

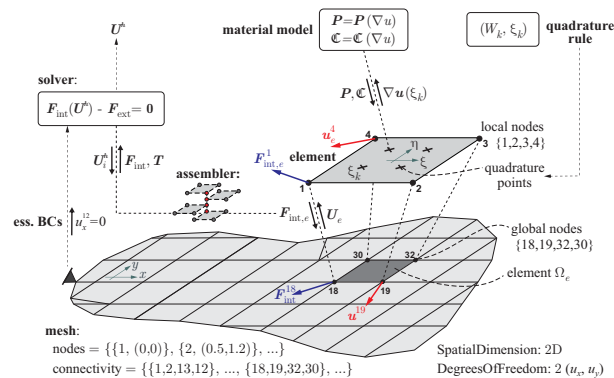
# Computational Solid Mechanics (Fall 2017)

Dennis M. Kochmann

Mechanics & Materials  
Department of Mechanical and Process Engineering  
ETH Zürich

course website:

[www.mm.ethz.ch/teaching.html](http://www.mm.ethz.ch/teaching.html)



Copyright © 2017 by Dennis M. Kochmann

These lecture notes are a *concise* collection of equations and comments. They are by no means a complete textbook or set of class notes that could replace lectures.

Therefore, you are strongly encouraged to take your own notes during lectures and to use this set of notes rather for reference.

## Contents

<b>1</b>	<b>Introduction: Continuum Mechanics and Notation</b>	<b>6</b>
1.1	An introductory example: heat conduction	6
1.2	A more advanced example: mechanical equilibrium	9
1.3	A special case: linearized kinematics	13
1.4	Summary and Looking Ahead	15
<b>2</b>	<b>Numerical Methods</b>	<b>16</b>
<b>3</b>	<b>Variational Calculus</b>	<b>18</b>
3.1	Functionals	18
3.2	Variations	18
3.3	Example: Hanging bar under its own weight	20
3.4	Example: Static Heat Conduction	21
<b>4</b>	<b>The weak form</b>	<b>24</b>
4.1	Classical and weak solutions	24
4.2	Equivalence of strong and weak forms	26
4.3	Approximate solutions	26
<b>5</b>	<b>The mechanical variational problem</b>	<b>28</b>
5.1	Linearized kinematics	28
5.2	Finite kinematics	30
5.3	Thermal problem revisited	31
5.4	A simple example: nonlinear springs	32
<b>6</b>	<b>Interpolation spaces</b>	<b>33</b>
<b>7</b>	<b>The Finite Element Method</b>	<b>35</b>
<b>8</b>	<b>Finite element spaces: polynomial shape functions in 1D</b>	<b>37</b>
8.1	One dimension	37
8.2	Example: linear elastic bar in 1D	37
8.3	Higher dimensions	39
<b>9</b>	<b>Simplicial elements</b>	<b>40</b>
9.1	Linear Triangle (T3)	40
9.2	Extension to three dimensions:	42
9.3	Finite element implementation	42
9.4	Higher-order triangles and tetrahedra:	42
<b>10</b>	<b>The bilinear quadrilateral element</b>	<b>44</b>
<b>11</b>	<b>Numerical quadrature</b>	<b>47</b>
11.1	Example: Riemann sums	47
11.2	Gauss quadrature	47
11.2.1	Gauss-Legendre quadrature	48
11.2.2	Other Gauss quadrature rules	50
11.3	Higher dimensions	51

11.4 Finite element implementation . . . . .	52
11.5 Quadrature error estimates . . . . .	52
11.6 Quadrature rules for simplicial elements: . . . . .	52
11.7 Which quadrature rule to use? . . . . .	53
<b>12 Generalization and implementation of the simplicial elements</b>	<b>54</b>
<b>13 Assembly</b>	<b>55</b>
<b>14 Overview: Numerical Implementation</b>	<b>56</b>
<b>15 Iterative solvers</b>	<b>58</b>
15.1 Netwon-Raphson (NR) method . . . . .	58
15.2 Damped Newton-Raphson (dNR) method . . . . .	59
15.3 Quasi-Newton (QN) method . . . . .	59
15.4 Line search method . . . . .	59
15.5 Gradient flow method . . . . .	60
15.6 Nonlinear Least Squares . . . . .	60
15.7 Conjugate Gradient (CG) method . . . . .	60
<b>16 Boundary conditions</b>	<b>62</b>
16.1 Neumann boundary conditions . . . . .	62
16.2 Examples of external forces . . . . .	63
16.3 Dirichlet boundary conditions . . . . .	64
16.4 Rigid body motion . . . . .	66
<b>17 Error estimates and adaptivity</b>	<b>67</b>
17.1 Finite element error analysis . . . . .	67
17.2 Smoothing and adaptivity . . . . .	68
<b>18 Element defects: shear locking and hourglassing</b>	<b>70</b>
<b>19 Dynamics</b>	<b>72</b>
19.1 Variational setting . . . . .	72
19.2 Free vibrations . . . . .	75
19.3 Modal decomposition . . . . .	77
19.4 Transient time-dependent solutions . . . . .	78
19.5 Explicit time integration . . . . .	79
19.6 A reinterpretation of finite differences . . . . .	79
19.7 Implicit time integration . . . . .	81
<b>20 Internal variables and inelasticity</b>	<b>84</b>
20.1 Inelastic material models . . . . .	84
20.2 Example: viscoelasticity, (visco)plasticity . . . . .	86
20.3 Example: viscoplasticity . . . . .	87
20.4 Example: linear viscoelasticity . . . . .	88
<b>A Introduction, Vector Spaces</b>	<b>93</b>
<b>B Function Spaces</b>	<b>94</b>

<b>C Approximation Theory</b>	<b>98</b>
C.1 Sobolev spaces . . . . .	101
C.2 Higher dimensions . . . . .	101
<b>D Operators</b>	<b>104</b>
<b>E Uniqueness</b>	<b>105</b>
<b>F Vainberg's theorem</b>	<b>106</b>
<b>G Energy norm</b>	<b>107</b>

# 1 Introduction: Continuum Mechanics and Notation

## 1.1 An introductory example: heat conduction

We describe a body  $\Omega \subset \mathbb{R}^d$  with boundary  $\partial\Omega$  as a collection of **material points**. Each point has a position  $\mathbf{x}$  in a Cartesian coordinate system  $(x_1, \dots, x_d)$  in  $d$  dimensions with origin  $\mathbf{O}$ .

Points are described by vectors defined by components in the Cartesian reference frame:

$$\mathbf{x} = \sum_{i=1}^d x_i \mathbf{g}_i = x_i \mathbf{g}_i. \quad (1.1)$$

Here and in the following we use Einstein's *summation convention* which implies summation over repeated indices. The usual index notation rules apply; e.g., the inner product is written as

$$\mathbf{a} \cdot \mathbf{b} = a_i \mathbf{g}_i \cdot b_j \mathbf{g}_j = a_i b_j \delta_{ij} = a_i b_i \quad \text{with Kronecker's delta} \quad \delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{else.} \end{cases} \quad (1.2)$$

This is used to define the length of a vector as

$$\|\mathbf{a}\| = \sqrt{\mathbf{a} \cdot \mathbf{a}} = \sqrt{a_i a_i} = \sqrt{\sum_{i=1}^d a_i a_i}. \quad (1.3)$$

Matrix-vector multiplication becomes

$$\mathbf{a} = \mathbf{M}\mathbf{b} \quad \Leftrightarrow \quad a_i = M_{ij} b_j \quad \text{and} \quad (\mathbf{M}^T)_{ij} = M_{ji}. \quad (1.4)$$

We use **mappings** to denote fields. For example, the temperature field in a static problem is described by a mapping

$$T(\mathbf{x}) : \Omega \rightarrow \mathbb{R}, \quad (1.5)$$

which assigns to each point  $\mathbf{x} \in \Omega$  a real number, the temperature. If the field is differentiable, one often introduces **kinematic variables** such as the temperature *gradient* field:

$$\boldsymbol{\beta} : \Omega \rightarrow \mathbb{R}^d \quad \text{and} \quad \boldsymbol{\beta} = \text{grad } T = \nabla T \quad \Leftrightarrow \quad \beta_i = \frac{\partial T}{\partial x_i} = T_{,i}. \quad (1.6)$$

Here and in the following, we use *comma indices* to denote partial derivatives.

For every kinematic variable, there is conjugate field (often called **flux**) like the heat flux  $\mathbf{q}$  in this thermal problem, which is also a mapping:

$$\mathbf{q} : \Omega \rightarrow \mathbb{R}^d. \quad (1.7)$$

The heat flux vector  $\mathbf{q}$  assigns to each point in  $\Omega$  a heat flux direction and magnitude. If we are interested, e.g., in the loss of heat through a point  $\mathbf{x} \in \partial\Omega$  on the surface of  $\Omega$  with outward unit normal  $\mathbf{n}(\mathbf{x})$ , then that amount of heat leaving  $\Omega$  is the projection  $\mathbf{q}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x})$ . Clearly, the components of  $\mathbf{q} = (q_1, \dots, q_d)^T$  imply the heat flux through surfaces perpendicular to each of the  $d$  Cartesian coordinate directions.

Next, **constitutive relations** link kinematic quantities to fluxes. For example, define the heat flux vector as  $\mathbf{q} = \mathbf{q}(\boldsymbol{\beta})$ . Fourier's law of heat conduction states, e.g.,

$$\mathbf{q} = -\mathbf{K}\boldsymbol{\beta} \quad \Leftrightarrow \quad q_i = -K_{ij}\beta_{,j}, \quad (1.8)$$

where  $\mathbf{K}$  denotes a conductivity *tensor*.  $K_{ij}$  are the components of the conductivity tensor; they form a  $d \times d$  matrix. Such a *second-order tensor* is a convenient way to store the conductivity properties in any arbitrary orientation in the form of a matrix, along with the coordinate basis in which it is defined.  $\mathbf{K}$  provides for each direction of the temperature gradient  $\boldsymbol{\beta} = \alpha \mathbf{n}$  the resulting normalized heat flux  $\mathbf{q}$ . To this end, one defines

$$\mathbf{K} = K_{ij} \mathbf{g}_i \otimes \mathbf{g}_j \quad \Rightarrow \quad \mathbf{K} \mathbf{n} = (K_{ij} \mathbf{g}_i \otimes \mathbf{g}_j) \mathbf{n} = K_{ij} \mathbf{g}_i (\mathbf{g}_j \cdot \mathbf{n}) = K_{ij} n_j \mathbf{g}_i. \quad (1.9)$$

Such a second-order tensor is hence a linear mapping of vectors onto vectors. Here, we defined the *dyadic product* (or *tensor product*), which produces a tensor according to

$$\mathbf{M} = \mathbf{a} \otimes \mathbf{b} \quad \Leftrightarrow \quad M_{ij} = a_i b_j. \quad (1.10)$$

A special tensor is the identity, which maps vectors onto themselves:

$$\mathbf{I} = \delta_{ij} \mathbf{g}_i \otimes \mathbf{g}_j. \quad (1.11)$$

To solve a thermal problem, we also need **balance laws**. Here, *conservation of energy* may be used, which states the change of **internal energy**  $E$  in a body over time  $t$  balances the inward and outward flux of energy and the energy being produced inside the body (by some heat source density  $\rho s$ ). Mathematically, this implies

$$\frac{d}{dt} E = \int_{\Omega} \rho s \, dV - \int_{\partial\Omega} \mathbf{q} \cdot \mathbf{n} \, dS. \quad (1.12)$$

We can use the **divergence theorem** to write

$$\int_{\partial\Omega} \mathbf{q} \cdot \mathbf{n} \, dS = \int_{\partial\Omega} q_i n_i \, dS = \int_{\Omega} q_{i,i} \, dV = \int_{\Omega} \operatorname{div} \mathbf{q} \, dV, \quad (1.13)$$

which defines the *divergence of a vector* as the scalar quantity

$$\operatorname{div}(\cdot) = (\cdot)_{i,i}. \quad (1.14)$$

Energy is an **extensive** variable, i.e., the total energy doubles when adding to bodies of the same energy. This is in contrast to **intensive** variables, such as temperature or pressure, when adding to bodies having the same, e.g., temperature. Since energy is an extensive variable, we can introduce an energy density  $e$  and write

$$E = \int_{\Omega} e \, dV \quad \Rightarrow \quad \frac{d}{dt} E = \int_{\Omega} \dot{e} \, dV. \quad (1.15)$$

Here and in the following, we use dots to denote *rates*, i.e., *time derivatives*. Note that we use a *Lagrangian* description and that, so far, there is no motion or deformation involved in our discussion, so the time derivative does not affect the volume integral. Rewriting the conservation of energy now yields

$$\int_{\Omega} \dot{e} \, dV = \int_{\Omega} \rho s \, dV - \int_{\Omega} \operatorname{div} \mathbf{q} \, dV. \quad (1.16)$$

This can be rewritten as

$$\int_{\Omega} (\dot{e} - \rho s + \operatorname{div} \mathbf{q}) \, dV = 0. \quad (1.17)$$

Since conservation of energy does not only have to hold for  $\Omega$  but for any subbody  $\omega \subset \Omega$ , we may conclude that the **local energy balance** equation is

$$\boxed{\dot{e} = \rho s - \operatorname{div} \mathbf{q}} \quad (1.18)$$

This is the local (i.e., pointwise) counterpart to the macroscopic energy balance and states that at each point  $\mathbf{x} \in \Omega$  the rate of energy change ( $\dot{e}$ ) is given by the local production of heat ( $\rho s$ ) minus the heat lost by outward fluxes  $\mathbf{q}$  away from the point.

Finally, exploiting that thermally stored energy gives  $e = \rho c_v T$  (with constant mass density  $\rho$  and specific heat capacity  $c_v$ ), and we insert Fourier's law of heat conduction to overall arrive at

$$\int_{\Omega} \rho c_v \dot{T} \, dV = \int_{\Omega} \rho s \, dV - \int_{\Omega} \operatorname{div}(-\mathbf{K}\boldsymbol{\beta}) \, dS = \int_{\Omega} \rho s \, dV - \int_{\Omega} (-K_{ij}T_{,j})_{,i} \, dS. \quad (1.19)$$

Note that the final term requires the use of the product rule since

$$(-K_{ij}T_{,j})_{,i} = -K_{ij,i}T_{,j} - K_{ij}T_{,ij}. \quad (1.20)$$

Let us assume a **homogeneous** body with  $\mathbf{K}(\mathbf{x}) = \mathbf{K} = \text{const.}$  Further, let us rewrite the above under a single integral:

$$\int_{\Omega} (\rho c_v \dot{T} - \rho s - K_{ij}T_{,ij}) \, dV = 0. \quad (1.21)$$

Again, by extension of energy conservation to arbitrary subbodies, we conclude the local energy balance equation

$$\rho c_v \dot{T} = K_{ij}T_{,ij} + \rho s \quad (1.22)$$

This is the heat equation in its *anisotropic form*. In the special case of **isotropy** (i.e., the conductivity is the same in all directions), we obtain the *Laplacian* since

$$K_{ij} = \kappa \delta_{ij} \quad \Rightarrow \quad K_{ij}T_{,ij} = \kappa \delta_{ij}T_{,ij} = \kappa T_{,ii} = \kappa \Delta T \quad \text{with} \quad \Delta(\cdot) = (\cdot)_{,ii}, \quad (1.23)$$

and we arrive at the well-known **heat equation** with sources:

$$\boxed{\rho c_v \dot{T} = \kappa \Delta T + \rho s} \quad (1.24)$$

Whenever we consider a *static* problem, we assume that the body is in equilibrium and the temperature field is constant, which reduces the above to *Poisson's equation*, viz.

$$\kappa \Delta T = -\rho s. \quad (1.25)$$



## 1.2 A more advanced example: mechanical equilibrium

The mechanics of solids (and fluids) generally describes deformable bodies. To this end, we label each material point by its position  $\mathbf{X}$  in a **reference configuration** (e.g., the configuration at time  $t = 0$ ). The current position  $\mathbf{x}$ , by contrast, is a function of  $\mathbf{X}$  and of time  $t$ :  $\mathbf{x} = \mathbf{x}(\mathbf{X}, t)$ . Fields in the reference and **current configuration** are generally referred to by upper- and lower-case characters (and the same applies to indices), e.g.,

$$\mathbf{x} = x_i \mathbf{g}_i, \quad \mathbf{X} = X_I \mathbf{G}_I. \quad (1.26)$$

Note that, to avoid complication, we will work with Cartesian coordinate systems only (see the tensor notes, e.g., for curvilinear coordinates).

The mechanics of a deformable body undergoing finite deformations is generally described by a **deformation mapping**

$$\varphi(\mathbf{X}, t) : \Omega \times \mathbb{R} \rightarrow \mathbb{R}^d \quad \text{such that} \quad \mathbf{x} = \varphi(\mathbf{X}, t). \quad (1.27)$$

Since it depends on time, we can take time derivatives to arrive at the **velocity** and **acceleration** fields, respectively:

$$\mathbf{V}(\mathbf{X}, t) = \frac{d}{dt} \mathbf{x}(\mathbf{X}, t) = \frac{d}{dt} \varphi(\mathbf{X}, t), \quad \mathbf{A}(\mathbf{X}, t) = \frac{d}{dt} \mathbf{V}(\mathbf{X}, t) = \frac{d^2}{dt^2} \varphi(\mathbf{X}, t). \quad (1.28)$$

Note that those are **Lagrangian** fields (one could also write those as functions of the current position  $\mathbf{x}$ , which results in the **Eulerian** counterparts; this is usually done in fluid mechanics).

Like in the thermal problem, we introduce *kinematics* by defining the **deformation gradient**

$$\mathbf{F} = \text{Grad } \varphi \quad \Leftrightarrow \quad F_{iJ} = \frac{\partial \varphi_i}{\partial X_J} = \varphi_{i,J}. \quad (1.29)$$

Note that this is a second-order, *two-point tensor* defined across both configurations. If one uses the same coordinate frame for the undeformed and deformed configurations, one may alternatively introduce the **displacement field**

$$\mathbf{u}(\mathbf{X}) = \mathbf{x} - \mathbf{X} \quad \text{and} \quad \mathbf{x} = \mathbf{X} + \mathbf{u}(\mathbf{X}), \quad (1.30)$$

so that

$$\mathbf{F} = \text{Grad } \varphi = \text{Grad}(\mathbf{X} + \mathbf{u}) = \mathbf{I} + \text{Grad } \mathbf{u} \quad \Leftrightarrow \quad F_{iJ} = \delta_{iJ} + u_{i,J}. \quad (1.31)$$

Note that in case of *no deformation*, we have  $\mathbf{x} = \mathbf{X}$  so that  $\mathbf{F} = \mathbf{I}$  (and *not*  $\mathbf{F} = \mathbf{0}$ ).

$\mathbf{F}$  contains plenty of information about the local deformation. For example, the **volume change** at a point is given by

$$\frac{dv}{dV} = J = \det \mathbf{F}, \quad (1.32)$$

and for physical reasons we must have  $J > 0$  (this ensures that the deformation mapping is injective; i.e., no two material points are mapped onto the same point in the current configuration). No volume change implies  $J = 1$ .

Similarly, the **stretch** in a direction defined by a unit vector  $\mathbf{N}$  is determined by

$$\lambda(\mathbf{N}) = \frac{ds}{dS} = \sqrt{\mathbf{N} \cdot \mathbf{C} \mathbf{N}} \quad \text{with} \quad \mathbf{C} = \mathbf{F}^T \mathbf{F} \quad (1.33)$$

the **right Cauchy-Green tensor**. As for  $\mathbf{F}$ , an undeformed point has  $\mathbf{C} = \mathbf{I}$ .

Next, we need a constitutive law that links the deformation gradient to a “flux”, which in mechanical problems we refer to as the **stress**. Heat flux may be defined as energy flow per oriented area,

$$q_i = \frac{dQ}{dA_i} \quad \text{with} \quad d\mathbf{A} = \mathbf{N} dA \quad \text{and} \quad \|\mathbf{N}\| = 1. \quad (1.34)$$

Analogously, stresses are force vector per oriented area:

$$P_{iJ} = \frac{dF_i}{dA_J}, \quad (1.35)$$

which defines the so-called **First Piola-Kirchhoff** (1st PK) **stress tensor**. This is a second-order tensor that captures the force components on any oriented area into all  $d$  coordinate directions. The resulting **traction** vector on a particular *infinitesimal* area with unit normal vector  $\mathbf{N}$  is

$$\mathbf{T} = \mathbf{P}\mathbf{N} \quad \Leftrightarrow \quad T_i = P_{iJ}N_J. \quad (1.36)$$

Notice that the thus defined tractions satisfy Newton’s law of action and reaction since

$$\mathbf{T}(\mathbf{N}) = \mathbf{P}\mathbf{N} = -\mathbf{P}(-\mathbf{N}) = -\mathbf{T}(-\mathbf{N}), \quad (1.37)$$

which we know well from inner forces in undergraduate mechanics. The total force acting on a surface  $A$  is hence

$$\mathbf{F}_{\text{tot}} = \int_A \mathbf{T} dS. \quad (1.38)$$

For a mechanical problem, the relevant balance laws are *conservation of linear momentum* and of *angular momentum*. As before, one can formulate those as macroscopic balance laws. For example, macroscopic linear momentum balance is nothing but the well-known equation  $\mathbf{F}_{\text{tot}} = m\mathbf{A}$  (*sum of all forces equals mass times mean acceleration*). To derive the local balance law of a continuous body, note that external forces include both surface tractions  $\mathbf{T}$  and *body forces*  $\rho_0\mathbf{B}$  – using  $\rho_0$  to denote the reference mass density. Overall, we obtain

$$\int_{\partial\Omega} \mathbf{T} dS + \int_{\Omega} \rho_0\mathbf{B} dV = \int_{\Omega} \rho_0\mathbf{A} dV \quad \Leftrightarrow \quad \int_{\partial\Omega} \mathbf{P}\mathbf{N} dS + \int_{\Omega} \rho_0\mathbf{B} dV = \int_{\Omega} \rho_0\mathbf{A} dV. \quad (1.39)$$

Practicing the divergence theorem once more, we see that

$$\int_{\partial\Omega} P_{iJ}N_J dS = \int_{\Omega} P_{iJ,J} dV \quad \Rightarrow \quad (\text{Div } \mathbf{P})_i = P_{iJ,J}, \quad (1.40)$$

which defines the *divergence of a second-order tensor*, which is a vector. Note that we use a *capitol operator* “Div” as opposed to “div” to indicate differentiation with respect to the undeformed coordinates.

When we again exploit that the above balance law must hold for all subbodies  $\omega \subset \Omega$ , we arrive at the *local statement of linear momentum balance*:

$$\boxed{\text{Div } \mathbf{P} + \rho_0\mathbf{B} = \rho_0\mathbf{A}} \quad \Leftrightarrow \quad \boxed{P_{iJ,J} + \rho_0 B_i = \rho_0 A_i} \quad (1.41)$$

Note that the special case of **quasistatics** assumes that inertial effects are negligible, so one solves the quasistatic linear momentum balance  $\text{Div } \mathbf{P} + \rho_0\mathbf{B} = \mathbf{0}$ . Except for gravity or electro/magnetomechanics, body forces also vanish in most cases, so that one simply arrives at  $\text{Div } \mathbf{P} = \mathbf{0}$ .

It is important to recall that stresses in finite deformations are not unique but we generally have different types of stress tensors. Above, we introduced the **first Piola-Kirchhoff stress tensor**  $\mathbf{P}$ , which implies *actual force per undeformed area*. Similarly, one can define the **Cauchy stress tensor**  $\boldsymbol{\sigma}$ , which denotes *actual force per deformed area*. The definition and link are given by

$$\sigma_{ij} = \frac{dF_i}{da_j}, \quad \boldsymbol{\sigma} = \frac{1}{J} \mathbf{P} \mathbf{F}^T. \quad (1.42)$$

Later on, it will be helpful to link stresses and deformation to energy. If the stored mechanical energy is characterized by the **strain energy density**  $W = W(\mathbf{F})$ , then one can show that

$$\mathbf{P} = \frac{\partial W}{\partial \mathbf{F}} \quad \Leftrightarrow \quad P_{iJ} = \frac{\partial W}{\partial F_{iJ}}. \quad (1.43)$$

Without knowing much about tensor analysis, we may interpret the above as a derivative of the energy density with respect to each component of  $\mathbf{F}$ , yielding the corresponding component of  $\mathbf{P}$ . This concept can also be extended to introduce a *fourth-order tensor*, the **incremental tangent modulus tensor**

$$\mathbb{C} = \frac{\partial \mathbf{P}}{\partial \mathbf{F}} \quad \Leftrightarrow \quad \mathbb{C}_{iJkL} = \frac{\partial P_{iJ}}{\partial F_{kL}}, \quad (1.44)$$

for which each component of  $\mathbf{P}$  is differentiated with respect to each component of  $\mathbf{F}$ . For further information on tensor analysis, see the tensor notes. Without further discussion, notice that a *stress-free* state implies that the energy attains an extremum.

Note that the dependence of  $W$  on  $\mathbf{F}$  (and consequently the constitutive law between  $\mathbf{P}$  and  $\mathbf{F}$ ) is generally strongly nonlinear, which limits opportunities for closed-form analytical solutions. Also, for *material frame indifference* we must in fact have  $W = W(\mathbf{C})$ , but that is a technical detail of minor importance here.

**Examples:** Switching between symbolic and index notation is often convenient (especially when taking derivatives) and should be practiced. Consider the following examples:

$$\bullet \quad \mathbf{a} = \mathbf{T} \mathbf{b} \quad \Leftrightarrow \quad a_i = T_{ij} b_j \quad (1.45)$$

$$\bullet \quad \mathbf{a} = \mathbf{T}^T \quad \Leftrightarrow \quad a_i = T_{ji} b_j \quad (1.46)$$

$$\bullet \quad \text{tr}(\mathbf{a} \otimes \mathbf{b}) \quad \Leftrightarrow \quad \text{tr}[a_i b_j] = a_i b_i = \mathbf{a} \cdot \mathbf{b} = b_i a_i = \text{tr}(\mathbf{b} \otimes \mathbf{a}) \quad (1.47)$$

$$\bullet \quad \text{tr}(\mathbf{R}^T \mathbf{T}) = \text{tr}[R_{ji} T_{jk}] = R_{ji} T_{ji} = \mathbf{R} \cdot \mathbf{T} \quad (1.48)$$

$$\bullet \quad \frac{\partial \mathbf{a}}{\partial \mathbf{a}} = \left[ \frac{\partial a_i}{\partial a_j} \right] = [\delta_{ij}] = \mathbf{I} \quad (1.49)$$

$$\bullet \quad \frac{\partial \text{tr} \mathbf{T}}{\partial \mathbf{T}} = \frac{\partial T_{kk}}{\partial \mathbf{T}} = \left[ \frac{\partial T_{kk}}{\partial T_{ij}} \right] = [\delta_{ki} \delta_{kj}] = [\delta_{ij}] = \mathbf{I} \quad (1.50)$$

$$\bullet \quad \frac{\partial \sqrt{\text{tr} \mathbf{T}}}{\partial \mathbf{T}} = \frac{\partial \sqrt{\text{tr} \mathbf{T}}}{\partial \text{tr} \mathbf{T}} \frac{\partial \text{tr} \mathbf{T}}{\partial \mathbf{T}} = \frac{\mathbf{I}}{2\sqrt{\text{tr} \mathbf{T}}} \quad (1.51)$$

$$\bullet \quad \frac{\partial \lambda(\mathbf{N})}{\partial \mathbf{F}} = \frac{\partial \sqrt{\mathbf{N} \cdot \mathbf{F}^T \mathbf{F} \mathbf{N}}}{\partial \mathbf{F}} = \frac{1}{2\lambda(\mathbf{N})} \left[ \frac{\partial N_K F_{jK} F_{jM} N_M}{\partial F_{iJ}} \right] \quad (1.52)$$

$$= \frac{1}{2\lambda(\mathbf{N})} [N_K \delta_{ij} \delta_{JK} F_{jM} N_M + N_K F_{jK} \delta_{ij} \delta_{JM} N_M] = \frac{[F_{iM} N_M N_J]}{\lambda(\mathbf{N})} \quad (1.53)$$

$$= \mathbf{F} \mathbf{M} \otimes \mathbf{N} / \lambda(\mathbf{N}) \quad (1.54)$$

**Examples:** Consider a special form of the compressible Neo-Hookean material model defined by

$$W(\mathbf{F}) = \frac{\mu}{2}(\text{tr} \bar{\mathbf{C}} - 3) + \frac{\kappa}{2}(J - 1)^2, \quad \bar{\mathbf{C}} = \bar{\mathbf{F}}^T \bar{\mathbf{F}} \quad \text{and} \quad \bar{\mathbf{F}} = \frac{\mathbf{F}}{J^{-1/3}}. \quad (1.55)$$

The first Piola-Kirchhoff stress tensor is computed as

$$\begin{aligned} \mathbf{P} &= \frac{\partial W}{\partial \mathbf{F}} = \frac{\partial}{\partial \mathbf{F}} \left[ \frac{\mu}{2} \left( \frac{\text{tr} \mathbf{F}^T \mathbf{F}}{J^{2/3}} - 3 \right) + \frac{\kappa}{2} (J - 1)^2 \right], \quad J = \det \mathbf{F} \\ &= \frac{\mu}{2} \left( \frac{1}{J^{2/3}} \frac{\partial \text{tr} \mathbf{F}^T \mathbf{F}}{\partial \mathbf{F}} + \text{tr} \mathbf{F}^T \mathbf{F} \frac{\partial}{\partial \mathbf{F}} \frac{1}{J^{2/3}} \right) + 2\kappa(J - 1) \frac{\partial J}{\partial \mathbf{F}} \\ &= \frac{\mu}{2} \left( \frac{2\mathbf{F}}{J^{2/3}} - \frac{2 \text{tr} \mathbf{F}^T \mathbf{F}}{3J^{5/3}} J \mathbf{F}^{-T} \right) + \kappa(J - 1) J \mathbf{F}^{-T} \\ &= \frac{\mu}{J^{2/3}} \mathbf{F} + \left[ \kappa(J - 1)J - \frac{\mu}{3J^{2/3}} \text{tr} \mathbf{F}^T \mathbf{F} \right] \mathbf{F}^{-T}, \end{aligned} \quad (1.56)$$

where we used

$$\frac{\partial J}{\partial \mathbf{F}} = \text{cof} \mathbf{F} = J \mathbf{F}^{-T} \quad \text{and} \quad \frac{\partial \text{tr} \mathbf{F}^T \mathbf{F}}{\partial \mathbf{F}} = 2\mathbf{F}. \quad (1.57)$$

Using index notation, the above identity is written as

$$P_{iJ} = \frac{\partial W}{\partial F_{iJ}} = \frac{\mu}{J^{2/3}} F_{iJ} + \left[ \kappa(J - 1)J - \frac{\mu}{3J^{2/3}} I_1 \right] F_{Ji}^{-1}, \quad I_1 = \text{tr} \mathbf{F}^T \mathbf{F} = F_{kL} F_{kL}. \quad (1.58)$$

The incremental stiffness tensor requires taking second derivatives:

$$\begin{aligned} \mathbb{C}_{iJkL} &= \frac{\partial P_{iJ}}{\partial F_{kL}} \\ &= \mu \left( F_{iJ} \frac{\partial}{\partial F_{kL}} \frac{1}{J^{2/3}} + \frac{1}{J^{2/3}} \frac{\partial F_{iJ}}{\partial F_{kL}} \right) + \kappa(2J - 1) \frac{\partial J}{\partial F_{kL}} F_{Ji}^{-1} + \kappa(J - 1) J \frac{\partial F_{Ji}^{-1}}{\partial F_{kL}} \\ &\quad - \frac{\mu}{3} \left[ \frac{\partial J^{-2/3}}{\partial F_{kL}} I_1 F_{Ji}^{-1} + \frac{1}{J^{2/3}} \frac{\partial I_1}{\partial F_{kL}} F_{Ji}^{-1} + \frac{I_1}{J^{2/3}} \frac{\partial F_{Ji}^{-1}}{\partial F_{kL}} \right] \\ &= \frac{\mu}{J^{2/3}} \delta_{ik} \delta_{JL} - \frac{2\mu}{3J^{2/3}} (F_{iJ} F_{Lk}^{-1} + F_{Ji}^{-1} F_{kL}) \\ &\quad + \left( \frac{2\mu}{9J^{2/3}} I_1 + \kappa(2J - 1)J \right) F_{Lk}^{-1} F_{Ji}^{-1} + \left( \frac{\mu}{3J^{2/3}} I_1 - \kappa(J - 1)J \right) F_{Jk}^{-1} F_{Li}^{-1}, \end{aligned} \quad (1.59)$$

where we rearranged the last equation by grouping terms and simplifying, and we used

$$\frac{\partial F_{Ji}^{-1}}{\partial F_{kL}} = -F_{Jk}^{-1} F_{Li}^{-1}. \quad (1.60)$$

To check the final answers, we verify that each side of the equation has the exact same free indices appearing only once.

### 1.3 A special case: linearized kinematics

Whenever only small deformation is expected, the above framework can be significantly simplified by using **linearized kinematics**. To this end, we assume that  $\|\text{Grad } \mathbf{u}\| \ll 1$  (“*small strains*”). Note that in this case it does not make a significant difference if we differentiate with respect to  $x_i$  or  $X_I$ , so that one generally uses only lower-case indices for simplicity.

In small strains, the *displacement field* is the key field to be determined (rather than the deformation mapping), i.e., we seek  $\mathbf{u} = \mathbf{u}(\mathbf{x}, t)$ .

Recall that

$$\mathbf{C} = \mathbf{F}\mathbf{F}^T = (\mathbf{I} + \text{Grad } \mathbf{u})(\mathbf{I} + \text{Grad } \mathbf{u}^T) = \mathbf{I} + \text{Grad } \mathbf{u} + (\text{Grad } \mathbf{u})^T + (\text{Grad } \mathbf{u})(\text{Grad } \mathbf{u})^T. \quad (1.61)$$

Now, the final term is dropped by a scaling argument ( $\|\text{Grad } \mathbf{u}\| \ll 1$ ). Therefore, we may introduce a kinematic relation like in the thermal problem:

$$\boldsymbol{\beta} = \text{Grad } \mathbf{u}, \quad (1.62)$$

and all important local deformation information is encoded in  $\boldsymbol{\beta}$ . Like a temperature gradient causes heat flux, a displacement gradient causes stresses (if displacements are constant everywhere, the body is undergoing **rigid body translation** and does not produce any stresses).

To make sure we also do not pick up **rigid body rotation**, one introduces the infinitesimal **strain tensor**

$$\boldsymbol{\varepsilon} = \frac{1}{2} [\text{Grad } \mathbf{u} + (\text{Grad } \mathbf{u})^T] = \frac{1}{2} (\boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T) \quad \Leftrightarrow \quad \varepsilon_{ij} = \frac{1}{2} (u_{i,j} + u_{j,i}). \quad (1.63)$$

Notice that, unlike in finite deformations, *no deformation* implies  $\boldsymbol{\varepsilon} = \mathbf{0}$ . Furthermore, by definition  $\boldsymbol{\varepsilon}$  is symmetric since  $\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}^T$  (not like  $\mathbf{F}$  which is asymmetric). The same applies for  $\boldsymbol{\sigma}$  and  $\mathbf{P}$  which are, respectively, symmetric and asymmetric.

As before, local deformation metrics are encoded into  $\boldsymbol{\varepsilon}$ . For example, volumetric deformation is characterized by the *trace* of  $\boldsymbol{\varepsilon}$ , viz.

$$\frac{dv}{dV} = 1 + \text{tr } \boldsymbol{\varepsilon} = 1 + \varepsilon_{ii}, \quad (1.64)$$

while stretches in the three coordinate directions are given by  $\varepsilon_{(ii)}$  (parentheses implying no summation over  $i$ ) and angle changes are identified as  $\gamma_{ij} = 2\varepsilon_{ij}$  with  $i \neq j$ .

In linearized kinematics, all three stress tensors coincide and one commonly uses only the Cauchy stress tensor  $\boldsymbol{\sigma}$  to define the constitutive relation  $\boldsymbol{\sigma} = \boldsymbol{\sigma}(\boldsymbol{\varepsilon})$ . In the simplest case of **linear elasticity**, those are linearly linked via

$$\boldsymbol{\sigma} = \mathbb{C} \boldsymbol{\varepsilon} \quad \Leftrightarrow \quad \sigma_{ij} = \mathbb{C}_{ijkl} \varepsilon_{kl} \quad (1.65)$$

with a fourth-order **elasticity tensor**  $\mathbb{C}$  linking each component of  $\boldsymbol{\sigma}$  to those of  $\boldsymbol{\varepsilon}$ . Alternatively, we can again encode the constitutive response in a *strain energy density*  $W = W(\boldsymbol{\varepsilon})$ , which, e.g., for the case of linear elasticity reads

$$W = \frac{1}{2} \boldsymbol{\varepsilon} \cdot \mathbb{C} \boldsymbol{\varepsilon} \quad \Leftrightarrow \quad W = \frac{1}{2} \varepsilon_{ij} \mathbb{C}_{ijkl} \varepsilon_{kl}, \quad (1.66)$$

so that

$$\boldsymbol{\sigma} = \frac{\partial W}{\partial \boldsymbol{\varepsilon}}, \quad \mathbb{C} = \frac{\partial \boldsymbol{\sigma}}{\partial \boldsymbol{\varepsilon}} = \frac{\partial^2 W}{\partial \boldsymbol{\varepsilon} \partial \boldsymbol{\varepsilon}}. \quad (1.67)$$

When taking derivatives with respect to  $\boldsymbol{\varepsilon}$ , caution is required since  $\boldsymbol{\varepsilon}$  is symmetric, so  $\varepsilon_{ij} = \varepsilon_{kl}$  and, consequently, derivatives with respect to  $\varepsilon_{ij}$  must also take into account those terms containing  $\varepsilon_{ji}$  (for  $i \neq j$ ). Therefore, the derivative should always be computed as

$$\frac{\partial}{\partial \varepsilon_{ij}} = \frac{1}{2} \left( \frac{\partial}{\partial \varepsilon_{ij}} + \frac{\partial}{\partial \varepsilon_{ji}} \right). \quad (1.68)$$

Alternatively, one may also use  $\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}^T$  and simply replace  $\boldsymbol{\varepsilon} = \frac{1}{2}(\boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T)$  before differentiating.

The traction vector on a surface with unit normal  $\mathbf{n}$  now becomes  $\mathbf{t} = \boldsymbol{\sigma} \mathbf{n}$ .

The local statement of linear momentum balance in linearized kinematics is

$$\boxed{\operatorname{div} \boldsymbol{\sigma} + \rho \mathbf{b} = \rho \mathbf{a}} \quad \Leftrightarrow \quad \boxed{\sigma_{ij,j} + \rho b_i = \rho a_i} \quad (1.69)$$

where the small-strain versions of density, body force density and acceleration field were introduced as  $\rho$ ,  $\mathbf{b}$  and  $\mathbf{a} = \ddot{\mathbf{u}}$ , respectively.

In small strains, we can insert the kinematic and linear elastic constitutive relations as well as the definition of the acceleration field into linear momentum balance to obtain

$$(\mathbb{C}_{ijkl} \varepsilon_{kl})_{,j} + \rho b_i = \rho \ddot{u}_i \quad \Leftrightarrow \quad (\mathbb{C}_{ijkl} u_{k,l})_{,j} + \rho b_i = \rho \ddot{u}_i \quad (1.70)$$

and in case of a *homogeneous* body with  $\mathbb{C}(\mathbf{x}) = \mathbb{C} = \text{const.}$  we finally arrive at

$$\mathbb{C}_{ijkl} u_{k,lj} + \rho b_i = \rho \ddot{u}_i, \quad (1.71)$$

which is known as **Navier's equation** to be solved for the unknown field  $\mathbf{u}(\mathbf{x}, t)$ .

Finally, the following will be helpful when implementing material models in our code. Note that we may use the relations

$$\varepsilon_{ij} = \frac{1}{2}(u_{i,j} + u_{j,i}) \quad \text{and} \quad F_{iJ} = \delta_{iJ} + u_{i,J} \quad (1.72)$$

to write (using the chain rule)

$$\frac{\partial W}{\partial u_{i,J}} = \frac{\partial W}{\partial F_{kL}} \frac{\partial F_{kL}}{\partial u_{i,J}} = P_{kL} \frac{\partial}{\partial u_{i,J}} (\delta_{kL} + u_{k,L}) = P_{kL} \delta_{ik} \delta_{JL} = P_{iJ} \quad (1.73)$$

and (exploiting the symmetry of  $\boldsymbol{\sigma}$ )

$$\frac{\partial W}{\partial u_{i,j}} = \frac{\partial W}{\partial \varepsilon_{kl}} \frac{\partial \varepsilon_{kl}}{\partial u_{i,j}} = \sigma_{kl} \frac{\partial}{\partial u_{i,j}} \frac{1}{2}(u_{k,l} + u_{l,k}) = \frac{1}{2} \sigma_{kl} (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}) = \frac{1}{2} (\sigma_{ij} + \sigma_{ji}) = \sigma_{ij}. \quad (1.74)$$

Therefore, we can use the alternative relations for our stresses

$$\boxed{P_{iJ} = \frac{\partial W}{\partial u_{i,J}}} \quad \text{and} \quad \boxed{\sigma_{ij} = \frac{\partial W}{\partial u_{i,j}}}. \quad (1.75)$$

The beauty in those relations is that the stress tensor definition is now *identical irrespective of whether we are working in linearized or finite kinematics*.

## 1.4 Summary and Looking Ahead

So far, we have seen how partial differential equations govern the thermal and mechanical behavior of solid bodies (and, of course, those two can be coupled as well to describe the thermo-mechanical behavior of deformable bodies). In order to solve a problem, we need an **initial boundary value problem** (IBVP), which furnishes the above equations with appropriate **boundary conditions** (BCs) and **initial conditions** (ICs).

To this end, we subdivide the boundary  $\partial\Omega$  of a body  $\Omega$  into

- $\partial\Omega_D \equiv$  **Dirichlet boundary**, prescribing the *primary field* ( $\varphi$ ,  $T$ , etc.):

$$\text{e.g.} \quad \mathbf{u}(\mathbf{x}, t) = \hat{\mathbf{u}}(\mathbf{x}, t) \text{ on } \partial\Omega_D \quad \text{or} \quad T(\mathbf{x}, t) = \hat{T}(\mathbf{x}, t) \text{ on } \partial\Omega_D. \quad (1.76)$$

- $\partial\Omega_N \equiv$  **Neumann boundary**, prescribing *derivatives of the primary field* ( $\mathbf{F}$ ,  $\beta$ , etc.):

$$\text{e.g.} \quad \mathbf{t}(\mathbf{x}, t) = \boldsymbol{\sigma}(\mathbf{x}, t)\mathbf{n}(\mathbf{x}, t) = \hat{\mathbf{t}}(\mathbf{x}, t) \text{ on } \partial\Omega_N \quad \text{or} \quad q(\mathbf{x}, t) = \hat{q}(\mathbf{x}, t) \text{ on } \partial\Omega_N, \quad (1.77)$$

Note that we may generally assume that

$$\partial\Omega_D \cup \partial\Omega_N = \partial\Omega \quad \text{and in most problems also} \quad \partial\Omega_D \cap \partial\Omega_N = \emptyset. \quad (1.78)$$

In addition, all time-dependent problems require **initial conditions**, e.g.,

$$\begin{aligned} T(\mathbf{x}, 0) &= T_0(\mathbf{x}) \quad \forall \mathbf{X} \in \Omega, \\ \text{or} \quad \boldsymbol{\varphi}(\mathbf{X}, 0) &= \mathbf{x}_0(\mathbf{X}) \quad \text{and} \quad \mathbf{V}(\mathbf{X}, 0) = \mathbf{V}_0(\mathbf{X}) \quad \forall \mathbf{X} \in \Omega. \end{aligned} \quad (1.79)$$

The number of required BCs/ICs depends on the **order of a PDE**, e.g.,

$$\rho c_v \dot{T} = \text{div}(\mathbf{K} \text{grad } T) + \rho s \quad (1.80)$$

is first-order in time and therefore requires one IC, e.g.,  $T(\mathbf{x}, 0) = T_0(\mathbf{x})$ . It is second-order in space and hence requires BCs along all  $\partial\Omega$  (e.g., two conditions per  $x$  and  $y$  coordinates).

