of constant angular momentum

$$\vec{\Pi} \cdot \vec{\Pi} = \text{constant.}$$

and the invariant ellipsoids of constant energy

$$\vec{\Pi} \cdot (\mathbf{I}^{-1} \vec{\Pi}) = \text{constant.}$$

If  $\mathbf{I}$  has distinct eigenvalues, this gives the picture shown in Figure 2.

To explain this picture in more detail, we write the rigid body equations in component form. Let  $\{\vec{e}_1, \vec{e}_2, \vec{e}_3\}$  be an orthonormal basis of eigenvectors, or principal axes, of  $\mathbf{I}$ . There is such a basis because  $\mathbf{I}$  is symmetric. We denote the corresponding eigenvalues, or principal moments of inertia, by  $I_j > 0$ , where

$$\mathbf{I} \vec{e}_j = I_j \vec{e}_j.$$

The eigenvalues are positive since  $\mathbf{I}$  is positive definite. (It also follows from (3.80) that if  $I_1 \leq I_2 \leq I_3$ , say, then  $I_3 \leq I_1 + I_2$ .)

We expand

$$\vec{\Pi}(t) = \sum_{j=1}^3 \Pi_j(t) \vec{e}_j$$

with respect to this principal axis basis. The component form of (3.81) is then

$$\begin{aligned} \dot{\Pi}_1 &= \left( \frac{1}{I_3} - \frac{1}{I_2} \right) \Pi_2 \Pi_3, \\ \dot{\Pi}_2 &= \left( \frac{1}{I_1} - \frac{1}{I_3} \right) \Pi_3 \Pi_1, \\ \dot{\Pi}_3 &= \left( \frac{1}{I_2} - \frac{1}{I_1} \right) \Pi_1 \Pi_2. \end{aligned}$$

Restricting this system to the invariant sphere

$$\Pi_1^2 + \Pi_2^2 + \Pi_3^2 = 1,$$

we see that there are three equilibrium points  $(1, 0, 0)$ ,  $(0, 1, 0)$ ,  $(0, 0, 1)$ , corresponding to steady rotations about each of the principle axes of the bodies. If  $I_1 < I_2 < I_3$ , then the middle equilibrium is an unstable saddle point, while the other two equilibria are stable centers (see Figure 2).

The instability of the middle equilibrium can be observed by stretching an elastic band around a book and spinning it around each of its three axes.

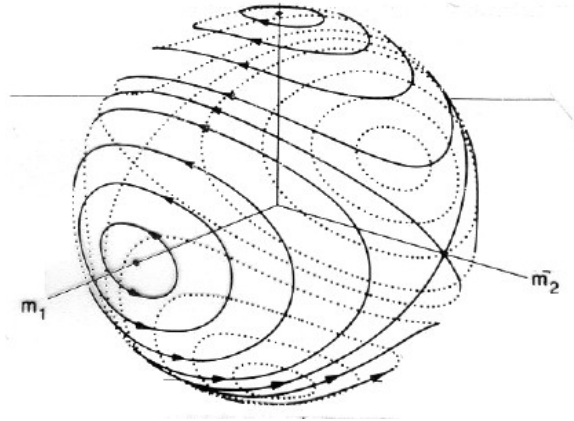


FIGURE 2. Phase portrait for rigid body rotation (from Bender and Orzag).

This rigid-body Poisson bracket has a geometrical interpretation as a Poisson bracket on  $\mathfrak{so}(3)^*$ , the dual of the Lie algebra of the three-dimensional rotation group  $SO(3)$ . Here, this dual space is identified with  $\mathbb{R}^3$  through the cross-product and the Euclidean inner product.

There is an analogous Poisson bracket, called a Lie-Poisson bracket, on the dual of any Lie algebra. Like the rigid-body bracket, it depends linearly on the coordinates of the dual Lie algebra. Arnold observed that the equations of incompressible, inviscid fluid flows may be interpreted as Lie-Poisson equations associated with the infinite-dimensional group of volume-preserving diffeomorphisms on the fluid domain.

