

## CHAPTER 11

# Finite element methods

### 1. Preliminaries

For a number of applications the restrictions imposed by finite difference or spectral methods with respect to the computational grid are too severe. This is especially the case in structural engineering where the elasticity equations are solved in domains of complicated geometry such as the interior of an automobile engine. A review of the finite difference and spectral methods would show that the reason relatively simple grids are required is that the differential form of the equation is used. Finite volume methods had no such restriction since they used an integral formulation, and indeed complicated geometries may be treated by finite volume methods. Another class of methods which are based upon an integral formulation are the finite element and closely related boundary element methods. We shall concentrate on finite element methods for now.

The basic idea behind the finite element methods is to employ a piecewise local approximation  $\tilde{q}$  of the unknown function  $q$  that satisfies some PDE of interest. The piecewise local approximation is defined over some general discretization of the domain of definition of  $q$  denoted by  $\Omega$ . Instead of directly using the piecewise local approximation in the PDE we employ a weighted residual formulation. There arises the significant question of how to best relate the integral formulation to the PDE of interest. Once the discretization, piecewise local approximation and integral formulation are determined a system of equations is obtained whose solution gives the complete approximation to the problem of interest. We shall look at each of these components in detail.

**1.1. Spatial discretizations.** A domain  $\Omega$  may be discretized into simple elements in very many ways. Nonetheless only a few are typically used in practice. General affine geometry furnishes some guidance for general discretization techniques. We know for instance that any  $d$ -dimensional domain may be expressed as a reunion of simplices

$$(1.1) \quad \Omega = \cup_k S_k .$$

Simplices are the simplest continuum geometric entities one can construct in a space of dimension  $d$ . For 1D spaces the simplices are line segments. In 2D they are triangles and in 3D they are tetrahedra. The measure of each of these elements is easily determined by the formulas:

$$(1.2) \quad \begin{aligned} & (1) \text{ line segment in 1D of nodes } \{x_1, x_2\} \\ & l = x_2 - x_1 = \left| \begin{array}{cc} 1 & 1 \\ x_1 & x_2 \end{array} \right| \end{aligned}$$

FIGURE 1. Example of the discretization into triangles of the domain between a circle and a NACA-0012 airfoil.

(2) triangle in 2D with nodes  $\{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}$

$$(1.3) \quad A = \frac{1}{2} \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{vmatrix}$$

(3) tetrahedron in 3D with nodes  $\{(x_1, y_1, z_1), (x_2, y_2, z_2), (x_3, y_3, z_3)\}$

$$(1.4) \quad V = \frac{1}{6} \begin{vmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \\ y_1 & y_2 & y_3 & y_4 \\ z_1 & z_2 & z_3 & z_4 \end{vmatrix}$$

In the above formulas the element measures are given with sign, the sign corresponding to orientation of the nodes. We understand that the positive value is to be taken whenever a true geometric measure (length, area, volume) is required. Simplices have many attractive theoretical properties, in particular there exists a definition of what an optimal discretization is for a number of PDE's of interest, especially elliptic PDE's such as the Poisson equation. Fig. 1 shows an example of such a discretization.

Another widely used discretization is into generalized polyhedra having  $2d$  sides, i.e. line segments in 1D, quadrilaterals in 2D, hexahedra in 3D. These have the advantage of enabling easier organization of programs since there is a natural ordering of the indices identifying each element. Thus discretizations which use these types of elements give rise to *structured computational grids*, similar to those encountered in finite difference methods whereas discretizations using simplices lead to *unstructured computational grids*.