In summary, we will have governing PDEs supplemented by ICs and BCs, as required (e.g., quasistatic problems, of course, do not require any initial conditions). Those need to be solved for the primary fields (e.g., temperature  $T$  or displacements  $\mathbf{u}$  or the deformation mapping  $\boldsymbol{\varphi}$ ).

Unfortunately, analytical solutions are hardly ever available – except for relatively simple problems involving

- simple geometries,
- simple material behavior,
- simple ICs/BCs.

For realistic geometries, materials and/or ICs/BCs, one usually requires numerical techniques to obtain *approximate solutions*.

## 2 Numerical Methods

To numerically solve such ODEs/PDEs, we generally have so-called direct and indirect methods.

**Direct methods** aim to solve the governing equations directly; for example, using **finite differences** (FD). Consider the isotropic heat equation discussed before, which for brevity we write in 1D as (absorbing the coefficient into  $k$  and  $r$ ).

$$\dot{T} = k T_{,xx} + r. \quad (2.1)$$

Introduce a regular  $(\Delta x, \Delta t)$ -grid with  $T_i^\alpha = T(x_i, t^\alpha)$  and use **Taylor expansions**, e.g., in space:

$$T(x_{i+1}, t^\alpha) = T_{i+1}^\alpha = T_i^\alpha + \Delta x \left. \frac{\partial T}{\partial x} \right|_{x_i, t^\alpha} + \frac{(\Delta x)^2}{2} \left. \frac{\partial^2 T}{\partial x^2} \right|_{x_i, t^\alpha} + \frac{(\Delta x)^3}{3!} \left. \frac{\partial^3 T}{\partial x^3} \right|_{x_i, t^\alpha} + O(\Delta x^4) \quad (2.2)$$

$$T(x_{i-1}, t^\alpha) = T_{i-1}^\alpha = T_i^\alpha - \Delta x \left. \frac{\partial T}{\partial x} \right|_{x_i, t^\alpha} + \frac{(\Delta x)^2}{2} \left. \frac{\partial^2 T}{\partial x^2} \right|_{x_i, t^\alpha} - \frac{(\Delta x)^3}{3!} \left. \frac{\partial^3 T}{\partial x^3} \right|_{x_i, t^\alpha} + O(\Delta x^4) \quad (2.3)$$

Addition of the two equations gives:

$$T_{i+1}^\alpha + T_{i-1}^\alpha = 2T_i^\alpha + (\Delta x)^2 \left. \frac{\partial^2 T}{\partial x^2} \right|_{x_i, t^\alpha} + O(\Delta x^4) \Rightarrow \boxed{\frac{\partial^2 T}{\partial x^2}(x_i, t^\alpha) = \frac{T_{i+1}^\alpha - 2T_i^\alpha + T_{i-1}^\alpha}{(\Delta x)^2} + O(\Delta x^2)} \quad (2.4)$$

This is the **second-order central difference** approximation.

Analogously, Taylor expansion in time and subtraction of the two equations gives:

$$T_i^{\alpha+1} - T_i^{\alpha-1} = 2\Delta t \left. \frac{\partial T}{\partial t} \right|_{x_i, t^\alpha} + O(\Delta t^3) \Rightarrow \boxed{\frac{\partial T}{\partial t}(x_i, t^\alpha) = \frac{T_i^{\alpha+1} - T_i^{\alpha-1}}{2\Delta t} + O(\Delta t^2)} \quad (2.5)$$

which is the **first-order central difference** approximation. Many other such finite-difference approximations of derivatives can be obtained in a similar fashion. For example, a simpler first-order stencil is obtained from the first Taylor equation (2.2) alone:

$$T_i^{\alpha+1} - T_i^\alpha = \Delta t \left. \frac{\partial T}{\partial t} \right|_{x_i, t^\alpha} + O(\Delta t^2) \Rightarrow \boxed{\frac{\partial T}{\partial t}(x_i, t^\alpha) = \frac{T_i^{\alpha+1} - T_i^\alpha}{\Delta t} + O(\Delta t)} \quad (2.6)$$

which is often referred to as the first-order **forward-Euler** approximation. Analogously, we can use the second Taylor equation (2.3) to obtain the **backward-Euler** approximation

$$T_i^\alpha - T_i^{\alpha-1} = \Delta t \left. \frac{\partial T}{\partial t} \right|_{x_i, t^\alpha} + O(\Delta t^2) \Rightarrow \boxed{\frac{\partial T}{\partial t}(x_i, t^\alpha) = \frac{T_i^\alpha - T_i^{\alpha-1}}{\Delta t} + O(\Delta t)} \quad (2.7)$$

In order to numerically solve a PDE directly, we choose suitable finite-difference approximations for all appearing derivatives. For example, using the second-order central-difference approximation for the spatial and the forward-Euler approximation for the temporal derivative in the heat equation, the *discretized governing equation* becomes

$$\frac{T_i^{\alpha+1} - T_i^\alpha}{\Delta t} = k \frac{T_{i+1}^\alpha - 2T_i^\alpha + T_{i-1}^\alpha}{(\Delta x)^2} + r(x_i, t^\alpha) + O(\Delta t, \Delta x^2), \quad (2.8)$$



which in the limit  $\Delta t, \Delta x \rightarrow 0$  is expected to converge towards the same solution as the governing equation (this is the requirement of *consistency* of the discretized equation).

Note that for known values  $T_i^\alpha$  at the current time, the above equation can easily be solved for  $T_i^{\alpha+1}$  at the new time. In fact, the right-hand side does not involve  $T_i^{\alpha+1}$ , which is why this finite-difference scheme is **explicit**.

By contrast, when using the backward-Euler approximation, we obtain

$$\frac{T_i^\alpha - T_i^{\alpha-1}}{\Delta t} = k \frac{T_{i+1}^\alpha - 2T_i^\alpha + T_{i-1}^\alpha}{(\Delta x)^2} + r(x_i, t^\alpha) + O(\Delta t, \Delta x^2), \quad (2.9)$$

which is a linear system to be solved for  $T_i^{\alpha+1}$  at the new time step and is therefore an **implicit** scheme.

Numerical solution can be interpreted via **stencils**, which may also reveal the required BCs/ICs.

**Problems** associated with direct methods include:

- a *regular grid* is required (which is fine for many fluidic mechanics problems but oftentimes problematic for complex solid geometries).
- variables are *defined only at grid points*, hence the error is minimized only at grid points (and we have no information about what happens between grid points; both the primary fields and their errors are undefined between grid points). This can be problematic when seeking approximate solutions that are “*globally optimal*”. Also, how to apply BCs/ICs in between grid points, how about moving BCs?
- *stability/efficiency issues* (probably known from fundamental computational mechanics classes: CFL-condition, von Neumann stability analysis, etc.). In a nutshell, the choice of  $\Delta t$  and  $\Delta x$  is not arbitrary but – aside from accuracy concerns – the stability of, especially explicit, finite-difference schemes dictates maximum step widths to be used.

As an **alternative**, **indirect methods** do not solve the ODEs/PDEs directly but search for optimal approximations, e.g.,  $\mathbf{u}^h(\mathbf{x}) \approx \mathbf{u}(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ , including BCs and ICs. To do so, we need to discuss a lot more.

Particular **questions** to be addressed include:

- How do we choose  $\mathbf{u}^h(\mathbf{x})$ ? For example, globally or locally defined functions? Which choices minimize the error?
- What is “*optimal*”? How do we quantify between error approximation and exact solution?
- What trial functions shall we use? For example, polynomial or Fourier series, piecewise defined or maybe even piece-constant?

We need a couple concepts to address those questions. Note that in the following, we will formulate most concepts in 1D with analogous generalizations possible for higher dimensions unless specifically mentioned.

## 3 Variational Calculus

### 3.1 Functionals

In order to understand the big picture of indirect methods, let us first discuss the energetics of deformable bodies and, in particular, introduce the concepts of functionals and variational calculus.

A **functional** is special type of mapping which maps from a function space  $\mathcal{U}$  to  $\mathbb{R}$ :

$$I : u \in \mathcal{U} \rightarrow I[u] \in \mathbb{R}. \quad (3.1)$$

Oftentimes, functionals impose constraints on the function space  $\mathcal{U}$ , usually differentiability/integrability constraint. As an example, consider

$$I[u] = \int_0^1 u^2(x) \, dx, \quad (3.2)$$

which is a *functional* requiring that  $u$  is square-integrable.

It will be important to restrict the space of admissible functions when seeking for "physically reasonable" approximations  $u^h(x)$ . For example, we may want to restrict how smooth a function is and whether or not it is allowed to have any poles or discontinuities, etc.

To this end, we quickly introduce the **Sobolev space**

$$H^k(\Omega) = \{u : \Omega \rightarrow \mathbb{R} \text{ such that } \|u\|_{H^k(\Omega)} < \infty\}, \quad (3.3)$$

with the **Sobolev norm** (in 1D)

$$\|u\|_{H^k(\Omega)}^2 = \int_{\Omega} u(x)^2 \, dx + \int_{\Omega} [u'(x)]^2 \, dx + \dots + \int_{\Omega} [u^{(k)}(x)]^2 \, dx. \quad (3.4)$$

Consequently,  $H^k(\Omega)$  denotes the space of all functions whose derivatives up to  $k$ th order are **square-integrable**. The above example in (3.2), e.g., requires  $u \in \mathcal{U} \subset H^0(0,1)$ . That is functions  $u(x)$  must be square-integrable on the interval  $(0,1)$ .

Our classical example will be the energy of a mechanical system which depends on the displacement field  $\mathbf{u}(\mathbf{x})$  and defines an energy  $I \in \mathbb{R}$ . Consider, e.g., the 1D strain energy of a bar of length  $L$  and with constant Young's modulus  $E$  and cross-sectional area  $A$ :

$$I[u] = \int_0^L \frac{E}{2} [u_{,x}(x)]^2 A \, dx, \quad (3.5)$$

where we used that  $W = \frac{E}{2} \varepsilon^2$  and  $\varepsilon = u_{,x}$  with  $u = u(x)$ . Here, we generally may want to impose the restriction  $u \in H^1(0,L)$ , unless dealing with discontinuities such as cracks.

Functionals are to be distinguished from *functions* such as  $f(\mathbf{x}) = \sqrt{x_1^2 + x_2^2} = |\mathbf{x}|$  with  $f : \mathbb{R}^2 \rightarrow \mathbb{R}_0^+$ . Unlike a *function* which is a mapping from  $\mathbb{R}^d \rightarrow \mathbb{R}$ , a *functional's* domain is generally a function space  $\mathcal{U}$  (e.g., all polynomial functions up to a certain degree, or all continuously differentiable functions, or all piecewise polynomial functions).

### 3.2 Variations

Consider a functional  $I : \mathcal{U} \rightarrow \mathbb{R}$  such as the potential energy. Analogous to the stationarity condition of classical optimization problems, a necessary condition for an extremum of  $I$  is that

the first variation of  $I$  vanishes, i.e.,  $\delta I[u] = 0$  (this is the **stationarity condition**). Like in functional analysis, we are going to compute the analog of a derivative to identify maxima and minima of the functional. To this end, we perturb the functional around a point by a small variation and verify if the functional increases or decreases. Complicating is the fact that a "point" is now in fact a function, and a variation must be a variation of that function. To this end, we define the following.

A **variation**  $\delta u$  is an arbitrary function that represents admissible changes of  $u$ . If  $\Omega \subset \mathbb{R}^d$  is the domain of  $u \in \mathcal{U}$  with boundary  $\partial\Omega$ , we seek solutions

$$u \in \mathcal{U} = \{u \in H^k(\Omega) : u = \hat{u} \text{ on } \partial\Omega_D\} \quad (3.6)$$

then the variation must satisfy

$$\delta u \in \mathcal{U}_0 = \{\delta u \in H^k(\Omega) : \delta u = 0 \text{ on } \partial\Omega_D\}. \quad (3.7)$$

$k$  can be determined from the specific form of  $I[u]$ , as will be discussed later. The fact that variations  $\delta u$  must vanish on the Dirichlet boundary  $\partial\Omega_D$  stems from the need for perturbations that allow the perturbed function  $u + \delta u$  to still satisfy the Dirichlet boundary conditions.

With this, we define the **first variation** of  $I$  (analog of a first derivative) as

$$\delta I[u] = \lim_{\epsilon \rightarrow 0} \frac{I[u + \epsilon \delta u] - I[u]}{\epsilon} = \left. \frac{d}{d\epsilon} I[u + \epsilon \delta u] \right|_{\epsilon \rightarrow 0} \quad (3.8)$$

and analogously higher-order variations via

$$\delta^k I[u] = \delta(\delta^{k-1} I) \quad \text{for} \quad k \geq 2 \quad (3.9)$$

Note that a *Taylor expansion of a functional*  $I$  can now be written as

$$I[u + \delta u] = I[u] + \delta I[u] + \frac{1}{2!} \delta^2 I[u] + \frac{1}{3!} \delta^3 I[u] + \dots \quad (3.10)$$

The following are helpful relations for  $u, v \in \mathcal{U}$ , further  $I_i : \mathcal{U} \rightarrow \mathcal{V} \subset \mathbb{R}$ , and constants  $\alpha_i \in \mathbb{R}$ :

- $\delta(\alpha_1 I_1 + \alpha_2 I_2) = \alpha_1 \delta I_1 + \alpha_2 \delta I_2$
- $\delta(I_1 I_2) = (\delta I_1) I_2 + I_1 (\delta I_2)$
- $\delta \frac{du}{dx} = \frac{d}{dx} \delta u$  (assuming differentiability of  $u$ , specifically  $u \in C^1$ )
- $\delta \int_{\Omega} u \, dx = \int_{\Omega} \delta u \, dx$  (assuming  $\Omega$  is independent of  $u$ )
- $\delta I[u, v, \dots] = \frac{d}{d\epsilon} I[u + \epsilon \delta u, v + \epsilon \delta v, \dots]_{\epsilon \rightarrow 0}$

### Example:

Let us consider

$$I[u] = \|u\|_{L_2(0,1)}^2 = \int_0^1 u^2 \, dx \quad \text{so that we seek} \quad u \in \mathcal{U} = H^0(0,1). \quad (3.11)$$

The variations follow as

$$\begin{aligned}\delta I &= \lim_{\epsilon \rightarrow 0} \frac{d}{d\epsilon} \int_0^1 (u + \epsilon \delta u)^2 dx = \lim_{\epsilon \rightarrow 0} \int_0^1 2(u + \epsilon \delta u) \delta u dx = 2 \int_0^1 u \delta u dx \\ \delta^2 I &= \lim_{\epsilon \rightarrow 0} \delta I[u + \epsilon \delta u] = \frac{d}{d\epsilon} \int_0^1 2(u + \epsilon \delta u) \delta u dx = 2 \int_0^1 (\delta u)^2 dx \\ \delta^k I &= 0 \quad \text{for all } k > 2.\end{aligned}\tag{3.12}$$

Notice that

$$\begin{aligned}I[u + \delta u] &= \int_0^1 (u + \delta u)^2 dx = \int_0^1 u^2 dx + \int_0^1 2u \delta u dx + \int_0^1 (\delta u)^2 dx \\ &= I[u] + \delta I[u] + \frac{1}{2} \delta^2 I[u].\end{aligned}\tag{3.13}$$

As a *practical hint*, note that for any integral

$$I[\mathbf{u}] = \int_{\Omega} f(u_i, u_{i,j}, \dots) dV\tag{3.14}$$

we can write

$$\frac{\partial I[\mathbf{u} + \epsilon \delta \mathbf{u}]}{\partial \epsilon} = \int_{\Omega} \left( \frac{\partial f}{\partial u_i} \delta u_i + \frac{\partial f}{\partial u_{i,j}} \delta u_{i,j} + \dots \right) dV\tag{3.15}$$

and omit the lengthy derivation of the first variation given above.

Here comes the reason we look into variations: some classes of partial differential equations possess a so-called **variational structure**; i.e., their solutions  $u \in \mathcal{U}$  can be interpreted as extremal points over  $\mathcal{U}$  of a functional  $I[u]$ .

### 3.3 Example: Hanging bar under its own weight

Consider a bar of length  $L$  (Young's modulus  $E$ , cross-sectional area  $A$ , density  $\rho$ ) that is hanging from the ceiling and deformed under its own weight (gravitational acceleration  $g$ ). Let us denote the 1D displacement field  $u(x)$  where  $x$  runs from top to bottom of the bar. The total potential energy is thus the total strain energy of the elastic bar minus the work done by the gravitational body forces:

$$I[u] = \int_0^L \frac{E}{2} u_{,x}^2(x) A dx - \int_0^L \rho g u(x) A dx \quad \text{and} \quad u(x) \in \mathcal{U} = \{u \in H^1(0, L) : u(0) = 0\}.\tag{3.16}$$

The first variation yields

$$\begin{aligned}\delta I[u] &= \int_0^L E u_{,x}(x) \delta u_{,x}(x) A dx - \int_0^L \rho g \delta u(x) A dx \\ &= - \int_0^L [E u_{,xx}(x) + \rho g] \delta u(x) A dx + E A u_{,x}(L) \delta u(L) - E A u_{,x}(0) \delta u(0) = 0,\end{aligned}\tag{3.17}$$

where we used integration by parts of the first integral to arrive at the final form. Noting that  $\delta u_x(0)$  because of boundary conditions, the last term vanishes. Finally, recall that the above variation must vanish for *all* variations  $\delta u \in \mathcal{U}_0$ . This implies that (3.17) implies we must have

$$E u_{,xx}(x) + \rho g = 0 \quad \text{and} \quad E A u_{,x}(L) = 0.\tag{3.18}$$

These are exactly the governing equations and traction-free boundary condition that the bar needs to satisfy. Hence, we have shown that minimizing (3.16) over all  $u(x) \in \mathcal{U}$  is equivalent to solving (3.18) with  $u(0) = 0$ . That is, we have to theoretical strategy to replace the solution of a differential equation by a minimization problem.

Note that we can easily find the analytical solution to the problem by integrating (3.18) twice for  $u(x)$  and inserting the boundary conditions, resulting in

$$u(x) = -\frac{\rho g}{2E}x(x - 2L). \quad (3.19)$$

By the way, this solution to the system of differential equations is called the **classical solution**.

One way to exploit the above variational structure is the so-called **Rayleigh-Ritz** approach, which introduces an approximation  $u^h(x) \approx u(x)$ , e.g., a polynomial series

$$u^h(x) = \sum_{i=0}^n c_i x^i \quad \text{where} \quad c_0 = 0 \text{ because of BCs} \quad (3.20)$$

with unknown coefficients  $c_i \in \mathbb{R}$ . Of course, any choice of ansatz functions is permissible including, e.g., Fourier series of cosine or sine terms, as long as they satisfy the essential boundary condition  $u(0) = 0$  and the differentiability/integrability requirements (e.g., a piecewise linear guess for  $u^h(x)$  would not be permissible as its derivatives are not square-integrable). Next, we insert  $u^h(x)$  into (3.16), i.e.,

$$I[u^h] = \int_0^L \frac{E}{2} (u_{,x}^h)^2(x) A dx - \int_0^L \rho g u^h(x) A dx, \quad (3.21)$$

which can be integrated to depend only on the unknown coefficients  $c_i$ . Finally, we find the solution by minimization with respect to the coefficients:

$$0 = \frac{\partial I[u^h]}{\partial c_i}, \quad (3.22)$$

which gives  $n$  equations for the  $n$  unknown coefficients. Note that the above form of the energy results in a linear system of equations for the unknown coefficients  $c_i$ , which can be solved numerically in an efficient manner.

If we use polynomial ansatz functions as in (3.20), then the exact solution (3.19) is contained in the solution space if  $n \geq 2$ . In fact, if one chooses  $n \geq 2$  and solves for the unknown coefficients, then one obtains

$$c_1 = \frac{\rho g}{E}L \quad c_2 = -\frac{\rho g}{2E}, \quad c_i = 0 \quad \text{for} \quad i > 2, \quad (3.23)$$

which is the exact solution (3.19).

### 3.4 Example: Static Heat Conduction

As an introductory example, let us review the *static heat conduction problem* in  $d$  dimensions, defined by ( $\mathbf{N} \in \mathbb{R}^d$  denoting the outward unit normal vector)

$$\begin{cases} \kappa \Delta T + \rho s = 0 & \text{in } \Omega, \\ T = \hat{T} & \text{on } \partial\Omega_D, \\ q = -\kappa \text{grad } T \cdot \mathbf{n} = \hat{q} & \text{on } \partial\Omega_N, \end{cases} \quad (3.24)$$

and we seek solutions  $T : \Omega \rightarrow \mathbb{R}$  that satisfy all of the above equations and meet the required differentiability conditions. Such solutions are called **classical solution**.

As an alternative to solving the above equations, consider the total potential energy defined by the functional  $I : \mathcal{U} \rightarrow \mathbb{R}$  with

$$I[T] = \int_{\Omega} \left( \frac{\kappa}{2} \|\text{grad } T\|^2 - \rho s T \right) dV + \int_{\partial\Omega_N} \hat{q} T dS. \quad (3.25)$$

The specific form shows that we need to seek solutions in the space

$$\mathcal{U} = \{T \in H^1(\Omega) : T = \hat{T} \text{ on } \partial\Omega_D\} \quad \text{and} \quad \mathcal{U}_0 = \{\delta T \in H^1(\Omega) : \delta T = 0 \text{ on } \partial\Omega_D\}. \quad (3.26)$$

Let us find extremal points  $T \in \mathcal{U}$  that render  $I[T]$  stationary.

The first variation follows as

$$\delta I[T] = \int_{\Omega} \left( \frac{\kappa}{2} 2T_{,i} \delta T_{,i} - \rho s \delta T \right) dV + \int_{\partial\Omega_N} \hat{q} \delta T dS = 0 \quad \text{for all } \delta T \in \mathcal{U}_0. \quad (3.27)$$

Application of the divergence theorem to the first term yields

$$\int_{\partial\Omega} \kappa T_{,i} n_i \delta T dS - \int_{\Omega} \kappa T_{,ii} \delta T dV - \int_{\Omega} \rho s \delta T dV + \int_{\partial\Omega_N} \hat{q} \delta T dS = 0 \quad \text{for all } \delta T \in \mathcal{U}_0. \quad (3.28)$$

Rearranging terms and using the fact that  $\delta T = 0$  on  $\partial\Omega_D$  leads to

$$- \int_{\Omega} (\kappa T_{,ii} + \rho s) \delta T dV + \int_{\partial\Omega_N} (\kappa T_{,i} n_i + \hat{q}) \delta T dS = 0 \quad \text{for all } \delta T \in \mathcal{U}_0. \quad (3.29)$$

This must hold for all admissible variations  $\delta T \in \mathcal{U}_0$ . Therefore, (3.29) is equivalent to stating

$$\kappa \Delta T + \rho s = 0 \text{ in } \Omega, \quad -\kappa(\text{grad } T)\mathbf{n} = \hat{q} \text{ on } \partial\Omega_N \quad \text{and} \quad T = \hat{T} \text{ on } \partial\Omega_D. \quad (3.30)$$

Ergo, extremal points  $T \in \mathcal{U}$  of (3.25) are guaranteed to satisfy the governing equations (3.24) and are thus classical solutions.

To see if it is a *maximizer or minimizer*, let us compute the second variation

$$\delta^2 I[T] = \int_{\Omega} \kappa \delta T_{,i} \delta T_{,i} dV = \int_{\Omega} \kappa \|\delta \text{grad } T\|^2 dV \geq 0. \quad (3.31)$$

Hence, the extremum is a minimizer, assuming that  $\kappa > 0$ . Otherwise, note that  $\kappa < 0$  leads to solutions being (unstable) energy maxima, which implies that  $\kappa > 0$  is a (necessary and sufficient) *stability condition*.

Notice that (assuming that  $\kappa = \text{const.}$ ) we can rewrite the energy functional for short as

$$I[T] = \frac{1}{2} \mathcal{B}(T, T) - \mathcal{L}(T), \quad (3.32)$$

where we introduced the **bilinear form**  $\mathcal{B}$  and the **linear form**  $\mathcal{L}$  as

$$\mathcal{B}(\cdot, \cdot) = \kappa \langle \text{grad } \cdot, \text{grad } \cdot \rangle_{\Omega} \quad \text{and} \quad \mathcal{L}(\cdot) = \langle \rho s, \cdot \rangle_{\Omega} - \langle \hat{q}, \cdot \rangle_{\partial\Omega_N} \quad (3.33)$$

with the inner product operator

$$\langle f, g \rangle_{\Omega} = \int_{\Omega} f g dV. \quad (3.34)$$

This is in fact a recipe for a more general class of variational problems: let us consider an energy functional of the general form

$$\begin{aligned} I[u] &= \frac{1}{2} \mathcal{B}(u, u) - \mathcal{L}_\Omega(u) - \mathcal{L}_{\partial\Omega}(u) \\ &= \frac{1}{2} \kappa \langle \text{grad } u, \text{grad } u \rangle_\Omega - \langle \rho s, u \rangle_\Omega - \langle \hat{q}, u \rangle_{\partial\Omega_N} \end{aligned} \quad (3.35)$$

with  $u \in \mathcal{U}$  being some (scalar- or vector-valued) mapping and  $\rho s$  and  $\hat{q}$  denoting, respectively, distributed body sources and surface fluxes. Now we have

$$\begin{aligned} \delta I[u] &= \mathcal{B}(u, \delta u) - \mathcal{L}(\delta u) \\ &= - \int_\Omega [\kappa \text{div}(\text{grad } u) + \rho s] \delta u \, dV - \int_{\partial\Omega_N} [\hat{q} - \kappa(\text{grad } u) \cdot \mathbf{n}] \delta u \, dS. \end{aligned} \quad (3.36)$$

Thus, the energy density (3.35) is generally suitable for quasistatic problems of the type

$$\begin{cases} \kappa \Delta u + \rho s = 0 & \text{in } \Omega \\ u = \hat{u} & \text{on } \partial\Omega_D \\ \kappa(\text{grad } u) \cdot \mathbf{n} = \hat{q} & \text{on } \partial\Omega_N \end{cases} \quad (3.37)$$

Note that (3.37) describes not only heat conduction but the general form also applies to electromagnetism, elasticity (to be discussed later), and various other fields. Notice that, while (3.35) required  $u \in H^1(\Omega)$  (highest derivatives are of first order), evaluating (3.37) in general requires that  $u \in C^2(\Omega) \cap C^0(\overline{\Omega})$  (second derivatives are required). We will get back to this point later.

For notational purposes, let us adopt the following notation found in various textbooks on finite elements: the first variation is usually abbreviated as an operator acting on both the unknown field  $u$  and its variation  $\delta u$ ; i.e., we write  $\mathcal{G} : \mathcal{U} \times \mathcal{U}_0 \rightarrow \mathcal{V} \subset \mathbb{R}$  with

$$\mathcal{G}(u, \delta u) = D_{\delta u} I[u] = \lim_{\epsilon \rightarrow 0} \frac{d}{d\epsilon} I[u + \epsilon \delta u] \quad (3.38)$$

One of the beauties of the above variational problem (3.37) is that a unique minimizer exists by the **Lax-Milgram theorem**, see Appendix E. Recall that for the linear heat problem above we already showed that the solution is a unique (global) minimizer if  $\kappa > 0$ .

## 4 The weak form

### 4.1 Classical and weak solutions

Consider a physical problem that is – as before – governed by the so-called **strong form**

$$\begin{cases} (\kappa u_{,i})_{,i} + s = 0 & \text{in } \Omega \\ u_i = \hat{u}_i & \text{on } \partial\Omega_D \\ \kappa u_{,i} n_i = \hat{q} & \text{on } \partial\Omega_N. \end{cases} \quad (4.1)$$

In order to describe the restrictions  $u$  must satisfy, let us define that a function  $u$  is of **class**  $C^k(\Omega)$  with an integer  $k \geq 0$  if it is  $k$  times continuously differentiable over  $\Omega$  (i.e.,  $u$  possesses derivatives up to the  $k$ th order and these derivatives are continuous functions).

**Examples:**

- Any  $k$ th-order polynomial  $u(x)$  with  $k \geq 0$  is generally  $C^\infty(\mathbb{R})$ .
- Consider a continuous, piecewise-linear function  $u : \Omega = (0, 2) \rightarrow \mathbb{R}$ . Function  $u$  is  $C^0(\Omega)$  but not  $C^1(\Omega)$ .
- The Heavyside function  $H(x)$  is said to be  $C^{-1}(\mathbb{R})$  since its “**zeroth derivative**” (i.e., the function itself) is not continuous.

If there are no discontinuities such as cracks, shocks, etc. (or discontinuities in the BCs/ICs) we usually assume that the classical solution fields are  $C^\infty(\Omega)$ , so we may take derivatives; otherwise, derivatives exist almost everywhere (**a.e.**)

It is convenient to define by  $C_0^k(\Omega)$  the space of all functions contained in  $C^k(\Omega)$  whose support is a *bounded subset* of  $\Omega$  (i.e.,  $u(x) \neq 0$  only on a finite subset of  $\Omega$ ). Then, notice that

$$C_0^k(\Omega) \subset H_0^k(\Omega) \quad (4.2)$$

and

$$C_0^\infty(\Omega) = \bigcap_{k \geq 0} C_0^k(\Omega). \quad (4.3)$$

Going back to the problem described by the strong form (4.1) above, we need to seek solutions

$$u \in C^2(\Omega) \cap C^0(\overline{\Omega}), \quad (4.4)$$

i.e., functions  $u$  must be twice continuously differentiable within  $\Omega$  and at least continuous up to the boundary  $\partial\Omega$ .

As we showed previously, the solution  $u$  can alternatively be found by using a variational approach, viz.

$$u = \arg \min \{ I[u] : u = \hat{u} \text{ on } \partial\Omega_D \} \quad (4.5)$$

whose stationarity condition is

$$\delta I[u] = \mathcal{G}(u, \delta u) = \mathcal{B}(u, \delta u) - \mathcal{L}(\delta u) = 0 \quad \text{for all } \delta u \in \mathcal{U}_0(\Omega). \quad (4.6)$$

Therefore, we can reformulate that problem (without, in principle, knowing anything about variational calculus) as:

$$\boxed{\text{find } u \in \mathcal{U} \quad \text{s.t.} \quad \mathcal{G}(u, v) = \mathcal{B}(u, v) - \mathcal{L}(v) = 0 \quad \text{for all } v \in \mathcal{U}_0(\Omega)} \quad (4.7)$$



This is called the **weak form** of the problem because we now seek solutions  $u \in \mathcal{U}$  where

$$\mathcal{U} = \{u \in \mathcal{H}^1(\Omega) : u = \hat{u} \text{ on } \partial\Omega_d\}, \quad (4.8)$$

that satisfy (4.7) for all  $v \in \mathcal{U}_0(\Omega)$ , and such a solution is called **weak solution**. There is one essential difference between the weak and strong form: solutions of the weak form are required to be in  $H^1(\Omega)$ , whereas the strong form required solutions to be in  $C^2(\Omega)$ . Thus, we have weakened/relaxed the conditions on the family of solutions, which is why the above is called the *weak form*.

Notice that, if  $v$  is interpreted as a virtual displacement field, then (4.7) is also referred to as the **principle of virtual work**.

Computationally, solving the weak form is usually preferable over the strong form. First,  $u \in H^1(\Omega)$  is simpler to satisfy than  $u \in C^2(\Omega)$  (e.g., piecewise linear interpolation is sufficient in the weak form but not in the strong form). Second, as we showed already for the *Rayleigh-Ritz* approach, the weak form boils down to solving a system of algebraic equations (rather than solving PDEs).

Let us show that we can also arrive at the weak form in an alternative fashion without the use of variational calculus. This is particularly helpful if no potential exists.

Let us take the first equation in (4.1), multiply it by some random trial function  $v \in \mathcal{U}_0(\Omega)$  that vanishes on  $\partial\Omega_D$ , and integrate over the entire domain. The result, which must still vanish due to (4.1), is

$$0 = - \int_{\Omega} [(\kappa u_{,i})_{,i} + s] v \, dV, \quad (4.9)$$

which must hold for all admissible  $v \in \mathcal{U}_0(\Omega)$ . This is the basis for the family of the so-called **methods of weighted residuals** (where one picks specific choices of  $v$ ).

Using the divergence theorem and the fact that  $v = 0$  on  $\partial\Omega_D$  reduces the above to

$$\begin{aligned} 0 &= \int_{\Omega} \kappa u_{,i} v_{,i} \, dV - \int_{\Omega} s v \, dV - \int_{\partial\Omega_N} \kappa u_{,i} n_i v \, dS \quad \text{for all } v \in \mathcal{U}_0(\Omega) \\ &= \int_{\Omega} \kappa u_{,i} v_{,i} \, dV - \int_{\Omega} s v \, dV - \int_{\partial\Omega_N} \hat{q} v \, dS \quad \text{for all } v \in \mathcal{U}_0(\Omega), \end{aligned} \quad (4.10)$$

where we used the Neumann boundary condition  $\kappa u_{,i} n_i = \hat{q}$  to transform the last term. The last equation in (4.10) is exactly identical to (4.7). In other words, we can find the weak form without the use of variational calculus, moreover, even without the existence of an energy functional by starting directly from the strong form. This is an important observation (even those problems that do not have a variational structure can thus be written in terms of a weak form). To decide whether or not a variational structure exists for a given problem, can be done by the use of Vainberg's theorem given in Appendix F.

## 4.2 Equivalence of strong and weak forms

We now have *two equivalent variational principles*:

Given a space

$$\mathcal{U} = \{u \in \mathcal{H}^k(\Omega) : u = \hat{u} \text{ on } \partial\Omega_d\}, \quad (4.11)$$

a functional  $I : \mathcal{U} \rightarrow \mathbb{R}$  and associated bilinear, continuous form  $\mathcal{B}(\cdot, \cdot)$  defined on  $\mathcal{U} \times \mathcal{U}$  and a continuous linear form  $\mathcal{L}(\cdot)$  defined on  $\mathcal{U}$ , we seek to

$$(A) \quad \boxed{\text{find } u \in \mathcal{U} \quad \text{s.t.} \quad u = \arg \min I[u]} \quad (4.12)$$

$$(B) \quad \boxed{\text{find } u \in \mathcal{U} \quad \text{s.t.} \quad \mathcal{B}(u, v) = \mathcal{L}(v) \quad \text{for all } v \in \mathcal{U}_0} \quad (4.13)$$

And we know that the two have a unique connection since  $\delta I = \mathcal{B}(u, \delta u) - \mathcal{L}(\delta u)$ . Thus, we know that

$$(A) \Leftrightarrow (B) \quad (4.14)$$

with a unique solution for this particular type of problem (if it is stable, e.g.,  $\kappa > 0$ ).

## 4.3 Approximate solutions

The idea of numerical approaches is to find an approximate solution: we replace the space  $\mathcal{U}$  by a *finite-dimensional* subspace

$$\mathcal{U}^h \subset \mathcal{U}, \quad (4.15)$$

in which we seek a solution  $u^h$ , where  $h$  stands for the discretization size.

An ***n-dimensional space***  $\mathcal{U}^h$  is defined by a set of  $n$  basis functions  $\{N_1, \dots, N_n\}$  and the approximation

$$\boxed{u^h(x) = \sum_{a=1}^n u^a N^a(x) \quad \text{and} \quad v^h(x) = \sum_{a=1}^n v^a N^a(x).} \quad (4.16)$$

Assume that the approximation space is chosen wisely, so the exact solution can be attained with infinite refinement; i.e., we assume that

$$\text{for all } u \in \mathcal{U} \quad \text{there exists } u_h(v) \in \mathcal{U}^h \quad \text{such that} \quad \lim_{h \rightarrow 0} \|u_h(v) - u\| = 0. \quad (4.17)$$

Then we can formulate the **discrete problem**

$$(C) \quad \boxed{\text{find } u^h \in \mathcal{U}^h \quad \text{s.t.} \quad \mathcal{B}(u^h, v^h) = \mathcal{L}(v^h) \quad \text{for all } v^h \in \mathcal{U}_0^h} \quad (4.18)$$

Next, let us insert the approximations (4.16) into (4.18) to obtain:

$$\mathcal{B}\left(\sum_{a=1}^n u^a N^a, \sum_{b=1}^n v^b N^b\right) = \mathcal{L}\left(\sum_{b=1}^n v^b N^b\right) \quad \text{for all } v^b \quad (4.19)$$

or, exploiting that  $\mathcal{B}$  is bilinear and  $\mathcal{L}$  is linear,

$$\sum_{b=1}^n v^b \left[ \sum_{a=1}^n u^a \mathcal{B}(N^a, N^b) - \mathcal{L}(N^b) \right] = 0 \quad \text{for all } v^b. \quad (4.20)$$

Since this must hold for all (admissible)  $v^b$ , we conclude that

$$\boxed{\sum_{a=1}^n u^a \mathcal{B}(N^a, N^b) = \mathcal{L}(N^b)} \quad \text{for } b = 1, \dots, n. \quad (4.21)$$

This is a *linear* system to be solved for  $u^a$  ( $a = 1, \dots, n$ ).

Let us define a vector of all unknown coefficients:

$$\mathbf{U}^h = \{u^1, \dots, u^n\}^T. \quad (4.22)$$

Further, we define a (symmetric) matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$  and vector  $\mathbf{F} \in \mathbb{R}^n$  with components

$$K_{ab} = \mathcal{B}(N^a, N^b), \quad F_b = \mathcal{L}(N^b). \quad (4.23)$$

Then, the linear system reads

$$\mathbf{K} \mathbf{U}^h = \mathbf{F} \quad \Leftrightarrow \quad K_{ba} U_a^h = F_b. \quad (4.24)$$

When we are using the same approximation space for  $u^h$  and  $v^h$ , this is the so-called the **Bubnov-Galerkin approximation**. Alternatively, one can choose different function spaces for the approximations  $u^h$  and  $v^h$ , which leads to the so-called **Petrov-Galerkin** method. The latter gains importance when solving *over/underconstrained problems* since it allows to control the number of equations by the choice of the dimension of the space of  $v^h$ .

## 5 The mechanical variational problem

### 5.1 Linearized kinematics

After all those precursors, let us analyze the mechanical variational problem and start with the simplest problem: *quasistatics in linearized kinematics*. Here, the **strong form** is

$$\begin{cases} \sigma_{ij,j} + \rho b_i = 0 & \text{in } \Omega, \\ u_i = \hat{u}_i & \text{on } \partial\Omega_D, \\ \sigma_{ij}n_j = \hat{t} & \text{on } \partial\Omega_N. \end{cases} \quad (5.1)$$

The associated total potential energy is

$$I[\mathbf{u}] = \int_{\Omega} W(\boldsymbol{\varepsilon}) \, dV - \int_{\Omega} \rho \mathbf{b} \cdot \mathbf{u} \, dV - \int_{\partial\Omega_N} \hat{\mathbf{t}} \cdot \mathbf{u} \, dS \quad (5.2)$$

and we seek displacement field solutions

$$\mathbf{u} = \arg \min \{ I[\mathbf{u}] : \mathbf{u} = \hat{\mathbf{u}} \text{ on } \partial\Omega_D \}. \quad (5.3)$$

We compute the first variation, defining  $\text{sym}(\cdot) = \frac{1}{2}(\cdot + \cdot^T)$ ,

$$\begin{aligned} \delta I[\mathbf{u}] &= \int_{\Omega} \frac{\partial W}{\partial \varepsilon_{ij}} \text{sym}(\delta u_{i,j}) \, dV - \int_{\Omega} \rho b_i \delta u_i \, dV - \int_{\partial\Omega_N} \hat{t}_i \delta u_i \, dS \\ &= \int_{\Omega} \sigma_{ij} \delta u_{i,j} \, dV - \int_{\Omega} \rho b_i \delta u_i \, dV - \int_{\partial\Omega_N} \hat{t}_i \delta u_i \, dS = 0 \quad \forall \quad \delta \mathbf{u} \in \mathcal{U}_0, \end{aligned} \quad (5.4)$$

where we used  $\sigma_{ij} = \partial W / \partial \varepsilon_{ij}$  and  $\sigma_{ij} = \sigma_{ji}$  (by angular momentum balance). Note that application of the divergence theorem shows the equivalence of the two forms since

$$\delta I[\mathbf{u}] = 0 = \int_{\partial\Omega_N} (\sigma_{ij}n_j - \hat{t}_i) \delta u_i \, dS - \int_{\Omega} (\sigma_{ij,j} + \rho b_i) \delta u_i \, dV \quad \forall \quad \delta \mathbf{u} \in \mathcal{U}_0. \quad (5.5)$$

We can use the first variation to define the weak form as

$$G(\mathbf{u}, \mathbf{v}) = \mathcal{A}(\mathbf{u}, \mathbf{v}) - \mathcal{L}(\mathbf{v}) = 0 \quad \text{for all adm. } \mathbf{v} \quad (5.6)$$

with

$$\mathcal{A}(\mathbf{u}, \mathbf{v}) = \int_{\Omega} \sigma_{ij} (\text{sym}(\nabla \mathbf{u})) v_{i,j} \, dV \quad \text{and} \quad \mathcal{L}(\mathbf{v}) = \int_{\Omega} \rho b_i v_i \, dV + \int_{\partial\Omega_N} \hat{t}_i v_i \, dS. \quad (5.7)$$

Notice that  $\mathcal{A}(\cdot, \cdot)$  is generally *not* a bilinear operator, while  $\mathcal{L}(\cdot)$  is a linear operator.

Next, we introduce the discrete weak form  $\mathcal{A}(\mathbf{u}^h, \mathbf{v}^h) - \mathcal{L}(\mathbf{v}^h) = 0$  with the Bubnov-Galerkin approximation

$$\mathbf{u}^h(\mathbf{x}) = \sum_{a=1}^n \mathbf{u}^a N^a(\mathbf{x}) \quad \text{and} \quad \mathbf{v}^h(\mathbf{x}) = \sum_{a=1}^n \mathbf{v}^a N^a(\mathbf{x}), \quad (5.8)$$

so that we arrive at (in component form)

$$\sum_{a=1}^n v_i^a \left[ \int_{\Omega} \sigma_{ij} (\text{sym}(\nabla \mathbf{u}^h)) N_{,j}^a \, dV - \int_{\Omega} \rho b_i N^a \, dV - \int_{\partial\Omega_N} \hat{t}_i N^a \, dS \right] = 0 \quad \text{for all adm. } \mathbf{v}^a \quad (5.9)$$

or

$$\boxed{\mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}} = \mathbf{0}} \quad \text{with} \quad \mathbf{U}^h = \{\mathbf{u}^1, \dots, \mathbf{u}^n\}^T \quad (5.10)$$

and

$$\boxed{F_{\text{int},i}^a = \int_{\Omega} \sigma_{ij}(\nabla \mathbf{u}^h) N_{,j}^a \, dV \quad \text{and} \quad F_{\text{ext},i}^a = \int_{\Omega} \rho b_i N^a \, dV + \int_{\partial\Omega_N} \hat{t}_i N^a \, dS} \quad (5.11)$$

For the special case of **linear elasticity** we have  $\sigma_{ij} = \mathbb{C}_{ijkl} u_{k,l}$  so that the *weak form* reads

$$G(\mathbf{u}, \mathbf{v}) = \mathcal{B}(\mathbf{u}, \mathbf{v}) - \mathcal{L}(\mathbf{v}) = 0 \quad \text{for all adm. } \mathbf{v} \quad (5.12)$$

with

$$\mathcal{B}(\mathbf{u}, \mathbf{v}) = \int_{\Omega} v_{i,j} \mathbb{C}_{ijkl} u_{k,l} \, dV \quad \text{and} \quad \mathcal{L}(\mathbf{v}) = \int_{\Omega} \rho b_i v_i \, dV + \int_{\partial\Omega_N} \hat{t}_i v_i \, dS, \quad (5.13)$$

so  $\mathcal{B}(\cdot, \cdot)$  is indeed a bilinear form. Inserting the approximate fields, (5.11) becomes

$$\begin{aligned} F_{\text{int},i}^a &= \int_{\Omega} \mathbb{C}_{ijkl} u_{k,l}^h N_{,j}^a \, dV = \sum_{b=1}^n \int_{\Omega} \mathbb{C}_{ijkl} u_k^b N_{,l}^b N_{,j}^a \, dV = \sum_{b=1}^n u_k^b \int_{\Omega} \mathbb{C}_{ijkl} N_{,j}^a N_{,l}^b \, dV \\ &= \sum_{b=1}^n K_{ik}^{ab} u_k^b \quad \text{with} \quad K_{ik}^{ab} = \int_{\Omega} \mathbb{C}_{ijkl} N_{,j}^a N_{,l}^b \, dV \\ &\Rightarrow \quad \mathbf{F}_{\text{int}} = \mathbf{K} \mathbf{U}^h \quad \Rightarrow \quad \mathbf{U}^h = \mathbf{K}^{-1} \mathbf{F}_{\text{ext}} \quad \text{if } \det \mathbf{K} \neq 0. \end{aligned} \quad (5.14)$$

That is, we arrive at a linear problem to be solved for the unknown coefficients  $\mathbf{U}^h = \{\mathbf{u}^1, \dots, \mathbf{u}^n\}$ .