### 13. Hamiltonian PDEs

The Euler-Lagrange equation of a variational PDE can be transformed into a canonical Hamiltonian PDE in an analogous way to ODEs.

For example, consider the wave equation (3.57) with Lagrangian  $\mathcal{L}(u, u_t)$  in (3.55). We define the momentum  $p(\cdot, t)$ , conjugate to the field variable  $u(\cdot, t)$  by

$$p = \frac{\delta \mathcal{L}}{\delta u_t}$$

For (3.55), we get

$$p = \rho_0 u_t.$$

We then define the Hamiltonian functional  $\mathcal{H}(u, p)$  by

$$\mathcal{H}(u, p) = \int p u_t dx - \mathcal{L}(u, u_t).$$

For (3.55), we get

$$\mathcal{H}(u, p) = \frac{1}{2} \int \left\{ \frac{p^2}{\rho_0} + k |\nabla u|^2 \right\} dx.$$

Hamilton's equations are

$$u_t = \frac{\delta \mathcal{H}}{\delta p}, \quad p_t = -\frac{\delta \mathcal{H}}{\delta u}.$$

For (3.55), we find that

$$u_t = \frac{p}{\rho_0}, \quad p_t = k\Delta u.$$

The elimination of  $p$  from this equation yields the wave equation (3.57).

The Poisson bracket of two functionals  $\mathcal{F}(u, p)$ ,  $\mathcal{G}(u, p)$  associated with these canonical variables is

$$\{\mathcal{F}, \mathcal{G}\} = \int \left( \frac{\delta \mathcal{F}}{\delta u} \frac{\delta \mathcal{G}}{\delta p} - \frac{\delta \mathcal{F}}{\delta p} \frac{\delta \mathcal{G}}{\delta u} \right) dx$$

Then, as before, for any functional  $\mathcal{F}(u, p)$ , evaluated on a solutions of Hamilton's equation, we have

$$\mathcal{F}_t = \{\mathcal{F}, \mathcal{H}\}.$$

### 13.1. Poisson brackets

One advantage of the Poisson bracket formulation is that it generalizes easily to PDE problems in which a suitable choice of canonical variables is not obvious.

Consider, for simplicity, an evolution equation that is first-order in time for a scalar-valued function  $u(x, t)$ . Suppose that  $\mathbf{J}(u)$  is a skew-symmetric, linear operator on functions, which may depend upon  $u$ . In other words, this means that

$$\int f(x) \mathbf{J}(u) [g(x)] dx = - \int \mathbf{J}(u) [f(x)] g(x) dx.$$

for all admissible functions  $f, g, u$ . Here, we choose the integration range as appropriate; for example, we take the integral over all of space if the functions are defined on  $\mathbb{R}$  or  $\mathbb{R}^n$ , or over a period cell if the functions are spatially periodic. We also assume that the boundary terms from any integration by parts can be neglected; for example, because the functions and their derivatives decay sufficiently rapidly at infinity, or by periodicity.

We then define a Poisson bracket of two spatial functionals  $\mathcal{F}(u)$ ,  $\mathcal{G}(u)$  of  $u$  by

$$\{\mathcal{F}, \mathcal{G}\} = \int \frac{\delta \mathcal{F}}{\delta u} \cdot \mathbf{J}(u) \left[ \frac{\delta \mathcal{G}}{\delta u} \right] dx$$

This bracket is a skew-symmetric derivation. If  $\mathbf{J}$  is a constant operator that is independent of  $u$ , then the bracket satisfies the Jacobi identity (3.67), but, in general, the Jacobi identity places severe restrictions on how  $\mathcal{J}(u)$  can depend on  $u$  (see [40]).