**1.2. Piecewise interpolations.** Once a discretization scheme for the geometric domain has been established the next step is to define a local approximation of  $q$  over the element  $E$ . Typically the approximation is an interpolation based upon values  $Q_j$  defined somewhere within the element  $E$ , but this is not obligatory and other approximations (spectral elements, Chebyshev elements) may be used. The position where the values  $Q_j$  are to be defined must be established. A simple choice is the element nodes but again this is not obligatory and the values may be positioned at other points within  $E$ . Finally an interpolation scheme must be established such as polynomial interpolation. Let us give some typical examples:

1.2.1. *Linear elements in 1D.* The element  $E$  has two nodes  $\{x_1, x_2\}$ ,  $x_2 > x_1$ . Values representing  $q(x)$  are defined at the nodes  $\{Q_1, Q_2\}$ . These define a linear polynomial approximation valid over  $E$

$$(1.5) \quad \tilde{q}(x) = \frac{(x - x_1)Q_2 + (x_2 - x)Q_1}{x_2 - x_1} = N_1(x)Q_1 + N_2(x)Q_2$$

The functions  $N_1(x)$ ,  $N_2(x)$  have properties reminiscent of the Dirac delta

$$(1.6) \quad N_1(x_1) = 1, N_1(x_2) = 0$$

$$(1.7) \quad N_2(x_1) = 0, N_2(x_2) = 1$$

FIGURE 2. Linear 1D form functions.

FIGURE 3. Quadratic element form functions in 1D.

and are called *form functions*. The particular ones used here are called the 1D linear form functions and are depicted in Fig. (2)

1.2.2. *Quadratic elements in 1D*. The element  $E$  has three nodes  $\{x_1, x_2, x_3\}$  and the local approximation is

$$(1.8) \quad \tilde{q}(x) = N_1(x)Q_1 + N_2(x)Q_2 + N_3(x)Q_3$$

with the form functions

$$(1.9) \quad N_1(x) = \frac{(x - x_2)(x - x_3)}{(x_1 - x_2)(x_1 - x_3)}$$

$$(1.10) \quad N_2(x) = \frac{(x - x_3)(x - x_1)}{(x_2 - x_3)(x_2 - x_1)}$$

$$(1.11) \quad N_3(x) = \frac{(x - x_1)(x - x_2)}{(x_3 - x_1)(x_3 - x_2)}$$

1.2.3. *Linear elements on triangles in 2D*. The element  $E$  has 3 nodes of coordinates  $\{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}$  at which the values  $Q_1, Q_2, Q_3$  are defined. The local approximation of  $q$  is given by

$$(1.12) \quad \tilde{q}(x, y) = N_1(x, y)Q_1 + N_2(x, y)Q_2 + N_3(x, y)Q_3$$

with the form functions

$$(1.13) \quad N_1(x, y) = \frac{1}{2A} \begin{vmatrix} 1 & 1 & 1 \\ x & x_2 & x_3 \\ y & y_2 & y_3 \end{vmatrix} = \frac{1}{2A} (xy_2 - yx_2 - xy_3 + yx_3 + x_2y_3 - x_3y_2)$$

$$(1.14) \quad N_2(x, y) = \frac{1}{2A} \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x & x_3 \\ y_1 & y & y_3 \end{vmatrix} = \frac{1}{2A} (yx_1 - xy_1 + xy_3 - yx_3 - x_1y_3 + y_1x_3)$$

$$(1.15) \quad N_3(x, y) = \frac{1}{2A} \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x \\ y_1 & y_2 & y \end{vmatrix} = \frac{1}{2A} (xy_1 - yx_1 - xy_2 + yx_2 + x_1y_2 - x_2y_1)$$

Notice how the properties of simplices enable the form functions to be easily determined.

1.2.4. *Linear along each direction elements on quadrilaterals in 2D*. The element  $E$  has 4 nodes  $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\}$ . It is convenient to introduce a local coordinate system  $(\xi, \eta)$  so that the nodes correspond to the local coordinates  $(\pm 1, \pm 1)$ . The local approximation is then given in the local coordinates

by

$$(1.16) \quad \tilde{q}(\xi, \eta) = \sum_{k=1}^4 N_k(\xi, \eta) Q_k$$

with the form functions

$$(1.17) \quad N_1(\xi, \eta) = \frac{1}{4}(1 + \xi)(1 + \eta)$$

$$(1.18) \quad N_2(\xi, \eta) = \frac{1}{4}(1 - \xi)(1 + \eta)$$

$$(1.19) \quad N_3(\xi, \eta) = \frac{1}{4}(1 - \xi)(1 - \eta)$$

$$(1.20) \quad N_4(\xi, \eta) = \frac{1}{4}(1 + \xi)(1 - \eta)$$

The local transformation of coordinates can also be written in terms of the form functions

$$(1.21) \quad x(\xi, \eta) = \sum_{k=1}^4 N_k(\xi, \eta) x_k, \quad y(\xi, \eta) = \sum_{k=1}^4 N_k(\xi, \eta) y_k$$