For computational purposes, notice that vectors  $\mathbf{U}^h$  and (internal or external)  $\mathbf{F}$ , e.g., in 3D are, respectively

$$\mathbf{U}^h = \begin{pmatrix} \mathbf{u}^1 \\ \vdots \\ \mathbf{u}^n \end{pmatrix} = \begin{pmatrix} u_1^1 \\ u_2^1 \\ u_3^1 \\ \vdots \\ u_1^n \\ u_2^n \\ u_3^n \end{pmatrix}, \quad \mathbf{F} = \begin{pmatrix} \mathbf{F}^1 \\ \vdots \\ \mathbf{F}^n \end{pmatrix} = \begin{pmatrix} F_1^1 \\ F_2^1 \\ F_3^1 \\ \vdots \\ F_1^n \\ F_2^n \\ F_3^n \end{pmatrix}. \quad (5.15)$$

If we use 0-index notation like in C++ (i.e., we sum over  $a = 0, \dots, n-1$  instead of  $a = 1, \dots, n$ ), then

$$u_i^a \text{ is the } (d \cdot a + i)\text{th component of vector } \mathbf{U}^h \text{ in } d \text{ dimensions.} \quad (5.16)$$

Similarly, we apply the same rule to the rows and columns of matrix  $\mathbf{K}$ , so that

$$K_{ik}^{ab} \text{ is the component at } (d \cdot a + i, d \cdot b + k) \text{ of matrix } \mathbf{K} \text{ in } d \text{ dimensions.} \quad (5.17)$$

There is a shortcut to computing the internal force vector via the **Rayleigh-Ritz** method. As introduced in Section 3.3, this technique inserts  $\mathbf{u}^h$  directly into the total potential energy and minimizes the latter with respect to the unknown coefficients, i.e., we must solve

$$\frac{\partial I[\mathbf{u}^h]}{\partial u^a} = \mathbf{0} \quad \forall \quad a = 1, \dots, n. \quad (5.18)$$

Note that

$$\begin{aligned} \frac{\partial I[\mathbf{u}^h]}{\partial u_i^a} = 0 &= \frac{\partial}{\partial u_i^a} \left[ \int_{\Omega} W(\boldsymbol{\varepsilon}^h) \, dV - \int_{\Omega} \rho \mathbf{b} \cdot \mathbf{u}^h \, dV - \int_{\partial\Omega_N} \hat{\mathbf{t}} \cdot \mathbf{u}^h \, dS \right] \\ &= \int_{\Omega} \frac{\partial W}{\partial \varepsilon_{kl}}(\boldsymbol{\varepsilon}^h) \frac{\partial \varepsilon_{kl}^h}{\partial u_i^a} \, dV - \int_{\Omega} \rho b_k \frac{\partial}{\partial u_i^a} \sum_{b=1}^n u_k^b N^b \, dV - \int_{\partial\Omega_N} \hat{t}_k \frac{\partial}{\partial u_i^a} \sum_{b=1}^n u_k^b N^b \, dS. \end{aligned} \quad (5.19)$$

where

$$\varepsilon_{kl}^h = \frac{1}{2}(u_{k,l}^h + u_{l,k}^h) = \sum_{b=1}^n \frac{1}{2}(u_k^b N_{,l}^b + u_l^b N_{,k}^b) \quad \Rightarrow \quad \frac{\partial \varepsilon_{kl}^h}{\partial u_i^a} = \frac{1}{2}(\delta_{ik} N_{,l}^a + N_{,k}^a \delta_{li}). \quad (5.20)$$

This is equivalent to

$$\begin{aligned} 0 &= \int_{\Omega} \sigma_{kl}(\boldsymbol{\varepsilon}^h) \frac{1}{2}(N_{,l}^a \delta_{ik} + N_{,k}^a \delta_{li}) \, dV - \int_{\Omega} \rho b_i N^a \, dV - \int_{\partial\Omega_N} \hat{t}_i N^a \, dS \\ &= \int_{\Omega} \sigma_{il}(\boldsymbol{\varepsilon}^h) N_{,l}^a \, dV - \int_{\Omega} \rho b_i N^a \, dV - \int_{\partial\Omega_N} \hat{t}_i N^a \, dS. \end{aligned} \quad (5.21)$$

By comparison, we see immediately that this yields  $F_{\text{int},i}^a$  and  $F_{\text{ext},i}^a$  directly. Thus, rather than resorting to variations, we can obtain the internal and external force vectors alternatively (and oftentimes much more simply) by the Rayleigh-Ritz approach, which inserts the approximation into the potential energy and then differentiates with respect to the unknown coefficients.

Also, notice that  $\mathbf{F}_{\text{ext}}$  is independent of the constitutive law and only depends on the applied body forces and surface tractions. That is,  $\mathbf{F}_{\text{ext}}$  is sufficiently general for arbitrary materials (in linearized kinematics), while the computation of  $\mathbf{F}_{\text{int}}$  depends on the particular material model.

## 5.2 Finite kinematics

The variational problem in finite-deformation quasistatics is quite similar:

$$I[\boldsymbol{\varphi}] = \int_{\Omega} W(\mathbf{F}) \, dV - \int_{\Omega} \rho_0 \mathbf{B} \cdot \boldsymbol{\varphi} \, dV - \int_{\partial\Omega_N} \hat{\mathbf{T}} \cdot \boldsymbol{\varphi} \, dS \quad (5.22)$$

and we seek solutions

$$\boldsymbol{\varphi} \in \mathcal{U} = \{ \boldsymbol{\varphi} \in H^1(\Omega) : \boldsymbol{\varphi} = \hat{\boldsymbol{\varphi}} \text{ on } \partial\Omega_D \} \quad \text{such that} \quad \boldsymbol{\varphi} = \arg \min I[\boldsymbol{\varphi}]. \quad (5.23)$$

In all our problems, we will assume that the undeformed and deformed coordinate systems coincide so that we write for convenience  $\boldsymbol{\varphi} = \mathbf{x} = \mathbf{X} + \mathbf{u}$  and we thus formulate the above problem in terms of the displacement field  $\mathbf{u} \in \mathcal{U}$ , like in the linear elastic case (note that this is a “notational crime” that we adopt here for convenience).

We may then write  $\mathbf{F} = \mathbf{I} + \text{Grad } \mathbf{u}$  and compute the first variation as

$$\begin{aligned} \delta I[\mathbf{u}] &= \int_{\Omega} \frac{\partial W}{\partial F_{iJ}} \delta u_{i,J} \, dV - \int_{\Omega} \rho_0 B_i \delta u_i \, dV - \int_{\partial\Omega_N} \hat{T}_i \delta u_i \, dS \\ &= \int_{\Omega} P_{iJ} \delta u_{i,J} \, dV - \int_{\Omega} \rho_0 B_i \delta u_i \, dV - \int_{\partial\Omega_N} \hat{T}_i \delta u_i \, dS = 0, \end{aligned} \quad (5.24)$$

where we used the first Piola-Kirchhoff stress tensor  $P_{iJ} = \partial W / \partial F_{iJ}$  (which is *not* symmetric).

Even though the form looks similar to (5.4), recall that  $\mathbf{P}(\mathbf{F})$  involves in a generally nonlinear relation between  $\mathbf{P}$  and the displacement gradient  $\text{Grad } \mathbf{u}$ . Therefore, the finite-deformation variational problem does *not* involve a *bilinear* form (even if the material is elastic).

As before, we produce a discrete approximation, e.g., with the **Bubnov-Galerkin** approximation

$$\mathbf{u}^h(\mathbf{X}) = \sum_{a=1}^n \mathbf{u}^a N^a(\mathbf{X}) \quad \text{and} \quad \mathbf{v}^h(\mathbf{X}) = \sum_{a=1}^n \mathbf{v}^a N^a(\mathbf{X}), \quad (5.25)$$

so that we again arrive at

$$\boxed{\mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}} = \mathbf{0}} \quad (5.26)$$

where now, by comparison,

$$\boxed{F_{\text{int},i}^a = \int_{\Omega} P_{iJ}(\nabla \mathbf{u}) N_{,J}^a \, dV \quad \text{and} \quad F_{\text{ext},i}^a = \int_{\Omega} \rho_0 B_i N^a \, dV + \int_{\partial\Omega_N} \hat{T}_i N^a \, dS} \quad (5.27)$$

In a nutshell, the linearized and finite elastic variational problems result in the same system of equations (5.26). For the special case of linear elasticity, that system is linear. Otherwise the problem is nonlinear and requires an iterative solution method.

Note that in both formulations of linearized and finite kinematics we assumed that  $\hat{\mathbf{T}} = \text{const.}$ , i.e., that the *externally applied forces are constant* and do not depend on deformation. Especially in finite deformations this is oftentimes not the case; e.g., consider pressure loading  $\hat{T} = p\mathbf{n}$  where  $\mathbf{n}$  is the *deformed* surface normal and  $\mathbf{n} \, ds = \mathbf{J} \mathbf{F}^{-T} \mathbf{N} \, dS$  by the Piola transform. In such cases, the above variational form does not apply and one needs to revise the external force terms appropriately. For example, for pressure loading we know the work done by pressure is  $p v$  with the deformed volume  $v$ , so we may use  $I_{\text{ext}} = \int_{\Omega} p J \, dV$ , which must replace the traction term in the above total potential energy (and the derivatives follow analogously).

### 5.3 Thermal problem revisited

For completeness, let us revisit the quasistatic thermal problem in an analogous fashion: we seek temperature values  $\mathbf{T}^h = \{T_1, \dots, T_n\}$  such that

$$\boxed{Q_{\text{int}}(\mathbf{T}^h) - Q_{\text{ext}} = 0} \quad (5.28)$$

with

$$\boxed{Q_{\text{int}}^a = \int_{\Omega} q_i(\nabla T) N_{,i}^a \, dV \quad \text{and} \quad Q_{\text{ext}}^a = \int_{\Omega} \rho s N^a \, dV + \int_{\partial\Omega_N} \hat{q} N^a \, dS} \quad (5.29)$$

where  $q_i = \partial W / \partial T_{,i}$ .

For *linear heat conduction*,  $\mathbf{q} = \kappa \, \text{grad } T$ , we obtain a linear system of equations since then

$$Q_{\text{int}}^a = \sum_{b=1}^n T^b \int_{\Omega} \kappa N_{,i}^a N_{,i}^b \, dV = \sum_{b=1}^n K^{ab} T^b \quad \text{with} \quad K^{ab} = \int_{\Omega} \kappa N_{,i}^a N_{,i}^b \, dV. \quad (5.30)$$

(Notice that for convenience we defined the flux vector  $\mathbf{q}$  without the negative sign, which results in the analogous forms as in linear elasticity.)

## 5.4 A simple example: nonlinear springs

Consider an axial spring that undergoes large deformation, i.e., its two end points move from  $(\mathbf{X}^0, \mathbf{X}^1)$  to  $(\mathbf{x}^0, \mathbf{x}^1)$  and  $\mathbf{x}^i = \mathbf{X}^i + \mathbf{u}^i$ . (Note that we use 0 and 1 instead of 1 and 2 to comply with C++ standard indexing which starts with 0.)

The bar stores strain energy upon stretching with an energy density  $W = W(\varepsilon)$  where  $\varepsilon$  is the axial bar strain:

$$\varepsilon = \frac{l - L}{L}, \quad \text{where} \quad \mathbf{l} = \mathbf{x}^1 - \mathbf{x}^0, \quad \mathbf{L} = \mathbf{X}^1 - \mathbf{X}^0, \quad \text{and} \quad l = |\mathbf{l}|, \quad L = |\mathbf{L}|. \quad (5.31)$$

The total energy of a bar with cross-sectional area  $A$  and initial length  $L$  is therefore

$$I = A L W(\varepsilon). \quad (5.32)$$

Without even defining interpolation functions, we can use the Rayleigh-Ritz shortcut to calculate the resulting *internal force* on node 0 as

$$\begin{aligned} \mathbf{F}_{\text{int}}^0 &= \frac{\partial I}{\partial \mathbf{u}^0} = A L \frac{\partial W}{\partial \varepsilon} \frac{\partial}{\partial \mathbf{u}^0} \frac{l - L}{L} \\ &= A \sigma(\varepsilon) \frac{\partial}{\partial \mathbf{u}^0} \sqrt{(\mathbf{X}^1 + \mathbf{u}^1 - \mathbf{X}^0 - \mathbf{u}^0) \cdot (\mathbf{X}^1 + \mathbf{u}^1 - \mathbf{X}^0 - \mathbf{u}^0)} \\ &= -A \sigma(\varepsilon) \frac{\mathbf{l}}{l}, \end{aligned} \quad (5.33)$$

where  $\sigma(\varepsilon) = \partial W / \partial \varepsilon$  is the axial stress in the bar. Analogously, the force on node 1 becomes

$$\mathbf{F}_{\text{int}}^1 = -\mathbf{F}_{\text{int}}^0 = A \sigma(\varepsilon) \frac{\mathbf{l}}{l}. \quad (5.34)$$

As expected, the force points along the (deformed) axis of the spring end points, and the forces on the two end points are of the same magnitude but of opposite sign.

Note that we did not specify whether or not the spring is linear elastic; i.e., the specific choice of  $W(\varepsilon)$  will determine the behavior of the spring and can be chosen to be quadratic, i.e.  $W(\varepsilon) = \frac{k}{2} \varepsilon^2$ , which results in a linear spring, but can also be more complex.

As a final remark, the assumption of a constant strain  $\varepsilon$  along the spring length tacitly implies that we assume a linear displacement profile along the spring. That is, the above formulation is equivalent to assuming linear shape functions  $N^0$  and  $N^1$  so that

$$\mathbf{u}^h(\mathbf{X}) = \mathbf{u}^0 N^0(\mathbf{X}) + \mathbf{u}^1 N^1(\mathbf{X}). \quad (5.35)$$

In our implementation, we will define two classes – one *material model* that defines the material point relations  $W = W(\varepsilon)$  and  $\sigma = \sigma(\varepsilon)$  as well as an nonlinear bar/spring *element* that defines  $I$  and  $\{\mathbf{F}_{\text{int}}^0, \mathbf{F}_{\text{int}}^1\}$ . Notice that the element will require the material model to compute its required quantities.



## 6 Interpolation spaces

So far, we have assumed approximations of the type

$$u^h(x) = \sum_{a=1}^n u^a N^a(x), \quad (6.1)$$

but we have not chosen particular spaces  $\mathcal{U}^h$  for the interpolation or **shape functions**  $N^a(x)$ .

In general, there are two possible choices:

- **global** shape functions that are defined everywhere in  $\Omega$ , i.e.,  $|\text{supp } N^a| \sim |\Omega|$ ,  
e.g., polynomials  $N^a(x) = x^{a-1}$  or trigonometric polynomials  $N^a(x) = \cos(\pi(a-1)x)$ .
- **local** shape functions that are defined only locally:  $|\text{supp } N^a| \ll |\Omega|$ ,  
e.g., piecewise linear shape functions,

where we introduced the **support** of a continuous function  $u$  defined on  $\Omega \in \mathbb{R}^d$  as the (closure in  $\Omega$  of the) set of all points where  $u(x) \neq 0$ , i.e.,

$$\text{supp } u = \{x \in \Omega : u(x) \neq 0\} \quad (6.2)$$

This means that  $u(x) = 0$  for  $x \in \Omega \setminus \text{supp } u$ .

For any set of shape functions, the following **shape function properties** must be satisfied:

- (1) for any  $x \in \Omega$  there is at least one  $a$  with  $1 \leq a \leq n$  such that  $N^a(x) \neq 0$  (i.e., the *whole domain must be covered*; otherwise, there is no approximation at all at certain points)
- (2) all  $N^a$  should allow to satisfy the *Dirichlet boundary conditions* if required.
- (3) *linear independence* of the shape functions:

$$\sum_{a=1}^n u^a N^a = 0 \quad \Leftrightarrow \quad u^a = 0 \text{ for all } a = 1, \dots, n. \quad (6.3)$$

In other words, given any function  $u^h \in \mathcal{U}^h$ , there exists a unique set of parameters  $\{u^1, \dots, u^n\}$  such that

$$u^h = \sum_{a=1}^n u^a N^a. \quad (6.4)$$

Then, functions  $\{N^1, \dots, N^n\}$  are a **basis** of  $\mathcal{U}^h$ . Linear independence is important since it avoids ill-posed problems.

For example, take  $\mathcal{U}^h = \mathbb{P}_2$  and  $\{N^1, N^2, N^3\} = \{1, x, x^2\}$  so that  $u^h = u^1 N^1 + u^2 N^2 + u^3 N^3$ . Hence, if, e.g.,  $u^h = a + bx + x^2$  then we immediately conclude that  $u^1 = a$ ,  $u^2 = b$ ,  $u^3 = c$  uniquely.

Otherwise, i.e., if there existed a set  $\{\alpha^1, \dots, \alpha^n\} \neq 0$  such that  $\sum_{a=1}^n \alpha^a N^a = 0$ , then this set of parameters could be added on top of any solution  $\{u^1, \dots, u^n\}$  such that

$$u^h = \sum_{a=1}^n u^a N^a = \sum_{a=1}^n (u^a + \alpha^a) N^a, \quad (6.5)$$

which means both  $\{u^a\}$  and  $\{u^a + \alpha^a\}$  are solutions (hence the problem is not well-posed).

- (4) The shape functions  $N^a$  must satisfy the *differentiability/integrability requirements* of the weak form (this depends on the problem to be solved and will be discussed later).
- (5) The shape functions must possess “*sufficient approximation power*”. In other words, consider  $u^h \in \mathcal{U}^h \subset \mathcal{U}$ : we should ensure that  $u^h = \sum_{a=1}^n u^a N^a \rightarrow u$  as  $n \rightarrow \infty$ .

Condition (5) is a crucial one. It tells us that for an approximation to converge, we must pick an approximate function space that gives the solution “a chance to converge”. For example, assume you aim to approximate a high-order polynomial  $u \in \mathbb{P}_n$  (with  $n \gg 1$ ) by an approximation  $u^h$  using shape functions  $\{1, x, x^2, x^3, \dots, x^n\}$ . This is expected to converge as  $n \rightarrow \infty$ , because the coefficients of  $u$  will approach the coefficients of  $u^h$ . But choosing shape functions  $\{1, x, x^3, \dots, x^n\}$  (notice the  $x^2$ -term is omitted) will never converge as  $n \rightarrow \infty$ . Polynomials do satisfy this requirement by the following theorem.

**Weierstrass approximation theorem:** Given a continuous function  $f : [a, b] \subset \mathbb{R} \rightarrow \mathbb{R}$  and any scalar  $\epsilon > 0$ , then there exists a polynomial

$$p_n(x) \in \mathbb{P}_\infty \quad \text{such that} \quad |f(x) - p_n(x)| < \epsilon \quad \text{for all } x \in [a, b]. \quad (6.6)$$

This means every continuous function  $u$  can be approximated by a polynomial function to within any level of accuracy.

Therefore,  $\{N^i\} = \{1, x, x^2, x^3, \dots\}$ , i.e., the polynomials in  $\mathbb{R}$ , is a suitable choice for the shape functions that satisfy the **completeness property** (and we have shown their linear independence).

Note that, as discussed above, one *cannot* omit any intermediate-order terms from the set

$$\{1, x, x^2, x^3, \dots\}. \quad (6.7)$$

If one omits a term, e.g., take  $\{1, x^2, x^3, \dots\}$ , then if  $u^h \in \mathcal{U}^h = \mathbb{P}_n$  then there is no set  $\{u^1, \dots, u^n\}$  such that  $u^h = \sum_{i=1}^n u^i N^i$ .

As an extension, the Weierstrass approximation theorem also applies to **trigonometric polynomials** (cf. Fourier series).

### completeness in higher dimensions:

A polynomial approximation in  $\mathbb{R}^d$  is **complete up to order  $q$** , if it contains independently all monomials  $\mathbf{x}^\alpha$  with  $|\alpha| = \alpha_1 + \dots + \alpha_d \leq q$ , i.e., using multi-indices we write

$$u^h = \sum_{\beta=0}^q \sum_{|\alpha|=\beta} c_\alpha \mathbf{x}^\alpha. \quad (6.8)$$

What does this mean in practice?

1D:  $\{1, x, x^2, x^3, \dots, x^q\}$  so that a polynomial of order  $q$  contains  $q + 1$  monomials

2D:  $q = 0$ :  $\{1\}$

$q = 1$ :  $\{1, x_1, x_2\}$

$q = 2$ :  $\{1, x_1, x_2, x_1^2, x_1 x_2, x_2^2\}$

$q = 3$ :  $\{1, x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3\}$

...

The number of independent monomials in 2D is hence  $(q + 1)(q + 2)/2$ .

## 7 The Finite Element Method

**Motivation:** we would like to define shape functions that are *local* and admit a simple way to enforce *Dirichlet BCs*.

**Idea:** we introduce a discretization  $\mathcal{T}_h$  that splits  $\Omega$  into subdomains  $\Omega_e$ , the so-called **elements**, such that

$$\Omega_e \subset \Omega, \quad \Omega = \bigcup_e \Omega_e, \quad \partial\Omega \subseteq \bigcup_e \partial\Omega_e. \quad (7.1)$$

$\mathcal{T}_h$  is defined by the collection of **nodes** and **elements** and is called a **mesh**.

Mathematically (and computationally), a **finite element** is an object that has

- (i) a FE subdomain  $\Omega_e$ .
- (ii) a (linear) space of shape functions  $N^i$  (restricted to  $\Omega_e$ , i.e.  $\text{supp } N^i = \overline{\Omega_e}$ ).
- (iii) a set of **degrees of freedom** (dofs), viz. the  $u^a$  associated with those  $N^i$ .

The **Finite Element Method** (FEM) defines continuous, piecewise-polynomial shape functions such that

$$N^i(x_j) = \delta_{ij} \quad \text{for all } i, j \in \{1, \dots, n\} \quad (7.2)$$

This is the defining relation that determines the shape functions. Notice that if we evaluate the approximation  $u^h(x)$  at one of the nodes  $x_j$ , then

$$u^h(x_j) = \sum_{a=1}^n u^a N^a(x_j) = \sum_{a=1}^n u^a \delta_{aj} = u^j. \quad (7.3)$$

That is, the coefficient  $u^j$  can now be identified as the value of approximate function  $u^h$  at node  $j$ . This makes for a very beneficial physical interpretation of the (yet to be determined) shape function coefficients.

Let us check the *requirements for shape functions*:

- (1) is automatically satisfied: if  $x \in \Omega$  then  $x \in \Omega_e$  for some  $e$ , then there are  $N^i(x) \neq 0$
- (2) can be satisfied by fixing degrees of freedom of the boundary nodes (*errors possible*)
- (3) Assume, by contradiction, that  $u^h(x) = 0$  for all  $x \in \Omega$  while some  $u^a \neq 0$ . Now, evaluate at a node  $x_j$ :

$$0 = u^h(x_j) = \sum_{a=1}^n u^a N^a(x_j) = u^j \quad \Rightarrow \quad u^j = 0, \quad (7.4)$$

which contradicts the assumption that some  $u_a^h \neq 0$ . Thus, we have *linear independence*.

- (4) *Integrability/differentiability* requirements depend on the variational problem to be solved and must be ensured. For example, for mechanics we have  $\mathcal{U}^h, \mathcal{V}^h \in H^1$ , i.e., first derivatives must be square-integrable. Note that this guarantees that displacements (0th derivatives) are continuous and thus *compatible* (no jumps in displacements).
- (5) *Completeness* requires  $u^h \rightarrow u$  (and thus  $\mathcal{U}^h \rightarrow \mathcal{U}$ ) to within desirable accuracy. In the FE method, one enriches  $\mathcal{U}^h$  by, e.g.,

- **h-refinement**: refining the discretization  $\mathcal{T}_h$  while keeping the polynomial order fixed.
- **p-refinement**: increasing the polynomial interpolation order within a fixed discretization  $\mathcal{T}_h$ .

- **hp-refinement**: combination of the two above.
- **r-refinement**: repositioning of nodes while keeping discretization/interpolation fixed.

A note on ensuring *sufficient approximation power*: consider the exact solution  $u(x)$  at a point  $x \in \Omega$  so

$$u(x+h) = u(x) + h u'(x) + \frac{1}{2} h^2 u''(x) + \dots + \frac{1}{q!} h^q u^{(q)}(x) + O(h^{q+1}) \quad (7.5)$$

Assume that  $\mathcal{U}^h$  contains all polynomials complete up to degree  $q$  (i.e.,  $u^h \in \mathbb{P}_q$ ), then there exists

$$u^h \in \mathcal{U}^h \quad \text{such that} \quad u(x) = u^h(x) + O(h^{q+1}). \quad (7.6)$$

Let  $p$  be the highest derivative in the weak form, then

$$\frac{d^p u}{dx^p} = \frac{d^p u^h}{dx^p} + O(h^{q+1-p}). \quad (7.7)$$

For the solution to converge as  $h \rightarrow 0$  we need  $q+1-p \geq 1$  so that we have at least order  $O(h)$ . Thus we must ensure that

$$\boxed{q \geq p} \quad (7.8)$$

## 8 Finite element spaces: polynomial shape functions in 1D

Let us start with the simplest of all choices: *continuous, piecewise-polynomial interpolation functions*. Note that we need  $q \geq 1$  since  $p = 1$  for the mechanical/thermal/electromagnetic variational problem; i.e., we need *at least linear interpolation* within elements.

### 8.1 One dimension

**Simplest example: 2-node bar element**

Interpolation with element dofs  $\{u_e^1, u_e^2\}$  so that

$$u_e^h(x) = N_e^1(x)u_e^1 + N_e^2(x)u_e^2 \quad (8.1)$$

and we must have  $u_e^h(0) = u_e^1$  and  $u_e^h(\Delta x) = u_e^2$ .

This gives the **element shape functions**:

$$N_e^1(x) = 1 - \frac{x}{\Delta x}, \quad N_e^2(x) = \frac{x}{\Delta x}. \quad (8.2)$$

Note that the interpolation space uses  $\{1, x\}$  which is complete up to  $q = 1$ .

### 8.2 Example: linear elastic bar in 1D

Recall that for a linear bar element we obtained

$$N_e^1(x) = 1 - \frac{x}{\Delta x}, \quad N_e^2(x) = \frac{x}{\Delta x} \quad \Rightarrow \quad N_{e,x}^1(x) = -\frac{1}{\Delta x}, \quad N_{e,x}^2(x) = \frac{1}{\Delta x}. \quad (8.3)$$

so that the (only non-zero) axial strain inside the bar element is *constant* since:

$$\varepsilon_{xx}^h = u_{,x}^h(x) = \sum_{i=1}^n u_e^i N_{e,x}^i(x) = u_e^1 N_{e,x}^1(x) + u_e^2 N_{e,x}^2(x) = \frac{u_e^2 - u_e^1}{\Delta x} = \text{const.} \quad (8.4)$$

If we assume a linear elastic bar, so  $\sigma = E\varepsilon$ , we obtain from (5.11)

$$F_{\text{int}}^a = \int_{\Omega} \sigma_{xx}(\varepsilon_{xx}^h) N_{,x}^a dV = \int_0^{\Delta x} E \frac{u_e^2 - u_e^1}{\Delta x} N_{,x}^a A dx = \frac{EA}{\Delta x} (u_e^2 - u_e^1) \int_0^{\Delta x} N_{,x}^a dx, \quad (8.5)$$

and inserting the constant shape function derivatives finally yields

$$F_{\text{int}}^1 = \frac{EA}{\Delta x} (u_e^2 - u_e^1) \left( -\frac{1}{\Delta x} \right) \Delta x = -\frac{EA}{\Delta x} (u_e^2 - u_e^1) \quad \text{and} \quad F_{\text{int}}^2 = -F_{\text{int}}^1. \quad (8.6)$$

The stiffness matrix follows as

$$\left[ \frac{\partial F_{\text{int}}^i}{\partial u_e^j} \right] = \frac{EA}{\Delta x} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}, \quad (8.7)$$

which has the typical form known from 1D assemblies of linear springs with stiffness  $k = EA/\Delta x$ .

## Lagrangian interpolation of higher order

Interpolation up to degree  $q$ , i.e.,  $\{1, x, x^2, \dots, x^q\}$  so that

$$u_e^h(x) = \sum_{a=1}^{q+1} N_e^a(x) u_e^a = a_0 + a_1 x + a_2 x^2 + \dots + a_q x^q. \quad (8.8)$$

In general, shape functions can be determined by solving the  $q + 1$  equations

$$u_e^h(x_i) = u^i \quad \text{for all} \quad i = 1, \dots, q + 1 \quad (8.9)$$

for the  $q + 1$  coefficients  $a_i$  ( $i = 0, \dots, q$ ). Then, rearranging the resulting polynomial allows to extract the shape functions  $N_e^a(x)$  by comparison of the coefficients of  $u_e^a$ .

Alternatively, we can solve

$$N_e^a(x_i) = \delta_{ai} \quad \text{for all nodes } i = 1, \dots, q + 1. \quad (8.10)$$

The solution is quite intuitive:

$$N_e^a(x) = \frac{(x - x_1) \cdot \dots \cdot (x - x_{a-1}) \cdot (x - x_{a+1}) \cdot \dots \cdot (x - x_{q+1})}{(x_a - x_1) \cdot \dots \cdot (x_a - x_{a-1}) \cdot (x_a - x_{a+1}) \cdot \dots \cdot (x_a - x_{q+1})} \quad (8.11)$$

One can readily verify that  $N_e^a(x_i) = \delta_{ai}$ . These are called **Lagrange polynomials**.

## Alternative: hierarchical interpolation

We can also construct higher-order interpolations based on lower-order shape functions. For example, start with a 2-node bar:

$$N_e^1(x) = 1 - \frac{x}{\Delta x}, \quad N_e^2(x) = \frac{x}{\Delta x}. \quad (8.12)$$

Let us enrich the interpolation to reach  $q = 2$ :

$$u^h(x) = N_e^1(x) u_e^1 + N_e^2(x) u_e^2 + \tilde{N}_e^3(x) \alpha_e \quad (8.13)$$

with

$$\tilde{N}_e^3(x) = a_0 + a_1 x + a_2 x^2. \quad (8.14)$$

We need to find the coefficients  $a_i$ . Note that we must have

$$\tilde{N}_e^3(0) = \tilde{N}_e^3(\Delta x) = 0 \quad \Rightarrow \quad \tilde{N}_e^3(x) = c \frac{x}{\Delta x} \left( 1 - \frac{x}{\Delta x} \right) \quad (8.15)$$

with some constant  $c \neq 0$ .

Note that  $\alpha_e$  does not have to be continuous across elements and can hence be determined *locally* (i.e., given  $u_e^1$  and  $u_e^2$ ,  $\alpha_e$  can be determined internally for each element, which allows for *condensation* of the  $\alpha_e$ -dof).

### Example of higher-order interpolation: 2-node beam element

Linear elastic Euler-Bernoulli beams are a most common structural element. From the variational form (or the strong form,  $EIw^{(4)}(x) = q(x)$  for statics, which is of 4th order in the deflection  $w(x)$ ) we know that  $p = 2$ . Therefore, we must have  $q \geq 2$ , and  $w \in H^2(\Omega)$ .

The simplest admissible interpolation is based on  $\{1, x, x^2, x^3\}$ , so

$$w_e^h(x) = c_0 + c_1x + c_2x^2 + c_3x^3. \quad (8.16)$$

We need *four dofs*, so we pick *two nodes* and assign to each node a deflection  $w$  and angle  $\theta = w'$ :

$$w_e^h(x) = \sum_{i=1}^2 [N_e^{i1}(x)w_e^i + N_e^{i2}(x)\theta_e^i] \quad (8.17)$$

and we must have

$$w_e^h(0) = w_e^1, \quad w_e^h(\Delta x) = w_e^2, \quad (w_e^h)'(0) = \theta_e^1, \quad (w_e^h)'(\Delta x) = \theta_e^2. \quad (8.18)$$

The resulting shape functions are known as **Hermitian polynomials** (and can be found in textbooks and the notes).

Note that this is only one possible choice; we could also define alternative nodes and nodal values. However, the above choice ensures that *both deflection and angle are continuous* across elements.

## 8.3 Higher dimensions

**Problem:** in higher dimensions it will be cumbersome to define polynomial shape functions on the actual shape of the element, unless one uses regular structured meshes (e.g. grids). In general, all elements have different shapes and it is beneficial to define shape functions independent of the specific element shape.

To this end, we introduce a (bijective) **isoparametric mapping**  $\phi$  from a *reference domain* (with reference coordinates  $\xi = \{\xi, \eta, \zeta\}$ ) onto the *physical domain* (with coordinates  $\mathbf{x} = \{x, y, z\}$ ) of an element  $e$ :

$$\mathbf{x} = \phi(\xi), \quad \text{i.e. in 3D:} \quad x = x(\xi, \eta, \zeta), \quad y = y(\xi, \eta, \zeta), \quad z = z(\xi, \eta, \zeta). \quad (8.19)$$

For simplicity, we reuse the interpolation concepts from before:

$$\begin{aligned} \text{so far:} \quad \mathbf{u}_e^h &= \sum_{i=1}^n N_e^i(\mathbf{x}) \mathbf{u}_e^i \\ \text{now:} \quad \mathbf{u}_e^h &= \sum_{i=1}^n N_e^i(\xi) \mathbf{u}_e^i \quad \text{and} \quad \mathbf{x} = \sum_{i=1}^m \tilde{N}_e^i(\xi) \mathbf{x}_e^i, \end{aligned} \quad (8.20)$$

where we have three options for the mapping:

- (i) **isoparametric:**  $n = m$  and  $N_e^i = \tilde{N}_e^i$  (same interpolation)
- (ii) **subparametric:**  $n > m$  (lower-order interpolation of positions)
- (iii) **superparametric:**  $n < m$  (higher-order interpolation of positions)

The strategy is now to define  $N_e^i(\xi)$  in the reference configuration. Let us first discuss this concept for simplicial elements in the next section, before proceeding to more general classes of polynomial element interpolations.

## 9 Simplicial elements

A **simplex** of order  $k$  is a  $k$ -dimensional polytope which is the convex hull of its  $k + 1$  vertices.

In plain English, a simplex of order  $k$  is a convex body made up of  $k + 1$  nodes:

- in 1D: a 2-node bar, interpolation basis is  $\{1, x\}$
- in 2D: a 3-node triangle (T3), interpolation basis is  $\{1, x, y\}$
- in 3D: a 4-node tetrahedron (T4), interpolation basis is  $\{1, x, y, z\}$

Overall, this shows that interpolation is of degree  $q = 1$  (linear) for all simplicial elements.

One usually uses special shape functions for simplices, based on barycentric coordinates.

### 9.1 Linear Triangle (T3)

In 2D, the 3-node triangular element uses the **barycentric coordinates**

$$l_e^i(\mathbf{x}) = \frac{A_e^i(\mathbf{x})}{A_e}, \quad (9.1)$$

where  $A_e = |\Omega_e|$  is the total triangle area, and  $A_e^i$  is the sub-area opposite from node  $i$ , so that  $\sum_{i=1}^3 A_e^i = A_e$ .

It is an easy check to see that  $0 \leq l_e^i \leq 1$  and  $\sum_{i=1}^3 l_e^i(\mathbf{x}) = 1$  for all  $\mathbf{x} \in \Omega_e$ , so the  $l_e^i$  qualify as shape functions for the T3 element.

For convenience we use  $r = l_e^1$  and  $s = l_e^2$  as *reference coordinates*, so that the shape functions become

$$N_e^1(r, s) = r, \quad N_e^2(r, s) = s, \quad N_e^3(r, s) = 1 - r - s. \quad (9.2)$$

The isoparametric mapping in 2D now means

$$x = \sum_{i=1}^3 N_e^i(r, s) x_e^i, \quad y = \sum_{i=1}^3 N_e^i(r, s) y_e^i. \quad (9.3)$$

There is one particular **difficulty with isoparametric elements**:

To compute force vectors, etc., we need shape function derivatives  $N_{,x}^i$  and  $N_{,y}^i$  but the shape functions were defined as  $N^i(r, s)$ , so only  $N_{,r}^i$  and  $N_{,s}^i$  are known.

Let us use  $x = x(r, s)$  and  $y = y(r, s)$  so the chain rule gives

$$\begin{pmatrix} u_{,r} \\ u_{,s} \end{pmatrix} = \begin{pmatrix} u_{,x}x_{,r} + u_{,y}y_{,r} \\ u_{,x}x_{,s} + u_{,y}y_{,s} \end{pmatrix} = \begin{pmatrix} x_{,r} & y_{,r} \\ x_{,s} & y_{,s} \end{pmatrix} \begin{pmatrix} u_{,x} \\ u_{,y} \end{pmatrix} = \mathbf{J} \begin{pmatrix} u_{,x} \\ u_{,y} \end{pmatrix} \quad (9.4)$$

with the **Jacobian** matrix

$$\mathbf{J} = \begin{pmatrix} x_{,r} & y_{,r} \\ x_{,s} & y_{,s} \end{pmatrix} \quad (9.5)$$

so that

$$J = \det \mathbf{J} = \frac{\partial x}{\partial r} \frac{\partial y}{\partial s} - \frac{\partial x}{\partial s} \frac{\partial y}{\partial r}. \quad (9.6)$$



Note that, like for the deformation mapping, for the isoparametric mapping to be invertible we need to have  $J > 0$ . This implies that *elements cannot be distorted* (inverted or non-convex).

Using our isoparametric mapping, we obtain

$$\mathbf{J} = \begin{pmatrix} \sum_{i=1}^3 N_{e,r}^i x_e^i & \sum_{i=1}^3 N_{e,r}^i y_e^i \\ \sum_{i=1}^3 N_{e,s}^i x_e^i & \sum_{i=1}^3 N_{e,s}^i y_e^i \end{pmatrix}. \quad (9.7)$$

This solves the problem. As discussed, we need to have  $J > 0$  so that we can invert  $\mathbf{J}$  to arrive at (by the [inverse function theorem](#))

$$\begin{pmatrix} u_{,x} \\ u_{,y} \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} u_{,r} \\ u_{,s} \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} \sum_{i=1}^3 N_{e,r}^i u_e^i \\ \sum_{i=1}^3 N_{e,s}^i u_e^i \end{pmatrix} \quad \text{but also} \quad \begin{pmatrix} u_{,x} \\ u_{,y} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^3 N_{e,x}^i u_e^i \\ \sum_{i=1}^3 N_{e,y}^i u_e^i \end{pmatrix}. \quad (9.8)$$

By equating these two and comparing the coefficients of  $u_e^i$  we thus obtain

$$\begin{pmatrix} N_{e,x}^i \\ N_{e,y}^i \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} N_{e,r}^i \\ N_{e,s}^i \end{pmatrix} \quad \text{and more generally:} \quad \boxed{\nabla_{\mathbf{x}} N_e^i = \mathbf{J}^{-1} \nabla_{\boldsymbol{\xi}} N_e^i} \quad (9.9)$$

with reference coordinates  $\boldsymbol{\xi} = \{r, s\}$ . This is generally applicable *for any isoparametric mapping*.

By inserting the above shape functions into the Jacobian matrix we find that

$$\mathbf{J} = \begin{pmatrix} x_{,r} & y_{,r} \\ x_{,s} & y_{,s} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^3 N_{e,r}^i x_e^i & \sum_{i=1}^3 N_{e,r}^i y_e^i \\ \sum_{i=1}^3 N_{e,s}^i x_e^i & \sum_{i=1}^3 N_{e,s}^i y_e^i \end{pmatrix} = \begin{pmatrix} x_e^1 - x_e^3 & y_e^1 - y_e^3 \\ x_e^2 - x_e^3 & y_e^2 - y_e^3 \end{pmatrix} \quad (9.10)$$

and

$$J = \det \mathbf{J} = (x_e^1 - x_e^3)(y_e^2 - y_e^3) - (x_e^2 - x_e^3)(y_e^1 - y_e^3) = 2A_e. \quad (9.11)$$

Notice that  $\mathbf{J}$  and hence  $J$  are constant and do not depend on  $(r, s)$ .

Further, we have

$$\begin{pmatrix} N_{e,x}^i \\ N_{e,y}^i \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} N_{e,r}^i \\ N_{e,s}^i \end{pmatrix}. \quad (9.12)$$

From (12.1) we see that all shape function derivatives are constant ( $-1$ ,  $+1$ , or  $0$ ), so that we may conclude that

$$N_{e,x}^i = \text{const.}, \quad N_{e,y}^i = \text{const.} \quad \text{and also} \quad dA = J \, dr \, ds = 2A_e \, dr \, ds. \quad (9.13)$$

The constant shape function derivatives indicate that *all strain components are constant* within the element since  $\boldsymbol{\varepsilon} = \text{sym}(\nabla \mathbf{u})$ . This is why the 3-node triangle element is also called **Constant Strain Triangle** or **CST**. This also has the important consequence that integration of force vectors or stiffness matrices can easily be performed *exactly* by a *single quadrature point* since the integrands are constant across the element.

In case of higher-order triangular elements, the following relation is helpful for evaluating integrals:

$$\int_{\Omega_e} r^\alpha s^\beta \, dA = \frac{\alpha! \beta!}{(\alpha + \beta + 2)!} 2A_e. \quad (9.14)$$

## 9.2 Extension to three dimensions:

The extension to three dimensions is straight-forward and results in the **4-node tetrahedron** (**constant strain tetrahedron**) with reference coordinates  $(r, s, t)$  and shape functions

$$N_e^1(r, s, t) = r, \quad N_e^2(r, s, t) = s, \quad N_e^3(r, s, t) = t, \quad N_e^4(r, s, t) = 1 - r - s - t. \quad (9.15)$$

Like in 2D, strains are constant in the linear tetrahedron, and  $dV = 6V \, dr \, ds \, dt$ .

Note the following important relation that is analogous to (9.14),

$$\int_{\Omega_e} r^\alpha s^\beta t^\gamma dA = \frac{\alpha! \beta! \gamma!}{(\alpha + \beta + \gamma + 2)!} 2A. \quad (9.16)$$

## 9.3 Finite element implementation

When using the above simplicial elements, the finite element implementation can make use of the isoparametric mapping. Specifically, we showed that

$$N_{e,j}^i = J_{j\xi}^{-1} N_{e,\xi}^i = \text{const.} \quad \text{and} \quad J_e = \det \mathbf{J} = 2A_e. \quad (9.17)$$

We further note that strains are constant within elements since

$$u_{i,j}^h = \sum_{a=1}^n u_i^a N_{e,j}^a = \sum_{a=1}^n u_i^a J_{j\xi}^{-1} N_{e,\xi}^a = \text{const.} \quad \Rightarrow \quad \sigma_{ij} = \sigma_{ij}(\nabla \mathbf{u}^h) = \text{const.} \quad (9.18)$$

The key integrals to be evaluated (e.g., in 2D with an element thickness  $t$ ) now become

$$F_{\text{int},i}^a = \int_{\Omega_e} \sigma_{ij} N_{e,j}^a dV = \int_{\Omega_e} \sigma_{ij} N_{e,j}^a t dA = \sigma_{ij} J_{j\xi}^{-1} N_{e,\xi}^a t \int_{\Omega_e} dA = \sigma_{ij} J_{j\xi}^{-1} N_{e,\xi}^a \frac{J_e t}{2}. \quad (9.19)$$

Note that if we pre-compute the constant quantities in the element constructor, e.g.,

$$w_e = \frac{J_e t}{2} = \text{const.} \quad \text{and} \quad N_{e,j}^a = J_{j\xi}^{-1} N_{e,\xi}^a = \text{const.}, \quad (9.20)$$

then the above reduces to

$$F_{\text{int},i}^a = \sigma_{ij}(\nabla \mathbf{u}^h) N_{e,j}^a w_e. \quad (9.21)$$

Similarly, we compute the element energy

$$I_e = \int_{\Omega_e} W(\nabla \mathbf{u}^h) dV = \int_{\Omega_e} W(\nabla \mathbf{u}^h) t dA = W(\nabla \mathbf{u}^h) w_e. \quad (9.22)$$

## 9.4 Higher-order triangles and tetrahedra:

One can also define higher-order triangular and tetrahedral elements.