As an example, let us consider the Hamiltonian formulation of the KdV equation

$$(3.84) \quad u_t + uu_x + u_{xxx} = 0.$$

We define a constant skew-symmetric operator

$$\mathbf{J} = \partial_x$$

and a Hamiltonian functional

$$\mathcal{H}(u) = \int \left\{ -\frac{1}{6}u^3 + \frac{1}{2}u_x^2 \right\}$$

Then

$$\frac{\delta \mathcal{H}}{\delta u} = -\frac{1}{2}u^2 - u_{xx}$$

and hence the KdV equation (3.84) may be written as

$$u_t = \mathbf{J} \left[ \frac{\delta \mathcal{H}}{\delta u} \right]$$

The associated Poisson bracket is

$$\{\mathcal{F}, \mathcal{G}\} = \int \frac{\delta \mathcal{F}}{\delta u} \partial_x \left[ \frac{\delta \mathcal{G}}{\delta u} \right] dx$$

The KdV equation is remarkable in that it has two different, but compatible, Hamiltonian formulations. This property is one way to understand the fact that the KdV equation is a completely integrable Hamiltonian PDE.

The second structure has the skew-symmetric operator

$$\mathbf{K}(u) = \frac{1}{3}(u\partial_x + \partial_x u) + \partial_x^3.$$

Note that the order of the operations here is important:

$$u\partial_x \cdot f = uf_x, \quad \partial_x u \cdot f = (uf)_x = uf_x + u_x f.$$

Thus, the commutator of the multiplication operator  $u$  and the partial derivative operator  $\partial_x$ , given by  $[u, \partial_x]f = -u_x f$ , is the multiplication operator  $-u_x$ .

The Poisson bracket associated with  $\mathbf{K}$  satisfies the Jacobi identity (this depends on a nontrivial cancelation). In fact, the Poisson bracket associated with  $\alpha\mathbf{J} + \beta\mathbf{K}$  satisfies the Jacobi identity for any constants  $\alpha, \beta$ , which is what it means for the Poisson structures to be compatible.

The KdV-Hamiltonian for  $\mathbf{K}$  is

$$\mathcal{P}(u) = -\frac{1}{2} \int u^2 dx,$$

with functional derivative

$$\frac{\delta \mathcal{P}}{\delta u} = -u.$$

The KdV equation may then be written as

$$u_t = \mathbf{K} \left[ \frac{\delta \mathcal{P}}{\delta u} \right].$$

#### 14. Path integrals

Feynman gave a remarkable formulation of quantum mechanics in terms of path integrals. The principle of stationary action for classical mechanics may be understood heuristically as arising from a stationary phase approximation of the Feynman path integral.

The method of stationary phase provides an asymptotic expansion of integrals with a rapidly oscillating integrand. Because of cancelation, the behavior of such integrals is dominated by contributions from neighborhoods of the stationary phase points where the oscillations are the slowest. Here, we explain the basic idea in the case of one-dimensional integrals. See Hormander [27] for a complete discussion.

### 14.1. Fresnel integrals

Consider the following Fresnel integral

$$(3.85) \quad I(\varepsilon) = \int_{-\infty}^{\infty} e^{ix^2/\varepsilon} dx.$$

This oscillatory integral is not defined as an absolutely convergent integral, since  $e^{ix^2/\varepsilon}$  has absolute value one, but it can be defined as an improper Riemann integral

$$I(\varepsilon) = \lim_{R \rightarrow \infty} \int_{-R}^R e^{ix^2/\varepsilon} dx.$$

The convergence follows from an integration by parts:

$$\int_1^R e^{ix^2/\varepsilon} dx = \left[ \frac{\varepsilon}{2ix} e^{ix^2/\varepsilon} \right]_1^R + \int_1^R \frac{\varepsilon}{2ix^2} e^{ix^2/\varepsilon} dx.$$

The integrand in (3.85) oscillates rapidly away from the stationary phase point  $x = 0$ , and these parts contribute terms that are smaller than any power of  $\varepsilon$  as  $\varepsilon \rightarrow 0$ , as we show below. The first oscillation near  $x = 0$ , where cancelation does not occur, has width of the order  $\varepsilon^{1/2}$ , and as a result  $I(\varepsilon) = O(\varepsilon^{1/2})$  as  $\varepsilon \rightarrow 0$ .