## 2. Variational derivation of weighted residual formulations

We now turn to the problem of how to obtain a measure of the error introduced in approximating the exact solution  $q$  to the PDE of interest with its piecewise approximation  $\tilde{q}$ . Some techniques were presented in the general presentation of weighted residual methods carried out in Chapter 2. For a wide class of problems of interest, especially elliptic problems there exist alternative formulations that lead to more efficient numerical algorithms. These are based upon variational and functional analysis and we shall consider the basics of the theory here.

**2.1. Variational calculus.** Consider the problem of determining the extremum of the integral

$$(2.1) \quad I(q) = \int_a^b f(x, q, q') dx$$

over all functions  $q : \mathbb{R} \rightarrow \mathbb{R}$  that belong to some class, for example piecewise continuous functions and that satisfy the boundary conditions  $q(x = a) = q_a$ ,  $q(x = b) = q_b$ .  $I(q)$  is called a functional in that it associates a scalar value to each element from a space of functions. We can consider small perturbations of the function  $q$  that we denote by  $\delta q$ . The perturbations maintain the boundary conditions, i.e.

$$(2.2) \quad \delta q(x = a) = 0, \quad \delta q(x = b) = 0.$$

The change in  $I$  is

$$(2.3) \quad \delta I = I(q + \delta q) - I(q) = \int_a^b f(x, q + \delta q, q' + \delta q') dx - \int_a^b f(x, q, q') dx.$$

We shall consider  $q, q'$  as independent variables in  $f$  and carry out series expansions to obtain

$$(2.4) \quad \delta I = \int_a^b \left[ \left( \frac{\partial f}{\partial q} \right) \delta q + \left( \frac{\partial f}{\partial q'} \right) \delta q' \right] dx.$$

We can interchange the  $\delta$  and  $d/dx$  operators in the second term and then integrate by parts

$$(2.5) \quad \int_a^b \left( \frac{\partial f}{\partial q'} \right) \delta q' dx = \int_a^b \left( \frac{\partial f}{\partial q'} \right) \delta \frac{dq}{dx} dx = \int_a^b \left( \frac{\partial f}{\partial q'} \right) \frac{d}{dx} (\delta q) dx =$$

$$(2.6) \quad = \left( \frac{\partial f}{\partial q'} \right) (\delta q) \Big|_{x=a}^{x=b} - \int_a^b \frac{d}{dx} \left( \frac{\partial f}{\partial q'} \right) (\delta q) dx$$

Applying the boundary conditions and then replacing the above result in (2.4) leads to

$$(2.7) \quad \delta I = \int_a^b \left[ \left( \frac{\partial f}{\partial q} \right) - \frac{d}{dx} \left( \frac{\partial f}{\partial q'} \right) \right] \delta q dx .$$

For  $I$  to be at an extremum  $\delta I$  must maintain the same sign under any perturbation of the extremum. This is only possible if the factor multiplying  $\delta q$  in the above integral is zero everywhere. If it were not then  $\delta q_1$  would give some value  $\delta I_1$  and  $-\delta q_1$  would lead to the opposite value  $-\delta I_1$  and  $I$  would not be at an extremum. We therefore have

$$(2.8) \quad \left( \frac{\partial f}{\partial q} \right) - \frac{d}{dx} \left( \frac{\partial f}{\partial q'} \right) = 0$$

as the condition for  $I$  to be at an extremum. This is known as the *Euler variational principle*. At the extremum we obviously have  $\delta I = 0$ .