For example, the **quadratic triangle** (T6) element has six nodes so that the interpolation function space is  $\{1, x, y, x^2, xy, y^2\}$  and therefore complete up to second order. Per convention, nodes 1-3 are the corners while 4-6 are at edge midpoints, and counting is counter-clockwise as before.

With the same reference coordinates  $(r, s)$  as for the T3 element, the shape functions can be found as

$$\begin{aligned} N_e^1(r, s) &= r(2r - 1), & N_e^2(r, s) &= s(2s - 1), & N_e^3(r, s) &= t(2t - 1), \\ N_e^4(r, s) &= 4rs, & N_e^5(r, s) &= 4st, & N_e^6(r, s) &= 4rt, \end{aligned} \quad (9.23)$$

where  $t = 1 - r - s$ .

Since shape function derivatives are linear, strains (and thus stresses) also vary linearly within the element which is therefore known as the **linear strain triangle** (LST). The quadratic interpolation implies that element edges of isoparametric elements can be curved.

Analogously, the **quadratic tetrahedron** (T10) has four corner nodes and six nodes at edge midpoints.

## 10 The bilinear quadrilateral element

The above simplicial element is simple to implement and its constant strains admit an exact integration of energy, stresses, and stiffness. However, the constant fields within CST elements may not be ideal. As an alternative in 2D, let us discuss the **4-node bilinear quadrilateral element** (also known as **Q4**).

By definition, the reference configuration is  $\Omega_e = [-1, 1]^2$  and node numbering starts in the bottom left corner and is counterclockwise.

Shape functions must satisfy  $N_e^i(\xi_j, \eta_j) = \delta_{ij}$  and can be obtained from the 2-node bar which has:

$$\text{2-node bar:} \quad \bar{N}_e^1(\xi) = \frac{1}{2}(1 - \xi), \quad \bar{N}_e^2(\xi) = \frac{1}{2}(1 + \xi). \quad (10.1)$$

It can easily be verified that  $\bar{N}_e^1(-1) = 1$ ,  $\bar{N}_e^1(1) = 0$ , and  $\bar{N}_e^2(-1) = 0$ ,  $\bar{N}_e^2(1) = 1$ .

The 2D element shape functions thus follow from combining the above in the  $\xi$  and  $\eta$  directions:

$$\begin{aligned} \text{Q4 element:} \quad N_e^1(\xi, \eta) &= \bar{N}_e^1(\xi)\bar{N}_e^1(\eta) = \frac{1}{4}(1 - \xi)(1 - \eta), \\ N_e^2(\xi, \eta) &= \bar{N}_e^2(\xi)\bar{N}_e^1(\eta) = \frac{1}{4}(1 + \xi)(1 - \eta), \\ N_e^3(\xi, \eta) &= \bar{N}_e^2(\xi)\bar{N}_e^2(\eta) = \frac{1}{4}(1 + \xi)(1 + \eta), \\ N_e^4(\xi, \eta) &= \bar{N}_e^1(\xi)\bar{N}_e^2(\eta) = \frac{1}{4}(1 - \xi)(1 + \eta). \end{aligned} \quad (10.2)$$

One can easily verify that  $N_e^i(\xi_j, \eta_j) = \delta_{ij}$  and  $\sum_{i=1}^4 N_e^i(\xi, \eta) = 1$  for all  $(\xi, \eta)$ .

The isoparametric mapping here implies

$$x = \sum_{i=1}^4 N_e^i(\xi, \eta) x_e^i, \quad y = \sum_{i=1}^4 N_e^i(\xi, \eta) y_e^i. \quad (10.3)$$

Notice that this implies *straight edges* (in the reference configuration) *remain straight* (in the actual mesh). For example, take the bottom edge ( $\eta = -1$ ), the interpolation along this edge is  $x = \sum_{i=1}^4 N_e^i(\xi, -1)x_e^i$  and  $y = \sum_{i=1}^4 N_e^i(\xi, -1)y_e^i$ , which are both linear in  $\xi$ . Therefore, this element has straight edges in physical space.

### Remarks:

- *completeness* up to only  $q = 1$  is given by  $\{1, \xi, \eta, \xi\eta\}$ . This means we must be able to represent solutions

$$u^h(x, y) = c_0 + c_1x + c_2y \quad (10.4)$$

exactly. Check:

$$\begin{aligned} u^h(x, y) &= \sum_{i=1}^4 N_e^i u_e^i = \sum_{i=1}^4 N_e^i u^h(x_i, y_i) = \sum_{i=1}^4 N_e^i (c_0 + c_1x_i + c_2y_i) \\ &= \left( \sum_{i=1}^4 N_e^i \right) c_0 + c_1 \sum_{i=1}^4 N_e^i x_i + c_2 \sum_{i=1}^4 N_e^i y_i = c_0 + c_1x + c_2y. \end{aligned} \quad (10.5)$$

- *integrability*: Note that this interpolation scheme ensures that  $u^h$  is continuous across elements. To see this, notice that on any element edge only those two shape functions are non-zero whose nodes are on that edge while the others are zero.

As for the simplicial element, computing shape function derivatives requires the use of the inverse function theorem.

In analogy to the simplicial element, we use  $x = x(\xi, \eta)$  and  $y = y(\xi, \eta)$  so the inverse function theorem yields

$$\begin{pmatrix} N_{e,x}^i \\ N_{e,y}^i \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} N_{e,\xi}^i \\ N_{e,\eta}^i \end{pmatrix} \quad \text{and more generally:} \quad \boxed{\nabla_{\mathbf{x}} N_e^i = \mathbf{J}^{-1} \nabla_{\boldsymbol{\xi}} N_e^i} \quad (10.6)$$

As a simple **example**, recall the 2-node bar element whose shape functions we computed as

$$N_e^1(x) = 1 - \frac{x}{\Delta x}, \quad N_e^2(x) = \frac{x}{\Delta x}. \quad (10.7)$$

For a reference bar element with nodes at  $\xi = \pm 1$ , the analogous shape functions in 1D read

$$N_e^1(\xi) = \frac{1 - \xi}{2}, \quad N_e^2(\xi) = \frac{1 + \xi}{2}. \quad (10.8)$$

Applying the above relations to the 1D problem gives

$$J = \frac{\partial x}{\partial \xi} = \frac{\partial N_e^1}{\partial \xi} x_e^1 + \frac{\partial N_e^2}{\partial \xi} x_e^2 = \frac{x_e^2 - x_e^1}{2} = \frac{\Delta x}{2}. \quad (10.9)$$

This confirms that indeed

$$N_{e,x}^i = J^{-1} N_{e,\xi}^i = \frac{2}{\Delta x} N_{e,\xi}^i. \quad (10.10)$$

A useful relation to evaluate area integrals in the following is ( $\mathbf{e}_1, \mathbf{e}_2$  being reference unit vectors)

$$d\mathbf{A} = d\mathbf{x} \times d\mathbf{y} = \left( \frac{\partial x}{\partial \xi} d\xi \mathbf{e}_1 + \frac{\partial x}{\partial \eta} d\eta \mathbf{e}_2 \right) \times \left( \frac{\partial y}{\partial \xi} d\xi \mathbf{e}_1 + \frac{\partial y}{\partial \eta} d\eta \mathbf{e}_2 \right) = J d\xi d\eta \mathbf{e}_1 \times \mathbf{e}_2 \quad (10.11)$$

so that

$$\boxed{dA = |d\mathbf{A}| = J d\xi d\eta} \quad (10.12)$$

### Extension to higher-order elements:

The **9-node quadratic quadrilateral (Q9)** derives its shape functions from applying the 1D shape functions of the 3-node bar (using Lagrangian interpolation) to the 2D case. For example,

$$N_1(\xi, \eta) = \frac{\xi(1 - \xi)\eta(1 - \eta)}{4}, \quad (10.13)$$

so that overall the interpolation includes monomials

$$\{1, \xi, \eta, \xi\eta, \xi^2, \eta^2, \xi^2\eta, \xi\eta^2, \xi^2\eta^2\}, \quad (10.14)$$

which is complete up to order  $q = 2$  (quadratic).

Since the above elements includes way more polynomial terms than required for quadratic interpolation, one can construct elements with less nodes, e.g., the **8-node quadratic quadrilateral (Q8)** also known as **serendipity element**.

Shape functions are constructed as follows:

$$\begin{aligned} N_5(\xi, \eta) &= \frac{(1-\eta)(1-\xi^2)}{4}, \\ N_6(\xi, \eta) &= \frac{(1-\eta^2)(1+\xi)}{4}, \\ N_1(\xi, \eta) &= \frac{(1-\eta)(1-\xi)}{4} - \frac{1}{2}(N_5 + N_6), \quad \text{etc.} \end{aligned} \tag{10.15}$$

### Extension to three dimensions:

The same procedure can be applied to 3D, resulting in the **8-node brick element**. The reference coordinates  $(\xi, \eta, \zeta)$  are linked to the physical coordinates  $(x, y, z)$  via the shape functions

$$N_e^1(\xi, \eta, \zeta) = \frac{1}{8}(1-\xi)(1-\eta)(1-\zeta), \quad \dots \tag{10.16}$$

again with

$$\boxed{\nabla_{\mathbf{x}} N_e^i = \mathbf{J}^{-1} \nabla_{\boldsymbol{\xi}} N_e^i} \quad \text{and} \quad \boxed{dV = J \, d\xi \, d\eta \, d\zeta} \tag{10.17}$$

## 11 Numerical quadrature

The finite element method frequently requires computing integrals, e.g., for the internal/external force vectors. Since these cannot be integrated analytically in general (with exceptions like the simplicial elements), we need **numerical integration** schemes.

Consider the integral

$$I[u] = \int_a^b u(x) \, dx. \quad (11.1)$$

For convenience, let us introduce the **shift**

$$\xi = 2 \frac{x-a}{b-a} - 1 \quad \Rightarrow \quad \xi(x=a) = -1, \quad \xi(x=b) = 1, \quad d\xi = \frac{2 \, dx}{b-a} \quad (11.2)$$

and  $x = \frac{\xi+1}{2}(b-a) + a$  so that

$$I[u] = \int_{-1}^1 f(\xi) \, d\xi \quad \text{where} \quad f(\xi) = \frac{b-a}{2} u\left(\frac{\xi+1}{2}(b-a) + a\right). \quad (11.3)$$

The goal is now to approximate the integral in (11.3) numerically.

### 11.1 Example: Riemann sums

Consider a *partition* with  $n+1$  nodes:

$$P = \{[\xi_0, \xi_1], [\xi_1, \xi_2], \dots, [\xi_{n-1}, \xi_n]\} \quad \text{such that} \quad -1 = \xi_0 < \xi_1 < \xi_2 < \dots < \xi_n = 1 \quad (11.4)$$

The **Riemann sum** makes use of the approximation

$$I \approx S = \sum_{i=1}^n f(\xi_i^*) (\xi_i - \xi_{i-1}) \quad \text{with} \quad \xi_{i-1} \leq \xi_i^* \leq \xi_i. \quad (11.5)$$

Different choices of  $\xi_i^*$ :

- (i) **left Riemann sum:**  $\xi_i^* = \xi_{i-1}$ :
- (ii) **right Riemann sum:**  $\xi_i^* = \xi_i$
- (iii) **middle Riemann sum:**  $\xi_i^* = \frac{1}{2}(\xi_{i-1} + \xi_i)$
- (iv) **trapezoidal sum:** average of left and right
- (v) **upper Riemann sum:**  $\xi_i^*$  s.t.  $g(\xi_i^*) = \sup_{\xi \in [\xi_{i-1}, \xi_i]} g(\xi)$
- (vi) **lower Riemann sum:**  $\xi_i^*$  s.t.  $g(\xi_i^*) = \inf_{\xi \in [\xi_{i-1}, \xi_i]} g(\xi)$

More refined formulations can be found in the so-called **Newton-Cotes** formulas (the trapezoidal rule is the Newton-Cotes formula of degree 1).

### 11.2 Gauss quadrature

A much more efficient alternative are **quadrature rules** (also called **cubature rules** in 3D) which are best suited for polynomial functions:

$$I[u] = \int_{-1}^1 f(\xi) \, d\xi \approx \sum_{i=0}^{n_{QP}-1} W_i f(\xi_i). \quad (11.6)$$

Now, we need to choose  $n_{QP}$ ,  $\{W_0, \dots, W_{n_{QP}-1}\}$  and  $\{\xi_0, \dots, \xi_{n_{QP}-1}\}$ . The choice should depend on the function to be integrated (more specifically, on its *smoothness*). Note that most functions of interest will be of polynomial type.

We say a quadrature rule is **exact of order  $q$**  if it integrates exactly all polynomial functions  $g \in \mathbb{P}_q([-1, 1])$ . **Gauss quadrature** generally chooses  $n_{QP}$  quadrature points and associated weights such that the quadrature rule is exact of order  $q = 2n_{QP} - 1$ .

### 11.2.1 Gauss-Legendre quadrature

**Gauss-Legendre quadrature** selects the nodes and weights such that the first  $2n_{QP}$  moments are computed exactly:

$$\mu_k = \int_{-1}^1 \xi^k d\xi \stackrel{!}{=} \sum_{i=0}^{n_{QP}-1} W_i \xi_i^k, \quad k = 0, 1, \dots, 2n_{QP} - 1. \quad (11.7)$$

These are  $2n_{QP}$  equations for the  $2n_{QP}$  free parameters  $(W_i, \xi_i)$  for  $i = 0, \dots, n_{QP} - 1$ . The equations are generally nonlinear and thus hard to solve analytically.

Let us compute the Gauss-Legendre weights and points for the lowest few orders in 1D:

- *a single quadrature point*, i.e.,  $n_{QP} = 1$ :

Two equations for two unknowns:

$$\int_{-1}^1 \xi^0 d\xi = 2 = W_0 \xi_0^0 = W_0 \quad \text{and} \quad \int_{-1}^1 \xi^1 d\xi = 0 = W_0 \xi_0^1 = W_0 \xi_0 \quad (11.8)$$

So that the first-order quadrature rule is given by

$$\boxed{W_0 = 2, \quad \xi_0 = 0} \quad (11.9)$$

Since linear functions are integrated exactly, this quadrature rule is *exact to order  $q = 1$* .

- *two quadrature points*, i.e.,  $n_{QP} = 2$ :

In close analogy, we now have

$$\begin{aligned} \int_{-1}^1 \xi^0 d\xi = 2 &= W_0 + W_1, & \int_{-1}^1 \xi^1 d\xi = 0 &= W_0 \xi_0 + W_1 \xi_1, \\ \int_{-1}^1 \xi^2 d\xi = \frac{1}{3} &= W_0 \xi_0^2 + W_1 \xi_1^2, & \int_{-1}^1 \xi^3 d\xi = 0 &= W_0 \xi_0^3 + W_1 \xi_1^3. \end{aligned} \quad (11.10)$$

A simple solution can be found for symmetric quadrature points  $\xi_0 = -\xi_1$ :

$$\boxed{W_0 = W_1 = 1, \quad \xi_0 = -\frac{1}{\sqrt{3}}, \quad \xi_1 = \frac{1}{\sqrt{3}}.} \quad (11.11)$$

This quadrature rule is *exact to order  $q = 3$*  (cubic polynomials are integrated exactly).

- *higher-order quadrature rules*:

Quadrature weights and points for arbitrary order can be obtained in analogous fashion and, most importantly, can be found in numerous look-up tables (see notes and textbooks). However, there is a better, systematic way to compute Gauss-Legendre quadrature weights and points.



Note that monomials  $\{1, \xi, \xi^2, \xi^3, \dots\}$ , although complete, are not orthogonal basis functions. We can turn them into orthogonal polynomials  $P_n(\xi)$  by, e.g., the **Gram-Schmidt orthogonalization** procedure. To this end, let us start with

$$P_0(\xi) = 1 \quad (11.12)$$

and obtain the next basis function by starting with the linear monomial  $\xi$  and computing

$$P_1(\xi) = \xi - \frac{\langle 1, \xi \rangle}{\langle 1, 1 \rangle} 1 = \xi, \quad (11.13)$$

where we used the inner product

$$\langle u, v \rangle = \int_{-1}^1 u(\xi) v(\xi) d\xi. \quad (11.14)$$

Similarly, the next higher basis function is obtained by starting from  $\xi^2$ , so that

$$P_2(\xi) = \xi^2 - \frac{\langle \xi, \xi^2 \rangle}{\langle \xi, \xi \rangle} \xi - \frac{\langle 1, \xi^2 \rangle}{\langle 1, 1 \rangle} 1 = \xi^2 - \frac{1}{3}. \quad (11.15)$$

Analogously, one finds

$$P_3(\xi) = \xi^3 - \frac{3}{5}\xi. \quad (11.16)$$

By continuing analogously, we create a *countably infinite set of orthogonal basis functions*  $P_n(\xi)$  such that

$$\int_{-1}^1 P_n(\xi) P_m(\xi) d\xi = 0 \quad \text{if } n \neq m. \quad (11.17)$$

These polynomials are known as **Legendre polynomials**. Note that they are defined only up to a constant, so one can renormalize them, which is commonly done by enforcing that  $P_n(1) = 1$  for all  $n$ . The result is the well known Legendre polynomials which can alternatively be defined via

$$P_n(\xi) = \frac{1}{2^n n!} \frac{d^n}{d\xi^n} [(\xi^2 - 1)^n] \quad (11.18)$$

These polynomials have another interesting feature, viz. by orthogonality with  $P_0(\xi) = 1$  we know that

$$\int_{-1}^1 P_n(\xi) d\xi = \langle 1, P_n \rangle = \begin{cases} 2, & \text{if } n = 0, \\ 0, & \text{else.} \end{cases} \quad (11.19)$$

$P_n(\xi)$  has exactly  $n$  roots in the interval  $[-1, 1]$ . Also, for  $n \neq 0$  we know that

$$\begin{aligned} P_n(\xi) &= -P_n(-\xi) & \text{for odd } n, \\ P_n(\xi) &= P_n(-\xi) & \text{for even } n. \end{aligned} \quad (11.20)$$

Moreover,  $P_n(0) = 0$  for odd  $n$ .

With this new set of basis functions, we can define the Gauss-Legendre quadrature rule to enforce

$$\int_{-1}^1 P_k(\xi) d\xi \stackrel{!}{=} \sum_{i=0}^{n_{QP}-1} W_i P_k(\xi_i), \quad k = 0, 1, \dots, 2n_{QP} - 1. \quad (11.21)$$

If  $n_{QP} = 1$ , then the solution is simple because the above equations simplify to

$$W_0 = 2 \quad \text{and} \quad 0 = W_0 P_1(\xi_0). \quad (11.22)$$

Therefore, the weight is, as before,  $W_0 = 2$  and the quadrature point is the root of  $P_1(\xi)$ , viz.  $\xi_0 = 0$ .

If  $n_{QP} = 2$ , then the four equations to be solved are

$$\begin{aligned} W_0 + W_1 &= 2 & \text{and} & & 0 &= W_0 P_1(\xi_0) + W_1 P_1(\xi_1), \\ 0 &= W_0 P_2(\xi_0) + W_1 P_2(\xi_1) & \text{and} & & 0 &= W_0 P_3(\xi_0) + W_1 P_3(\xi_1). \end{aligned} \quad (11.23)$$

By analogy, we choose the quadrature points to be the roots of  $P_2(\xi)$ , so that

$$P_2(\xi_0) = P_2(\xi_1) = 0 \quad \Rightarrow \quad \xi_0 = \frac{1}{\sqrt{3}}, \quad \xi_1 = -\frac{1}{\sqrt{3}}. \quad (11.24)$$

Using  $P_n(\xi) = -P_n(-\xi)$ , the above equations reduce to

$$W_0 = W_1 = 1. \quad (11.25)$$

Further, note that  $P_n(0) = 0$  for odd  $n$ . Therefore, the same procedure can be continued as follows. For an arbitrary number of quadrature points,  $n_{QP}$ , the **Gauss-Legendre quadrature** points and associated weights are defined by

$$\boxed{P_{n_{QP}}(\xi_i) = 0, \quad w_i = \frac{2}{(1 - \xi_i^2)[P'_{n_{QP}}(\xi_i)]^2} \quad i = 0, \dots, n_{QP} - 1.} \quad (11.26)$$

As a check, take, e.g.,  $n_{QP} = 1$  so that  $P_1(\xi) = \xi$  with root  $\xi_0 = 0$ . The weight is computed as

$$w_0 = \frac{2}{(1 - \xi_0^2)[P'_{n_{QP}}(\xi)]^2} = 2, \quad (11.27)$$

as determined above. Similarly, for  $n_{QP} = 2$  we have  $P_2(\xi) = \frac{1}{2}(3\xi^2 - 1)$  with the above roots of  $\pm 1/\sqrt{3}$ . The associated weights are computed as

$$w_0 = \frac{2}{(1 - \xi_0^2)[P'_2(\xi_0)]^2} = \frac{2}{(1 - \xi_0^2)[3\xi_0]^2} = \frac{2}{\frac{2}{3}\frac{3^2}{3}} = 1 = w_1, \quad (11.28)$$

which agrees with our prior solution.

### 11.2.2 Other Gauss quadrature rules

Note that if, for general functions  $f$ , one can sometimes find a decomposition  $f(\xi) = w(\xi)g(\xi)$  where  $w(\cdot)$  is a known weighting function and  $g(\xi)$  is (approximately) polynomial, so that a more suitable quadrature rule may be found via

$$I[u] = \int_{-1}^1 f(\xi) d\xi = \int_{-1}^1 w(\xi) g(\xi) d\xi \approx \sum_{i=0}^{n_{QP}-1} w(\xi_i) g(\xi_i). \quad (11.29)$$

Examples of such Gaussian quadrature rules include those of **Gauss-Chebyshev** type, which are obtained from a weighting function  $w(\xi) = (1 - \xi^2)^{-1/2}$ , and the quadrature points are the roots of Chebyshev polynomials. **Gauss-Hermite** quadrature uses a weighting function

$w(\xi) = \exp(-\xi^2)$  (and the integral is taken over the entire real axis). Gauss-Legendre quadrature is included as the special case  $w(\xi) = 1$ .

Another popular alternative (less for FE though) is **Gauss-Lobatto quadrature** which includes the interval end points as quadrature points and is accurate for polynomials up to degree  $2n_{QP} - 3$ , viz.

$$I[u] = \int_{-1}^1 f(\xi) d\xi \approx \frac{2}{n_{QP}(n_{QP} - 1)} [f(-1) + f(1)] + \sum_{i=2}^{n_{QP}-1} W_i f(\xi_i). \quad (11.30)$$

### 11.3 Higher dimensions

Like the polynomial shape functions, the above quadrature rules can easily be extended to 2D and 3D, e.g.,

$$\begin{aligned} \int_{-1}^1 \int_{-1}^1 f(\xi, \eta) d\xi d\eta &= \int_{-1}^1 \left[ \sum_{i=0}^{N-1} W_i f(\xi_i, \eta) \right] d\eta = \sum_{j=0}^{N-1} W_j \left[ \sum_{i=0}^{N-1} W_i f(\xi_i, \eta_j) \right] \\ &= \sum_{k=0}^{n_{QP}-1} W_k^* f(\xi_k, \eta_k) \end{aligned} \quad (11.31)$$

with the combined weights  $W_k^* = W_i W_j$  and points  $(\xi_k, \eta_k) = (\xi_i, \eta_j)$  obtained from the individual quadrature rules in each direction. By symmetry we choose  $N = \sqrt{n_{QP}}$  so that  $N^2 = n_{QP}$ .

For example, consider the **Q4 element**. By reusing the 1D Gauss-Legendre weights and points, we now have:

- *first-order quadrature* ( $q = 1$ ), as in 1D, has only a single quadrature point ( $n_{QP} = 1$ ):

$$\boxed{W_0 = 1 \quad \text{and} \quad (\xi_0, \eta_0) = (0, 0)} \quad (11.32)$$

Bilinear functions (at most linear in  $\xi$  and  $\eta$ ) are integrated exactly with this rule.

- *third-order quadrature* ( $q = 3$ ), now has four quadrature points ( $n_{QP} = 2^2 = 4$ ):

$$\boxed{W_0 = W_1 = W_2 = W_3 = 1 \quad \text{and} \quad (\xi_i, \eta_i) = \left( \pm \frac{1}{\sqrt{3}}, \pm \frac{1}{\sqrt{3}} \right)} \quad (11.33)$$

Bicubic polynomial functions (at most cubic in  $\xi$  and  $\eta$ ) are integrated exactly.

Similarly, the **brick element** in 3D uses Gauss-Legendre quadrature as follows:

- *first-order quadrature* ( $q = 1$ ) still has a single quadrature point ( $n_{QP} = 1$ ):

$$\boxed{W_0 = 1 \quad \text{and} \quad (\xi_0, \eta_0) = (0, 0)} \quad (11.34)$$

- *third-order quadrature* ( $q = 3$ ), now has four quadrature points ( $n_{QP} = 2^3 = 8$ ):

$$\boxed{W_i = 1 \quad \text{and} \quad (\xi_i, \eta_i, \zeta_i) = \left( \pm \frac{1}{\sqrt{3}}, \pm \frac{1}{\sqrt{3}}, \pm \frac{1}{\sqrt{3}} \right)}. \quad (11.35)$$

## 11.4 Finite element implementation

The key integrals to be evaluated numerically are (e.g., in 2D)

$$I^e = \int_{\Omega_e} W(\nabla \mathbf{u}^h) dV = \int_{-1}^1 \int_{-1}^1 W(\nabla \mathbf{u}^h(\boldsymbol{\xi})) J(\boldsymbol{\xi}) t_e d\xi d\eta \approx \sum_{k=0}^{n_{QP}-1} W_k W(\nabla \mathbf{u}^h(\boldsymbol{\xi}_k)) J(\boldsymbol{\xi}_k) t_e,$$

$$F_{\text{int},i}^a = \int_{\Omega_e} \sigma_{ij} N_{,j}^a dV = \int_{-1}^1 \int_{-1}^1 \sigma_{ij}(\boldsymbol{\xi}) N_{,j}^a(\boldsymbol{\xi}) J(\boldsymbol{\xi}) t_e d\xi d\eta \approx \sum_{k=0}^{n_{QP}-1} W_k \sigma_{ij}(\boldsymbol{\xi}_k) N_{,j}^a(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k) t_e$$
(11.36)

and the tangent matrix

$$T_{ik}^{ab} = \int_{\Omega_e} \mathbb{C}_{ijkl} N_{,j}^a N_{,l}^b dV = \int_{-1}^1 \int_{-1}^1 \mathbb{C}_{ijkl}(\boldsymbol{\xi}) N_{,j}^a(\boldsymbol{\xi}) N_{,l}^b(\boldsymbol{\xi}) J(\boldsymbol{\xi}) t_e d\xi d\eta$$

$$\approx \sum_{k=0}^{n_{QP}-1} W_k \mathbb{C}_{ijkl}(\boldsymbol{\xi}_k) N_{,j}^a(\boldsymbol{\xi}_k) N_{,l}^b(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k) t_e$$
(11.37)

Notice for implementation purposes that  $\sigma_{ij} N_{,j}^a$  is a simple matrix-vector multiplication so that, using the isoparametric mapping of Section 8.3,

$$\mathbf{F}_{\text{int}}^a \approx \sum_{k=0}^{n_{QP}-1} W_k J(\boldsymbol{\xi}_k) t_e \boldsymbol{\sigma}(\boldsymbol{\xi}_k) \nabla_{\mathbf{x}} N^a(\boldsymbol{\xi}_k) \quad \text{with} \quad \nabla_{\mathbf{x}} N^a = \mathbf{J}^{-1} = \nabla_{\boldsymbol{\xi}} N^a. \quad (11.38)$$

## 11.5 Quadrature error estimates

Using numerical quadrature to approximate an exact integral introduces a **quadrature error**, which is bounded as follows. For a function  $f(\boldsymbol{\xi}) \in C^{2n_{QP}}(\Omega)$  we have that (without proof here)

$$\left| \int_{-1}^1 f(\boldsymbol{\xi}) d\xi - \sum_{k=0}^{n_{QP}-1} W_k f(\boldsymbol{\xi}_k) \right| \leq C h^{2n_{QP}} \max_{\boldsymbol{\xi} \in [-1,1]} |f^{(2n_{QP})}(\boldsymbol{\xi})| \quad (11.39)$$

with a constant  $C > 0$ . This shows that, as can be expected, the quadrature error decreases with decreasing mesh size  $h$  and with increasing smoothness of function  $f$ . The rate of convergence under  $h$ -refinement also depends on the smoothness of the function.

The exact error depends on the chosen quadrature rule. For example, for Gauss-Legendre quadrature an error estimate is given by

$$e = \frac{2^{2n_{QP}+1} (n_{QP}!)^4}{(2n_{QP}+1) [(2n_{QP})!]^3} \max_{\boldsymbol{\xi} \in [-1,1]} \|f^{(2n_{QP})}(\boldsymbol{\xi})\|. \quad (11.40)$$

## 11.6 Quadrature rules for simplicial elements:

As discussed above, simplicial elements (bar, triangle, tetrahedron) produce *constant strains* and thus *constant stresses* within elements. Hence, a **single quadrature point** at an **arbitrary location** inside the element is sufficient. Usually one chooses the point to be located at the element center, which gives

$$W_0 = 1, \quad r_0 = s_0 = \frac{1}{3} \quad (\text{in 2D}) \quad \text{and} \quad r_0 = s_0 = t_0 = \frac{1}{4} \quad (\text{in 3D}). \quad (11.41)$$

If higher-order quadrature rules are required (e.g., for triangular and tetrahedral elements of higher interpolation order as discussed next), the same concepts as for Gauss-Legendre quadrature can be applied here, resulting in **simplicial quadrature** rules whose weights and quadrature point locations can be found in look-up tables.

## 11.7 Which quadrature rule to use?

Stresses, shape function derivatives, and Jacobians are not necessarily smooth polynomials. Thus, rather than finding the exact integration order for each element and constitutive model, we introduce a minimum required integration order for a particular element type.

Our minimum requirement is that an undistorted elastic element is integrated exactly. Thus, we define **full integration** as the order needed to integrate an *undistorted, homogeneous, linear elastic element exactly*. An element is undistorted if element angles are preserved or, in other words, if  $J = \text{const.}$

For **example**, the 4-node quadrilateral (Q4) is undistorted if the physical element has the shape of a rectangle (side lengths  $a$  and  $b$ ), so that

$$J = \frac{ab}{4} = \text{const.} \quad (11.42)$$

Then, for a linear elastic Q4 element we have

$$F_{\text{int},i}^a = \int_{-1}^1 \int_{-1}^1 \mathbb{C}_{ijkl} \varepsilon_{kl}^h(\boldsymbol{\xi}) N_{,j}^a(\boldsymbol{\xi}) \frac{ab}{4} t_e d\xi d\eta = \frac{abt_e}{4} \mathbb{C}_{ijkl} \int_{-1}^1 \int_{-1}^1 \varepsilon_{kl}^h(\boldsymbol{\xi}) N_{,j}^a(\boldsymbol{\xi}) d\xi d\eta, \quad (11.43)$$

where  $\varepsilon^h = \text{sym}(\nabla_{\mathbf{x}} \mathbf{u}^h)$  is at most linear (since the interpolation of  $\mathbf{u}^h$  is bilinear),  $\nabla_{\mathbf{x}} N^a$  is also at most linear for the same reason. Overall, the integrand is at most a quadratic polynomial, so that we need integration order  $q \geq 2$ .

Recall that in 1D we showed that  $q = 2n_{QP} - 1$ , so that full integration of the Q4 element in 2D requires  $n_{QP} = 2 \times 2 = 4$  quadrature points.

Analogously, full integration of the quadratic 2D elements Q8/Q9 requires  $n_{QP} = 3^2 = 9$  quadrature points. Full integration of the 8-node brick element requires  $n_{QP} = 8$  quadrature points.

For simplicial elements, we showed that a single quadrature point ( $n_{QP} = 1$ ) is always sufficient to integrate exactly, since stresses and strains within elements are constant. For the quadratic triangle (T6), full integration requires order  $q = 2$ , which corresponds to *three quadrature points*.

Note that not only does full integration guarantee that the internal force vector of an undistorted, elastic element is integrated exactly. By reviewing the element energy and the element tangent matrix, we make the same observation (i.e., those are integrated exactly as well):

$$\begin{aligned} I_e &= \int_{\Omega_e} W dV = \frac{abt_e}{4} \frac{1}{2} \mathbb{C}_{ijkl} \int_{-1}^1 \int_{-1}^1 \varepsilon_{ij}^h(\boldsymbol{\xi}) \varepsilon_{kl}^h(\boldsymbol{\xi}) d\xi d\eta, \\ T_{ij}^{ab} &= \frac{abt_e}{4} \mathbb{C}_{ijkl} \int_{-1}^1 \int_{-1}^1 N_{,j}^a(\boldsymbol{\xi}) N_{,l}^b(\boldsymbol{\xi}) d\xi d\eta. \end{aligned} \quad (11.44)$$

Using an integration rule less than full integration is called **under-integration**; the opposite is called **over-integration**. Which integration order to use depends very much on the element, material model, etc. Sometimes under-integration can be beneficial (e.g., to avoid locking). We will discuss some of these issues later.

## 12 Generalization and implementation of the simplicial elements

For a general simplicial element in  $d$  dimensions having  $n = d + 1$  nodes, we have shape functions

$$N_e^1(r_1, \dots, r_d) = r_1, \quad N_e^2(r_1, \dots, r_d) = r_2, \quad \dots \quad N_e^n(r_1, \dots, r_d) = 1 - \sum_i^d r_i, \quad (12.1)$$

where  $\boldsymbol{\xi} = \{r_1, \dots, r_d\}$  denote the  $d$  *barycentric coordinates*. Then, the Jacobian  $\mathbf{J}$  is given by

$$J_{ij} = \frac{\partial X_j}{\partial r_i} = \sum_{a=1}^n X_j^a \frac{\partial N_e^a}{\partial r_i} \quad \text{or} \quad \mathbf{J} = \sum_{a=1}^n \nabla_{\boldsymbol{\xi}} N_e^a \otimes \mathbf{X}^a, \quad J = \det \mathbf{J}, \quad (12.2)$$

so that shape function derivatives in *physical coordinates* follow as

$$\begin{pmatrix} N_{e,X_1}^a \\ \vdots \\ N_{e,X_d}^a \end{pmatrix} = \mathbf{J}^{-1} \begin{pmatrix} N_{e,r_1}^a \\ \vdots \\ N_{e,r_d}^a \end{pmatrix} \quad \text{or} \quad \nabla_{\mathbf{X}} N_e^a = \mathbf{J}^{-1} \nabla_{\boldsymbol{\xi}} N_e^a. \quad (12.3)$$

Using numerical quadrature with weights  $W_k$  and points  $\boldsymbol{\xi}_k$  (in *reference coordinates*), the element energy is given by

$$I_e \approx \sum_{k=1}^{n_{QP}} W_k W(\mathbf{F}(\boldsymbol{\xi}_k)) J(\boldsymbol{\xi}_k) t_e. \quad (12.4)$$

Nodal forces are approximated by

$$F_{\text{int},i}^a \approx \sum_{k=1}^{n_{QP}} W_k P_{iJ}(\mathbf{F}(\boldsymbol{\xi}_k)) N_{,J}^a(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k) t_e. \quad (12.5)$$

Finally, the element stiffness matrix components become

$$T_{il}^{ab} \approx \sum_{k=1}^{n_{QP}} W_k \mathbb{C}_{iJlL}(\mathbf{F}(\boldsymbol{\xi}_k)) N_{,J}^a(\boldsymbol{\xi}_k) N_{,L}^b(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k) t_e. \quad (12.6)$$

Here,  $t_e$  is an element constant (e.g., the cross-sectional area  $A$  in 1D, the thickness  $t$  in 2D, and simply 1 in 3D).

The deformation gradient,  $\mathbf{F} = \mathbf{I} + \nabla_{\mathbf{X}} \mathbf{u}$ , and strain tensor,  $\boldsymbol{\varepsilon} = \text{sym}(\nabla_{\mathbf{X}} \mathbf{u})$ , can be obtained directly from  $\nabla_{\mathbf{X}} \mathbf{u}$ . Also, recall that

$$\nabla_{\mathbf{X}} \mathbf{u}(\boldsymbol{\xi}) = \sum_{a=1}^n \mathbf{u}_e^a \otimes \nabla_{\mathbf{X}} N^a(\boldsymbol{\xi}). \quad (12.7)$$

Note that, in *linearized kinematics*, the above equations hold analogously with  $P_{iJ}$  replaced by the Cauchy stress tensor, upper-case coordinates  $\mathbf{X}$  by lower-case  $\mathbf{x}$ , etc.

### 13 Assembly

So far, we have defined **local** element vectors and matrices. The solution of any problem requires the assembly of **global** vectors and matrices. Specifically, recall that shape functions associated with a node  $a$  are non-zero only in elements adjacent to node  $a$  (this is the principle of *locality*). When computing, e.g., the total energy of a body discretized into  $n_e$  elements, we exploit that

$$I = \int_{\Omega} W(\nabla \mathbf{u}^h) dV = \sum_{e=1}^{n_e} \int_{\Omega_e} W(\nabla \mathbf{u}^h) dV = \sum_{e=1}^{n_e} I^e. \quad (13.1)$$

When computing force vectors and incremental stiffness matrices, the situation is more complex, which is why we introduce the **assembly operator**, viz.

$$\mathbf{F}_{\text{int}} = \mathcal{A} \mathbf{F}_{\text{int},e}, \quad \mathbf{T}_{\text{int}} = \mathcal{A} \mathbf{T}_{\text{int},e}, \quad (13.2)$$

which loops over all  $n_e$  elements  $e$  and *adds* their respective contributions to the global quantities. This requires careful book-keeping to keep track of the correspondence between local and global node numbering (and is the reason our implemented elements store the node IDs).

Similarly, an inverse assignment operator extracts element quantities from a global vector:

$$\mathbf{U}_e^h = \mathcal{A}_e^{-1}(\mathbf{U}^h). \quad (13.3)$$

For practical purposes, it is helpful to recall that entry  $U_i^a$  is located at position  $a \cdot d + i$  in the assembled vector  $\mathbf{U}^h$  (using C++ notation with 0-indexed lists, in  $d$  dimensions). For example, consider a 2D problem using 2-node bar elements. If an element connecting nodes 1 and 3 computes the nodal force vectors  $\mathbf{F}_{\text{int},e}^1$  and  $\mathbf{F}_{\text{int},e}^3$ , then they are to be added into the global force vector as

$$\mathbf{F}_{\text{int}} = \begin{pmatrix} \cdot \\ \cdot \\ F_{\text{int},e,1}^1 \\ F_{\text{int},e,2}^1 \\ \cdot \\ \cdot \\ F_{\text{int},e,1}^3 \\ F_{\text{int},e,2}^3 \\ \cdot \\ \cdot \end{pmatrix} \quad (13.4)$$

Analogously, if the element connecting nodes 0 and 2 computes a stiffness matrix  $\mathbf{K}_e^{02}$ , then its components must be added onto the global tangent stiffness matrix as

$$\mathbf{T} = \begin{pmatrix} K_{11}^{02} & K_{12}^{02} & \cdot & \cdot & K_{13}^{02} & K_{14}^{02} & \cdot & \cdot \\ K_{21}^{02} & K_{22}^{02} & \cdot & \cdot & K_{23}^{02} & K_{24}^{02} & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ K_{31}^{02} & K_{32}^{02} & \cdot & \cdot & K_{33}^{02} & K_{34}^{02} & \cdot & \cdot \\ K_{41}^{02} & K_{42}^{02} & \cdot & \cdot & K_{43}^{02} & K_{44}^{02} & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix} \quad (13.5)$$

Note that because nodes in FEM have non-zero shape functions only in adjacent elements, the above matrix is generally sparse (which comes with significant computational advantages).

## 14 Overview: Numerical Implementation

The above theoretical concepts are implemented in our c++ finite element code, whose general structure is summarized in the following. Details can be extracted from the source files. The code structure is schematically shown in Fig. 1.

### material model:

The material model computes

- $W = W(\nabla \mathbf{u})$
- $P_{iJ} = P_{iJ}(\nabla \mathbf{u})$  or  $\sigma_{ij} = \sigma_{ij}(\nabla \mathbf{u})$
- $\mathbb{C}_{iJkL} = \mathbb{C}_{iJkL}(\nabla \mathbf{u})$  or  $\mathbb{C}_{ijkl}(\nabla \mathbf{u})$

Depending on the (finite/linearized) model, the strain is

$$\mathbf{F} = \mathbf{I} + \nabla \mathbf{u} \quad \text{or}$$

$$\boldsymbol{\varepsilon} = \frac{1}{2}(\nabla \mathbf{u} + \nabla \mathbf{u}^T)$$

Also, internal variables must be updated<sup>1</sup>:

$$\mathbf{z}_{\alpha+1} = \arg \inf \mathcal{F}_{\mathbf{z}_\alpha}(\nabla \mathbf{u})$$

The `MaterialModel` class implements:

- `computeEnergy`( $\nabla \mathbf{u}, \mathbf{z}_\alpha, t$ )  $\rightarrow W^*$
- `computeStresses`( $\nabla \mathbf{u}, \mathbf{z}_\alpha, t$ )  $\rightarrow P_{iJ}$  or  $\sigma_{ij}$
- `computeTangentMatrix`( $\nabla \mathbf{u}, \mathbf{z}_\alpha, t$ )  $\rightarrow \mathbb{C}_{iJkL}$

The respective strain tensor is provided by

- `computeStrain`( $\nabla \mathbf{u}$ )  $\rightarrow F_{iJ}$  or  $\varepsilon_{ij}$

which can, of course, also be called within the element to conveniently compute stresses, etc.

Internal variables are updated by<sup>1</sup>

- `updateInternalVariables`( $\nabla \mathbf{u}, \mathbf{z}_\alpha, t$ )  
 $\rightarrow \mathbf{z}_{\alpha+1}$  (if no update, simply return  $\mathbf{z}_\alpha$ )

### element:

The element computes

- $I_e \cong \sum_{k=1}^{n_{QP}} W_k W(\nabla \mathbf{u}(\boldsymbol{\xi}_k)) J(\boldsymbol{\xi}_k)$
- $(F_{\text{int},e})_i^a \cong \sum_{k=1}^{n_{QP}} W_k P_{iJ}(\boldsymbol{\xi}_k) N_{,J}^a(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k)$
- $(T_{\text{int},e})_{ik}^{ab} \cong \sum_{k=1}^{n_{QP}} W_k \mathbb{C}_{iJkL}(\boldsymbol{\xi}_k) N_{,J}^a(\boldsymbol{\xi}_k) N_{,L}^b(\boldsymbol{\xi}_k) \times J(\boldsymbol{\xi}_k)$

In addition, internal variables must be updated<sup>1</sup>:

$$\mathbf{z}_{\alpha+1} = \arg \inf \mathcal{F}_{\mathbf{z}_\alpha}(\nabla \mathbf{u}, t)$$

The `Element` class implements:

- `computeEnergy`( $\mathbf{U}_e^h, t$ )  $\rightarrow I_e$   
where  $\mathbf{U}_e^h = \{\mathbf{u}_e^1, \dots, \mathbf{u}_e^n\}$
- `computeForces`( $\mathbf{U}_e^h, t$ )  $\rightarrow \{\mathbf{F}_{\text{int},e}^1, \dots, \mathbf{F}_{\text{int},e}^n\}$
- `computeStiffnessMatrix`( $\mathbf{U}_e^h, t$ )  $\rightarrow (T_{\text{int},e})_{ik}^{ab}$

Note that the element has a `MaterialModel`, which is used to compute  $W$ ,  $P_{iJ}$ , and  $\mathbb{C}_{iJkL}$  from  $\nabla \mathbf{u}$ ,  $\mathbf{z}_\alpha$ , and time  $t$ .