Using contour integration, and changing variables  $x \mapsto e^{i\pi/4}s$  if  $\varepsilon > 0$  or  $x \mapsto e^{-i\pi/4}s$  if  $\varepsilon < 0$ , one can show that

$$\int_{-\infty}^{\infty} e^{ix^2/\varepsilon} dx = \begin{cases} e^{i\pi/4} \sqrt{2\pi|\varepsilon|} & \text{if } \varepsilon > 0, \\ e^{-i\pi/4} \sqrt{2\pi|\varepsilon|} & \text{if } \varepsilon < 0. \end{cases}$$

### 14.2. Stationary phase

Next, we consider the integral

$$(3.86) \quad I(\varepsilon) = \int_{-\infty}^{\infty} f(x) e^{i\varphi(x)/\varepsilon} dx,$$

where  $f : \mathbb{R} \rightarrow \mathbb{C}$  and  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  are smooth functions. A point  $x = c$  is a stationary phase point if  $\varphi'(c) = 0$ . We call the stationary phase point nondegenerate if  $\varphi''(c) \neq 0$ .

Suppose that  $I$  has a single stationary phase point at  $x = c$ , and it is nondegenerate. If there are several such points, we simply add together the contributions from each one. Then, using the idea that only the part of the integrand near the stationary phase point  $x = c$  contributes significantly, we Taylor expand the function  $f$  and the phase  $\varphi$  to approximate  $I(\varepsilon)$  as follows:

$$\begin{aligned} I(\varepsilon) &\sim \int f(c) \exp \frac{i}{\varepsilon} \left[ \varphi(c) + \frac{1}{2} \varphi''(c) (x - c)^2 \right] dx \\ &\sim f(c) e^{i\varphi(c)/\varepsilon} \int \exp \left[ \frac{i\varphi''(c)}{2\varepsilon} s^2 \right] ds \\ &\sim \sqrt{\frac{2\pi\varepsilon}{|\varphi''(c)|}} f(c) e^{i\varphi(c)/\varepsilon + i\sigma\pi/4}, \end{aligned}$$

where

$$\sigma = \text{sign } \varphi''(c).$$

### 14.3. The Feynman path integral

Consider a single, non-relativistic quantum mechanical particle of mass  $m$  in a potential  $V(\vec{x})$ . Suppose that the particle is located at  $\vec{x}_0$  at time  $t_0$ , and we observe the location of the particle at time  $t_1$ . We would like to calculate the probability of finding the particle in some specific region  $\Omega \subset \mathbb{R}^n$ .

According to Feynman's formulation of quantum mechanics [21], every event has an associated complex number,  $\Psi$ , called its amplitude. If an event can occur in a number of different independent ways, the amplitude of the event is obtained by adding together the amplitudes of the different subevents. Finally, the probability of observing an event when some measurement is made is the modulus of the amplitude squared  $|\Psi|^2$ .

The fact that amplitudes add, not probabilities, leads to the interference effects characteristic of quantum mechanics. For example, consider an event (like the observation of an electron in the 'double slit' experiment) which can occur in two different ways with equal probability. If the two amplitudes have opposite phase, then the probability of the event is zero, while if they have the same phase, then the probability of the event is four times the probability of the separate subevents.