The importance of the Euler variational principle for numerical solution of PDE's rests upon the link it furnishes between an integral formulation  $I(q)$  and a differential equation (2.8). We can write down specific forms of  $f$  that lead to PDE's of great practical interest. For example replacing

$$(2.9) \quad f(x, q, q') = \frac{1}{2} \left( \frac{dq}{dx} \right)^2 - gq$$

in (2.8) leads to the differential equation

$$(2.10) \quad q'' = g$$

with the boundary conditions  $q(x=a) = q_a$ ,  $q(x=b) = q_b$ . This is the standard 2 point boundary problem for a second order ODE. Recall that this can be solved by either direct discretization leading to the linear system of equations

$$(2.11) \quad Q_{j-1} - 2Q_j + Q_{j+1} = h^2 g_j, \quad j = 1, \dots, N-1$$

or by using a shooting method combined with an initial value solve in which we seek  $z = q'(x=a)$  that leads to  $q(x=b; z) = q_b$ . The variational formulation above suggests a third approach. Instead of directly solving the ODE we can seek  $q$  that minimizes  $I(q)$  with  $f$  given by (2.9). This is extremely useful in constructing finite element approximations as we will see below.

Other important expressions of the Euler variational principle can be derived for various situations. Let us consider the ones most often encountered.

(1) Functional of two functions in 1D. The functional is

$$(2.12) \quad I(p, q) = \int_a^b f(x, p, p', q, q') dx$$

and the Euler variational principle leads to

$$(2.13) \quad \left(\frac{\partial f}{\partial p}\right) + \left(\frac{\partial f}{\partial q}\right) - \frac{d}{dx} \left(\frac{\partial f}{\partial p'}\right) - \frac{d}{dx} \left(\frac{\partial f}{\partial q'}\right) = 0$$

(2) Functional of a 2D function.

$$(2.14) \quad I(q) = \int_c^d \int_a^b f(x, y, q, q_x, q_y) dx dy$$

$$(2.15) \quad \left(\frac{\partial f}{\partial q}\right) - \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial q_x}\right) - \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial q_y}\right) = 0$$

(3) Functional involving second order derivatives in 1D.

$$(2.16) \quad I(p, q) = \int_a^b f(x, p, p', p'') dx$$

$$(2.17) \quad \left(\frac{\partial f}{\partial q}\right) - \frac{d}{dx} \left(\frac{\partial f}{\partial q'}\right) + \frac{d^2}{dx^2} \left(\frac{\partial f}{\partial q''}\right) = 0$$

**2.2. Ritz methods.** In the Ritz formulation of the finite element method we seek a piecewise approximation that minimizes the functional associated with the PDE of interest. The piecewise local approximation can be expressed as

$$(2.18) \quad \tilde{q}(x) = \sum_e \sum_k Q_k^e N_k^e(x)$$

where the  $e$  sum is over all elements and the  $k$  sum is over all nodes within an element. The unknowns of the problem are the nodal values  $Q_k^e$ . The form functions  $N_k^e(x)$  correspond to some chosen approximation scheme. Let  $f$  be associated with the PDE we are interested in solving. The problem reduces to finding  $\{Q_k^e\}$  that minimizes

$$(2.19) \quad I(\tilde{q}) = \int_a^b f(x, \tilde{q}, \tilde{q}_x) dx .$$

This can be solved by finding the solution to the system of equations

$$(2.20) \quad \frac{\partial}{\partial Q_k^e} I(\tilde{q}) = 0$$

with  $e$  going over all elements and  $k$  over all element nodes.

Note that the entire procedure rests upon the ability to determine a function  $f$  that corresponds to a PDE of practical interest. In many situations we have physical guidance that such a variational principle formulation exists. The basic underpinning is furnished by analytical mechanics and the physical principle of least action which finds various expressions in different disciplines. The principle of least action asserts that of all the generalized trajectories  $(p, q) = \{q_k(t), p_k(t) \mid k = 1, \dots, 3N\}$  of a system of  $N$  particles, the one actually followed minimizes the action  $S$

$$(2.21) \quad S = \int_{t_0}^{t_1} \mathcal{L}(t, p, q) dt$$

with  $\mathcal{L}$  being the Lagrangean of the system. Here  $q_k$  denote generalized coordinates and  $p_k$  generalized momenta. Though not always immediately apparent this leads to other expressions typically called *minimum energy functionals*. These can be

written for systems with no dissipative effects. Here are some examples of functions  $f$  linked to important PDE's:

(1) Poisson equation in 2D

$$(2.22) \quad f = \frac{1}{2} (q_x^2 + q_y^2) - gq$$

for which (2.15) gives

$$(2.23) \quad q_{xx} + q_{yy} = g$$

(2) Poisson equation in 3D

$$(2.24) \quad f = \frac{1}{2} (q_x^2 + q_y^2 + q_z^2) - gq$$

for which the Euler variational principle

$$(2.25) \quad \left( \frac{\partial f}{\partial q} \right) - \frac{\partial}{\partial x} \left( \frac{\partial f}{\partial q_x} \right) - \frac{\partial}{\partial y} \left( \frac{\partial f}{\partial q_y} \right) - \frac{\partial}{\partial z} \left( \frac{\partial f}{\partial q_z} \right) = 0$$

gives

$$(2.26) \quad q_{xx} + q_{yy} + q_{zz} = g$$

**2.3. Galerkin methods.** The Ritz formulation typically leads to a system of equations which has nice numerical properties. However there are many systems for which a variational formulation is not possible typically because the system has dissipative behavior. In such situations we can again use an integral reformulation of the PDE of interest based upon the concept of a weak solution already introduced in the study of hyperbolic problems. Suppose we're looking for a solution to the problem

$$(2.27) \quad \mathcal{A}q = g$$

with  $\mathcal{A}$  some differential operator. A function  $q$  that directly satisfies (2.27) is called a *classical solution*. Consider now some space of test functions  $v$  and a scalar product defined for the functions  $q$  and  $v$ . From (2.27) we can derive

$$(2.28) \quad (\mathcal{A}q, v) = (g, v)$$

where  $(\cdot, \cdot)$  denotes the scalar product, e.g.

$$(2.29) \quad (u, v) = \int_a^b u(x)v(x)dx .$$

In (2.28) we can apply integration by parts to obtain

$$(2.30) \quad (q, \mathcal{A}^*v) = (g, v)$$

where  $\mathcal{A}^*$  is the adjoint operator of  $\mathcal{A}$ . This typically enables us to avoid differentiating functions  $q$  that might be discontinuous. We can now use (2.30) to determine the unknown coefficients of a finite element approximation

$$(2.31) \quad \tilde{q}(x) = \sum_e \sum_k Q_k^e N_k^e(x)$$

by requiring

$$(2.32) \quad \sum_e \sum_k Q_k^e (N_k^e(x), \mathcal{A}^*v) = (g, v) .$$

The only piece missing is how we choose the test functions  $v$ . In a Galerkin formulation these are chosen to be the form functions themselves leading to

$$(2.33) \quad \sum_e \sum_k Q_k^e (N_k^e(x), \mathcal{A}^* N_j^e(x)) = (g, N_j^e(x)) ,$$

thus defining a linear system

$$(2.34) \quad \mathbf{A}Q = b$$

$$(2.35) \quad A_{jk} = (N_k^e(x), \mathcal{A}^* N_j^e(x)) .$$

**2.4. A detailed example.** Let us now carry out the steps involved in solving a Poisson equation in 2D using a Ritz formulation and quadrilateral elements. The mathematical statement of the problem is

$$(2.36) \quad \begin{cases} q_{xx} + q_{yy} = g & (x, y) \in \Omega \\ q = b & (x, y) \in \partial\Omega \end{cases}$$

with the domain  $\Omega = [a, b] \times [c, d]$  and  $\partial\Omega$  denoting the boundary of  $\Omega$  on which Dirichlet conditions are given. The element form functions are given by (1.17)-(1.20) and the function  $f$  is given by (2.22). The function  $I(\tilde{q})$  is

$$(2.37) \quad I(\tilde{q}) = \int_c^d \int_a^b f(x, y, \tilde{q}, \tilde{q}_x, \tilde{q}_y) dx dy = \int_c^d \int_a^b \left[ \frac{1}{2} (\tilde{q}_x^2 + \tilde{q}_y^2) - g\tilde{q} \right] dx dy .$$