The element stores and updates<sup>1</sup>  $\mathbf{z}$ :

- `updateInternalVariables`( $\mathbf{U}_e^h, t$ ):  
update  $\mathbf{z}_\alpha \leftarrow \mathbf{z}_{\alpha+1}$

by calling the `MaterialModel`.

<sup>1</sup>Internal variables will be discussed later in the context of inelasticity. They are included here for completeness.



## assembler:

The assembly procedure calculates

$$\begin{aligned} \bullet \quad I(\mathbf{U}^h) &= \sum_{e=1}^{n_e} I_e(\mathbf{U}_e^h) \\ \bullet \quad \mathbf{F}_{\text{int}}(\mathbf{U}^h) &= \sum_{e=1}^{n_e} \mathbf{F}_{\text{int},e}(\mathbf{U}_e^h) \\ \bullet \quad \mathbf{T}_{\text{int}}(\mathbf{U}^h) &= \sum_{e=1}^{n_e} \mathbf{T}_{\text{int},e}(\mathbf{U}_e^h) \end{aligned}$$

These are the *global* quantities derived from *local* element quantities.

The **Assembler** class implements:

- `assembleEnergy( $\mathbf{U}^h, t$ )`  $\rightarrow I$   
where  $\mathbf{U}^h = \{\mathbf{u}^1, \dots, \mathbf{u}^n\}$
- `assembleForces( $\mathbf{U}^h, t$ )`  $\rightarrow \mathbf{F}_{\text{int}} - \mathbf{F}_{\text{ext}}$
- `assembleStiffnessMatrix( $\mathbf{U}^h, t$ )`  $\rightarrow \mathbf{T}_{\text{int}}$

The assembler calls all elements to request their contributions, then assembles those into the global matrices.

The assembler also implements<sup>1</sup>

- `updateInternalVariables( $\mathbf{U}^h, t$ )`

which asks each element to update its internal variables (usually called at the end of a converged solver step; this will be discussed later).

Finally, note that material models and elements use **vector notation** for all stiffness matrices.

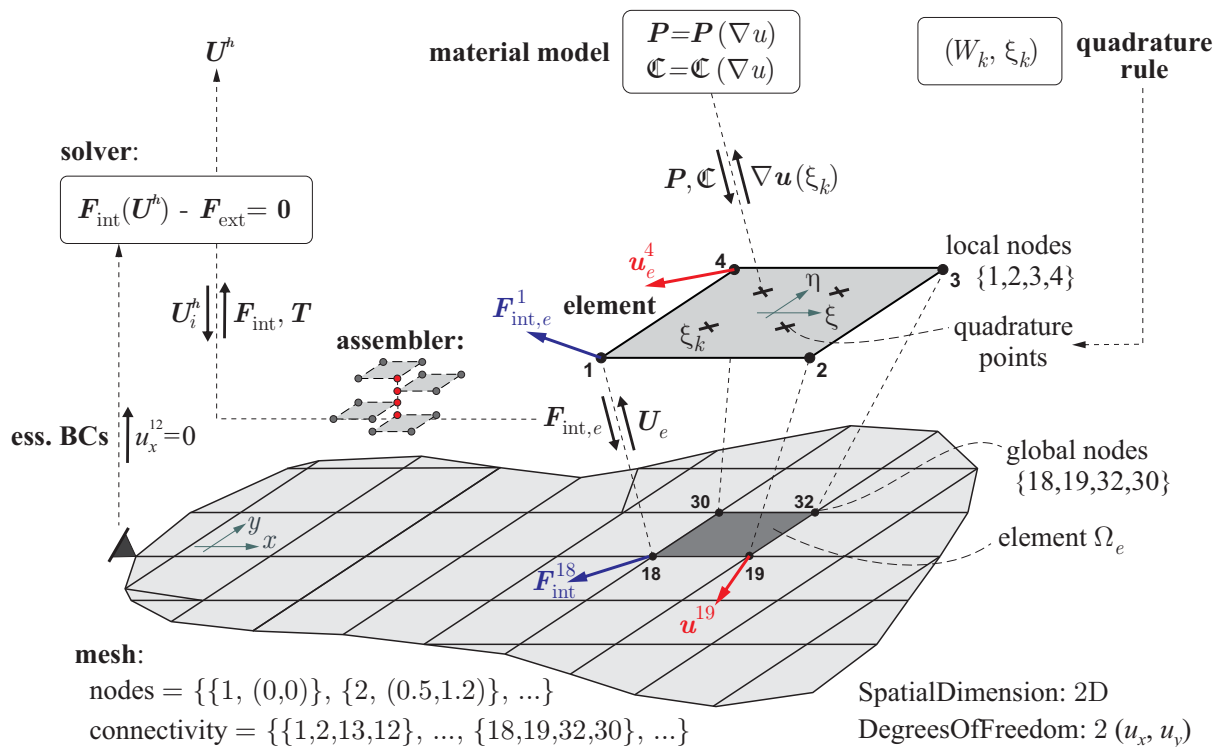


Figure 1: Illustration of the overall code structure.

## 15 Iterative solvers

In the linear elastic or thermal problem, the solution is simple to obtain from a linear system of equations, as discussed before. In case of nonlinear problems, an iterative solution method is required. Here, we discuss a few common examples.

The problem to be solved has the general form

$$\mathbf{f}(\mathbf{U}^h) = \mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}} = \mathbf{0}. \quad (15.1)$$

All iterative solvers start with an initial guess  $\mathbf{U}_0^h$ , which is then corrected in a multitude of ways to find

$$\mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \Delta \mathbf{U}_n^h. \quad (15.2)$$

The iterative scheme converges if  $\Delta \mathbf{U}_n^h \rightarrow 0$  as  $n \rightarrow \infty$ , or equivalently  $\mathbf{f}(\mathbf{U}_n^h) \rightarrow \mathbf{0}$  as  $n \rightarrow \infty$ .

### 15.1 Newton-Raphson (NR) method

The **Newton-Raphson method** (introduced by Newton in 1669 and generalized by Raphson in 1690) starts with a Taylor expansion:

$$\begin{aligned} 0 &= \mathbf{f}(\mathbf{U}_{n+1}^h) = \mathbf{f}(\mathbf{U}_n^h + \Delta \mathbf{U}_n^h) \\ &= \mathbf{f}(\mathbf{U}_n^h) + \frac{\partial \mathbf{f}}{\partial \mathbf{U}}(\mathbf{U}_n^h) \Delta \mathbf{U}_n^h + O(\Delta \mathbf{U}_n^h). \end{aligned} \quad (15.3)$$

If we neglect higher-order terms, then the above can be solved for the increment

$$\Delta \mathbf{U}_n^h = -[\mathbf{T}(\mathbf{U}_n^h)]^{-1} \mathbf{f}(\mathbf{U}_n^h) \quad \text{with} \quad \mathbf{T}(\mathbf{U}_n^h) = \frac{\partial \mathbf{f}}{\partial \mathbf{U}}(\mathbf{U}_n^h) \quad (15.4)$$

being the **tangent matrix**.

For the mechanical problem, this gives (e.g., in finite deformations)

$$\begin{aligned} T_{ik}^{ab}(\mathbf{U}_n^h) &= \frac{\partial F_i^a}{\partial U_k^b}(\mathbf{U}_n^h) = \frac{\partial}{\partial U_k^b} \left( \int_{\Omega} P_{iJ} N_{,J}^a dV - F_{\text{ext},i}^a \right) \\ &= \int_{\Omega} \frac{\partial P_{iJ}}{\partial F_{kL}} N_{,L}^b N_{,J}^a dV - \frac{\partial F_{\text{ext},i}^a}{\partial U_k^b} \\ &= \int_{\Omega} \mathbb{C}_{iJkL} N_{,L}^b N_{,J}^a dV - \frac{\partial F_{\text{ext},i}^a}{\partial U_k^b}, \end{aligned} \quad (15.5)$$

where  $\mathbb{C}_{iJkL}$  is the incremental stiffness tensor (in linearized kinematics, the expression is the same with  $\mathbb{C}_{ijkl}$  the linearized stiffness tensor).

Note that the solver requires that  $\det \mathbf{T} \neq 0$ , which is guaranteed in linear elasticity if no rigid body mode exists (i.e., the linearized system has no zero-energy mode so that  $\mathbf{U}^h \cdot \mathbf{T} \mathbf{U}^h \neq 0$  for all admissible  $\mathbf{U}^h$ ). The Newton-Raphson solver displays *quadratic convergence*.

Note that, if the problem is linear as in linear elasticity, the solver converges in one step since

$$\begin{aligned} \mathbf{U}_{n+1}^h &= \mathbf{U}_n^h + \Delta \mathbf{U}_n^h = \mathbf{U}_n^h - \mathbf{K}^{-1} [\mathbf{F}_{\text{int}} - \mathbf{F}_{\text{ext}}] \\ &= \mathbf{U}_n^h - \mathbf{K}^{-1} [\mathbf{K} \mathbf{U}_n^h - \mathbf{F}_{\text{ext}}] \\ &= \mathbf{U}_n^h - \mathbf{U}_n^h + \mathbf{K}^{-1} \mathbf{F}_{\text{ext}} \\ &= \mathbf{K}^{-1} \mathbf{F}_{\text{ext}}. \end{aligned} \quad (15.6)$$

## 15.2 Damped Newton-Raphson (dNR) method

A slight modification of the Newton-Raphson method, the **damped Newton-Raphson method** is beneficial, e.g., when the NR method tends to overshoot (e.g., in case of oscillatory energy landscapes or multiple minima such as in finite-deformation elasticity).

The iterative scheme is identical to the classical NR method except that

$$\mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \alpha \Delta \mathbf{U}_n^h \quad \text{with} \quad \alpha \in (0, 1). \quad (15.7)$$

The damping parameter  $\alpha$  can be chosen constant or adjusted based on convergence.

## 15.3 Quasi-Newton (QN) method

The **Quasi-Newton method** is the same as the classical NR method with the exception that one does not use the actual tangent matrix  $\mathbf{T}$  for computational simplicity or efficiency.

Motivation is thus to avoid the computation of  $\mathbf{T}(\mathbf{U}^h)$  and its inversion at each iteration step. Instead one uses a matrix  $\mathbf{B}_n$  and updates its inverse directly.

The general algorithm is as follows:

- (1) start with an initial guess  $\mathbf{U}_0^h$  and  $\mathbf{B}_0 = \mathbf{T}(\mathbf{U}_0^h)$
- (2) compute  $\Delta \mathbf{U}_n^h = -\mathbf{B}_n^{-1} \mathbf{f}(\mathbf{U}_n^h)$  and  $\mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \Delta \mathbf{U}_n^h$  and

$$\mathbf{B}_{n+1}^{-1} = \mathbf{B}_n^{-1} - \frac{(\mathbf{B}_n^{-1} \mathbf{z}_n - \Delta \mathbf{U}_n^h) \otimes \Delta \mathbf{U}_n^h \mathbf{B}_n^{-1}}{\Delta \mathbf{U}_n^h \cdot \mathbf{B}_n^{-1} \mathbf{z}_n} \quad \text{with} \quad \mathbf{z}_n = \mathbf{f}(\mathbf{U}_{n+1}^h) - \mathbf{f}(\mathbf{U}_n^h). \quad (15.8)$$

We omit the full derivation of the update for  $\mathbf{B}_{n+1}$  here for brevity. The idea is that  $\mathbf{B}_{n+1}^{-1}$  and  $\mathbf{B}_n^{-1}$  are approximately rank-one-connected using the *Sherman-Morrison formula*. The added benefit is that not only does  $\mathbf{T}$  not have to be recomputed exactly but also can the inversion or linear solver be skipped since the updated inverse is computed explicitly.

## 15.4 Line search method

The **line search method** can be used as an improvement for other nonlinear iterative solvers. Similar to the Quasi-Newton schemes, updates are made according to

$$\mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \beta \Delta \mathbf{U}_n^h, \quad (15.9)$$

where now  $\beta$  is not a constant but chosen such that  $\mathbf{f}(\mathbf{U}_{n+1}^h) = \mathbf{0}$ . For example, we can find  $\beta$  from solving

$$\Delta \mathbf{U}_n^h \cdot \mathbf{f}(\mathbf{U}_n^h + \beta \Delta \mathbf{U}_n^h) = 0. \quad (15.10)$$

This is generally a nonlinear but *scalar* problem that can be solved by bisection, regula falsi, secant, and other methods.

Notice that (15.10) is in fact the stationarity condition of the minimization problem

$$\beta = \arg \inf \left\| \mathbf{f}(\mathbf{U}_n^h + \beta \Delta \mathbf{U}_n^h) \right\|^2, \quad (15.11)$$

which is the motivation for the *nonlinear least-squares method* described below.

## 15.5 Gradient flow method

Although not with a proper physical meaning, the **gradient flow method** (also known as *gradient descent*) has become popular as an iterative solver for quasistatic problems.

The idea is to replace the equation

$$\mathbf{0} = \mathbf{f}(\mathbf{U}_{n+1}^h) \quad (15.12)$$

by a dynamic evolution equation:

$$\mathbf{C} \dot{\mathbf{U}}_{n+1/2}^h = -\mathbf{f}(\mathbf{U}_n^h) \quad \text{and} \quad \mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \Delta t \dot{\mathbf{U}}_{n+1/2}^h \quad (15.13)$$

with, e.g.,  $\mathbf{C} = c\mathbf{I}$  with  $c > 0$ . It is obvious that as  $\mathbf{f} \rightarrow \mathbf{0}$  we have  $\dot{\mathbf{U}}_{n+1/2}^h \rightarrow \mathbf{0}$  and thus the method converges. Although there is no guarantee to reach an extremum, the method is popular because it does not require a tangent matrix and is quite robust.

For example, using a simple backward-Euler discretization for the time derivative and  $\mathbf{C} = c\mathbf{I}$ , we obtain

$$\mathbf{U}_{n+1}^h = \mathbf{U}_n^h - \frac{1}{c} \mathbf{f}(\mathbf{U}_n^h). \quad (15.14)$$

## 15.6 Nonlinear Least Squares

The family of methods based on **nonlinear least squares** aim to minimize

$$r(\mathbf{U}^h) = \|\mathbf{f}(\mathbf{U}^h)\|^2 = \mathbf{f}(\mathbf{U}^h) \cdot \mathbf{f}(\mathbf{U}^h) \quad \text{so we solve} \quad \frac{\partial r}{\partial \mathbf{U}^h} = \mathbf{0}. \quad (15.15)$$

This approach is helpful, e.g., in case of over-constrained systems. Application of Newton-Raphson to this nonlinear system of equations leads to

$$\begin{aligned} \Delta \mathbf{U}^h &= - \left[ \frac{\partial}{\partial \mathbf{U}^h} \frac{\partial r}{\partial \mathbf{U}^h} \right]_{\mathbf{U}_n^h}^{-1} \frac{\partial r}{\partial \mathbf{U}^h}(\mathbf{U}_n^h) \\ &= - \left[ \mathbf{T}^T(\mathbf{U}_n^h) \frac{\partial \mathbf{f}}{\partial \mathbf{U}^h}(\mathbf{U}_n^h) \mathbf{f}(\mathbf{U}_n^h) + \frac{\partial \mathbf{T}^T}{\partial \mathbf{U}^h}(\mathbf{U}_n^h) \mathbf{f}(\mathbf{U}_n^h) \right]^{-1} \mathbf{T}^T(\mathbf{U}_n^h) \mathbf{f}(\mathbf{U}_n^h). \end{aligned} \quad (15.16)$$

If updates are small, we can neglect the second term in brackets (which requires higher derivatives than what is commonly computed in FEM), which gives

$$\Delta \mathbf{U}^h = - [\mathbf{T}^T(\mathbf{U}_n^h) \mathbf{T}(\mathbf{U}_n^h)]^{-1} \mathbf{T}^T(\mathbf{U}_n^h) \mathbf{f}(\mathbf{U}_n^h). \quad (15.17)$$

This is known as the **Gauss-Newton method**. Note that this reduces to Newton-Raphson for standard problems (i.e., those with as many equations as unknowns). However, Gauss-Newton can also be applied to overdetermined systems (i.e., more equations than unknowns).

## 15.7 Conjugate Gradient (CG) method

The **conjugate gradient** method follows the idea of iterating into the direction of steepest descent in order to minimize the total potential energy (as a variation, it can also be applied to the nonlinear least squares problem).

Here, the update is

$$\mathbf{U}_{n+1}^h = \mathbf{U}_n^h + \alpha_n \mathbf{S}_n, \quad (15.18)$$

where both the direction  $\mathbf{S}_n$  and increment  $\alpha_n$  are determined in an optimal way as follows.

The conjugate direction is updated according to

$$\mathbf{S}_n = -\mathbf{f}(\mathbf{U}_n^h) + \beta_n \mathbf{S}_{n-1} \quad (15.19)$$

with  $\beta$  computed from the current solution  $\mathbf{U}_n^h$  and the previous solution  $\mathbf{U}_n^h$  according to one of several options (Polak-Ribière, Fletcher-Reeves, etc.)

Then, the scalar increment  $\alpha_n$  is obtained from a *line search* to find

$$\alpha_n = \arg \min r(\mathbf{U}_n^h + \alpha \mathbf{S}_n). \quad (15.20)$$

A benefit of the conjugate gradient technique is that, as for gradient flow, *no tangent matrix is required*.

A variation of this scheme, originally developed for atomistics but also applicable to the FE method (and oftentimes converging faster than CG) is the so-called **Fast Inertial Relaxation Engine** (FIRE).

## 16 Boundary conditions

### 16.1 Neumann boundary conditions

Recall that the nonlinear system of equations to solve,

$$\mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}} = \mathbf{0}, \quad (16.1)$$

requires computation of the external force vector. The latter includes both *body forces* and *surface tractions*.

Now that we have introduced isoparametric mappings and numerical quadrature rules, we can apply those to the external force terms.

Body forces  $\rho \mathbf{b}$  produce the external force on node  $a$  in direction  $i$ , e.g., in 2D for a single element

$$F_{\text{ext},i}^a = \int_{\Omega_e} \rho b_i N^a dV = \int_{-1}^1 \int_{-1}^1 \rho b_i(\boldsymbol{\xi}) N_e^a(\boldsymbol{\xi}) J(\boldsymbol{\xi}) d\xi d\eta \approx \sum_{k=0}^{n_{QP}-1} W_k \rho b_i(\boldsymbol{\xi}_k) N_e^a(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k). \quad (16.2)$$

Surface tractions  $\hat{\mathbf{t}}$  result in an external force on node  $a$  in direction  $i$ , again in 2D and for a single element,

$$F_{\text{ext},i}^a = \int_{\partial\Omega_{N,e}} \hat{t}_i N^a dV = \int_{-1}^1 \hat{t}_i(\boldsymbol{\xi}) N_e^a(\boldsymbol{\xi}) J(\boldsymbol{\xi}) d\xi d\eta \approx \sum_{k=0}^{n_{QP}-1} W_k \hat{t}_i(\boldsymbol{\xi}_k) N_e^a(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k). \quad (16.3)$$

Note that the surface traction term integrates over the boundary of the element, so in  $d$  dimensions we can use a quadrature rule for  $d - 1$  dimensions (e.g., for a 2D element we use 1D quadrature on the element edges).

#### Implementation:

The variational formulation allows us to implement **external force elements** just like regular finite elements: we have defined their energy,

$$I_e = - \int_{\Omega_e} \rho \mathbf{b} \cdot \mathbf{u} dV - \int_{\Omega_e} \hat{\mathbf{t}} \cdot \mathbf{u} dS. \quad (16.4)$$

The corresponding force vectors  $\mathbf{F}_{\text{ext}}$  are given above, and the tangent matrices are obtained from

$$\mathbf{T} = \frac{\partial \mathbf{F}_{\text{ext}}}{\partial \mathbf{U}_e^h}. \quad (16.5)$$

Notice that if  $\hat{\mathbf{t}}$  and  $\rho \mathbf{b}$  are independent of displacements (as in most cases), we have  $\mathbf{T} = \mathbf{0}$ . An exception is, e.g., the application of *pressure* to the boundary *in finite deformations* (since the resulting force depends on the surface area which, in turn, depends on the sought deformation).

## 16.2 Examples of external forces

### Constant force:

Consider a constant force  $\mathbf{P}$  applied to a particular node  $i$  at deformed position  $\mathbf{x}_i$ . The element energy of this external force vector is

$$I_e = -\mathbf{P} \cdot \mathbf{u}^i, \quad (16.6)$$

so that the resulting force vector is

$$\mathbf{F}^a = \frac{\partial I_e}{\partial \mathbf{u}^a} = -\mathbf{P} \delta_{ia}, \quad (16.7)$$

i.e., an external force is applied only to node  $i$ . The stiffness matrix vanishes since

$$\mathbf{T}^{ab} = \frac{\partial \mathbf{F}^a}{\partial \mathbf{u}^b} = \mathbf{0}. \quad (16.8)$$

### Linear spring:

Next, consider a linear elastic spring (stiffness  $k$ ) attached to a node  $i$  (deformed position  $\mathbf{x}_i$ , undeformed position  $\mathbf{X}_i$ ) and anchored at a constant position  $\mathbf{x}_0 = \mathbf{X}_0$ . The element energy in this case is simply the spring energy

$$I_e = \frac{k}{2} \left( \|\mathbf{x}_i - \mathbf{x}_0\| - \|\mathbf{X}_i - \mathbf{X}_0\| \right)^2, \quad (16.9)$$

and force vector and stiffness matrix follow by differentiation.

### Indenter:

When simulating indentation tests, it is oftentimes convenient to apply the indenter forces via an external potential rather than by modeling contact. Consider a spherical (in 3D) or circular (in 2D) indenter of radius  $R$  whose center is located at the constant, known point  $\mathbf{x}_0$ . Here, one may use potentials of the type

$$I_e = C \left[ \|\mathbf{x}_0 - \mathbf{x}_i\| - R \right]_-^n \quad (16.10)$$

with a force constant  $C > 0$  and integer exponent  $n \geq 2$  (a common choice is  $n = 3$ ). The bracket  $[\cdot]_- = \min(\cdot, 0)$  implies that the energy is only non-zero if the term in bracket is negative (i.e., if the point enters the indenter radius. Again, forces and stiffness matrix follow by differentiation.

### Pressure in linearized kinematics:

Applying a constant pressure  $p$  to a surface can be accomplished rather easily in linearized kinematics via (16.3). Specifically, we have for a single element surface

$$I_e = - \int_{\partial\Omega_e} \hat{\mathbf{t}} \cdot \mathbf{u}^h dS = - \int_{\partial\Omega_e} (-p) \mathbf{n} \cdot \sum_{a=1}^n \mathbf{u}_e^a N_e^a dS = \sum_{a=1}^n \int_{\partial\Omega_e} p \mathbf{n} \cdot \mathbf{u}_e^a N_e^a dS. \quad (16.11)$$

Integration can be carried out using numerical quadrature on the element boundary (the element boundary normal  $\mathbf{n}$  can be computed from the nodal locations). Again, forces and stiffness matrix follow by differentiation (forces are constant, and the stiffness matrix vanishes).

### Pressure in finite kinematics:

Applying a constant pressure  $p$  to a surface in finite kinematics is more complicated as the element boundary undergoes finite changes during deformation, resulting in nodal forces that depend on deformation. Here, we start with the work done by the pressure (which we assume *constant*), viz.

$$\begin{aligned} I_e &= -pv = - \int_{\varphi(\Omega_e)} p dv = - \int_{\varphi(\Omega_e)} p \frac{1}{d} x_{i,i} dv \\ &= -\frac{p}{d} \int_{\partial\varphi(\Omega_e)} \varphi_i n_i ds = -\frac{p}{d} \int_{\partial\Omega_e} \varphi_i J F_{J_i}^{-1} N_J dS, \end{aligned} \quad (16.12)$$

where we used that  $x_{i,i} = \delta_{ii} = d$  in  $d$  dimensions as well as the Piola transform  $n_i ds = J F_{J_i}^{-1} N_J dS$ . Thus the numerical integration uses

$$I_e = -\frac{p}{d} \int_{\partial\Omega_e} \mathbf{u}^h \cdot \mathbf{F}^{-T} \mathbf{N} J dS, \quad (16.13)$$

which can be evaluated by numerical quadrature. As before, forces and stiffness matrix follow by differentiation but are non-zero and rather complex.

### 16.3 Dirichlet boundary conditions

Essential boundary conditions require us to replace individual equations in the nonlinear system by  $\mathbf{u}^a = \hat{\mathbf{u}}^a$ . This is accomplished computationally usually in one of two ways.

#### Substitution:

First, brute-force substitution is a simple method: one replaces the respective equation for  $\Delta u_i^a$  in the linearized system  $\mathbf{T} \Delta \mathbf{U}^h = \mathbf{F}$  by  $\Delta u_i^a = \Delta \hat{u}_i^a$ ; e.g.,

$$\begin{pmatrix} T_{11} & T_{12} & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & T_{55} \end{pmatrix} \begin{pmatrix} \Delta u^1 \\ \Delta u^2 \\ \Delta u^3 \\ \Delta u^4 \\ \Delta u^5 \end{pmatrix} = \begin{pmatrix} F^1 \\ F^2 \\ \Delta \hat{u}^3 \\ F^4 \\ F^5 \end{pmatrix} \quad (16.14)$$

For example, when using iterative solvers, one may want to choose the initial guess as  $\hat{\mathbf{u}}^a$  and subsequently enforce  $\Delta \mathbf{u}^a = \mathbf{0}$  in iterations by using the above condensation method.

#### Constraints:

The same method discussed above can also impose other types of boundary conditions, e.g., *periodic boundary conditions* of the type

$$\mathbf{u}^+ = \mathbf{u}^- \quad (16.15)$$

for some opposite nodes  $(+, -)$ . Also, constraints of the general type

$$f(\mathbf{u}^i, \mathbf{u}^j, \dots) = 0 \quad (16.16)$$

can be implemented in a similar fashion (e.g., rigid links between nodes).



In case of linear constraints on boundary nodes, it is sometimes convenient to introduce a *reduced* set of degrees of freedom  $\mathbf{U}_{\text{red}}^h$ , which contains only the free degrees of freedom, from which the full set of nodal degrees of freedom can be obtained via

$$\mathbf{U}^h = \mathbf{A} \mathbf{U}_{\text{red}}^h, \quad (16.17)$$

where  $\mathbf{A}$  is a matrix defined by the constraints. For example, consider a node  $a$  being constrained to move along a particular direction  $\mathbf{n} \in \mathbb{R}^d$ . That is, its degrees of freedom  $\mathbf{u}^a = \{u_1^a, u_2^a\}$  in 2D are reduced<sup>2</sup> to  $\mathbf{u}^a = \chi^a \mathbf{n}$  with the only unknown being  $\chi^a \in \mathbb{R}$ . E.g., in 2D we have

$$\mathbf{U}^h = \begin{pmatrix} \dots \\ \cdot & 1 & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & n_1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & n_2 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & n_3 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & 1 & \cdot \\ \dots \end{pmatrix} \begin{pmatrix} \vdots \\ u_1^{a-1} \\ u_2^{a-1} \\ \chi^a \\ u_1^{a+1} \\ u_2^{a+1} \\ \vdots \end{pmatrix} = \mathbf{A} \mathbf{U}_{\text{red}}^h \quad (16.18)$$

The governing equation now becomes

$$\mathbf{F}_{\text{int}}(\mathbf{A} \mathbf{U}_{\text{red}}^h) - \mathbf{F}_{\text{ext}} = \mathbf{0}. \quad (16.19)$$

Note that in the special case of a linear problem we know that  $\mathbf{F}_{\text{int}} = \mathbf{T} \mathbf{U}^h$ , so that (16.19) form can be transformed into a symmetric problem by pre-multiplying by  $\mathbf{A}^T$ , i.e.,

$$\mathbf{A}^T \mathbf{T} \mathbf{A} \mathbf{U}_{\text{red}}^h - \mathbf{A}^T \mathbf{F}_{\text{ext}} = \mathbf{0} \quad \Leftrightarrow \quad \mathbf{T}_{\text{red}} \mathbf{U}_{\text{red}}^h - \mathbf{F}_{\text{ext,red}} = \mathbf{0}, \quad (16.20)$$

which is to be solved for the reduced degrees of freedom. While efficient due to the reduced number of degrees of freedom, the pre- and post-multiplication by  $\mathbf{A}$  can add considerable numerical expenses.

Notice a convenient relation: if matrix  $\mathbf{A}$  embeds  $\mathbf{U}_{\text{red}}^h$  according to (16.17), then the total forces acting on the reduced degrees of freedom are given by  $\mathbf{A}^T \mathbf{F}$ ; e.g., in the above example:

$$\mathbf{F}_{\text{red}} = \mathbf{A}^T \mathbf{F} = \begin{pmatrix} \dots \\ \cdot & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & n_1 & n_2 & n_3 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 1 & \cdot & \cdot \\ \dots \end{pmatrix} \mathbf{F} = \begin{pmatrix} \vdots \\ \mathbf{F}^{a-1} \\ \mathbf{F}^a \cdot \mathbf{n} \\ \mathbf{F}^{a+1} \\ \vdots \end{pmatrix}. \quad (16.21)$$

Recall that in the nonlinear case, one may still solve for the reduced the linearized equation to compute the Newton-Raphson updates is

$$\Delta \mathbf{U}_{\text{red}}^h = -(\mathbf{T}^*)^{-1}(\mathbf{U}_{\text{red}}^h) [\mathbf{F}_{\text{int}}(\mathbf{A} \mathbf{U}_{\text{red}}^h) - \mathbf{F}_{\text{ext}}] \quad \text{where} \quad \mathbf{T}^* = \mathbf{A}^T \frac{\partial \mathbf{F}_{\text{int}}}{\partial \mathbf{U}^h}(\mathbf{A} \mathbf{U}_{\text{red}}^h). \quad (16.22)$$

There are many alternative ways to introduce constraints such as using the penalty method which appends to the potential energy each constraint with a penalty (Lagrange-type) multiplier. For further information, please refer to the literature.

<sup>2</sup>This is an *embedding* operation: node  $a$  is constraint to lie on a line manifold and is embedded into 3D space via embedding. The inverse operation is known as *submerging*.

## Condensation:

The substitution method introduced above for essential boundary conditions is simple to implement. However, it is quite expensive since the number of equations remains the same when imposing essential boundary conditions. The **condensation method** removes from the linear system to be solved those equations imposing essential boundary conditions.

Let us rewrite the linearized system by moving the third column to the right-hand side:

$$\begin{pmatrix} T_{11} & T_{12} & 0 & T_{14} & \cdot \\ \cdot & \cdot & 0 & \cdot & \cdot \\ \cdot & \cdot & 0 & \cdot & \cdot \\ \cdot & \cdot & 0 & \cdot & \cdot \\ \cdot & \cdot & 0 & \cdot & T_{55} \end{pmatrix} \begin{pmatrix} \Delta u^1 \\ \Delta u^2 \\ \Delta u^3 \\ \Delta u^4 \\ \Delta u^5 \end{pmatrix} = \begin{pmatrix} F^1 \\ F^2 \\ F^3 \\ F^4 \\ F^5 \end{pmatrix} - \begin{pmatrix} \Delta T_{13} \\ \Delta T_{23} \\ \Delta T_{33} \\ \Delta T_{43} \\ \Delta T_{53} \end{pmatrix} \Delta \hat{u}^3 \quad (16.23)$$

so that we can eliminate the third row and column from the system:

$$\begin{pmatrix} T_{11} & T_{12} & T_{14} & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & T_{55} \end{pmatrix} \begin{pmatrix} \Delta u^1 \\ \Delta u^2 \\ \Delta u^4 \\ \Delta u^5 \end{pmatrix} = \begin{pmatrix} F^1 \\ F^2 \\ F^4 \\ F^5 \end{pmatrix} - \begin{pmatrix} \Delta T_{13} \\ \Delta T_{23} \\ \Delta T_{43} \\ \Delta T_{53} \end{pmatrix} \Delta \hat{u}^3, \quad (16.24)$$

and solve for the remaining unknowns along with  $\Delta u^3 = \Delta \hat{u}^3$ .

The clear advantage of this method is the *reduction in size* of the system to be solved. The disadvantage is that it is computationally more involved (can be even more expensive for a small number of essential boundary conditions applied to large systems).

## 16.4 Rigid body motion

Consider the composite deformation mapping  $\mathbf{x} = \varphi^*(\varphi(\mathbf{X}))$  with  $\mathbf{x} = \varphi(\mathbf{X})$  being an admissible deformation mapping and  $\varphi^*(\mathbf{x}) = \mathbf{R}\mathbf{x} + \mathbf{c}$  denoting rigid body motion, i.e.,  $\mathbf{R} \in SO(d)$  and  $\mathbf{c} \in \mathbb{R}^d$  is a constant vector. Note that the combined deformation gradient is given by  $\mathbf{F}^* = \mathbf{R}\mathbf{F}$  where  $\mathbf{F} = \text{Grad } \varphi$ . Recall that the weak form read

$$\mathcal{G}(u, v) = \int_{\Omega} P_{iJ}(\mathbf{F}) v_{i,J} dV - \int_{\Omega} R B_i v_i dV - \int_{\partial\Omega_N} \hat{T}_i v_i dS = 0. \quad (16.25)$$

Insertion of the composite mapping yields

$$\mathcal{G}(u^*, v) = \int_{\Omega} P_{iJ}(\mathbf{R}\mathbf{F}) v_{i,J} dV - \int_{\Omega} R B_i v_i dV - \int_{\partial\Omega_N} \hat{T}_i v_i dS = 0. \quad (16.26)$$

However, note that material frame indifference requires that  $\mathbf{P} = \mathbf{P}(\mathbf{C})$  and  $\mathbf{C}^* = (\mathbf{F}^*)^T \mathbf{F}^* = \mathbf{F}^T \mathbf{F} = \mathbf{C}$ , so that the rotation has no affect on the weak form (neither does the translation  $\mathbf{c}$ ). Therefore, rigid body motion can be superimposed onto any admissible solution and must be suppressed by appropriate essential boundary conditions to ensure uniqueness of solutions.

For the linear elastic case, the tangent matrix  $\mathbf{T}$  therefore has as many zero eigenvalues as it has **rigid-body modes**, or **zero-energy modes**  $\mathbf{U}^*$  such that

$$\mathbf{U}^* \cdot \mathbf{T} \mathbf{U}^* = 0. \quad (16.27)$$

This implies that  $\mathbf{T}$  has zero eigenvalues and is thus not invertible.

The remedy is to suppress rigid body modes via appropriate essential boundary conditions; specifically in  $d$  dimensions we need  $d(d-1)/2$  such essential boundary conditions.

## 17 Error estimates and adaptivity

### 17.1 Finite element error analysis

Solving systems of PDEs by the finite element method introduces numerous sources of errors that one should be aware of:

- (i) The **discretization error** (also known as the **first fundamental error**) arises from discretizing the domain into elements of finite size  $h$ . As a result, the body  $\Omega$  is not represented correctly and the model (e.g., the outer boundary) may not match the true boundary  $\partial\Omega$  (e.g., think of approximating a circular domain  $\Omega$  by CST or Q4 elements). This error can be reduced by mesh **refinement** (and we discussed  $r$ -refinement,  $h$ -refinement,  $p$ -refinement, and  $hp$ -refinement).
- (ii) The **numerical integration error** results from the application of numerical quadrature:

$$\int_{\Omega_e} f(\xi) d\xi \cong \sum_{k=1}^{n_{QP}} W_k f(\xi_k) \quad (17.1)$$

We discussed that for  $f \in C^{k+1}(\Omega)$  (the extension to higher dimensions is analogous)

$$\left| \int_{-1}^1 f(\xi) d\xi - \sum_{q=1}^{n_{QP}} W_i f(\xi_i) \right| \leq C \|\Omega\| h^{k+1} \max_{\substack{\xi \in [-1,1] \\ |\alpha|=k+1}} \|D^\alpha f(\xi)\|. \quad (17.2)$$

Hence, the numerical integration error depends on the smoothness of the integrand and calls for a proper choice of the integration order.

- (iii) The **solution error** stems from numerically solving linear systems  $\mathbf{T}\mathbf{U}^h = \mathbf{F}$ . In general, the accuracy of the solution depends on the **condition number** of the matrix,

$$\kappa = \|\mathbf{T}\| \cdot \|\mathbf{T}^{-1}\| = \left| \frac{\lambda_{\max}}{\lambda_{\min}} \right| \quad (17.3)$$

with  $\lambda_{\max}$  ( $\lambda_{\min}$ ) being the largest (smallest) eigenvalue of  $\mathbf{T}$ . The higher the condition number, the larger the numerical error.

A practical consequence is the guideline to choose wisely the units of model parameters (such as material constants, domain size features, etc.). For example, when performing a linear elastic simulation, it is advisable to normalize elastic constants by 1 GPa instead of assigning, e.g.,  $E = 210 \cdot 10^9$  (instead, use  $E = 2.1$  and know that your results will be in 100 GPa's).

- (iv) An **approximation error** is introduced by approximating the functional space  $\mathcal{U}$  (in which to find the solution  $u(x)$ ) by a finite-dimensional subspace  $\mathcal{U}^h \subset \mathcal{U}$ .

As an example, consider an exact solution  $u(x)$  in 1D which is approximated by a piecewise linear polynomial function  $u^h(x)$ . The bar hanging from the ceiling under its own weight was such an example, for which the solution was found to be exact at the nodes with an error  $e(x) = u(x) - u^h(x)$  arising within elements. Therefore, we can find a point  $z$  within each element such that

$$\frac{\partial e}{\partial x}(z) = 0 \quad \text{for} \quad x_i \leq x \leq x_{i+1}. \quad (17.4)$$

Consequently, we can expand the error to find the solution at a node as

$$e(x_i) = 0 = e(z) + (x_i - z) \frac{\partial e}{\partial x}(z) + \frac{(x_i - z)^2}{2} \frac{\partial^2 e}{\partial x^2}(z) + O((x_i - z)^3). \quad (17.5)$$

Using (17.4), we find that

$$e(z) = -\frac{(x_i - z)^2}{2} \frac{\partial^2 e}{\partial x^2}(z) + O(h^3). \quad (17.6)$$

Note that

$$(x_i - z)^2 \leq \left( \frac{x_{i+1} - x_i}{2} \right)^2 = \frac{h^2}{4}, \quad (17.7)$$

where  $h$  denotes the nodal spacing. Altogether, we have thus shown that the maximum error in an element is bounded by

$$|e(x)|_{\max} \leq \frac{h^2}{8} \max_{x_i \leq x \leq x_{i+1}} \left| \frac{\partial^2 u}{\partial x^2} \right| \quad (17.8)$$

As shown in Appendix C, the above error bound can be significantly generalized. for an interpolation of order  $k$  and  $u \in H^{k+1}(\Omega)$ , we have

$$|u_h - u|_{H^1(\Omega)} \leq \frac{h^k}{\pi^k} |u|_{H^{k+1}(\Omega)} \quad \text{and} \quad \|u_h - u\|_{H^1(\Omega)} \leq c h^k |u|_{H^{k+1}(\Omega)}, \quad (17.9)$$

using Sobolev norms. Thus the error is again determined by the smoothness of the function to be interpolated; and it is expected to decrease with decreasing element size (as  $h \rightarrow 0$ ) – the faster the higher the interpolation order.

Note that special caution is required if stress concentrations of any kind are to be represented (e.g., imagine a linear elastic fracture problem and the issues arising from using linear elements to capture the  $1/r^n$ -type stress concentration near the crack tip).

- (v) A **truncation error** is made by every computer when storing and operating numeric values with only a finite number of digits (e.g., floats, doubles, etc.). This is unavoidable and one should be aware of what this error is (especially when choosing, e.g., solver tolerances). Choosing a solver tolerance in itself produces truncation error because we contend with a solution  $\mathbf{U}^h$  that satisfies  $\mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}} = \text{tol.}$  (instead of being zero).
- (vi) Finally, no simulation is free of **modeling error**, which refers to the large collection of errors made by the selection of the analytical model to be solved (before starting any numerical approximation). For example, we make choices for an appropriate material model, choose material parameters, boundary conditions, and geometric simplifications including reductions to lower dimensions (e.g., plane strain or plane stress instead of a 3D simulation).

The sum of all of the above error sources makes up the numerical error inherent in every simulation.

## 17.2 Smoothing and adaptivity

Both for postprocessing of approximate solutions and for mesh adaptivity (discussed below) it is convenient to introduce a **smoothing** scheme that takes piecewise-defined solutions (e.g., the stress and strain fields in case of simplicial elements) and computes smoothed nodal values of the discontinuous quantities.

In case of simplicial elements, the quantities of interest are constant within elements. Here one can define nodal quantities, e.g., as

$$(u^a)^* = \frac{\sum_{j=1}^{n_{\text{nb}}} u_{e,j} / V_{e,j}}{\sum_{j=1}^{n_{\text{nb}}} 1 / V_{e,j}}, \quad (17.10)$$

where the element quantities  $u_e$  of all  $n_{nb}$  neighboring elements meeting at a node  $a$  are weighted by the respective element volume  $V_e$ . This weighting results in smaller elements (whose quadrature points are closer to the node) having a larger weight than large elements (where quadrature points are far from the node).

In case of higher-order elements such as, e.g., the Q4 element, one can extrapolate element quantities that are defined at the quadrature points  $\xi_k$  to the nodes by invoking the definition that

$$u(\xi_k) = \sum_{a=1}^n (u^a)^* N^a(\xi_k). \quad (17.11)$$

For example for the Q4 element, there are four nodes and four quadrature points, so that the four equations can be solved for the four nodal values  $u^a$ . Once the nodal values are known, one again uses a smoothing relation like (17.10) with the element quantities  $u_{e,j}$  replaced by the nodal value  $u_e^a$  from the respective element, and element volume  $V_e$  replaced by the nodal weight (obtained from extrapolating the quadrature point weights  $W_k J_k t$  to the nodes).

When performing adaptive mesh refinement, we need an *error norm* where, as an example, we discuss the **ZZ error estimator** named after its inventors, Zienkiewicz and Zhu. If we use a smoothing scheme like the above, we can define a smoothed, continuous displacement gradient  $\nabla \mathbf{u}^*$  which is contrasted with the approximate solution  $\nabla \mathbf{u}^h$ , so that one may define the error per element by using the energy norm, viz.

$$\|e(\nabla \mathbf{u}, \nabla \mathbf{u}^*)\|_e^2 = \int_{\Omega_e} W(\nabla \mathbf{u}^* - \nabla \mathbf{u}) dV. \quad (17.12)$$

Note that this definition is not unique and one could also use the  $L_2$ -norm

$$\|e\|_e^2 = \int_{\Omega_e} \|\nabla \mathbf{u}^* - \nabla \mathbf{u}\| dV. \quad (17.13)$$

or any other sensible norm.

In order to define an error estimate, it makes sense to introduce the normalization

$$\eta_e = \frac{\|e(\nabla \mathbf{u}, \nabla \mathbf{u}^*)\|_e}{\|e(\nabla \mathbf{u}, \mathbf{0})\|_e}, \quad (17.14)$$

i.e., the error is divided by the energy in the element so as to not over- or underestimate element errors based on element sizes as well as in case of vast differences in element strain energies. The **mesh refinement criterion** states that an element is refined if

$$\eta_e > \eta_{tol}. \quad (17.15)$$

with some tolerance  $\eta_{tol}$ , determined as a compromise between accuracy and efficiency. Based on the refinement criterion, a set of elements can be identified after each equilibration which is flagged for refinement.

Needed next is a **mesh refinement algorithm**. For example, a frequent choice for triangular elements is **longest edge bisection**, which identifies the longest edge in an element to be refined and inserts a new node at this edge's midpoint, followed by an update of the element connectivity (removing two existing elements and creating four new elements). Note that this involves some book-keeping since adjacent elements flagged for refinement interfere, so that one needs an algorithmic decision about which elements to refine first and how to handle adjacent elements identified for refinement, etc.