To apply this formulation to the motion of a quantum mechanical particle, we take as the basic subevents the possible paths  $\vec{x}(t)$  of the particle from  $\vec{x}_0$  at time  $t_0$  to  $\vec{x}_1$  at time  $t_1$ . The amplitude of a path  $\vec{x}$  is proportional to  $e^{i\mathcal{S}(\vec{x})/\hbar}$  where  $\mathcal{S}(\vec{x})$  is the action of the path

$$\mathcal{S}(\vec{x}) = \int_{t_0}^{t_1} \left\{ \frac{1}{2} m |\dot{\vec{x}}|^2 - V(\vec{x}) \right\} dt,$$

and  $\hbar$  is Planck's constant. Like the action, Planck's constant has the dimension of energy · time, or momentum · length; its approximate value is  $\hbar = 1.054 \times 10^{-34}$  Js.

Thus the action, which is a somewhat mysterious quantity in classical mechanics, corresponds to a phase, measured in units of  $\hbar$ , in quantum mechanics.

The amplitude  $\psi(\vec{x}_1, t_1; \vec{x}_0, t_0)$  of the particle moving from  $\vec{x}_0$  at time  $t_0$  to  $\vec{x}_1$  at time  $t_1$  is then obtained formally by summing the amplitudes of each path over 'all' possible paths

$$\mathcal{P}(\vec{x}_1, t_1; \vec{x}_0, t_0) = \{ \vec{x} \mid \vec{x}(t) : [t_0, t_1] \rightarrow \mathbb{R}^n \text{ is continuous, } \vec{x}(t_0) = \vec{x}_0, \vec{x}(t_1) = \vec{x}_1 \}.$$

This gives

$$(3.87) \quad \psi(\vec{x}_1, t_1; \vec{x}_0, t_0) = \int_{\mathcal{P}(\vec{x}_1, t_1; \vec{x}_0, t_0)} e^{i\mathcal{S}(\vec{x})/\hbar} D\vec{x},$$

where  $D\vec{x}$  is supposed to be a measure on the path space that weights all paths equally, normalized so that  $|\psi|^2$  is a probability density.

This argument has great intuitive appeal, but there are severe difficulties in making sense of the result. First, there is no translation-invariant 'flat' measure  $D\vec{x}$  on an infinite-dimensional path space, analogous to Lebesgue measure on  $\mathbb{R}^n$ , that weights all paths equally. Second, for paths  $\vec{x}(t)$  that are continuous but not differentiable, which include the paths one expects to need, the action  $\mathcal{S}(\vec{x})$  is undefined, or, at best, infinite. Thus, in qualitative terms, the Feynman path integral in the expression for  $\psi$  in fact looks something like this:

$$“ \int_{\mathcal{P}} e^{i\mathcal{S}(\vec{x})/\hbar} D\vec{x} = \int e^{i\infty/\hbar} D? ”.$$

Nevertheless, there are ways to make sense of (3.87) as providing the solution  $\psi(\vec{x}, t; \vec{x}_0, t_0)$  of the Schrödinger equation

$$\begin{aligned} i\hbar\psi_t &= -\frac{\hbar^2}{2m}\Delta\psi + V(\vec{x})\psi, \\ \psi(\vec{x}, t_0; \vec{x}_0, t_0) &= \delta(\vec{x} - \vec{x}_0). \end{aligned}$$

For example, the Trotter product formula gives an expression for  $\psi$  as a limit of finite dimensional integrals over  $\mathbb{R}^N$  as  $N \rightarrow \infty$ , which may be taken as a definition of the path integral in (3.87) (see (3.92)–(3.94) below).

After seeing Feynman's work, Kac (1949) observed that an analogous formula for solutions of the heat equation with a lower-order potential term (the 'imaginary time' version of the Schrödinger equation),

$$(3.88) \quad u_t = \frac{1}{2}\Delta u - V(x)u,$$

can be given rigorous sense as a path integral with respect to Wiener measure, which describes the probability distribution of particle paths in Brownian motion.