## 18 Element defects: shear locking and hourglassing

Elements make specific approximations about how displacements, stresses and strains vary inside each element, and these approximations may introduce errors. Aside from the above general numerical errors stemming from the interpolation order, the quadrature rule, and geometric errors, etc., particular elements can have defects, some of which will be discussed in the following.

The **bilinear Q4 element** displays a defect when used to simulate beam bending. Consider a rectangular element having side lengths  $a$  and  $b$  in the  $x$  and  $y$  directions. If bending moments  $M$  are applied on the two vertical edges, then the element is expected to undergo bending.

In an actual (slender) beam, we know from elastic beam theory that the stresses in the beam vary linearly across the thickness, so we may write

$$\sigma_{xx} = -\frac{y}{b}\sigma_{\max}, \quad \sigma_{yy} = \sigma_{xy} = 0, \quad (18.1)$$

where  $\sigma_{\max}$  is the maximum tensile stress reached on the surface of the beam. From Hooke's law we obtain

$$\varepsilon_{xx} = -\frac{\sigma_{xx}}{E} = -\frac{y\sigma_{\max}}{bE}, \quad \varepsilon_{yy} = -\nu\frac{\sigma_{xx}}{E} = \nu\frac{y\sigma_{\max}}{bE}, \quad \varepsilon_{xy} = 0. \quad (18.2)$$

Let us try to find this solution using a Q4 element. To this end, we apply horizontal forces to all four nodes of the element, obtained from lumping the distributed beam stresses to the nodes:

$$\begin{aligned} F_{\text{ext},1} &= \int_{\partial\Omega_e} \sigma_{xx} N_e^1 dS = - \int_{-1}^1 \left( -\frac{y}{b}\sigma_{\max} \right) N_e^1(\xi, \eta) tb d\eta \\ &= \int_{-1}^1 \eta \sigma_{\max} \frac{(1-\xi)(1-\eta)}{4} tb d\eta = -\frac{\sigma_{\max} tb}{3} = -F, \end{aligned} \quad (18.3)$$

where we introduced the abbreviation  $F = \frac{\sigma_{\max} tb}{3}$  for convenience. Analogously, one obtains

$$F_{\text{ext},2} = F, \quad F_{\text{ext},3} = -F, \quad F_{\text{ext},4} = F. \quad (18.4)$$

That is, the bottom edge is stretched, while the top edge is compressed (as in beam bending).

Solving the problem for the case of linear elasticity reduces to a simple linear system of equations:

$$\mathbf{K}_e \mathbf{U}_e^h = \mathbf{F}_{\text{ext}} = \begin{pmatrix} -F \\ 0 \\ F \\ 0 \\ -F \\ 0 \\ F \\ 0 \end{pmatrix} \Rightarrow \mathbf{U}_e^h = \begin{pmatrix} -U \\ 0 \\ U \\ 0 \\ 0 \\ -U \\ 0 \\ U \\ 0 \end{pmatrix} \quad \text{with} \quad U = \frac{2a\sigma_{\max}}{E} \frac{(1+\nu)(1-2\nu)}{2(1-\nu) + (an)^2(1-2\nu)}, \quad (18.5)$$

where we evaluated all element forces and stiffness matrices exactly (without quadrature errors).

The resulting strains are evaluated as

$$\varepsilon_{xx} = -\frac{U}{a}\eta, \quad \varepsilon_{yy} = 0, \quad \varepsilon_{xy} = -\frac{U}{b}\xi. \quad (18.6)$$

While the axial strain component is reasonable, the presence of a shear stress is problematic. This so-called **parasitic shear** is not physical but a numerical artifact (elastic Euler beams produce no shear stresses). This parasitic shear contributes to the energy of the element, so that – when minimizing the potential energy in a boundary value problem – this shear component introduces artificial stiffness. This becomes apparent when calculating, e.g., the angle of bending  $\theta$  of the beam in the above problem. Comparing the exact one (obtained from beam theory) to the approximate one (obtained from the above construction), one finds

$$\frac{\theta_{\text{approx}}}{\theta_{\text{exact}}} = \frac{1 - \nu^2}{1 + \frac{1 - \nu}{2} \left( \frac{a}{b} \right)^2}. \quad (18.7)$$

Notice that, as the Q4 element becomes more and more slender ( $a/b \rightarrow \infty$ ), the numerically obtained angle approaches 0. That is, the element shows what is known as **locking** or, more specifically, **shear locking**: in case of very slender elements ( $a/b \ll 1$ ) the high shear strain and associated elastic energy prevents the element from deforming, it “locks”.

Next, let us use numerical quadrature to evaluate the element quantities of interest. For example, for the energy we have

$$I_e = \int_{\Omega_e} \frac{1}{2} \boldsymbol{\varepsilon}^h \cdot \mathbb{C} \boldsymbol{\varepsilon}^h dV \approx \frac{t}{2} \sum_{k=0}^{n_{QP}-1} W_k \boldsymbol{\varepsilon}^h(\boldsymbol{\xi}_k) \cdot \mathbb{C} \boldsymbol{\varepsilon}^h(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k). \quad (18.8)$$

We showed before that full integration requires  $n_{QP} = 2 \times 2$  quadrature points. If instead we use reduced integration with only a single quadrature point ( $n_{QP} = 1$ ) located at  $\boldsymbol{\xi}_0 = \mathbf{0}$ , then notice that strain fields of the form (18.6) – which vanish at  $\xi = \eta = 0$  – produce no energy. In other words, any strain field of the form (18.6) can appear without causing any energy and therefore also without causing any resistance. This is a **zero-energy mode** of the under-integrated Q4 element and a serious defect. The resulting deformation of elements undergoing alternating strains (18.6) is a numerical artifact and often referred to as **hourglass mode** or **chicken-wire mode** because of its appearance.

Note that not only the underintegrated Q4 element has such a zero-energy mode. For example, the Q8/Q9 elements have a similar zero-energy mode with curved element edges.

Finally, it is possible to use **selective integration**, which applies different quadrature rules for different energy contributions. For example, integrating the Q4 element with a  $2 \times 2$  quadrature rule for the nominal strains  $\varepsilon_{xx}$  and  $\varepsilon_{yy}$  while using reduced integration with a single quadrature point at  $\boldsymbol{\xi}_k = \mathbf{0}$  removes the spurious shear contribution while maintaining the correct stiffness against axial deformation. Notice that, however, selective integration is harder to implement, not directly applicable beyond linear elasticity, and should be used with caution.

## 19 Dynamics

### 19.1 Variational setting

The mechanical problems so far have been limited to quasistatic conditions. Let us consider the extension to dynamic problems where inertial effects matter and the strong form reads (for simplicity stated in linearized kinematics; the finite-deformation setting is analogous)

$$\begin{cases} \sigma_{ij,j} + \rho b_i = \rho a_i & \text{in } \Omega \\ u_i(\mathbf{x}, t) = \hat{u}_i(\mathbf{x}, t) & \text{on } \partial\Omega_D \\ \sigma_{ij}n_j(\mathbf{x}, t) = \hat{t}_i(\mathbf{x}, t) & \text{on } \partial\Omega_N \end{cases} \quad (19.1)$$

Now, we have  $u : \Omega \times \mathbb{R} \rightarrow \mathbb{R}^d$  with sufficient differentiability in both space and time.

The variational approach here makes use of the so-called **action principle** which uses the **action**

$$\mathcal{A}[u] = \int_{t_1}^{t_2} \mathcal{L}[u] dt \quad \text{with} \quad \mathcal{L}[u] = T[u] - I[u], \quad (19.2)$$

where  $I$  is the potential energy functional from before and  $T$  the kinetic energy functional:

$$T[u] = \int_{\Omega} \frac{\rho}{2} |\dot{\mathbf{u}}|^2 dV. \quad (19.3)$$

The **action principle** states that the solution  $\mathbf{u}(\mathbf{x}, t)$  renders  $\mathcal{A}$  stationary with  $\mathbf{u}(\mathbf{x}, t_1) = \mathbf{u}_1(\mathbf{x})$  and  $\mathbf{u}(\mathbf{x}, t_2) = \mathbf{u}_2(\mathbf{x})$ .

For variational material models (with  $W$  replaced by  $W^*$  for inelasticity), we have

$$\mathcal{A}[u] = \int_{t_1}^{t_2} \left[ \int_{\Omega} \left( \frac{\rho}{2} |\dot{\mathbf{u}}|^2 - W(\varepsilon) \right) dV + \int_{\Omega} \rho \mathbf{b} \cdot \mathbf{u} dV + \int_{\partial\Omega_N} \hat{\mathbf{t}} \cdot \mathbf{u} dS \right] dt. \quad (19.4)$$

Taking the first variation (with the divergence theorem and the same assumptions from before):

$$\delta \mathcal{A}[u] = 0 = \int_{t_1}^{t_2} \left[ \int_{\Omega} (\rho \dot{u}_i \delta \dot{u}_i - \sigma_{ij} \delta u_{i,j}) dV + \int_{\Omega} \rho b_i \delta u_i dV + \int_{\partial\Omega_N} \hat{t}_i \delta u_i dS \right] dt. \quad (19.5)$$

The **weak form** is thus obtained as

$$\mathcal{G}(u, v) = \int_{t_1}^{t_2} \left[ \int_{\Omega} (\rho \dot{u}_i \dot{v}_i - \sigma_{ij} v_{i,j}) dV + \int_{\Omega} \rho b_i v_i dV + \int_{\partial\Omega_N} \hat{t}_i v_i dS \right] dt = 0 \quad (19.6)$$

with

$$\mathbf{v} \in \{ \mathbf{v} \in H^1(\Omega) : \mathbf{v} = \mathbf{0} \text{ on } \partial\Omega_D \text{ and at } t = t_1 \text{ or } t = t_2 \}. \quad (19.7)$$

To avoid the time derivative in the variation, let us integrate by parts in time (the “boundary term” vanishes since  $\mathbf{v} = 0$  at  $t = t_1$  and  $t = t_2$ ):

$$\mathcal{G}(u, v) = - \int_{t_1}^{t_2} \left[ \int_{\Omega} (\rho \ddot{u}_i v_i + \sigma_{ij} v_{i,j}) dV - \int_{\Omega} \rho b_i v_i dV - \int_{\partial\Omega_N} \hat{t}_i v_i dS \right] dt = 0. \quad (19.8)$$

Note that, without the first term, we recover the elastostatic formulation.



Since in the dynamic problem the displacement field depends on time, we introduce a **semi-discretization**, i.e., we discretize the solution in space but not in time:

$$\mathbf{u}^h(\mathbf{x}, t) = \sum_{a=1}^n \mathbf{u}^a(t) N^a(\mathbf{x}) \quad \text{and} \quad \mathbf{v}^h(\mathbf{x}, t) = \sum_{a=1}^n \mathbf{v}^a(t) N^a(\mathbf{x}), \quad (19.9)$$

so that

$$\dot{\mathbf{u}}^h(\mathbf{x}, t) = \sum_{a=1}^n \dot{\mathbf{u}}^a(t) N^a(\mathbf{x}) \quad \text{and} \quad \ddot{\mathbf{u}}^h(\mathbf{x}, t) = \sum_{a=1}^n \ddot{\mathbf{u}}^a(t) N^a(\mathbf{x}). \quad (19.10)$$

Insertion into the weak form results in Galerkin's **discrete weak form**:

$$\begin{aligned} \mathcal{G}(\mathbf{u}^h, \mathbf{v}^h) = & - \int_{t_1}^{t_2} \sum_{a=1}^n \sum_{b=1}^n \left[ \ddot{u}_i^a v_i^b \int_{\Omega} \rho N^a N^b dV + v_i^b \int_{\Omega} \sigma_{ij} N_{,j}^b dV \right. \\ & \left. - v_i^b \int_{\Omega} \rho b_i N^b dV - v_i^b \int_{\partial\Omega_N} \hat{t}_i N^b dS \right] dt = 0 \end{aligned} \quad (19.11)$$

for all  $\mathbf{v}^b(t)$  histories that vanish at  $t_1$  and  $t_2$ .

Analogously to before, we now write

$$\mathbf{U}^h(t) = \{\mathbf{u}^1(t), \dots, \mathbf{u}^n(t)\}, \quad (19.12)$$

so that solving (19.11) for all  $\mathbf{v}^b(t)$  is equivalent to solving

$$\mathbf{M} \ddot{\mathbf{U}}^h + \mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}}(t) = \mathbf{0} \quad (19.13)$$

with

$$M_{ij}^{ab} = \delta_{ij} \int_{\Omega} \rho N^a N^b dV, \quad F_{\text{int},i}^b = \int_{\Omega} \sigma_{ij} N_{,j}^b dV, \quad F_{\text{ext},i}^b = \int_{\Omega} \rho b_i N^b dV + \int_{\partial\Omega_N} \hat{t}_i N^b dS. \quad (19.14)$$

Matrix  $\mathbf{M}$  is called the **consistent mass matrix**.

### Examples:

The consistent mass matrix for a *two-node bar element* is computed from shape functions

$$N_1(\xi) = \frac{1-\xi}{2}, \quad N_2(\xi) = \frac{1+\xi}{2}. \quad (19.15)$$

Specifically, we have (with  $m = \rho A L_e$ )

$$M^{ab} = \int_{\Omega} \rho N^a N^b dV = \int_{-1}^1 \rho N^a N^b A \frac{L_e}{2} dxi \quad \Rightarrow \quad \mathbf{M}_{1D} = \frac{m}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}. \quad (19.16)$$

Note that this is the consistent mass matrix for 1D motion. If each node moves has two degrees of freedom ( $u_1, u_2$ ) in the plane, then each pair of dof is linked by the above mass matrix, so that the total consistent mass matrix becomes

$$\mathbf{M}_{2D} = \frac{m}{6} \begin{pmatrix} 2 & 0 & 1 & 0 \\ 0 & 2 & 0 & 1 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 0 & 2 \end{pmatrix}. \quad (19.17)$$

Similarly, the consistent mass matrix of the CST is computed by integration of the shape functions, resulting for 1D motion in

$$\mathbf{M}_{\text{CST}} = \frac{m}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} \quad \text{with} \quad m = \rho A t. \quad (19.18)$$

As before, the corresponding mass matrix for 2D motion is obtained by applying the above matrix for each dof independently.

For other (non-simplicial) elements, the stiffness matrix can be evaluated analogously or computed by *numerical quadrature*:

$$M^{ab} = \int_{\Omega} \rho N^a N^b dV \approx \sum_{k=0}^{n_{QP}-1} W_k \rho N^a(\boldsymbol{\xi}_k) N^b(\boldsymbol{\xi}_k) J(\boldsymbol{\xi}_k). \quad (19.19)$$

In summary, the dynamic problem is quite analogous to the quasistatic one. The key difference is the first term in (19.13), which requires a strategy to obtain numerical solutions that are time-dependent. Note that the above formulation in linearized kinematics can easily be adopted for **finite deformations** (the final matrix equations are the same with internal/external force vectors replaced by the finite-deformation counterparts).

We note that various references call this dynamic variational principle the “*principle of least action*” or “*principle of minimum action*”, which is in fact not correct since the solution must not necessarily be a minimizer of  $\mathcal{A}$  (it is merely guaranteed to be an extremizer).

The consistent mass matrix above is *dense*, which may be inconvenient for numerical solutions. Therefore, one often resorts to the so-called *lumped mass matrix* which is an approximation that is diagonal. For example, by using *particle-mass lumping* for a 2-node bar element, one distributes the mass evenly to the two end points, resulting in

$$\mathbf{M}_{1D} = \frac{m}{2} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{M}_{2D} = \frac{m}{2} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (19.20)$$

Note that a comparison of the kinetic energies (e.g., in 1D) reveals

$$T_{\text{lumped}} = \frac{1}{2} \dot{\mathbf{U}}_e^h \cdot \mathbf{M}_{\text{lumped}} \dot{\mathbf{U}}_e^h = \frac{m}{4} [(\dot{u}_e^1)^2 + (\dot{u}_e^2)^2] \quad (19.21)$$

and

$$T_{\text{consistent}} = \frac{1}{2} \dot{\mathbf{U}}_e^h \cdot \mathbf{M}_{\text{consistent}} \dot{\mathbf{U}}_e^h = \frac{m}{6} [(\dot{u}_e^1)^2 + (\dot{u}_e^2)^2 + \dot{u}_e^1 \dot{u}_e^2]. \quad (19.22)$$

Hence, the lumped mass matrix yields to an approximation of the kinetic energy (which is generally not exact).

Particle-mass lumping can be extended to arbitrary elements having  $n$  nodes by defining

$$\mathbf{M} = \frac{m}{n} \mathbf{I}. \quad (19.23)$$

Finally, note that in *structural dynamics*, one often includes velocity-proportional damping through a **damping matrix**  $\mathbf{C}$  such that (19.13) turns into

$$\mathbf{M} \ddot{\mathbf{U}}^h + \mathbf{C} \dot{\mathbf{U}}^h + \mathbf{F}_{\text{int}}(\mathbf{U}^h) - \mathbf{F}_{\text{ext}}(t) = \mathbf{0} \quad (19.24)$$

with, oftentimes, mass- and stiffness-proportional damping via

$$\mathbf{C} = \alpha \mathbf{M} + \beta \mathbf{K}, \quad \alpha, \beta \in \mathbb{R}_+. \quad (19.25)$$

The choice of  $\alpha > 0$  controls low-frequency vibration attenuation, while  $\beta > 0$  suppresses high-frequency vibrations.

## 19.2 Free vibrations

Free vibrations of infinitesimal amplitude are a frequent case of interest. Starting with the general, nonlinear equations of motion,

$$\mathbf{M} \ddot{\mathbf{U}}^h + \mathbf{F}_{\text{int}}(\mathbf{U}^h) = \mathbf{F}_{\text{ext}}(t), \quad (19.26)$$

we linearize about a stable equilibrium configuration  $\mathbf{U}_0^h$  with constant forces  $\mathbf{F}_{\text{ext}} = \mathbf{F}_{\text{ext},0}$  such that  $\mathbf{F}_{\text{int}}(\mathbf{U}_0^h) = \mathbf{F}_{\text{ext},0}$ . Consider now a small time-varying perturbation  $\delta \mathbf{U}^h(t)$  such that  $\mathbf{U}^h = \mathbf{U}_0^h + \delta \mathbf{U}^h$  gives to leading order

$$\mathbf{M} \delta \ddot{\mathbf{U}}^h + \mathbf{F}_{\text{int}}(\mathbf{U}_0^h) + \frac{\partial \mathbf{F}_{\text{int}}}{\partial \mathbf{U}^h}(\mathbf{U}_0^h) \delta \mathbf{U}^h + \text{h.o.t.} = \mathbf{F}_{\text{ext},0} \quad (19.27)$$

or, invoking equilibrium,

$$\mathbf{M} \delta \ddot{\mathbf{U}}^h + \mathbf{T}(\mathbf{U}_0^h) \delta \mathbf{U}^h + \text{h.o.t.} = \mathbf{0}. \quad (19.28)$$

Similarly, when considering free vibrations about the undeformed ground state of an elastic body, we would have arrived immediately at

$$\mathbf{M} \ddot{\mathbf{U}}^h + \mathbf{K} \mathbf{U}^h = \mathbf{0}. \quad (19.29)$$

Thus, (19.28) is the generalized form of (19.29) for equilibrated systems. In both cases the form of the equation of motion governing free vibrations without external forcing and without damping is

$$\mathbf{M} \ddot{\mathbf{U}}^h + \mathbf{T} \mathbf{U}^h = \mathbf{0}. \quad (19.30)$$

The solution to the above ODE is of the separable form

$$\mathbf{U}^h = \hat{\mathbf{U}}^h \exp(i\omega t). \quad (19.31)$$

Insertion into (19.30) leads to (exploiting that the equation must hold for all times  $t$ )

$$(\mathbf{T} - \omega^2 \mathbf{M}) \hat{\mathbf{U}}^h = \mathbf{0}, \quad (19.32)$$

which is an **eigenvalue problem** with **eigenfrequency**  $\omega$  and **eigenvector**  $\hat{\mathbf{U}}^h$ . The latter is also referred to as the **eigenmode**.

For a FE discretization with  $n_{\text{dof}} = dn_n$  degrees of freedom ( $n_n$  nodes in  $d$  dimensions), the eigenvalue problem has  $n_{\text{dof}}$  eigenfrequencies with  $n_{\text{dof}}$  associated distinct eigenmodes. Each rigid body mode of the FE problem corresponds to a zero eigenfrequency. Therefore, a 2D (3D) free vibrational problem without essential BCs has 3 (6) *zero eigenfrequencies*.

### Example: free vibration of a bar

We can use the example of a freely vibrating bar to assess the influence of the different mass matrices. Consider a 2-node bar element with only axial displacements, so that each node has only a single degree of freedom  $u_i$ . The mass matrices and the stiffness matrix for this case were derived previously as

$$\mathbf{M}_{\text{consistent}} = \frac{m}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}, \quad \mathbf{M}_{\text{lumped}} = \frac{m}{2} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{K} = \frac{EA}{L} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}. \quad (19.33)$$

In both cases of the lumped and consistent mass matrices, we compute the two eigenfrequencies and eigenmodes by solving the eigenvalue problem

$$(\mathbf{K} - \omega^2 \mathbf{M}) \hat{\mathbf{U}}^h = \mathbf{0} \quad \Rightarrow \quad \det(\mathbf{K} - \omega^2 \mathbf{M}) = 0. \quad (19.34)$$

Insertion of the stiffness and consistent mass matrix results in the two solutions

$$\omega_0^{\text{consistent}} = 0, \quad \omega_1^{\text{consistent}} = \sqrt{12 \frac{EA}{mL}} \approx 3.464 \sqrt{\frac{EA}{mL}}. \quad (19.35)$$

The corresponding eigenvectors follow from insertion of the eigenfrequencies into the eigenvalue problem, giving

$$\hat{\mathbf{U}}_0^{\text{consistent}} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \hat{\mathbf{U}}_1^{\text{consistent}} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}. \quad (19.36)$$

As expected, we have one zero eigenfrequency associated with rigid body translation. When repeating the above procedure with the lumped mass matrix, we instead obtain

$$\omega_0^{\text{lumped}} = 0, \quad \omega_1^{\text{lumped}} = 2 \sqrt{\frac{EA}{mL}}. \quad (19.37)$$

and

$$\hat{\mathbf{U}}_0^{\text{lumped}} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \hat{\mathbf{U}}_1^{\text{lumped}} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}. \quad (19.38)$$

Hence, the two cases yield the same eigenmodes but significantly differ in the fundamental frequency. For comparison, let us compute the exact solution by studying the free vibration of a continuous, homogeneous, linear elastic bar. Linear momentum balance, i.e.,

$$Eu_{,xx} = \rho \ddot{u}, \quad (19.39)$$

admits the separable solution

$$u(x, t) = \hat{u}(x) \exp(i\omega t) \quad \Rightarrow \quad E\hat{u}_{,xx}(x) = -\omega^2 \rho \hat{u}(x) \quad (19.40)$$

Free-end boundary conditions imply that

$$\hat{u}_{,x}(0) = \hat{u}_{,x}(L) \quad \Rightarrow \quad \hat{u}(x) = A \cos\left(\frac{n\pi x}{L}\right) \quad (19.41)$$

with an integer  $n$  and a constant  $A \in \mathbb{R}$ . Insertion into (19.39) leads to

$$-E \left(\frac{n\pi}{L}\right)^2 = -\omega^2 \rho \quad \Rightarrow \quad \omega_n = n \pi \sqrt{\frac{EA}{mL}}. \quad (19.42)$$

Consequently, the eigenfrequencies are obtained as

$$\omega_0 = 0, \quad \omega_1 = \pi \sqrt{\frac{EA}{mL}}. \quad (19.43)$$

Comparison with the above two cases results in

$$\omega_1^{\text{lumped}} \leq \omega_1 \leq \omega_1^{\text{consistent}}. \quad (19.44)$$

Note that the eigenfrequency obtained from the consistent mass matrix is an upper bound on the eigenfrequency, while the lumped mass matrix is usually a lower bound (the latter is not rigorous, though, since it depends on the choice of the lumped mass matrix).

### 19.3 Modal decomposition

Let us introduce a few important features of the eigenmodes discussed above. By considering the eigenvalue problem for two distinct eigenmodes/-frequencies, we may write

$$\begin{aligned} (\mathbf{T} - \omega_i^2 \mathbf{M}) \hat{\mathbf{U}}_i^h = \mathbf{0} & \Rightarrow \hat{\mathbf{U}}_j^h \cdot (\mathbf{T} - \omega_i^2 \mathbf{M}) \hat{\mathbf{U}}_i^h = 0, \\ (\mathbf{T} - \omega_j^2 \mathbf{M}) \hat{\mathbf{U}}_j^h = \mathbf{0} & \Rightarrow \hat{\mathbf{U}}_i^h \cdot (\mathbf{T} - \omega_j^2 \mathbf{M}) \hat{\mathbf{U}}_j^h = 0, \end{aligned} \quad (19.45)$$

where we simply pre-multiplied each of the two equations by the respective other eigenvector. Subtraction of the two equations (using that  $\mathbf{T}$  is by definition symmetric) results in

$$(\omega_i^2 - \omega_j^2) \hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_j^h = 0. \quad (19.46)$$

This implies that either  $\omega_i = \omega_j$  (considering only positive eigenfrequencies) or  $\hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_j^h = 0$ . If we eliminate duplicated eigenfrequencies by, e.g., Gram-Schmid orthonormalization, then we may conclude that

$$\hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_j^h = 0 \quad \text{if} \quad i \neq j. \quad (19.47)$$

We can normalize the eigenvectors in the following fashion:

$$\hat{\mathbf{U}}_i^h \leftarrow \frac{\hat{\mathbf{U}}_i^h}{\sqrt{\hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_i^h}} \Rightarrow \hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_i^h = 1, \quad (19.48)$$

so that overall

$$\hat{\mathbf{U}}_i^h \cdot \mathbf{M} \hat{\mathbf{U}}_j^h = \delta_{ij}. \quad (19.49)$$

If we now consider

$$(\mathbf{T} - \omega_i^2 \mathbf{M}) \hat{\mathbf{U}}_i^h = \mathbf{0} \Rightarrow \hat{\mathbf{U}}_j^h \cdot (\mathbf{T} - \omega_i^2 \mathbf{M}) \hat{\mathbf{U}}_i^h = 0, \quad (19.50)$$

then we may also conclude that, by invoking (19.49),

$$\hat{\mathbf{U}}_j^h \cdot \mathbf{T} \hat{\mathbf{U}}_i^h = 0 \quad \text{if} \quad i \neq j. \quad (19.51)$$

Similarly, by pre-multiplying by the respective same eigenvector, i.e.,

$$\hat{\mathbf{U}}_i^h \cdot (\mathbf{T} - \omega_i^2 \mathbf{M}) \hat{\mathbf{U}}_i^h = 0, \quad (19.52)$$

we solve for  $\omega_i^2$  and obtain **Rayleigh's quotient**, which here simplifies due to the normalization:

$$\omega_i^2 = \frac{\hat{\mathbf{U}}_{(i)}^h \cdot \mathbf{T} \hat{\mathbf{U}}_{(i)}^h}{\hat{\mathbf{U}}_{(i)}^h \cdot \mathbf{M} \hat{\mathbf{U}}_{(i)}^h} = \hat{\mathbf{U}}_{(i)}^h \cdot \mathbf{T} \hat{\mathbf{U}}_{(i)}^h. \quad (19.53)$$

This all forms the basis for the method known as **modal decomposition**.

The starting point for modal decomposition of *linearized, elastic* system is the Fourier representation

$$\mathbf{U}^h(t) = \sum_{i=1}^n z_i(t) \hat{\mathbf{U}}_i^h, \quad (19.54)$$

where  $\{\hat{\mathbf{U}}_1^h, \dots, \hat{\mathbf{U}}_n^h\}$  are the  $n$  eigenmodes of the system. That is, we pre-compute the eigenvectors and seek a solution as a linear superposition of all eigenvectors with some unknown scalar coefficients that are continuous functions of time (maintaining the *semi*-discretization).

We substitute (19.54) into the linearized equations of motion with external forces,  $\mathbf{M}\ddot{\mathbf{U}}^h + \mathbf{T}\mathbf{U}^h = \mathbf{F}_{\text{ext}}$ , and pre-multiply the system of equations by  $\hat{\mathbf{U}}_i^h$ , so we arrive at

$$\sum_{i=1}^n [\ddot{z}_i(t) \mathbf{M} \hat{\mathbf{U}}_i^h + z_i(t) \mathbf{T} \hat{\mathbf{U}}_i^h] = \mathbf{F}_{\text{ext}}(t). \quad (19.55)$$

Pre-multiplying by  $\hat{\mathbf{U}}_j^h$  and exploiting the above orthogonality gives

$$\sum_{i=1}^n [\ddot{z}_i(t) \hat{\mathbf{U}}_j^h \cdot \mathbf{M} \hat{\mathbf{U}}_i^h + z_i(t) \hat{\mathbf{U}}_j^h \cdot \mathbf{T} \hat{\mathbf{U}}_i^h] = \hat{\mathbf{U}}_j^h \cdot \mathbf{F}_{\text{ext}}(t) \quad \Leftrightarrow \quad \ddot{z}_j(t) + \omega_j^2 z_j(t) = \hat{\mathbf{U}}_j^h \cdot \mathbf{F}_{\text{ext}}(t), \quad (19.56)$$

where no summation over  $j$  is implied. This equation is quite remarkable as it is a *scalar, linear* ODE for the unknown function  $z_j(t)$ . Moreover, the system of equations for  $n$  unknown equations  $z_i(t)$  (with  $i = 1, \dots, n$ ) has decoupled into  $n$  uncoupled scalar ODEs to be solved independently for the  $z_i(t)$ . The strategy is thus to first pre-compute all eigenfrequencies  $\omega_i$  and all eigenmodes  $\hat{\mathbf{U}}_j^h$ , so that the coefficients and right-hand sides are known in

$$\ddot{z}_i(t) + \omega_i^2 z_i(t) = \hat{\mathbf{U}}_i^h \cdot \mathbf{F}_{\text{ext}}(t), \quad (19.57)$$

which can relatively inexpensively be solved for functions  $z_i(t)$ .

For many practical problems, only a limited number of modes are important (i.e., taking the first  $m < n$  modes only). Therefore, numerical efficiency can be gained by truncating the Fourier sum, which is often referred to as **order reduction**:

$$\mathbf{U}^h(t) = \sum_{i=1}^m z_i(t) \hat{\mathbf{U}}_i^h, \quad m < n. \quad (19.58)$$

## 19.4 Transient time-dependent solutions

The time-dependent solution for the nodal variables  $\mathbf{U}^h(t)$  is obtained either in a continuous manner (e.g., by modal decomposition or in case of free vibrations as discussed above) or in a time-discretized fashion (e.g., by using finite-difference approximations in time). The turns the *semi*-discretization into a proper discretization in both space and time.

## 19.5 Explicit time integration

We discretize the solution in time, e.g., we assume constant time increments  $\Delta t > 0$  and write

$$\mathbf{u}^a(t^\alpha) = \mathbf{u}^{a,\alpha}, \quad \Delta t = t^{\alpha+1} - t^\alpha. \quad (19.59)$$

By using *central-difference approximations*, we obtain

$$\dot{\mathbf{U}}^h(t^\alpha) = \frac{\mathbf{U}^{h,\alpha+1} - \mathbf{U}^{h,\alpha-1}}{2\Delta t}, \quad \ddot{\mathbf{U}}^h(t^\alpha) = \frac{\mathbf{U}^{h,\alpha+1} - 2\mathbf{U}^{h,\alpha} + \mathbf{U}^{h,\alpha-1}}{(\Delta t)^2}. \quad (19.60)$$

Insertion into (19.24) leads to

$$\mathbf{M} \frac{\mathbf{U}^{h,\alpha+1} - 2\mathbf{U}^{h,\alpha} + \mathbf{U}^{h,\alpha-1}}{(\Delta t)^2} + \mathbf{C} \frac{\mathbf{U}^{h,\alpha+1} - \mathbf{U}^{h,\alpha-1}}{2\Delta t} + \mathbf{F}_{\text{int}}(\mathbf{U}^{h,\alpha}) - \mathbf{F}_{\text{ext}}(t^\alpha) = \mathbf{0}, \quad (19.61)$$

which can be reorganized into

$$\left[ \frac{\mathbf{M}}{(\Delta t)^2} + \frac{\mathbf{C}}{2\Delta t} \right] \mathbf{U}^{\alpha+1} = \frac{2\mathbf{M}}{(\Delta t)^2} \mathbf{U}^\alpha + \left[ \frac{\mathbf{C}}{2\Delta t} - \frac{\mathbf{M}}{(\Delta t)^2} \right] \mathbf{U}^{\alpha-1} - \mathbf{F}_{\text{int}}(\mathbf{U}^{h,\alpha}) + \mathbf{F}_{\text{ext}}(t^\alpha). \quad (19.62)$$

This is an update rule for  $\mathbf{U}^{\alpha+1}$ , using **explicit time integration**. Note that here and in the following we drop the superscript  $h$  for simplicity and to avoid confusion with the time step superscript  $\alpha$ .

Note that stability limits the time step of the explicit scheme. Specifically, we must ensure that

$$\Delta t \leq \Delta t_{\text{cr}} = \frac{2}{\omega_{\text{max}}}, \quad (19.63)$$

with  $\omega_{\text{max}}$  being the highest eigenfrequency. Note that the highest eigenfrequency scales inversely proportional to the smallest element size. For example, recall the two-node bar for which the eigenfrequency was of the form  $\omega_1 \propto \frac{1}{L} \sqrt{E/\rho}$  with  $L$  being the element size.

## 19.6 A reinterpretation of finite differences

Instead of introducing finite differences as done above, we could alternatively use interpolation functions in both space and time. For example, consider a discretization in time which evaluates  $\mathbf{U}^h$  at discrete time intervals  $\Delta t$  and then uses a quadratic interpolation in time, i.e., we define

$$\mathbf{U}^h(t) = \mathbf{U}^{\alpha+1} N^{\alpha+1}(t) + \mathbf{U}^\alpha N^\alpha(t) + \mathbf{U}^{\alpha-1} N^{\alpha-1}(t) \quad \text{for } t \in [t^\alpha - \frac{\Delta t}{2}, t^\alpha + \frac{\Delta t}{2}]. \quad (19.64)$$

where we dropped the superscript  $h$  for simplicity. Shape functions  $N(t)$  interpolate in time. The chosen quadratic interpolation and the range of validity ensure that  $\mathbf{U}^h(t)$  is twice differentiable (as needed for the acceleration) if

$$N^{\alpha+1}(t) = \frac{(t - t^\alpha)(t - t^{\alpha-1})}{2(\Delta t)^2}, \quad N^\alpha(t) = \frac{(t^{\alpha+1} - t)(t - t^{\alpha-1})}{(\Delta t)^2}, \quad N^{\alpha-1}(t) = \frac{(t^{\alpha+1} - t)(t^\alpha - t)}{2(\Delta t)^2}. \quad (19.65)$$

This choice ensures that

$$\frac{\partial}{\partial t} N^{\alpha+1}(t^\alpha - \frac{\Delta t}{2}) = 0 \quad \text{and} \quad \frac{\partial}{\partial t} N^{\alpha-1}(t^\alpha + \frac{\Delta t}{2}) = 0, \quad (19.66)$$

so that  $\dot{\mathbf{U}}^h(t)$  is indeed continuous. The acceleration evaluates to

$$\begin{aligned}\ddot{\mathbf{U}}^h(t) &= \mathbf{U}^{\alpha+1} \ddot{N}^{\alpha+1}(t) + \mathbf{U}^\alpha \ddot{N}^\alpha(t) + \mathbf{U}^{\alpha-1} \ddot{N}^{\alpha-1}(t) \\ &= \frac{\mathbf{U}^{\alpha+1}}{(\Delta t)^2} - 2 \frac{\mathbf{U}^\alpha}{(\Delta t)^2} + \frac{\mathbf{U}^{\alpha-1}}{(\Delta t)^2} = \frac{\mathbf{U}^{\alpha+1} - 2\mathbf{U}^\alpha + \mathbf{U}^{\alpha-1}}{(\Delta t)^2}.\end{aligned}\quad (19.67)$$

Thus, we recover the second-order central-difference scheme in time.

The full, discretized displacement field is now

$$\mathbf{u}^h(\mathbf{x}, t) = \sum_{a=1}^n \sum_{\gamma=\alpha-1}^{\alpha+1} \mathbf{u}^{a,\gamma} N^\gamma(t) N^a(\mathbf{x}) \quad \text{for } t \in [t^\alpha - \frac{\Delta t}{2}, t^\alpha + \frac{\Delta t}{2}], \quad (19.68)$$

where  $\mathbf{u}^{a,\gamma}$  is the displacement at node  $a$  at time  $t^\gamma$ . The acceleration field follows as the piecewise-constant approximation

$$\ddot{\mathbf{u}}^h(\mathbf{x}, t) = \sum_{a=1}^n \frac{\mathbf{u}^{a,\alpha+1} - 2\mathbf{u}^{a,\alpha} + \mathbf{u}^{a,\alpha-1}}{(\Delta t)^2} N^a(\mathbf{x}) \quad \text{for } t \in [t^\alpha - \frac{\Delta t}{2}, t^\alpha + \frac{\Delta t}{2}]. \quad (19.69)$$

Next, one may wish to evaluate information at the  $n_t$  discrete time steps, motivating the choice of the trial function as

$$\mathbf{v}^h(\mathbf{x}, t) = \sum_{b=1}^n \mathbf{v}^b(t) N^b(\mathbf{x}) = \sum_{b=1}^n \sum_{\alpha=1}^{n_t} \mathbf{v}^{b,\alpha} N^b(\mathbf{x}) \delta(t - t^\alpha) \quad (19.70)$$

When this choice of  $\mathbf{v}^h$  along with  $\mathbf{u}^h$  into the weak form, which derived earlier as

$$\mathcal{G}(u, v) = \int_{t_1}^{t_2} \left[ \int_{\Omega} (-\rho \ddot{u}_i v_i - \sigma_{ij} v_{i,j}) \, dV + \int_{\Omega} \rho b_i v_i \, dV + \int_{\partial\Omega_N} \hat{t}_i v_i \, dS \right] dt = 0, \quad (19.71)$$

then we may integrate explicitly with respect to time, resulting in

$$\begin{aligned}\sum_{b=1}^{n_t} v_i^b \left[ \int_{\Omega} \left( -\rho \frac{u_i^{a,\alpha+1} - 2u_i^{a,\alpha} + u_i^{a,\alpha-1}}{(\Delta t)^2} N^a N^b - \sigma_{ij} (\nabla \mathbf{u}^\alpha) N_{,j}^b \right) dV \right. \\ \left. + \int_{\Omega} \rho b_i(t^\alpha) N^b \, dV + \int_{\partial\Omega_N} \hat{t}_i(t^\alpha) N^b \, dS \right] = 0\end{aligned}\quad (19.72)$$

for each time step  $t^\alpha$  and for all admissible choices of  $v_i^b$ . The latter implies that we must have

$$\begin{aligned}\int_{\Omega} \rho N^a N^b \, dV \frac{u_i^{a,\alpha+1} - 2u_i^{a,\alpha} + u_i^{a,\alpha-1}}{(\Delta t)^2} + \int_{\Omega} \sigma_{ij} (\nabla \mathbf{u}^\alpha) N_{,j}^b \, dV \\ - \int_{\Omega} \rho b_i(t^\alpha) N^b \, dV - \int_{\partial\Omega_N} \hat{t}_i(t^\alpha) N^b \, dS = 0.\end{aligned}\quad (19.73)$$

This, however, is equivalent to

$$\mathbf{M} \frac{\mathbf{U}^{h,\alpha+1} - 2\mathbf{U}^{h,\alpha} + \mathbf{U}^{h,\alpha-1}}{(\Delta t)^2} + \mathbf{F}_{\text{int}}(\mathbf{U}^{h,\alpha}) - \mathbf{F}_{\text{ext}}(t^\alpha) = \mathbf{0}, \quad (19.74)$$

which was derived above for a second-order finite-difference time discretization, see (19.61).



## 19.7 Implicit time integration

Next, **implicit time integration** uses the same discretization in time but requires solving a (non)linear system of equations for  $\mathbf{U}^{\alpha+1}$ .

The most popular scheme for mechanical problems, is the so-called **Newmark- $\beta$**  method which is a combination of the linear acceleration and average acceleration schemes. Specifically, one assumes

$$\mathbf{U}^{\alpha+1} = \mathbf{U}^\alpha + \Delta t \dot{\mathbf{U}}^\alpha + \frac{(\Delta t)^2}{2} [2\beta \ddot{\mathbf{U}}^{\alpha+1} + (1-2\beta)\ddot{\mathbf{U}}^\alpha] \quad (19.75a)$$

$$\dot{\mathbf{U}}^{\alpha+1} = \dot{\mathbf{U}}^\alpha + \Delta t [\gamma \ddot{\mathbf{U}}^{\alpha+1} + (1-\gamma)\ddot{\mathbf{U}}^\alpha] \quad (19.75b)$$

with parameters  $0 \leq \beta \leq \gamma/2$  and  $0 \leq \gamma \leq 1$ , often chosen as

- $\beta = \frac{1}{4}$ ,  $\gamma = \frac{1}{2}$  in the **average acceleration** scheme, which is implicit and unconditionally stable (the underlying approximation is that the acceleration within the time interval  $\Delta t$  remains constant).
- $\beta = \frac{1}{6}$ ,  $\gamma = \frac{1}{2}$  in the **linear acceleration** scheme,
- $\beta = \gamma = 0$  returns to the explicit central-difference scheme discussed above.

Most popular is the average accelerations scheme which is unconditionally stable for arbitrary  $\Delta t > 0$ .

For linear structural dynamics problems, the method is *unconditionally stable* if  $2\beta \geq \gamma \geq 1/2$ . It is *conditionally stable* if  $\gamma < 1/2$ . For  $\gamma = 1/2$  the scheme is at least second-order accurate, while being first-order accurate otherwise.

For implementation purposes, let us solve (19.75) for the acceleration at the new time  $t^{\alpha+1}$ :

$$\ddot{\mathbf{U}}^{\alpha+1} = \frac{1}{\beta(\Delta t)^2} (\mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha - \Delta t \dot{\mathbf{U}}^\alpha) - \frac{1-2\beta}{2\beta} \ddot{\mathbf{U}}^\alpha. \quad (19.76)$$

Insertion into (19.75b) yields

$$\begin{aligned} \dot{\mathbf{U}}^{\alpha+1} &= \dot{\mathbf{U}}^\alpha + \Delta t(1-\gamma)\ddot{\mathbf{U}}^\alpha + \frac{\gamma}{\beta \Delta t} (\mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha + \Delta t \dot{\mathbf{U}}^\alpha) - \gamma \Delta t \frac{1-2\beta}{2\beta} \ddot{\mathbf{U}}^\alpha \\ &= \left(1 - \frac{\gamma}{\beta}\right) \dot{\mathbf{U}}^\alpha + \frac{\gamma}{\beta \Delta t} (\mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha) - \Delta t \left(\frac{\gamma}{2\beta} - 1\right) \ddot{\mathbf{U}}^\alpha \end{aligned} \quad (19.77)$$

Next, inserting both velocity (19.76) and acceleration (19.77) into the equation of motion at the new time  $t^{\alpha+1}$ ,

$$\mathbf{M} \ddot{\mathbf{U}}^{\alpha+1} + \mathbf{C} \dot{\mathbf{U}}^{\alpha+1} + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^{\alpha+1}) = \mathbf{0}, \quad (19.78)$$

leads after some rearrangement to

$$\begin{aligned} &\left( \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} \right) \mathbf{U}^{\alpha+1} + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^{\alpha+1}) \\ &= \mathbf{M} \left[ \frac{1}{\beta(\Delta t)^2} \mathbf{U}^\alpha + \frac{1}{\beta \Delta t} \dot{\mathbf{U}}^\alpha + \left( \frac{1}{2\beta} - 1 \right) \ddot{\mathbf{U}}^\alpha \right] \\ &\quad + \mathbf{C} \left[ \frac{\gamma}{\beta \Delta t} \mathbf{U}^\alpha + \left( \frac{\gamma}{\beta} - 1 \right) \dot{\mathbf{U}}^\alpha + \Delta t \left( \frac{\gamma}{2\beta} - 1 \right) \ddot{\mathbf{U}}^\alpha \right]. \end{aligned} \quad (19.79)$$

The right-hand side is fully known as it only involves  $\mathbf{U}^\alpha$ ,  $\dot{\mathbf{U}}^\alpha$ , and  $\ddot{\mathbf{U}}^\alpha$  from the previous time step. The left-hand side is generally nonlinear and requires an iterative solver.

Note that the implementation is quite similar to that of the quasistatic Newton-Raphson solver discussed before. The right-hand side is significantly longer (but contains only known information). The left-hand side is extended by the first term in (19.79), which is linear in  $\mathbf{U}^{\alpha+1}$  and the tangent matrix used for iterations is simply

$$\mathbf{T}^* = \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} + \mathbf{T}. \quad (19.80)$$

Boundary conditions can be implemented in the same fashion as for the quasistatic solvers discussed above.

As an **alternative derivation** of the implementation, we may introduce the changes in displacements, velocities and accelerations during a time step as

$$\Delta \mathbf{U}^\alpha = \mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha, \quad \Delta \dot{\mathbf{U}}^\alpha = \dot{\mathbf{U}}^{\alpha+1} - \dot{\mathbf{U}}^\alpha, \quad \Delta \ddot{\mathbf{U}}^\alpha = \ddot{\mathbf{U}}^{\alpha+1} - \ddot{\mathbf{U}}^\alpha. \quad (19.81)$$

We can rearrange the first equation of (19.75) into

$$\mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha = \Delta t \dot{\mathbf{U}}^\alpha + \frac{(\Delta t)^2}{2} [2\beta(\ddot{\mathbf{U}}^{\alpha+1} - \ddot{\mathbf{U}}^\alpha) + \ddot{\mathbf{U}}^\alpha] \quad (19.82)$$

or

$$\Delta \mathbf{U}^\alpha = \Delta t \dot{\mathbf{U}}^\alpha + \frac{(\Delta t)^2}{2} [2\beta \Delta \ddot{\mathbf{U}}^\alpha + \ddot{\mathbf{U}}^\alpha] \quad \Leftrightarrow \quad \Delta \ddot{\mathbf{U}}^\alpha = \frac{1}{\beta(\Delta t)^2} \Delta \mathbf{U}^\alpha - \frac{1}{\beta \Delta t} \dot{\mathbf{U}}^\alpha - \frac{1}{2\beta} \ddot{\mathbf{U}}^\alpha. \quad (19.83)$$

Similarly, the second equation of (19.75) gives

$$\dot{\mathbf{U}}^{\alpha+1} - \dot{\mathbf{U}}^\alpha = \Delta t [\gamma(\ddot{\mathbf{U}}^{\alpha+1} - \ddot{\mathbf{U}}^\alpha) + \ddot{\mathbf{U}}^\alpha] \quad \Leftrightarrow \quad \Delta \dot{\mathbf{U}}^\alpha = \Delta t \gamma \Delta \ddot{\mathbf{U}}^\alpha + \Delta t \ddot{\mathbf{U}}^\alpha, \quad (19.84)$$

so that insertion of the final result in (19.83) results in

$$\begin{aligned} \Delta \dot{\mathbf{U}}^\alpha &= \Delta t \gamma \left[ \frac{1}{\beta(\Delta t)^2} \Delta \mathbf{U}^\alpha - \frac{1}{\beta \Delta t} \dot{\mathbf{U}}^\alpha - \frac{1}{2\beta} \ddot{\mathbf{U}}^\alpha \right] + \Delta t \ddot{\mathbf{U}}^\alpha \\ &= \frac{\gamma}{\beta \Delta t} \Delta \mathbf{U}^\alpha - \frac{\gamma}{\beta} \dot{\mathbf{U}}^\alpha + \Delta t \left( 1 - \frac{\gamma}{2\beta} \right) \ddot{\mathbf{U}}^\alpha. \end{aligned} \quad (19.85)$$

The equations of motion (19.24) in their general time-discretized form are

$$\mathbf{M} \ddot{\mathbf{U}}^{\alpha+1} + \mathbf{C} \dot{\mathbf{U}}^{\alpha+1} + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^{\alpha+1}) = \mathbf{0}. \quad (19.86)$$

Subtracting the above equation for two consecutive time steps results in the incremental form

$$\mathbf{M} \Delta \ddot{\mathbf{U}}^\alpha + \mathbf{C} \Delta \dot{\mathbf{U}}^\alpha + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) - \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha) - [\mathbf{F}_{\text{ext}}(t^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^\alpha)] = \mathbf{0}. \quad (19.87)$$

Now inserting (19.83) and (19.85) leads to

$$\begin{aligned} \mathbf{M} \left[ \frac{1}{\beta(\Delta t)^2} \Delta \mathbf{U}^\alpha - \frac{1}{\beta \Delta t} \dot{\mathbf{U}}^\alpha - \frac{1}{2\beta} \ddot{\mathbf{U}}^\alpha \right] + \mathbf{C} \left[ \frac{\gamma}{\beta \Delta t} \Delta \mathbf{U}^\alpha - \frac{\gamma}{\beta} \dot{\mathbf{U}}^\alpha + \Delta t \left( 1 - \frac{\gamma}{2\beta} \right) \ddot{\mathbf{U}}^\alpha \right] \\ + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) - \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha) - [\mathbf{F}_{\text{ext}}(t^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^\alpha)] = \mathbf{0}, \end{aligned} \quad (19.88)$$

which can be regrouped into

$$\begin{aligned} \left[ \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} \right] \Delta \mathbf{U}^\alpha + \mathbf{F}_{\text{int}}(\mathbf{U}^{\alpha+1}) &= \mathbf{F}_{\text{ext}}(t^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^\alpha) + \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha) \\ - \left[ -\frac{1}{2\beta} \mathbf{M} + \Delta t \left( 1 - \frac{\gamma}{2\beta} \right) \mathbf{C} \right] \ddot{\mathbf{U}}^\alpha &- \left[ -\frac{1}{\beta \Delta t} \mathbf{M} - \frac{\gamma}{\beta} \mathbf{C} \right] \dot{\mathbf{U}}^\alpha, \end{aligned} \quad (19.89)$$

whose right-hand side is entirely known from the previous time step (and the external forces at the current time). Noting that  $\Delta \mathbf{U}^\alpha = \mathbf{U}^{\alpha+1} - \mathbf{U}^\alpha$ , this is a generally nonlinear system of equations to be solved for  $\Delta \mathbf{U}^\alpha$  or  $\mathbf{U}^{\alpha+1}$ . For example, we can use a *Newton-Raphson iterative approach* that starts with an initial guess  $\Delta \mathbf{U}_0^\alpha$  and incrementally finds

$$\Delta \mathbf{U}_{i+1}^\alpha = \Delta \mathbf{U}_i^\alpha + \delta \mathbf{U}_i^\alpha, \quad (19.90)$$

so that we need to solve

$$\left[ \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} \right] (\Delta \mathbf{U}_i^\alpha + \delta \mathbf{U}_i^\alpha) + \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha + \Delta \mathbf{U}_i^\alpha + \delta \mathbf{U}_i^\alpha) = \text{RHS}^\alpha \quad (19.91)$$

with

$$\text{RHS}^\alpha = \mathbf{F}_{\text{ext}}(t^{\alpha+1}) - \mathbf{F}_{\text{ext}}(t^\alpha) + \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha) + \left[ \frac{1}{2\beta} \mathbf{M} - \Delta t \left( 1 - \frac{\gamma}{2\beta} \right) \mathbf{C} \right] \ddot{\mathbf{U}}^\alpha + \left[ \frac{1}{\beta \Delta t} \mathbf{M} + \frac{\gamma}{\beta} \mathbf{C} \right] \dot{\mathbf{U}}^\alpha. \quad (19.92)$$

Linearization about  $\delta \mathbf{U}_i^\alpha = \mathbf{0}$  finally gives

$$\begin{aligned} &\left[ \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} + \mathbf{T}(\mathbf{U}^\alpha + \Delta \mathbf{U}_i^\alpha) \right] \delta \mathbf{U}_i^\alpha \\ &\approx \text{RHS}^\alpha - \left[ \frac{1}{\beta(\Delta t)^2} \mathbf{M} + \frac{\gamma}{\beta \Delta t} \mathbf{C} \right] \Delta \mathbf{U}_i^\alpha - \mathbf{F}_{\text{int}}(\mathbf{U}^\alpha + \Delta \mathbf{U}_i^\alpha). \end{aligned} \quad (19.93)$$

After a time step  $\alpha$ , we know  $\mathbf{U}^\alpha$ ,  $\dot{\mathbf{U}}^\alpha$  and  $\ddot{\mathbf{U}}^\alpha$  as well as  $\mathbf{F}_{\text{ext}}(t^{\alpha+1})$  and  $\mathbf{F}_{\text{ext}}(t^\alpha)$ . To solve for the new solution at  $t^{\alpha+1}$ , we choose an initial guess  $\Delta \mathbf{U}_0^\alpha$  (e.g.,  $\Delta \mathbf{U}_0^\alpha = \mathbf{0}$ ), so that the entire right-hand side in (19.93) as well as the matrix in front of the left-hand side in (19.93) are known. Next, the linear system in (19.93) is solved for  $\delta \mathbf{U}_i^\alpha$ , followed by the update  $\Delta \mathbf{U}_{i+1}^\alpha = \Delta \mathbf{U}_i^\alpha + \delta \mathbf{U}_i^\alpha$ , until convergence is achieved when  $\|\delta \mathbf{U}_i^\alpha\| \rightarrow 0$ . Once  $\Delta \mathbf{U}^\alpha = \lim_{i \rightarrow \infty} \Delta \mathbf{U}_i^\alpha$  has been found, we can use (19.83) to compute  $\Delta \ddot{\mathbf{U}}^\alpha$  and  $\ddot{\mathbf{U}}^{\alpha+1} = \ddot{\mathbf{U}}^\alpha + \Delta \ddot{\mathbf{U}}^\alpha$  as well as (19.85) to compute  $\Delta \dot{\mathbf{U}}^\alpha$  and  $\dot{\mathbf{U}}^{\alpha+1} = \dot{\mathbf{U}}^\alpha + \Delta \dot{\mathbf{U}}^\alpha$ . Then, the procedure restarts for the next time step.

Note that, if essential boundary conditions of the type  $\mathbf{U}^\alpha = \hat{\mathbf{U}}^\alpha$  are to be imposed, this implies that  $\Delta \mathbf{U}^\alpha = \hat{\mathbf{U}}^{\alpha+1} - \hat{\mathbf{U}}^\alpha$  is known for all times.

## 20 Internal variables and inelasticity

### 20.1 Inelastic material models

Inelastic material models describe a variety of phenomena, e.g.,

- **viscoelasticity**, i.e., time- and rate-dependent reversible behavior (the stress–strain relation depends on the loading rate; stresses and strains evolve over time); internal variables, e.g., for viscoelasticity with  $n$  Maxwell elements are the inelastic strain contributions  $\mathbf{z} = \{\mathbf{e}_p^1, \dots, \mathbf{e}_p^n\}$ .
- **plasticity**, i.e., history-dependent irreversible behavior (the stress–strain relation depends on the loading history); internal variables are usually the plastic strains, accumulated plastic strains, and possibly further history variables:  $\mathbf{z} = \{\mathbf{e}_p, \epsilon_p, \dots\}$
- **viscoplasticity**, i.e., history- and time-dependent irreversible behavior; internal variables are similar to the case of plasticity above.
- **damage**, i.e., irreversible degradation of the elastic stiffness with loading; internal variable is a damage parameter, e.g., a scalar measure  $z = d$ .
- **ferroelectricity**, i.e., irreversible electro-mechanical-coupling; internal variable can be, e.g., the polarization field  $\mathbf{z} = \mathbf{p}$ .

All these phenomena can be described by the same underlying principles.

The general description of an inelastic (variational) material model starts with a strain energy density

$$W = W(\boldsymbol{\varepsilon}, \mathbf{z}), \quad (20.1)$$

where  $\mathbf{z}$  denotes a collection of internal variables. While the stress tensor and linear momentum balance remain untouched, i.e., (in linearized kinematics)

$$\boldsymbol{\sigma} = \frac{\partial W}{\partial \boldsymbol{\varepsilon}} \quad \text{and} \quad \operatorname{div} \boldsymbol{\sigma} + \rho \mathbf{b} = \mathbf{0}, \quad (20.2)$$

the internal variables evolve according to a kinetic law

$$\frac{\partial W}{\partial \mathbf{z}} + \frac{\partial \phi^*}{\partial \dot{\mathbf{z}}} \ni 0, \quad (20.3)$$

where  $\phi^*$  denotes the **dual (dissipation) potential** and we assume that such a potential exists. The differential inclusion is replaced by an equality in case of rate-dependent models. (20.3) can alternatively be cast into an effective variational problem:

$$\dot{\mathbf{z}} = \arg \inf \{ \dot{W} + \phi^* \}. \quad (20.4)$$

Let us introduce a discretization in time:  $t^\alpha = \alpha \Delta t$ , where we assume constant time steps  $\Delta t = t^{\alpha+1} - t^\alpha$  and, for conciseness, we write  $\Delta(\cdot) = (\cdot)^{\alpha+1} - (\cdot)^\alpha$ , where  $(\cdot)^\alpha$  denotes a quantity at time  $t^\alpha$ . Using simple backward-Euler rules then gives  $\dot{W} = (W^{\alpha+1} - W^\alpha)/\Delta t$  and  $\dot{\mathbf{z}} = (\mathbf{z}^{\alpha+1} - \mathbf{z}^\alpha)/\Delta t = \Delta \mathbf{z}/\Delta t$ .

Thus,

$$\frac{\mathbf{z}^{\alpha+1} - \mathbf{z}^\alpha}{\Delta t} = \arg \inf \left\{ \frac{W^{\alpha+1} - W^\alpha}{\Delta t} + \phi^* \left( \frac{\mathbf{z}^{\alpha+1} - \mathbf{z}^\alpha}{\Delta t} \right) \right\}. \quad (20.5)$$

Multiplication by  $\Delta t$  and omitting  $W^\alpha$  (since it does not depend on  $\mathbf{z}^{\alpha+1}$ ) leads to

$$\mathbf{z}^{\alpha+1} = \arg \inf \left\{ W^{\alpha+1}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) + \Delta t \phi^* \left( \frac{\mathbf{z}^{\alpha+1} - \mathbf{z}^\alpha}{\Delta t} \right) \right\}, \quad (20.6)$$

where the right-hand side defines an **effective incremental potential**:

$$\mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) = W^{\alpha+1}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) + \Delta t \phi^* \left( \frac{\mathbf{z}^{\alpha+1} - \mathbf{z}^\alpha}{\Delta t} \right) \quad (20.7)$$

Notice that  $\boldsymbol{\sigma}^{\alpha+1} = \partial W / \partial \boldsymbol{\varepsilon}^{\alpha+1}$  so that the effective potential can be used to replace the classical strain energy density in the total potential energy:

$$I_{\mathbf{z}^\alpha}[\mathbf{u}^{\alpha+1}, \mathbf{z}^{\alpha+1}] = \int_{\Omega} \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) dV - \int_{\Omega} \rho \mathbf{b} \cdot \mathbf{u}^{\alpha+1} dV - \int_{\partial\Omega_N} \hat{\mathbf{t}} \cdot \mathbf{u}^{\alpha+1} dS. \quad (20.8)$$

By the subscripts  $\mathbf{z}^\alpha$  we denote that those potentials do depend on  $\mathbf{z}^\alpha$  (i.e., the internal variables at the previous time step) but that those fields are known when evaluating the respective quantities.

The solution can now be found from

$$\{\mathbf{u}^{\alpha+1}, \mathbf{z}^{\alpha+1}\} = \arg \inf I_{\mathbf{z}^\alpha}[\mathbf{u}^{\alpha+1}, \mathbf{z}^{\alpha+1}]. \quad (20.9)$$

We can exploit that only the internal energy term depends on the internal variables and further assume that the energy density is *local* in the internal variables (this does not apply, e.g., for phase field models whose energy involves gradients of the internal variables). Then,

$$\begin{aligned} \inf_{\mathbf{u}^{\alpha+1}} \inf_{\mathbf{z}^{\alpha+1}} \int_{\Omega} \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) dV - \dots &= \inf_{\mathbf{u}^{\alpha+1}} \int_{\Omega} \inf_{\mathbf{z}^{\alpha+1}} \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) dV - \dots \\ &= \inf_{\mathbf{u}^{\alpha+1}} \int_{\Omega} W_{\mathbf{z}^\alpha}^*(\boldsymbol{\varepsilon}^{\alpha+1}) dV - \dots, \end{aligned} \quad (20.10)$$

where

$$W_{\mathbf{z}^\alpha}^*(\boldsymbol{\varepsilon}^{\alpha+1}) = \inf_{\mathbf{z}^{\alpha+1}} \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}^{\alpha+1}) = \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}_*^{\alpha+1}), \quad \mathbf{z}_*^{\alpha+1} = \arg \inf \mathcal{F}_{\mathbf{z}^\alpha}(\boldsymbol{\varepsilon}^{\alpha+1}, \cdot) \quad (20.11)$$

is often referred to as the **condensed energy density** (the internal variables have been “condensed out”). Notice that the omitted terms (...) do not depend on the internal variables.

This finally leads to the **incremental variational problem**

$$I_{\mathbf{z}^\alpha}[\mathbf{u}^{\alpha+1}, \mathbf{z}^{\alpha+1}] = \int_{\Omega} W_{\mathbf{z}^\alpha}^*(\boldsymbol{\varepsilon}^{\alpha+1}) dV - \int_{\Omega} \rho \mathbf{b} \cdot \mathbf{u}^{\alpha+1} dV - \int_{\partial\Omega_N} \hat{\mathbf{t}} \cdot \mathbf{u}^{\alpha+1} dS, \quad (20.12)$$

which has the same structure as before. This concept of introducing internal variables into the variational framework is also known as **variational constitutive updates** and goes back to Ortiz and Stainier (2000).

Note that – for numerical implementation purposes – evaluation of  $W_{\mathbf{z}^\alpha}^*(\boldsymbol{\varepsilon}^{\alpha+1})$  always requires us to compute the updated internal variables  $\mathbf{z}_*^{\alpha+1}$  based on  $\boldsymbol{\varepsilon}^{\alpha+1}$ , before the energy, stresses, or incremental stiffness matrix can be evaluated.

The stresses are now

$$\begin{aligned} \boldsymbol{\sigma}^{\alpha+1} &= \frac{d}{d\boldsymbol{\varepsilon}^{\alpha+1}} W_{\mathbf{z}^\alpha}^*(\boldsymbol{\varepsilon}^{\alpha+1}) \\ &= \frac{\partial \mathcal{F}}{\partial \boldsymbol{\varepsilon}^{\alpha+1}}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}_*^{\alpha+1}) + \frac{\partial \mathcal{F}}{\partial \mathbf{z}^{\alpha+1}}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}_*^{\alpha+1}) \cdot \frac{\partial \mathbf{z}_*^{\alpha+1}}{\partial \boldsymbol{\varepsilon}^{\alpha+1}} \\ &= \frac{\partial W^*}{\partial \boldsymbol{\varepsilon}^{\alpha+1}}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{z}_*^{\alpha+1}), \end{aligned} \quad (20.13)$$

where the second term vanished because  $\mathbf{z}_*^{\alpha+1}$  renders  $\mathcal{F}$  stationary by definition.

The same trick unfortunately does not apply when computing the incremental stiffness matrix,

$$\mathbb{C}_{ijkl}^{\alpha+1} = \frac{d\sigma_{ij}^{\alpha+1}}{d\varepsilon_{kl}^{\alpha+1}} = \frac{\partial\sigma_{ij}^{\alpha+1}}{\partial\varepsilon_{ij}^{\alpha+1}} + \frac{\partial\sigma_{ij}^{\alpha+1}}{\partial\mathbf{z}^{\alpha+1}} \cdot \frac{\partial\mathbf{z}_*^{\alpha+1}}{\partial\varepsilon_{ij}^{\alpha+1}}, \quad (20.14)$$

where the second term does *not* vanish in general. It requires calculating the sensitivity of the internal variables with respect to the strain tensor components.

## 20.2 Example: viscoelasticity, (visco)plasticity

Viscoelasticity and (visco)plasticity all start with the same fundamental structure (here presented in linearized kinematics). The total strain tensor *decomposes additively* into elastic and inelastic (or plastic) contributions:

$$\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}_e + \boldsymbol{\varepsilon}_p, \quad (20.15)$$

where  $\boldsymbol{\varepsilon}_p$  belongs to the set of internal variables. (Note that in finite deformations, the decomposition is multiplicative:  $\mathbf{F} = \mathbf{F}_e \mathbf{F}_p$ .)

In case of history dependence, one introduces additional internal history variables. For example for **von Mises plasticity**, the *accumulated plastic strain*  $\epsilon_p$  captures the history of plastic strains  $\boldsymbol{\varepsilon}_p$  through the coupling

$$\dot{\epsilon}_p = \|\dot{\boldsymbol{\varepsilon}}_p\|_{\text{vM}} = \sqrt{\frac{2}{3}} \dot{\boldsymbol{\varepsilon}}_p \cdot \dot{\boldsymbol{\varepsilon}}_p. \quad (20.16)$$

With internal variables  $\mathbf{z} = \{\boldsymbol{\varepsilon}_p, \epsilon_p\}$ , the Helmholtz energy density decomposes into elastic and plastic energy:

$$W(\boldsymbol{\varepsilon}, \boldsymbol{\varepsilon}_p, \epsilon_p) = W_{\text{el}}(\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_p) + W_{\text{pl}}(\epsilon_p). \quad (20.17)$$

In case of reversible behavior (viscoelasticity), we have no plastic energy storage, i.e.,  $W_{\text{pl}} = 0$ .

The dual dissipation potential can be defined, e.g., by the general power-law structure

$$\phi^*(\dot{\epsilon}_p) = \sigma_0 |\dot{\epsilon}_p| + \frac{\tau_0 \dot{\epsilon}_0}{m+1} \left( \frac{\dot{\epsilon}_p}{\dot{\epsilon}_0} \right)^{m+1} \quad (20.18)$$

with positive constants  $\sigma_0$  (initial yield stress),  $\tau_0$  (hardening rate),  $\dot{\epsilon}_0$  (reference rate), and  $m$  (rate sensitivity). In case of viscoelasticity we choose  $\sigma_0 = 0$ . By contrast, for rate-independent plasticity one chooses  $\tau_0 = 0$ .

The differential inclusion in (20.3) is required because of the first term,  $\sigma_0 |\dot{\epsilon}_p|$ , whose derivative is not defined at the origin (for  $\dot{\epsilon}_p = 0$ ). Here, a **subdifferential** is required:

$$\frac{\partial}{\partial \dot{\epsilon}_p} \sigma_0 |\dot{\epsilon}_p| = \begin{cases} \sigma_0, & \text{if } \dot{\epsilon}_p > 0, \\ -\sigma_0, & \text{if } \dot{\epsilon}_p < 0, \\ (-\sigma_0, \sigma_0) & \text{if } \dot{\epsilon}_p = 0. \end{cases} \quad (20.19)$$

In the first two cases (i.e. for *plastic* flow), the kinetic law becomes an equality; in the third case (i.e., for *elastic* loading such that  $|\sigma| < \sigma_0$ ), the differential inclusion is required.

The specific forms of  $W_{\text{el}}$ ,  $W_{\text{pl}}$ , and  $\phi^*$  depend on the particular material model.

### 20.3 Example: viscoplasticity

Since plastic/viscous deformation is observed to be isochoric, one commonly assumes

$$\text{tr } \boldsymbol{\varepsilon}_p = 0 \quad \text{so that} \quad \boldsymbol{\varepsilon}_p = \mathbf{e}_p. \quad (20.20)$$

Similarly, only the deviatoric stress tensor  $\mathbf{s}$  should cause plastic deformation. Here and in the following, we denote the **deviatoric** tensors by

$$\text{dev}(\cdot) = (\cdot) - \frac{1}{3} \text{tr}(\cdot) \mathbf{I}, \quad \text{so that} \quad \mathbf{e} = \text{dev } \boldsymbol{\varepsilon}, \quad \mathbf{s} = \text{dev } \boldsymbol{\sigma}. \quad (20.21)$$

Using the above definitions of energy density and dual potential, we obtain

$$\mathcal{F}_{\{\mathbf{e}_p^\alpha, \epsilon_p^\alpha\}}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{e}_p^{\alpha+1}, \epsilon_p^{\alpha+1}) = W^{\alpha+1}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{e}_p^{\alpha+1}, \epsilon_p^{\alpha+1}) + \Delta t \phi^* \left( \frac{\epsilon_p^{\alpha+1} - \epsilon_p^\alpha}{\Delta t} \right) \quad (20.22)$$

Minimization with respect to  $\mathbf{e}_p^{\alpha+1}$  gives

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \mathbf{e}_p^{\alpha+1}} &= \frac{\partial W_{\text{el}}^{\alpha+1}}{\partial \mathbf{e}_p^{\alpha+1}} + \frac{\partial W_{\text{pl}}^{\alpha+1}}{\partial \epsilon_p^{\alpha+1}} \frac{\partial \epsilon_p^{\alpha+1}}{\partial \mathbf{e}_p^{\alpha+1}} + \Delta t \frac{\partial \phi^*}{\partial \epsilon_p^{\alpha+1}} \frac{1}{\Delta t} \frac{\partial \epsilon_p^{\alpha+1}}{\partial \mathbf{e}_p^{\alpha+1}} \\ &= -\text{dev } \boldsymbol{\sigma}^{\alpha+1} + [\tau(\epsilon_p^{\alpha+1}) + \tau^*(\epsilon_p^{\alpha+1})] \frac{\partial \epsilon_p^{\alpha+1}}{\partial \mathbf{e}_p^{\alpha+1}} \ni 0 \end{aligned} \quad (20.23)$$

with back-stresses

$$\tau(\epsilon_p^{\alpha+1}) = \frac{\partial W_{\text{pl}}^{\alpha+1}}{\partial \epsilon_p^{\alpha+1}}, \quad \tau^*(\epsilon_p^{\alpha+1}) = \frac{\partial \phi^*}{\partial \epsilon_p^{\alpha+1}}. \quad (20.24)$$

Further, note that

$$\frac{\partial \epsilon_p^{\alpha+1}}{\partial \mathbf{e}_p^{\alpha+1}} = \frac{\partial \Delta \epsilon_p^{\alpha+1}}{\partial \Delta \mathbf{e}_p^{\alpha+1}} = \frac{\partial}{\partial \Delta \mathbf{e}_p^{\alpha+1}} \sqrt{\frac{2}{3} \Delta \mathbf{e}_p^{\alpha+1} \cdot \Delta \mathbf{e}_p^{\alpha+1}} = \frac{2 \Delta \mathbf{e}_p^{\alpha+1}}{3 \Delta \epsilon_p^{\alpha+1}}. \quad (20.25)$$

Altogether, this results in

$$-\mathbf{s}^{\alpha+1} + [\tau(\epsilon_p^{\alpha+1}) + \tau^*(\epsilon_p^{\alpha+1})] \frac{2 \Delta \mathbf{e}_p^{\alpha+1}}{3 \Delta \epsilon_p^{\alpha+1}} \ni 0 \quad (20.26)$$

or

$$\mathbf{s}^{\alpha+1} \in [\tau(\epsilon_p^{\alpha+1}) + \tau^*(\epsilon_p^{\alpha+1})] \frac{2 \Delta \mathbf{e}_p^{\alpha+1}}{3 \Delta \epsilon_p^{\alpha+1}}. \quad (20.27)$$

Let us consider the *elastically isotropic* case (i.e.,  $\boldsymbol{\sigma} = \kappa(\text{tr } \boldsymbol{\varepsilon}) \mathbf{I} + 2\mu \mathbf{e}_e$ ). Assume that  $\Delta \epsilon_p^{\alpha+1} > 0$ , so we have

$$2\mu (\mathbf{e}^{\alpha+1} - \mathbf{e}_p^{\alpha+1}) = \left[ \tau(\epsilon_p^{\alpha+1}) + \sigma_0 + \tau_0 \left( \frac{\Delta \epsilon_p^{\alpha+1}}{\dot{\epsilon}_0 \Delta t} \right)^m \right] \frac{2 \Delta \mathbf{e}_p^{\alpha+1}}{3 \Delta \epsilon_p^{\alpha+1}}. \quad (20.28)$$

and we introduce an **elastic predictor** (i.e., strain if the plastic strain remained unaltered)

$$\mathbf{e}_{\text{pre}}^{\alpha+1} = \mathbf{e}^{\alpha+1} - \mathbf{e}_p^\alpha \quad \text{such that} \quad \mathbf{e}^{\alpha+1} - \mathbf{e}_p^{\alpha+1} = \mathbf{e}_{\text{pre}}^{\alpha+1} - \Delta \mathbf{e}_p^{\alpha+1}. \quad (20.29)$$

That gives

$$2\mu \mathbf{e}_{\text{pre}}^{\alpha+1} = 2\mu \Delta \mathbf{e}_{\text{p}}^{\alpha+1} + \frac{2}{3\Delta \epsilon_{\text{p}}^{\alpha+1}} \left[ 3\mu \Delta \epsilon_{\text{p}}^{\alpha+1} + \tau(\epsilon_{\text{p}}^{\alpha+1}) + \sigma_0 + \tau_0 \left( \frac{\Delta \epsilon_{\text{p}}^{\alpha+1}}{\dot{\epsilon}_0 \Delta t} \right)^m \right] \Delta \mathbf{e}_{\text{p}}^{\alpha+1}. \quad (20.30)$$

or

$$\frac{2\mu \mathbf{e}_{\text{pre}}^{\alpha+1}}{2\mu + \frac{2}{3\Delta \epsilon_{\text{p}}^{\alpha+1}} \left( \tau(\epsilon_{\text{p}}^{\alpha+1}) + \sigma_0 + \tau_0 \left( \frac{\Delta \epsilon_{\text{p}}^{\alpha+1}}{\dot{\epsilon}_0 \Delta t} \right)^m \right)} = \Delta \mathbf{e}_{\text{p}}^{\alpha+1}. \quad (20.31)$$

Now, let us use that

$$\Delta \epsilon_{\text{p}}^{\alpha+1} = \sqrt{\frac{2}{3} \Delta \mathbf{e}_{\text{p}}^{\alpha+1} \cdot \Delta \mathbf{e}_{\text{p}}^{\alpha+1}} = \frac{2\mu \|\mathbf{e}_{\text{pre}}^{\alpha+1}\|_{\text{vM}}}{2\mu + \frac{2}{3\Delta \epsilon_{\text{p}}^{\alpha+1}} \left( \tau(\epsilon_{\text{p}}^{\alpha+1}) + \sigma_0 + \tau_0 \left( \frac{\Delta \epsilon_{\text{p}}^{\alpha+1}}{\dot{\epsilon}_0 \Delta t} \right)^m \right)}. \quad (20.32)$$

This is a scalar equation to be solved for the increment  $\Delta \epsilon_{\text{p}}^{\alpha+1}$ , which is then inserted into (20.31) in order to determine  $\Delta \mathbf{e}_{\text{p}}^{\alpha+1}$ . This completes the calculation of the updated internal variables.

Notice that the above equation is equivalent to saying (introducing the **von Mises stress**  $\sigma_{\text{vM}}$ )

$$\sigma_{\text{vM}} = \sqrt{\frac{3}{2} \mathbf{s}^{\alpha+1} \cdot \mathbf{s}^{\alpha+1}} = \tau(\epsilon_{\text{p}}^{\alpha+1}) + \tau^*(\epsilon_{\text{p}}^{\alpha+1}). \quad (20.33)$$

Analogous relations can be obtained for  $\Delta \epsilon_{\text{p}}^{\alpha+1} < 0$ .

Note that rate-independent plasticity (including the special case of **von Mises plasticity**) assumes that  $\tau_0 = 0$ ; i.e., the dissipation potential only provides the initial yield threshold  $\sigma_0$ .

## 20.4 Example: linear viscoelasticity

In the simplest viscoelastic case, we take  $W_{\text{pl}} = 0$  and thus  $\tau(\epsilon_{\text{p}}^{\alpha+1}) = 0$ , i.e., the material has no “memory”. Also, the yield threshold is removed by choosing  $\sigma_0 = 0$ , and the accumulated plastic stress  $\epsilon_{\text{p}}$  is of no interest nor needed. The dissipation is reformulated as

$$\phi^*(\dot{\epsilon}_{\text{p}}) = \frac{\eta}{2} \|\dot{\epsilon}_{\text{p}}\|^2, \quad (20.34)$$

which agrees with the plastic case above with velocity-proportional damping ( $m = 1$ ), viscosity  $\eta = \tau_0/\dot{\epsilon}_0$ , and the von Mises norm replaced by the classical vector norm.

In this case, (20.28) reduces to

$$2\mu (\mathbf{e}_{\text{pre}}^{\alpha+1} - \Delta \mathbf{e}_{\text{p}}) = \tau_0 \frac{\Delta \mathbf{e}_{\text{p}}}{\dot{\epsilon}_0 \Delta t} = \eta \frac{\Delta \mathbf{e}_{\text{p}}}{\Delta t}, \quad (20.35)$$

which can be rearranged to

$$2\mu \mathbf{e}_{\text{pre}}^{\alpha+1} = \eta \frac{\Delta \mathbf{e}_{\text{p}}}{\Delta t} + 2\mu \Delta \mathbf{e}_{\text{p}} \quad (20.36)$$

giving

$$2\mathbf{e}_{\text{pre}}^{\alpha+1} = \left( \frac{\eta/\mu}{\Delta t} + 2 \right) \Delta \mathbf{e}_{\text{p}}. \quad (20.37)$$



If we define the relaxation time  $\tau = \eta/\mu$ , then

$$\Delta \mathbf{e}_p = \frac{2}{2 + \tau/\Delta t} \mathbf{e}_{\text{pre}}^{\alpha+1} \quad (20.38)$$

or

$$\mathbf{e}_p^{\alpha+1} = \mathbf{e}_p^\alpha + \frac{2}{2 + \tau/\Delta t} (\mathbf{e}^{\alpha+1} - \mathbf{e}_p^\alpha). \quad (20.39)$$

This can be extended to the **generalized Maxwell model**. Let us assume isotropic elasticity with shear and bulk modulus  $\mu_\infty$  and  $\kappa_\infty$ , respectively, while the  $n$  viscoelastic branches are characterized by shear stiffnesses  $\mu_i$  and viscosities  $\eta_i$  (for  $i = 1, \dots, n$ ). The effective incremental energy density now becomes

$$\begin{aligned} \mathcal{F}_{\{\mathbf{e}_p^{i,\alpha}\}}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{e}_p^{1,\alpha+1}, \dots, \mathbf{e}_p^{n,\alpha+1}) &= \frac{\kappa_\infty}{2} (\text{tr } \boldsymbol{\varepsilon}^{\alpha+1})^2 + \mu_\infty \mathbf{e}^{\alpha+1} \cdot \mathbf{e}^{\alpha+1} + \sum_{i=1}^n \mu_i \|\mathbf{e}^{\alpha+1} - \mathbf{e}_p^{i,\alpha+1}\|^2 \\ &\quad + \sum_{i=1}^n \frac{\eta_i}{2\Delta t} \|\mathbf{e}_p^{i,\alpha+1} - \mathbf{e}_p^{i,\alpha}\|^2, \end{aligned} \quad (20.40)$$

whose minimization with respect to the new internal variables yields

$$\mathbf{e}_p^{i,\alpha+1} = \mathbf{e}_p^{i,\alpha} + \frac{2}{2 + \tau_i/\Delta t} (\mathbf{e}^{\alpha+1} - \mathbf{e}_p^{i,\alpha}) \quad (20.41)$$

with relaxation times  $\tau_i = \eta_{(i)}/\mu_{(i)}$ . Note that this agrees with (20.39) for a single Maxwell element.

Insertion and differentiation leads to the stress tensor

$$\boldsymbol{\sigma}^{\alpha+1}(\boldsymbol{\varepsilon}^{\alpha+1}, \mathbf{e}_p^{1,\alpha}, \dots, \mathbf{e}_p^{n,\alpha}) = \kappa_\infty (\text{tr } \boldsymbol{\varepsilon}^{\alpha+1}) \mathbf{I} + 2\mu_\infty \mathbf{e}^{\alpha+1} + \sum_{i=1}^n 2\mu_i \frac{\tau_i/\Delta t}{2 + \tau_i/\Delta t} (\mathbf{e}^{\alpha+1} - \mathbf{e}_p^{i,\alpha}). \quad (20.42)$$

Similarly, the consistent incremental tangent matrix can be computed by differentiating  $\boldsymbol{\sigma}^{\alpha+1}$ :

$$\begin{aligned} \mathbb{C}_{ijkl}^{\alpha+1} = \frac{\partial \boldsymbol{\sigma}^{\alpha+1}}{\partial \boldsymbol{\varepsilon}^{\alpha+1}} &= \left[ \kappa_\infty - \frac{2}{3} \left( \mu_\infty + \sum_{i=1}^n \mu_i \frac{\tau_i/\Delta t}{2 + \tau_i/\Delta t} \right) \right] \delta_{ij} \delta_{kl} \\ &\quad + \left( \mu_\infty + \sum_{i=1}^n \mu_i \frac{\tau_i/\Delta t}{2 + \tau_i/\Delta t} \right) (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}). \end{aligned} \quad (20.43)$$

Finally, the condensed incremental energy density  $W^*$  can be computed analytically by inserting (20.41) into (20.40).

## Index

- $L_2$ -norm, 95
- $n$ -dimensional space, 26
- (proper) subset, 93
- 2-node bar element, 37
- 2-node beam element, 39
- 4-node bilinear quadrilateral, 44
- 4-node tetrahedron, 42
- 8-node brick element, 46
- 8-node quadratic quadrilateral, 46
- 9-node quadratic quadrilateral, 45
  
- a.e., 24, 94
- acceleration, 9
- action, 72
- action principle, 72
- approximation error, 67
- assembly operator, 55
- average acceleration, 81
  
- backward-Euler, 16
- balance laws, 7
- Banach space, 96
- barycentric coordinates, 40
- basis, 33
- bijective, 93
- bilinear form, 22
- bilinear operator, 104
- boundary, 93
- boundary conditions, 15
- Bubnov-Galerkin approximation, 27
  
- Cauchy stress tensor, 11
- chicken-wire mode, 71
- class  $C^k(\Omega)$ , 24, 94
- classical solution, 21, 22
- closure, 96
- complete polynomial, 102
- complete space, 96
- complete up to order  $q$ , 34
- completeness property, 34
- condensation method, 66
- condensed energy density, 85
- condition number, 67
- conjugate gradient, 60
- consistent mass matrix, 73
- Constant Strain Triangle, 41
- constitutive relations, 6
- continuous at a point  $x$ , 94
- continuous over  $\Omega$ , 94
- convergence, 96
- CST, 41
- cubature rules, 47
- current configuration, 9
  
- damped Newton-Raphson method, 59
- damping matrix, 74
- deformation gradient, 9
- deformation mapping, 9
- degree, 101
- degrees of freedom, 35
- deviatoric, 87
- Direct methods, 16
- Dirichlet boundary, 15
- Dirichlet-Poincaré inequality, 98
- discrete problem, 26
- discrete weak form, 73
- discretization error, 67
- displacement field, 9
- distance, 95, 96
- divergence theorem, 7
- domain, 93
- dual (dissipation) potential, 84
  
- effective incremental potential, 85
- eigenfrequency, 75
- eigenmode, 75
- eigenvector, 75
- elasticity tensor, 13
- elements, 35
- energy norm, 107
- essential supremum, 95
- Euclidean norm, 95
- Eulerian, 9
- explicit, 17
- explicit time integration, 79
- extensive, 7
- external force elements, 62
  
- Fast Inertial Relaxation Engine, 61
- finite differences, 16
- finite element, 35
- Finite Element Method, 35
- first fundamental error, 67
- First Piola-Kirchhoff, 10
- first Piola-Kirchhoff stress tensor, 11
- first variation, 19

- first-order central difference, 16
- flux, 6
- forward-Euler, 16
- full integration, 53
- function, 93
- functional, 18
  
- Gauss quadrature, 48
- Gauss-Chebyshev, 50
- Gauss-Hermite, 50
- Gauss-Legendre quadrature, 48, 50
- Gauss-Lobatto quadrature, 51
- Gauss-Newton method, 60
- generalized Maxwell model, 89
- global, 33
- global error estimate, 98, 100
- gradient flow method, 60
- Gram-Schmidt orthogonalization, 49
  
- h-refinement, 35
- heat equation, 8
- Hermitian polynomials, 39
- hierarchical interpolation, 38
- higher-order interpolation, 100
- Hilbert space, 96
- homogeneous, 8
- hourglass mode, 71
- hp-refinement, 36
  
- identity mapping, 93
- implicit, 17
- implicit time integration, 81
- incremental tangent modulus tensor, 11
- incremental variational problem, 85
- indirect methods, 17
- initial boundary value problem, 15
- initial conditions, 15
- injective, 93
- inner product, 94
- inner product space, 94
- intensive, 7
- internal energy, 7
- inverse function theorem, 41
- isomorphism, 93
- isoparametric, 39
- isoparametric mapping, 39
- isotropy, 8
  
- Jacobian, 40
  
- kinematic variables, 6
  
- $L_2$ -inner product, 94
  
- $L_2$ -space of functions, 96
- $L_p$ -norm, 95
- Lagrange polynomials, 38
- Lagrangian, 9
- Lagrangian interpolation, 38
- Lax-Milgram theorem, 23, 105
- Legendre polynomials, 49
- line search method, 59
- linear, 104
- linear acceleration, 81
- linear elasticity, 13, 29
- linear form, 22
- linear strain triangle, 43
- linear subspace, 93
- linear/vector space, 93
- linearized kinematics, 13
- local, 33
- local energy balance, 7
- local error estimate, 98
- locking, 71
- longest edge bisection, 69
  
- mapping, 93
- mappings, 6
- material points, 6
- mesh, 35
- mesh refinement algorithm, 69
- mesh refinement criterion, 69
- methods of weighted residuals, 25
- modal decomposition, 78
- modeling error, 68
- monomial, 101
- multi-index, 101
  
- Navier's equation, 14
- neighborhood, 96
- Neumann boundary, 15
- Neumann-Poincaré, 98
- Newmark- $\beta$ , 81
- Newton-Cotes, 47
- Newton-Raphson method, 58
- nodes, 35
- nonlinear least squares, 60
- norm, 95
- normed linear space, 95
- numerical integration, 47
- numerical integration error, 67
  
- one-to-one, 93
- onto, 93
- open, 96
- open set, 93

- operator, 104
- order of a PDE, 15
- order reduction, 78
- ordered triad, 93
- orthogonal, 94
- over-integration, 53
  
- p-refinement, 35
- parasitic shear, 71
- Petrov-Galerkin, 27
- Poincaré inequalities, 98
- positive, 104
- predictor, 87
- principle of virtual work, 25
  
- Q4, 44
- Q8, 46
- Q9, 45
- quadratic tetrahedron, 43
- quadratic triangle, 42
- quadrature error, 52
- quadrature rules, 47
- Quasi-Newton method, 59
- quasistatics, 10
  
- r-refinement, 36
- range, 93
- Rayleigh's quotient, 78
- Rayleigh-Ritz, 21, 29
- reference configuration, 9
- Riemann sum, 47
- right Cauchy-Green tensor, 10
- rigid body rotation, 13
- rigid body translation, 13
- rigid-body modes, 66
  
- second-order central difference, 16
- selective integration, 71
- semi-discretization, 73
- semi-norm, 99
- serendipity element, 46
- set, 93
- shape function properties, 33
- shape functions, 33
- shear locking, 71
- simplex, 40
- simplicial quadrature, 53
- smoothing, 68
- Sobolev norm, 18, 100
- Sobolev semi-norm, 99
- Sobolev space, 18, 101, 102
- solution error, 67
  
- space of all second-order polynomial functions, 93
- square-integrable, 18, 96, 101
- stationarity condition, 19
- stencils, 17
- strain energy density, 11
- strain tensor, 13
- stress, 10
- stretch, 9
- strong form, 24, 28
- subdifferential, 86
- subparametric, 39
- superparametric, 39
- support, 33, 103
- surjective, 93
- symmetric, 104
  
- tangent matrix, 58
- Taylor expansions, 16
- traction, 10
- triangle inequality, 95
- truncation error, 68
  
- under-integration, 53
  
- Vainberg's theorem, 106
- variation, 19
- variational constitutive updates, 85
- variational structure, 20
- velocity, 9
- volume change, 9
- von Mises plasticity, 88
- von Mises stress, 88
  
- weak form, 25, 72
- weak solution, 25
- Weierstrass approximation theorem, 34
  
- zero-energy mode, 71
- zero-energy modes, 66
- zeroth derivative, 24, 94
- ZZ error estimator, 69

## A Introduction, Vector Spaces

We define a **set**  $\Omega$  as a collection of points  $\mathbf{X} \in \Omega$ .