Explicitly, for sufficiently smooth potential functions  $V(x)$ , the Green's function of (3.88), with initial data  $u(x, t_0) = \delta(x - x_0)$ , is given by the Feynman-Kac formula

$$u(x, t; x_0, t_0) = \int_{\mathcal{P}(x, t; x_0, t_0)} e^{\int_{t_0}^t V(x(s)) ds} dW(x).$$

Here,  $\mathcal{P}(x, t; x_0, t_0)$  denotes the space of all continuous paths  $x(s)$ , with  $t_0 \leq s \leq t$  from  $x_0$  at  $t_0$  to  $x$  at  $t$ . The integral is taken over  $\mathcal{P}$  with respect to Wiener measure. Formally, we have

$$(3.89) \quad dW(x) = e^{-\int_{t_0}^t \frac{1}{2}|\dot{x}|^2 ds} Dx.$$

However, neither the 'flat' measure  $Dx$ , nor the exponential factor on the right-hand side of this equation are well-defined. In fact, the Wiener measure is supported on continuous paths that are almost surely nowhere differentiable (see Section 3.4).

#### 14.4. The Trotter product formula

To explain the idea behind the Trotter product formula, we write the Schrödinger equation as

$$(3.90) \quad i\hbar\psi_t = H\psi,$$

where the Hamiltonian operator  $H$  is given by

$$H = T + V$$

and the kinetic and potential energy operators  $T$  and  $V$ , respectively, are given by

$$T = -\frac{\hbar^2}{2m}\Delta, \quad V = V(\vec{x}).$$

Here,  $V$  is understood as a multiplication operator  $V : \psi \mapsto V(\vec{x})\psi$ .

We write the solution of (3.90) as

$$\psi(t) = e^{-it\hbar^{-1}H}\psi_0$$

where  $\psi(t) = \psi(\cdot, t) \in L^2(\mathbb{R}^n)$  denotes the solution at time  $t$ ,  $\psi_0 = \psi(0)$  is the initial data, and

$$e^{-it\hbar^{-1}H} : L^2(\mathbb{R}^n) \rightarrow L^2(\mathbb{R}^n)$$

is the one-parameter group of solution operators (or flow).

Assuming that  $V(\vec{x})$  is not constant, the operators  $T, V$  do not commute:

$$[T, V] = TV - VT = -\frac{\hbar^2}{2m} (2\nabla V \cdot \nabla + \Delta V).$$

(Here,  $\Delta V$  denotes the operation of multiplication by the function  $\Delta V$ .) Thus, the flows  $e^{-it\hbar^{-1}T}, e^{-it\hbar^{-1}V}$  do not commute, and  $e^{-it\hbar^{-1}H} \neq e^{-it\hbar^{-1}V} e^{-it\hbar^{-1}T}$ .

For small times  $\Delta t$ , however, we have

$$\begin{aligned} e^{-i\Delta t\hbar^{-1}H} &= I - \frac{i\Delta t}{\hbar} H - \frac{\Delta t^2}{2\hbar^2} H^2 + O(\Delta t^3) \\ &= I - \frac{i\Delta t}{\hbar} (T + V) - \frac{\Delta t^2}{2\hbar^2} (T^2 + TV + VT + V^2) + O(\Delta t^3), \\ e^{-i\Delta t\hbar^{-1}T} &= I - \frac{i\Delta t}{\hbar} T - \frac{\Delta t^2}{2\hbar^2} T^2 + O(\Delta t^3) \\ e^{-i\Delta t\hbar^{-1}V} &= I - \frac{i\Delta t}{\hbar} V - \frac{\Delta t^2}{2\hbar^2} V^2 + O(\Delta t^3). \end{aligned}$$

Thus,

$$e^{-i\Delta t\hbar^{-1}H} = e^{-i\Delta t\hbar^{-1}V} e^{-i\Delta t\hbar^{-1}T} - \frac{\Delta t^2}{2\hbar^2} [T, V] + O(\Delta t^3),$$

and we can obtain a first-order accurate approximation for the flow associated with  $H$  by composing the flows associated with  $V$  and  $T$ .

The numerical implementation of this idea is the fractional step method. We solve the evolution equation

$$u_t = (A + B)u$$

by alternately solving the equations

$$u_t = Au, \quad u_t = Bu$$

over small time-steps  $\Delta t$ . In this context, the second-order accurate approximation in  $\Delta t$

$$e^{\Delta t(A+B)} = e^{\frac{1}{2}\Delta t A} e^{\Delta t B} e^{\frac{1}{2}\Delta t A}$$

is called ‘Strang splitting.’

To obtain the solution of (3.90) at time  $t$ , we take  $N$  time-steps of length  $\Delta t = t/N$ , and let  $N \rightarrow \infty$ , which gives the Trotter product formula

$$(3.91) \quad \psi(t) = \lim_{N \rightarrow \infty} \left[ e^{-it(\hbar N)^{-1}V} e^{-it(\hbar N)^{-1}T} \right]^N \psi_0.$$

Under suitable assumptions on  $V$ , the right-hand side converges strongly to  $\psi(t)$  with respect to the  $L^2(\mathbb{R}^n)$ -norm.

The flows associated with  $V, T$  are easy to find explicitly. The solution of

$$i\hbar\psi_t = V\psi$$

is given by the multiplication operator

$$\psi_0 \mapsto e^{-it\hbar^{-1}V} \psi_0.$$

The solution of

$$i\hbar\psi_t = T\psi$$

may be found by taking the spatial Fourier transform and using the convolution theorem, which gives

$$\begin{aligned}\psi(\vec{x}, t) &= \int e^{\{-it|\vec{p}|^2/(2\hbar m) + i\vec{p}\cdot\vec{x}/\hbar\}} \hat{\psi}_0(\vec{p}) d\vec{p} \\ &= \left(\frac{m}{2\pi i\hbar t}\right)^{n/2} \int e^{im|\vec{x}-\vec{y}|^2/(2\hbar t)} \psi_0(\vec{y}) d\vec{y}.\end{aligned}$$

Using these results in the Trotter product formula (3.91), writing the spatial integration variable at time  $t_k = kt/N$  as  $\vec{x}_k$ , with  $\vec{x}_N = \vec{x}$ , and assuming that  $\psi(\vec{x}, 0) = \delta(\vec{x} - \vec{x}_0)$ , we get, after some algebra, that

$$(3.92) \quad \psi(\vec{x}, t) = \lim_{N \rightarrow \infty} C_{N,t} \int e^{iS_{N,t}(\vec{x}_0, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_{N-1}, \vec{x}_N)/\hbar} d\vec{x}_1 d\vec{x}_2, \dots, d\vec{x}_{N-1}$$

where the normalization factor  $C_{N,t}$  is given by

$$(3.93) \quad C_{N,t} = \left(\frac{mN}{2\pi i\hbar t}\right)^{n(N-1)/2}$$

and the exponent  $S_{N,t}$  is a discretization of the classical action functional

$$(3.94) \quad S_{N,t}(\vec{x}_0, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_{N-1}, \vec{x}_N) = \sum_{k=1}^{N-1} \frac{t}{N} \left[ \frac{m}{2} \left| \frac{\vec{x}_{k+1} - \vec{x}_k}{t/N} \right|^2 - V(\vec{x}_k) \right].$$

Equations (3.92)–(3.94) provide one way to interpret the path integral formula (3.87).

#### 14.5. Semiclassical limit

One of the most appealing features of the Feynman path integral formulation is that it shows clearly the connection between classical and quantum mechanics. The phase of the quantum mechanical amplitude is the classical action, and, by analogy with the method of stationary phase for finite-dimensional integrals, we expect that for semi-classical processes whose actions are much greater than  $\hbar$ , the amplitude concentrates on paths of stationary phase. Again, however, it is difficult to make clear analytical sense of this argument while maintaining its simple intuitive appeal.