We further say  $\Omega$  is a **(proper) subset** of all space if  $\Omega \subseteq \mathbb{R}^d$  in  $d$  dimensions (*proper* if  $\Omega \subset \mathbb{R}^d$ ).

We usually take  $\Omega$  to be an **open set**, i.e.,  $\Omega \cap \partial\Omega = \emptyset$  with **boundary**  $\partial\Omega$ .

Deformation and motion are described by a **mapping**

$$\varphi : \mathbf{X} \in \Omega \rightarrow \varphi(\mathbf{X}) \in \mathbb{R}^d \quad \text{or} \quad \varphi : \Omega \rightarrow \mathbb{R}^d, \quad (\text{A.1})$$

where  $\Omega$  is the **domain** and  $\mathbb{R}^d$  the **range** of  $\varphi$ . The mapped (current) configuration of  $\Omega$  is  $\varphi(\Omega)$ .

Every **function**  $f(x) : \mathbb{R} \rightarrow \mathbb{R}$  is a mapping from  $\mathbb{R}$  to  $\mathbb{R}$ .

We call a mapping **injective** (or **one-to-one**) if for each  $\mathbf{x} \in \varphi(\Omega)$  there is one unique  $\mathbf{X} \in \Omega$  such that  $\mathbf{x} = \varphi(\mathbf{X})$ . In other words, no two points  $\mathbf{X} \in \Omega$  are mapped onto the same position  $\mathbf{x}$ . A mapping is **surjective** (or **onto**) if the entire set  $\Omega$  is mapped onto the entire set  $\varphi(\Omega)$ ; i.e., for every  $\mathbf{X} \in \Omega$  there exists at least one  $\mathbf{x} \in \varphi(\Omega)$  such that  $\mathbf{x} = \varphi(\mathbf{X})$ . If a mapping is both injective and surjective (or one-to-one and onto) we say it is **bijective**. A bijective mapping is also called an **isomorphism**.

For example,  $\varphi : \Omega \rightarrow \mathbb{R}^d$  is injective, whereas  $\varphi : \Omega \rightarrow \varphi(\Omega)$  would be bijective.

For time-dependent problems, we have  $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}^d$  and  $\mathbf{x} = \varphi(\mathbf{X}, t)$ . This describes a family of configurations  $\varphi(\Omega, t)$ , from which we arbitrarily define a reference configuration  $\Omega$  for which  $\varphi = \text{id}$  (the **identity mapping**).

A **linear/vector space**  $\{\Omega, +; \mathbb{R}, \cdot\}$  is defined by the following identities. For any  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \Omega$  and  $\alpha, \beta \in \mathbb{R}$  it holds that

- (i) *closure*:  $\alpha \cdot \mathbf{u} + \beta \cdot \mathbf{v} \in \Omega$
- (ii) *associativity w.r.t. +*:  $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$
- (iii) *null element*:  $\exists \mathbf{0} \in \Omega$  such that  $\mathbf{u} + \mathbf{0} = \mathbf{u}$
- (iv) *negative element*: for all  $\mathbf{u} \in \Omega \quad \exists -\mathbf{u} \in \Omega$  such that  $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$
- (v) *commutativity*:  $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
- (vi) *associativity w.r.t.  $\cdot$* :  $(\alpha\beta) \cdot \mathbf{u} = \alpha(\beta \cdot \mathbf{u})$
- (vii) *distributivity w.r.t.  $\mathbb{R}$* :  $(\alpha + \beta) \cdot \mathbf{u} = \alpha \cdot \mathbf{u} + \beta \cdot \mathbf{u}$
- (viii) *distributivity w.r.t.  $\Omega$* :  $\alpha \cdot (\mathbf{u} + \mathbf{v}) = \alpha \cdot \mathbf{u} + \alpha \cdot \mathbf{v}$
- (ix) *identity*:  $1 \cdot \mathbf{u} = \mathbf{u}$

**examples:**

- $\mathbb{R}^d$  is a vector space. By contrast,  $\Omega \subset \mathbb{R}^d$  is *not* a vector space, since – in general – it violates, e.g., (i) closure and (iv) null element.
- $\mathbb{P}_2 = \{ax^2 + bx + c : a, b, c \in \mathbb{R}\}$  is the **space of all second-order polynomial functions**, or an **ordered triad**  $(a, b, c) \in \mathbb{R}^3$ . More generally,  $\mathbb{P}_k(\Omega)$  is the space of all  $k$ th-order polynomial functions defined on  $\Omega$ .  $\mathbb{P}_k(\Omega)$  is a *linear space*.

We call  $\mathbb{P}_2$  a **linear subspace** of  $\mathbb{P}_k$  with  $k \geq 2$ , and we write  $\mathbb{P}_2 \subseteq \mathbb{P}_k$ .

## B Function Spaces

Consider a function  $u(x) : \Omega \rightarrow \mathbb{R}$  and  $\Omega \subset \mathbb{R}$ .

$u$  is **continuous at a point  $x$**  if, given any scalar  $\epsilon > 0$ , there is a  $r(\epsilon) \in \mathbb{R}$  such that

$$|u(y) - u(x)| < \epsilon \quad \text{provided that} \quad |y - x| < r. \quad (\text{B.1})$$

A function  $u$  is **continuous over  $\Omega$**  if it is continuous at all points  $x \in \Omega$ .

$u$  is of **class  $C^k(\Omega)$**  with an integer  $k \geq 0$  if it is  $k$  times continuously differentiable over  $\Omega$  (i.e.,  $u$  possesses derivatives up to the  $k$ th order and these derivatives are continuous functions).

**Examples:**

- Functions  $u(x) \in \mathbb{P}_k$  with  $k \geq 0$  are generally  $C^\infty(\mathbb{R})$ .
- Consider a continuous, piecewise-linear function  $u : \Omega = (0, 2) \rightarrow \mathbb{R}$ . Function  $u$  is  $C^0(\Omega)$  but not  $C^1(\Omega)$ .
- The Heaviside function  $H(x)$  is said to be  $C^{-1}(\mathbb{R})$  since its “**zeroth derivative**” (i.e., the function itself) is not continuous.

If there are no discontinuities such as cracks, shocks, etc. (or discontinuities in the BCs/ICs) we usually assume the solution fields are  $C^\infty(\Omega)$ , so we may take derivatives; otherwise, derivatives exist almost everywhere (**a.e.**)

To evaluate the global errors of functions, we need *norms*.

Consider a linear space  $\{\mathcal{U}, +, \mathbb{R}, \cdot\}$ . A mapping  $\langle \cdot, \cdot \rangle : \mathcal{U} \times \mathcal{U} \rightarrow \mathbb{R}$  is called **inner product** on  $\mathcal{U} \times \mathcal{U}$  if for all  $u, v, w \in \mathcal{U}$  and  $\alpha \in \mathbb{R}$ :

- (i)  $\langle u + v, w \rangle = \langle u, w \rangle + \langle v, w \rangle$
- (ii)  $\langle u, v \rangle = \langle v, u \rangle$
- (iii)  $\langle \alpha \cdot u, v \rangle = \alpha \langle u, v \rangle$
- (iv)  $\langle u, u \rangle \geq 0$  and  $\langle u, u \rangle = 0 \Leftrightarrow u = 0$

A linear space  $\mathcal{U}$  endowed with an inner product is called an **inner product space**.

**Examples:**

- $\langle \mathbf{u}, \mathbf{v} \rangle = u_i v_i = \mathbf{u} \cdot \mathbf{v}$  defines an inner product for  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ .
- The  **$L_2$ -inner product** for functions  $u, v \in \mathcal{U}$  with domain  $\Omega$ :

$$\langle u, v \rangle_{L_2(\Omega)} = \int_{\Omega} u(x) v(x) dx \quad \text{and often just} \quad \langle u, v \rangle = \langle u, v \rangle_{L_2(\Omega)}. \quad (\text{B.2})$$

Note that if  $\langle u, v \rangle = 0$  we say  $u$  and  $v$  are **orthogonal**.

**Examples:**

- *Legendre polynomials*:

$$p_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n \quad \text{so that} \quad p_0 = 1, \quad p_1 = x, \quad p_2 = \frac{1}{2}(3x^2 - 1), \dots \quad (\text{B.3})$$

orthogonality on  $\Omega = (-1, 1)$ :

$$\int_{-1}^1 p_n(x) p_m(x) dx = \frac{2}{2n + 1} \delta_{mn} \quad (\text{B.4})$$

- *trigonometric functions:*

$$p_n(x) = \cos\left(\frac{\pi n x}{L}\right) \quad (\text{B.5})$$

orthogonality on  $\Omega = (-L, L)$

$$\int_{-L}^L p_n(x) p_m(x) dx = \begin{cases} 2L, & \text{if } m = n = 0 \\ L, & \text{if } m = n \neq 0 \\ 0, & \text{else} \end{cases} \quad (\text{B.6})$$

Now we are in place to define the **distance** between  $x_1$  and  $x_2$ :

$$d(x_1, x_2) = \sqrt{\langle x_1 - x_2, x_1 - x_2 \rangle} \quad (\text{B.7})$$

We need this concept not only for points in space but also to define the closeness or proximity of functions.

Consider a linear space  $\{\mathcal{U}, +; \mathbb{R}, \cdot\}$ . A mapping  $\|\cdot\| : \mathcal{U} \rightarrow \mathbb{R}_+$  is called a **norm** on  $\mathcal{U}$  if for all  $u, v \in \mathcal{U}$  and  $\alpha \in \mathbb{R}$ :

- (i)  $\|u + v\| \leq \|u\| + \|v\|$  (**triangle inequality**)
- (ii)  $\|\alpha \cdot u\| = |\alpha| \|u\|$
- (iii)  $\|u\| \geq 0$  and  $\|u\| = 0 \Leftrightarrow u = 0$ .

A linear space  $\Omega$  endowed with a norm is called a **normed linear space** (NLS).

**Examples** of norms:

- Consider the  $d$ -dimensional Euclidean space, so  $\mathbf{x} = \{x_1, \dots, x_d\}^T$ . Then we define

- the 1-norm:  $\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$
- the 2-norm:  $\|\mathbf{x}\|_2 = \left(\sum_{i=1}^d |x_i|^2\right)^{1/2}$  (**Euclidean norm**)
- the  $n$ -norm:  $\|\mathbf{x}\|_n = \left(\sum_{i=1}^d |x_i|^n\right)^{1/n}$
- the  $\infty$ -norm:  $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

- Now turning to functions, the  **$L_p$ -norm** of a function  $u : \Omega \rightarrow \mathbb{R}$ :

$$\|u\|_{L_p(\Omega)} = \left(\int_{\Omega} u^p dx\right)^{1/p} \quad (\text{B.8})$$

The most common norm is the  **$L_2$ -norm**:

$$\|u\|_{L_2(\Omega)} = \langle u, u \rangle_{L_2(\Omega)}^{1/2} = \left(\int_{\Omega} u^2(x) dx\right)^{1/2}. \quad (\text{B.9})$$

Furthermore, notice that

$$\|u\|_{L_\infty(\Omega)} = \operatorname{ess\,sup}_{x \in \Omega} |u(x)|, \quad (\text{B.10})$$

where we introduced the **essential supremum**

$$\operatorname{ess\,sup}_{x \in \Omega} |u(x)| = M \quad \text{with the smallest } M \text{ that satisfies } |u(x)| \leq M \text{ for a.e. } x \in \Omega. \quad (\text{B.11})$$

Now, that we have norms, we can generalize our definition of the distance. If  $u_n, u \in \mathcal{U}$  equipped with a norm  $\|\cdot\| : \mathcal{U} \rightarrow \mathbb{R}$ , then we define the **distance** as

$$d(u_n, u) = \|u_n - u\|. \quad (\text{B.12})$$

Now, we are in place to define the **convergence** of a sequence of functions  $u_n$  to  $u$  in  $\mathcal{U}$ : we say  $u_n \rightarrow u \in \mathcal{U}$  if for all  $\epsilon > 0$  there exists  $N(\epsilon)$  such that  $d(u_n, u) < \epsilon$  for all  $n > N$ .

### Examples:

- Consider  $u_n \in \mathcal{U} = \mathbb{P}_2(\Omega)$  with  $L_2$ -norm and  $\Omega \subset \mathbb{R}$

$$u_n(x) = \left(1 + \frac{1}{n}\right)x^2 \quad \rightarrow \quad u(x) = x^2 \quad \text{since} \quad d(u_n - u) = \frac{1}{n} \int_{\Omega} x^2 dx \quad (\text{B.13})$$

with  $u \in \mathcal{U} = \mathbb{P}_2(\Omega)$ . For example, for  $d(u_n - u) < \epsilon$  we need  $n > N = \int_{\Omega} x^2 dx / \epsilon$ .

- Fourier series:*

$$u(x) = \sum_{i=0}^{\infty} c_i x^i \quad \Rightarrow \quad u_n(x) = \sum_{i=0}^n c_i x^i \quad \text{such that} \quad u_n \rightarrow u \text{ as } n \rightarrow \infty. \quad (\text{B.14})$$

Given a point  $u$  in a normed linear space  $\mathcal{U}$ , a **neighborhood**  $\mathcal{N}_r(u)$  of radius  $r > 0$  is defined as the set of points  $v \in \mathcal{U}$  for which  $d(u, v) < r$ . Now, we can define sets properly:

A subset  $\mathcal{V} \subset \mathcal{U}$  is called **open** if, for each point  $u \in \mathcal{V}$ , there exists a neighborhood  $\mathcal{N}_r(u)$  which is fully contained in  $\mathcal{V}$ . The complement  $\widetilde{\mathcal{V}}$  of an open set  $\mathcal{V}$  is, by definition a closed set. The **closure**  $\overline{\mathcal{V}}$  of an open set  $\mathcal{V}$  is the smallest closed set that contains  $\mathcal{V}$ . In simple terms, a closed set is defined as a set which contains all its limit points. Therefore, note that

$$\sup_{x \in \Omega} |u(x)| = \max_{x \in \overline{\Omega}} |u(x)|. \quad (\text{B.15})$$

For example,  $(0, 1)$  is an open set in  $\mathbb{R}$ .  $[0, 1]$  is a closed set, and  $[0, 1]$  is the closure of  $(0, 1)$ .

A linear space  $\mathcal{U}$  is a **complete space** if every sequence  $u_n$  in  $\mathcal{U}$  converges to  $u \in \mathcal{U}$ . In simple terms, the space must contain all *limit points*.

A *complete normed linear space* is called a **Banach space**; i.e.,  $\{\mathcal{U}, +; \mathbb{R}, \cdot\}$  with a norm  $\|\cdot\|$  and  $u_n \rightarrow u \in \mathcal{U}$ . A *complete inner product space* is called a **Hilbert space**.

Note that  $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$  defines a norm. Hence, every Hilbert space is also a Banach space (but not the other way around).

As an example, consider  $\mathcal{U} = \mathbb{P}_n$  (the space of all polynomial functions of order  $n \in \mathbb{N}$ ). This is a *linear space* which we equip with a norm, e.g., the  $L_2$ -norm. It is *complete* since  $(a_n, b_n, c_n, \dots) \rightarrow (a, b, c, \dots)$  for  $a, b, c, \dots \in \mathbb{R}$ . And an *inner product* is defined via  $\langle u, v \rangle = \int_{\Omega} uv \, dx$ . With all these definitions,  $\mathcal{U}$  is a Hilbert space.

We can use these norms to define function spaces, e.g., the  **$L_2$ -space of functions**:

$$L_2(\Omega) = \left\{ u : \Omega \rightarrow \mathbb{R} : \int_{\Omega} u^2 \, dx < \infty \right\} \quad (\text{B.16})$$

We say  $L_2(\Omega)$  contains all functions that are **square-integrable** on  $\Omega$ .



**Examples:**

- $u : \Omega \rightarrow \mathbb{R}$  with  $u \in \mathbb{P}_k(\Omega)$  and  $\text{ess sup}_{x \in \Omega} |u(x)| < \infty$ . Then,  $u \in L_2(\Omega)$ .
- $f : \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = x^{-2}$  is not in  $L_2(\Omega)$  if  $0 \in \Omega$ .
- Piecewise constant functions  $u$  (with  $\text{ess sup}_{x \in \Omega} |u(x)| < \infty$ ) are square-integrable and thus in  $L_2$ .

Note that we can write alternatively

$$L_2(\Omega) = \left\{ u : \Omega \rightarrow \mathbb{R} : \|u\|_{L_2(\Omega)} < \infty \right\}. \quad (\text{B.17})$$

## C Approximation Theory

**Motivation:** in computational mechanics, we seek approximate solutions  $u_h(x) = \sum_{a=1}^N u_a N_a(x)$ , e.g., a linear combination of basis functions  $N_a(x)$  with amplitudes  $u_a \in \mathbb{R}$ .

**Questions:** How does  $u_h(x)$  converge to  $u(x)$ , if at all? Can we find an error estimate  $\|u_h - u\|$ ? What is the rate of convergence (how fast does it converge, cf. the truncation error arguments for grid-based direct methods)?

Fundamental tools for estimating errors are the **Poincaré inequalities**:

(i) **Dirichlet-Poincaré inequality**:

$$\int_0^h |v(x)|^2 dx \leq c_h \int_0^h [v'(x)]^2 dx \quad \text{if} \quad v(0) = v(h) = 0. \quad (\text{C.1})$$

with a constant  $c_h > 0$  that depends on the interval size  $h$ .

(ii) **Neumann-Poincaré** (or Poincaré-Wirtinger) inequality:

$$\int_0^h |v(x) - \bar{v}|^2 dx \leq c_h \int_0^h [v'(x)]^2 dx \quad \text{with} \quad \bar{v} = \frac{1}{h} \int_0^h u(x) dx \quad (\text{C.2})$$

In 1D an optimal constant can be found:  $c_h = h^2/\pi^2$ .

(iii) extension:

$$\int_0^h |v(x)|^2 dx \leq \frac{h^2}{\pi^2} \left[ \int_0^h [v'(x)]^2 dx + |v(x_0)|^2 \right] \quad \text{with} \quad x_0 \in [0, h]. \quad (\text{C.3})$$

Now, let us use those inequalities to find error bounds. Suppose a general function  $u(x)$  is approximated by a *piecewise linear* approximation  $u_h(x)$ . Let's first find a *local* error estimate.

Consider  $v(x) = u'_h(x) - u'(x)$  and note that by *Rolle's theorem*

$$u'_h(x_0) - u'(x_0) = 0 \quad \text{for some} \quad x_0 \in (0, h). \quad (\text{C.4})$$

Next, use inequality (iii):

$$\int_0^h |u'_h(x) - u'(x)|^2 dx \leq \frac{h^2}{\pi^2} \int_0^h |u''_h(x) - u''(x)|^2 dx, \quad (\text{C.5})$$

but since  $u_h(x)$  is piecewise linear, we have  $u''_h(x) = 0$ , so that we arrive at the **local error estimate**

$$\int_0^h |u'_h(x) - u'(x)|^2 dx \leq \frac{h^2}{\pi^2} \int_0^h |u''(x)|^2 dx. \quad (\text{C.6})$$

Now, let's seek a **global error estimate** by using

$$\int_a^b (\cdot) dx = \sum_{i=0}^n \int_{x_i}^{x_{i+1}} (\cdot) dx \quad \text{with} \quad x_0 = a, \quad x_{n+1} = b, \quad x_{i+1} = x_i + h \quad (\text{C.7})$$

so that

$$\boxed{\int_a^b |u'_h(x) - u'(x)|^2 dx \leq \frac{h^2}{\pi^2} \int_a^b |u''(x)|^2 dx} \quad (\text{C.8})$$

Taking square roots, we see that for  $\Omega = (a, b)$

$$\|u'_h - u'\|_{L_2(\Omega)} \leq \frac{h}{\pi} \|u''\|_{L_2(\Omega)} \quad (\text{C.9})$$

and hence that  $\|u'_h - u'\|_{L_2(\Omega)} \rightarrow 0$  as  $h \rightarrow 0$  linearly in  $h$ .

We want to write this a bit more concise. Let us define the **Sobolev semi-norm**:

$$\boxed{|u|_{H^k(\Omega)} = \left[ \int_{\Omega} |D^k u|^2 \, dx \right]^{1/2}} \quad \text{or for short} \quad |u|_k = |u|_{H^k(\Omega)} \quad (\text{C.10})$$

where in 1D  $D^k u = u^{(k)}$ . A **semi-norm** in general must satisfy the following conditions:

- (i)  $\|u + v\| \leq \|u\| + \|v\|$  (like for a *norm*)
- (ii)  $\|\alpha \cdot u\| = |\alpha| \|u\|$  (like for a *norm*)
- (iii)  $\|u\| \geq 0$  (a *norm* also requires  $\|u\| = 0$  iff  $u = 0$ , not so for a semi-norm).

**Examples** in 1D:

- from before:

$$|u|_{H^1(a,b)} = \left[ \int_a^b |u'(x)|^2 \, dx \right]^{1/2} \quad (\text{C.11})$$

- analogously:

$$|u|_{H^2(a,b)} = \left[ \int_a^b |u''(x)|^2 \, dx \right]^{1/2} \quad (\text{C.12})$$

so that we can write (C.8) as

$$|u_h - u|_{H^1(a,b)}^2 \leq \frac{h^2}{\pi^2} |u|_{H^2(a,b)}^2 \quad \Rightarrow \quad \boxed{|u_h - u|_{H^1(a,b)} \leq \frac{h}{\pi} |u|_{H^2(a,b)}} \quad (\text{C.13})$$

Hence, the topology of convergence is bounded by the regularity of  $u$ . Convergence with  $h$ -refinement is linear.

- note the special case

$$|u|_{H^0}^2 = \int_{\Omega} u(x)^2 \, dx = \|u\|_{L^2}^2 \quad (L_2\text{-norm}) \quad (\text{C.14})$$

We can extend this to **higher-order interpolation**. For example, use a piecewise quadratic interpolation  $u_h$ . From Poincaré:

$$\int_0^h |u'_h - u'|^2 dx \leq \frac{h^2}{\pi^2} \int_0^h |u''_h - u''|^2 dx \leq \frac{h^4}{\pi^4} \int_0^h |u'''_h - u'''|^2 dx = \frac{h^4}{\pi^4} \int_0^h |u''|^2 dx \quad (\text{C.15})$$

Extension into a **global error estimate** with quadratic  $h$ -convergence:

$$|u_h - u|_{H^1(a,b)} \leq \frac{h^2}{\pi^2} |u|_{H^3(a,b)}. \quad (\text{C.16})$$

For a **general interpolation of order  $k$** :

$$\boxed{|u_h - u|_{H^1(a,b)} \leq \frac{h^k}{\pi^k} |u|_{H^{k+1}(a,b)}} \quad (\text{C.17})$$

Why is the Sobolev semi-norm not a norm? Simply consider the example  $u(x) = c > 0$ . All higher derivatives vanish on  $\mathbb{R}$ , so that  $|u|_{H^k(\Omega)} = 0$  for  $\Omega \subset \mathbb{R}$  and  $k \geq 1$ . However, that does not imply that  $u = 0$  (in fact, it is not).

Let us introduce the **Sobolev norm** (notice the double norm bars)

$$\|u\|_{H^k(\Omega)} = \left( \sum_{m=0}^k |u|_{H^m(\Omega)}^2 \right)^{1/2} \quad \text{or for short} \quad \|u\|_k = \|u\|_{H^k(\Omega)} \quad (\text{C.18})$$

For example, in one dimension

$$\|u\|_{H^1(\Omega)}^2 = |u|_0^2 + |u|_1^2 = \int_{\Omega} u(x)^2 dx + \int_{\Omega} [u'(x)]^2 dx = \|u\|_{L_2(\Omega)}^2 + |u|_{H^1(\Omega)}^2 \quad (\text{C.19})$$

Note that this also shows that, more generally,

$$\|u\|_{H^k(\Omega)}^2 \geq |u|_{H^k(\Omega)}^2. \quad (\text{C.20})$$

Let us derive a final global error estimate, one that involves proper norms – here for the example of a piecewise linear  $u_h$ . Start with Poincaré inequality (i):

$$\|u_h - u\|_{L_2(a,b)}^2 = \int_a^b |u_h - u|^2 dx \leq c_h \int_a^b |u'_h - u'|^2 dx = c_h |u_h - u|_{H^1(a,b)}^2 \quad (\text{C.21})$$

and from (C.19):

$$\begin{aligned} \|u_h - u\|_{H^1(\Omega)}^2 &= \|u_h - u\|_{L_2(\Omega)}^2 + |u_h - u|_{H^1(\Omega)}^2 \\ &\leq (1 + c_h) |u_h - u|_{H^1(a,b)}^2 \leq c^* h |u|_{H^2(a,b)}^2 \end{aligned} \quad (\text{C.22})$$

which along with (C.13) gives

$$\boxed{\|u_h - u\|_{H^1(a,b)} \leq ch |u|_{H^2(a,b)} \leq ch \|u\|_{H^2(a,b)}} \quad (\text{C.23})$$

## Summary and Extension of Norms:

$$\begin{aligned}
 L_p\text{-norm: } \|u\|_{L_p(\Omega)} &= \left( \int_{\Omega} u^p \, dx \right)^{1/p} \\
 \text{Sobolev semi-norm: } |u|_{H^k(\Omega)} &= \left[ \int_{\Omega} |D^k u|^2 \, dx \right]^{1/2} = |u|_k \\
 \text{Sobolev norm: } \|u\|_{H^k(\Omega)} &= \left( \sum_{m=0}^k |u|_{H^m(\Omega)}^2 \right)^{1/2} = \|u\|_k \\
 \text{generalization: } |u|_{W^{k,p}(\Omega)} &= \left[ \int_{\Omega} |D^k u|^p \, dx \right]^{1/p} = |u|_{k,p} \\
 \|u\|_{W^{k,p}(\Omega)} &= \left( \sum_{m=0}^k |u|_{W^{m,p}(\Omega)}^p \right)^{1/p} = \|u\|_{k,p}
 \end{aligned}$$

## C.1 Sobolev spaces

The Sobolev norm is used to define a **Sobolev space**:

$$H^k(\Omega) = \{u : \Omega \rightarrow \mathbb{R} \text{ such that } \|u\|_k < \infty\}, \quad (\text{C.24})$$

which includes all functions whose  $k$ th-order derivatives are **square-integrable**.

Examples:

- Consider a piecewise linear function  $u \in C^0$  defined on  $\Omega = (0, 2)$ . Then  $u \in H^1(\Omega)$  since the first derivative is piecewise constant and therefore square-integrable.
- Consider the Heavyside step function  $H(x) \in C^{-1}$  defined on  $\mathbb{R}$ . Then, e.g.,  $h \in H^0(\Omega)$  with  $\Omega = (-1, 1)$  since the first derivative (the Dirac delta function) is not square-integrable over  $(-1, 1)$ .

Overall, note that the above examples imply that

$$\boxed{H^m(\Omega) \subset C^k(\Omega) \quad \text{with} \quad m > k.} \quad (\text{C.25})$$

For example, if a function has a  $k$ th continuous derivative, then the  $(k+1)$ th derivative is defined piecewise and therefore square-integrable.

## C.2 Higher dimensions

To extend the above concepts to higher dimensions, we need multi-indices. A **multi-index** is an array of non-negative integers:

$$\alpha = (\alpha_1, \dots, \alpha_n) \in (\mathbb{Z}_0^+)^n \quad (\text{C.26})$$

The **degree** of a multi-index is defined as

$$|\alpha| = \alpha_1 + \dots + \alpha_n. \quad (\text{C.27})$$

This can be used to define a **monomial** for  $\mathbf{x} \in \mathbb{R}^n$ :

$$\mathbf{x}^\alpha = x_1^{\alpha_1} \cdot \dots \cdot x_n^{\alpha_n} \quad (\text{C.28})$$

For example, we can now extend our definition of polynomials to higher dimensions:

$$p(\mathbf{x}) \in \mathbb{P}_k(\mathbb{R}^2) \quad \Rightarrow \quad p(\mathbf{x}) = \sum_{\beta=0}^k \sum_{|\alpha|=\beta} a_{\alpha} \mathbf{x}^{\alpha} \quad (\text{C.29})$$

Specifically, the monomials above for  $\mathbf{x} \in \mathbb{R}^2$  are

$$\begin{aligned} \text{for } |\alpha| = 0: \quad & \{x^0 y^0\} = \{1\} \\ \text{for } |\alpha| = 1: \quad & \{x^1 y^0, x^0 y^1\} = \{x, y\} \\ \text{for } |\alpha| = 2: \quad & \{x^2 y^0, x^1 y^1, x^0 y^2\} = \{x^2, xy, y^2\} \quad \dots \end{aligned} \quad (\text{C.30})$$

so that

$$p(\mathbf{x}) = a_{(0,0)} + a_{(1,0)}x_1 + a_{(0,1)}x_2 + a_{(2,0)}x_1^2 + a_{(1,1)}x_1x_2 + a_{(0,2)}x_2^2 + \dots \quad (\text{C.31})$$

Note that this defines a **complete polynomial** of degree  $k$ .

Now we can use multi-indices to define partial derivatives via

$$\boxed{D^{\alpha}u = \frac{\partial^{|\alpha|}u}{\partial x_1^{\alpha_1} \cdots \partial x_n^{\alpha_n}} \quad \text{and} \quad D^0u = u} \quad (\text{C.32})$$

A common notation is

$$\sum_{|\alpha|=\beta} D^{\alpha}u = \sum_{\substack{\alpha_1, \dots, \alpha_n \\ \text{s.t. } |\alpha|=\beta}} \frac{\partial^{|\alpha|}u}{\partial x_1^{\alpha_1} \cdots \partial x_n^{\alpha_n}} \quad (\text{C.33})$$

With the above derivatives, we may redefine the *inner product*

$$\langle u, v \rangle_{H^m(\Omega)} = \int_{\Omega} \sum_{\beta=0}^m \sum_{|\alpha|=\beta} D^{\beta}u D^{\beta}v \, dx \quad (\text{C.34})$$

and the *Sobolev norm*

$$\|u\|_{H^m(\Omega)} = \langle u, u \rangle_{H^m(\Omega)}^{1/2} = \left[ \sum_{\beta=0}^m \sum_{|\alpha|=\beta} \int_{\Omega} (D^{\alpha}u)^2 \, dx \right]^{1/2} = \left[ \sum_{\beta=0}^m \sum_{|\alpha|=\beta} \|D^{\alpha}u\|_{L_2(\Omega)}^2 \right]^{1/2} \quad (\text{C.35})$$

Let's look at some examples; e.g., consider  $\Omega = \mathbb{R}^2$  and  $m = 1$ . Then we have

$$D^0u = u \quad \text{and} \quad D^1u = \left\{ \frac{\partial u}{\partial x_1}, \frac{\partial u}{\partial x_2} \right\} \quad (\text{C.36})$$

so that

$$\langle u, v \rangle_{H^1(\mathbb{R}^2)} = \int_{\mathbb{R}^2} \left( uv + \frac{\partial u}{\partial x_1} \frac{\partial v}{\partial x_1} + \frac{\partial u}{\partial x_2} \frac{\partial v}{\partial x_2} \right) dx_1 \, dx_2 \quad (\text{C.37})$$

and

$$\|u\|_{H^1(\mathbb{R}^2)}^2 = \int_{\mathbb{R}^2} \left[ u^2 + \left( \frac{\partial u}{\partial x_1} \right)^2 + \left( \frac{\partial u}{\partial x_2} \right)^2 \right] dx_1 \, dx_2. \quad (\text{C.38})$$

Altogether we can now properly define a **Sobolev space** in arbitrary dimensions:

$$\boxed{H^m(\Omega) = \left\{ u : \Omega \rightarrow \mathbb{R} : D^{\alpha}u \in L_2(\Omega) \, \forall \, \alpha \leq m \right\}} \quad (\text{C.39})$$

This is the set of all functions whose derivatives up to  $m$ th order all exist and are square-integrable.

As an example,  $u \in H^1(\Omega)$  implies that  $u$  and all its first partial derivatives must be square-integrable over  $\Omega$  because (C.38) must be finite.

Let us look at the example  $u(x) = |x|$  and take  $\Omega = (-1, 1)$ . Then, we have  $u'(x) = H(x)$  (the Heaviside jump function) and  $u''(x) = \delta(x)$  (the Dirac delta function). Therefore,

$$\begin{aligned} \int_a^b u^2(x) \, dx < \infty & \Rightarrow u \in L_2(\Omega) = H^0(\Omega) \\ \int_a^b \left( \frac{\partial u}{\partial x} \right)^2 \, dx = \int_a^b H^2(x) \, dx < \infty & \Rightarrow \frac{\partial u}{\partial x} \in L_2(\Omega) \quad \text{and} \quad u \in H^1(\Omega) \quad (\text{C.40}) \\ \int_a^b \left( \frac{\partial^2 u}{\partial x^2} \right)^2 \, dx = \int_a^b \delta^2(x) \, dx = \infty & \Rightarrow \frac{\partial^2 u}{\partial x^2} \notin L_2(\Omega) \quad \text{and} \quad u \notin H^2(\Omega) \end{aligned}$$

Note that one usually indicates the *highest order*  $k$  that applies (since this is what matters for practical purposes), so here we thus conclude that  $u \in H^1(\Omega)$ .

From the above, we also see that

$$H^\infty \subset \dots \subset H^2 \subset H^1 \subset H^0 = L_2. \quad (\text{C.41})$$

Notice that even though polynomials  $u \in \mathbb{P}_k(\Omega)$  are generally in  $H^\infty(\Omega)$  for any bounded  $\Omega \subset \mathbb{R}^d$ , they are not square-integrable over  $\Omega = \mathbb{R}^d$ . Luckily, in practical problems we usually consider only finite bodies  $\Omega$ . To more properly address this issue, let us introduce the **support** of a continuous function  $u$  defined on the open set  $\Omega \subset \mathbb{R}^d$  as the closure in  $\Omega$  of the set of all points where  $u(x) \neq 0$ , i.e.,

$$\text{supp } u = \overline{\{x \in \Omega : u(x) \neq 0\}} \quad (\text{C.42})$$

This means that  $u(x) = 0$  for  $x \in \Omega \setminus \text{supp } u$ . We may state, e.g., that functions  $u : \Omega \rightarrow \mathbb{R}$  with a *finite support*  $\Omega \subset \mathbb{R}^d$  and  $\text{ess sup}_\Omega < \infty$  are square-integrable over  $\Omega$ .

Finally, let us define by  $C_0^k(\Omega)$  the set of all functions contained in  $C^k(\Omega)$  whose support is a bounded subset of  $\Omega$ . Also, notice that

$$C_0^k(\Omega) \subset H_0^k(\Omega) \quad (\text{C.43})$$

and

$$C_0^\infty(\Omega) = \bigcap_{k \geq 0} C_0^k(\Omega). \quad (\text{C.44})$$

## D Operators

Energies are often defined via operators. Generally, we call  $\mathcal{A}$  an **operator** if

$$\mathcal{A} : u \in \mathcal{U} \rightarrow \mathcal{A}(u) \in \mathcal{V}, \quad (\text{D.1})$$

where both  $\mathcal{U}$  and  $\mathcal{V}$  are function spaces.

A simple example is

$$\mathcal{A}(u) = c \frac{du}{dx}, \quad (\text{D.2})$$

which is a (linear differential) *operator* requiring  $u \in C^1$ .

An operator  $\mathcal{A} : \mathcal{U} \rightarrow \mathcal{V}$  is **linear** if for all  $u_1, u_2 \in \mathcal{U}$  and  $\alpha, \beta \in \mathbb{R}$

$$\mathcal{A}(\alpha \cdot u_1 + \beta \cdot u_2) = \alpha \cdot \mathcal{A}(u_1) + \beta \cdot \mathcal{A}(u_2). \quad (\text{D.3})$$

For example,  $\mathcal{L}$  is a linear operator in

$$a u_{,xx} + b u_{,x} = c \quad \Leftrightarrow \quad \mathcal{L}(u) = c \quad \text{with} \quad \mathcal{L}(\cdot) = a(\cdot)_{,xx} + b(\cdot)_{,x}. \quad (\text{D.4})$$

Operators (such as the inner product operator) can also act on more than one function. Consider, e.g., an operator  $\mathcal{B} : \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$  where  $\mathcal{U}, \mathcal{V}$  are Hilbert spaces.  $\mathcal{B}$  is called a **bilinear operator** if for all  $u, u_1, u_2 \in \mathcal{U}$  and  $v, v_1, v_2 \in \mathcal{V}$  and  $\alpha, \beta \in \mathbb{R}$

- (i)  $\mathcal{B}(\alpha \cdot u_1 + \beta \cdot u_2, v) = \alpha \cdot \mathcal{B}(u_1, v) + \beta \cdot \mathcal{B}(u_2, v)$
- (ii)  $\mathcal{B}(u, \alpha \cdot v_1 + \beta \cdot v_2) = \alpha \cdot \mathcal{B}(u, v_1) + \beta \cdot \mathcal{B}(u, v_2)$

An example of a bilinear operator is the inner product  $\langle \cdot, \cdot \rangle : \mathcal{U} \times \mathcal{U} \rightarrow \mathbb{R}$  for some Hilbert space  $\mathcal{U}$ .

An operator  $\mathcal{A} : \mathcal{U} \rightarrow \mathcal{V}$  is called **symmetric** if

$$\langle \mathcal{A}(u), v \rangle = \langle u, \mathcal{A}(v) \rangle \quad \text{for all} \quad u, v \in \mathcal{U}. \quad (\text{D.5})$$

Furthermore, the operator is **positive** if

$$\langle \mathcal{A}(u), u \rangle \geq 0 \quad \text{for all} \quad u \in \mathcal{U}. \quad (\text{D.6})$$

An example of a symmetric operator is  $\mathcal{A}(u) = \mathbf{M}u$  with  $u \in \mathbb{R}^d$  and  $\mathbf{M} \in \mathbb{R}^{d \times d}$ , which is positive if  $\mathbf{M}$  is positive-semidefinite.



## E Uniqueness

One of the beauties of the above variational problem (3.37) is that a unique minimizer exists by the **Lax-Milgram theorem**. This is grounded in (assuming  $|\Omega| < \infty$  and  $u, v \in \mathcal{U}$  with some Hilbert space  $\mathcal{U}$ ):

- *boundedness* of the bilinear form:

$$|\mathcal{B}(u, v)| \leq C \|u\| \|v\| \quad \text{for some} \quad C > 0. \quad (\text{E.1})$$

For a bilinear form  $\mathcal{B}(u, v) = \langle \text{Grad } u, \text{Grad } v \rangle$ , this is satisfied by the *Cauchy-Schwarz inequality* (using  $L_2$ -norms):

$$|\mathcal{B}(u, v)| \leq C \|\text{Grad } u\|_{L_2(\Omega)} \|\text{Grad } v\|_{L_2(\Omega)} \leq C \|\text{Grad } u\|_{H^1(\Omega)} \|\text{Grad } v\|_{H^1(\Omega)} \quad (\text{E.2})$$

- *coercivity* of the bilinear form (*ellipticity*):

$$\mathcal{B}(u, u) \geq c \|u\|^2 \quad \text{for some} \quad c > 0. \quad (\text{E.3})$$

Again, for a bilinear form  $\mathcal{B}(u, v) = \langle \text{Grad } u, \text{Grad } v \rangle$  this is satisfied by *Poincaré's inequality*:

$$\mathcal{B}(u, u) = \|\text{Grad } u\|_{L_2(\Omega)}^2 \geq c \|u\|_{L_2(\Omega)}^2 \quad (\text{E.4})$$

These two requirements imply the well-posedness of the variational problem and thus imply the existence of a unique solution (or, that the potential has a unique global minimizer). In simple terms, the two conditions that the functional has sufficient *growth properties* (i.e., the bilinear form has upper and lower bounds).

## F Vainberg's theorem

The question arises whether or not a general form like (3.35) always exists for any set of PDEs/ODEs as governing equations. Vainberg's theorem helps us answer this question. Consider a weak form

$$\mathcal{G}[u, v] = 0 \quad \forall \quad v \in \mathcal{U}_0(\Omega). \quad (\text{F.1})$$

Let us see if  $\mathcal{G}$  derives from a potential  $I$  via its first variation. That would imply that

$$\mathcal{G}(u, \delta u) = D_{\delta u} I[u] = \lim_{\epsilon \rightarrow 0} \frac{d}{d\epsilon} I[u + \delta u]. \quad (\text{F.2})$$

Now recall from calculus that for any continuously differentiable function  $f(x, y)$  we must have by *Schwartz' theorem*

$$\frac{\partial}{\partial y} \frac{\partial f}{\partial x} = \frac{\partial}{\partial x} \frac{\partial f}{\partial y}. \quad (\text{F.3})$$

We can use the same strategy to formulate whether or not a weak form derives from a potential. Specifically, we can take one more variation and state that

$$\boxed{D_{\delta u_2} \mathcal{G}(u, \delta u_1) = D_{\delta u_1} \mathcal{G}(u, \delta u_2) \quad \text{if and only if } I[u] \text{ exists}} \quad (\text{F.4})$$

This is known as **Vainberg's theorem**.

We can easily verify this for the general form given in (3.35):

$$\begin{aligned} \mathcal{G}(u, \delta u_1) &= D_{\delta u_1} I[u] = \int_{\Omega} [\lambda \text{Grad } u \text{ Grad } \delta u_1 - S \delta u_1] \, dx - \int_{\partial\Omega_N} \hat{Q} \delta u_1 \, dx \\ \Rightarrow D_{\delta u_1} \mathcal{G}(u, \delta u_2) &= \int_{\Omega} \lambda \text{Grad } \delta u_2 \text{ Grad } \delta u_1 \, dx = D_{\delta u_2} \mathcal{G}(u, \delta u_1) \end{aligned} \quad (\text{F.5})$$

In simple terms (and not most generally), Vainberg's theorem holds if the potential obeys symmetry. This in turn implies that the governing PDE contains derivatives of *even* order (such as linear momentum balance which is second-order in both spatial and temporal derivatives, or the equilibrium equations of beams which are fourth-order in space and second-order in time). If the PDEs are of *odd* order (such as, e.g., the time-dependent diffusion or heat equation), then no direct potential  $I$  exists – there are work-arounds using so-called *effective potentials* that will be discussed later in the context of internal variables.

Of course, knowing that a variational structure exists is beneficial but it does not reveal anything about the actual solution  $u$  which will be obtained by solving the above system of equations.

## G Energy norm

For many purposes it will be convenient to introduce the so-called **energy norm**

$$|\cdot|_E = \sqrt{\mathcal{B}(\cdot, \cdot)}. \quad (\text{G.1})$$

For example, if in subsequent sections we find an approximate solution  $T^h$  for the temperature field, then the error can be quantified by

$$|T - T^h|_E = \sqrt{\mathcal{B}(T - T^h, T - T^h)} = \sqrt{\int_{\Omega} \lambda \|\text{Grad}(T - T^h)\|^2 \, dV}. \quad (\text{G.2})$$

Notice that in this case  $\|T - T^h\|$  does not necessarily imply  $T - T^h = 0$  so that, strictly speaking, the above energy norm is only a *semi-norm*. To turn it into a proper norm, we need to exclude solutions corresponding to *rigid-body motion* (i.e., solutions that imply uniform translations or rotations of the  $T$ -field but without producing gradients in  $T$ ). If we assume that essential boundary conditions are chosen appropriately to suppress rigid body motion by seeking solutions

$$T \in \mathcal{U} = \{T \in H^1 : T = \hat{T} \text{ on } \partial\Omega_D\}, \quad (\text{G.3})$$

then, for this space,  $\|\cdot\|_E$  is indeed a proper norm.