The finite element approximation is determined by the chosen form functions and the nodal values  $\{Q_k^e\}$ . The extremum of  $I(\tilde{q})$  is attained when

$$(2.38) \quad \frac{\partial}{\partial Q_k^e} I(\tilde{q}) = 0$$

which leads to

$$(2.39) \quad \int_c^d \int_a^b \left[ \left( \tilde{q}_x \frac{\partial \tilde{q}_x}{\partial Q_k^e} + \tilde{q}_y \frac{\partial \tilde{q}_y}{\partial Q_k^e} \right) - g \frac{\partial \tilde{q}}{\partial Q_k^e} \right] dx dy = 0$$

Note that

$$(2.40) \quad \frac{\partial \tilde{q}}{\partial Q_k^e} = N_k^e, \quad \frac{\partial \tilde{q}_x}{\partial Q_k^e} = \frac{\partial N_k^e}{\partial x}, \quad \frac{\partial \tilde{q}_y}{\partial Q_k^e} = \frac{\partial N_k^e}{\partial y}$$

so these derivatives no longer contain the unknowns  $\{Q_k^e\}$ . We thus obtain

$$(2.41) \quad \sum_e \sum_j \left[ \iint \left( \frac{\partial N_j^e}{\partial x} \frac{\partial N_k^e}{\partial x} + \frac{\partial N_j^e}{\partial y} \frac{\partial N_k^e}{\partial y} \right) dx dy \right] Q_k^e = \sum_e \iint g N_k^e dx dy$$

with  $k$  going over all the element nodes. The sum over the elements is typically known as an assembly operation, leading to the computation of the matrix elements

$$(2.42) \quad A_{jk} = \sum_e \iint \left( \frac{\partial N_j^e}{\partial x} \frac{\partial N_k^e}{\partial x} + \frac{\partial N_j^e}{\partial y} \frac{\partial N_k^e}{\partial y} \right) dx dy$$

known as the system *stiffness matrix*. We can easily compute the elements of this matrix. Analytical computation is possible as in

$$(2.43) \quad \frac{\partial N_k^e}{\partial x} = \frac{\partial N_k^e}{\partial \xi} \frac{\partial \xi}{\partial x} + \frac{\partial N_k^e}{\partial \eta} \frac{\partial \eta}{\partial x}$$



$$(2.44) \quad \frac{\partial \xi}{\partial x} = \frac{\frac{D(\xi, y)}{D(\xi, \eta)}}{\frac{D(x, y)}{D(\xi, \eta)}} = \frac{1}{J} \begin{vmatrix} 1 & 0 \\ y_\xi & y_\eta \end{vmatrix} = \frac{y_\eta}{J}$$

$$(2.45) \quad \frac{\partial \eta}{\partial x} = \frac{\frac{D(\eta, y)}{D(\xi, \eta)}}{\frac{D(x, y)}{D(\xi, \eta)}} = \frac{1}{J} \begin{vmatrix} 0 & 1 \\ y_\xi & y_\eta \end{vmatrix} = -\frac{y_\xi}{J}$$

$$(2.46) \quad J = \begin{vmatrix} x_\xi & x_\eta \\ y_\xi & y_\eta \end{vmatrix} = x_\xi y_\eta - x_\eta y_\xi$$

The  $x(\xi, \eta)$  and  $y(\xi, \eta)$  dependencies are given by (1.21) so we obtain

$$(2.47) \quad \frac{\partial x}{\partial \xi} = \sum_{k=1}^4 \frac{\partial N_k}{\partial \xi} x_k, \quad \frac{\partial y}{\partial \xi} = \sum_{k=1}^4 \frac{\partial N_k}{\partial \xi} y_k$$

$$(2.48) \quad \frac{\partial x}{\partial \eta} = \sum_{k=1}^4 \frac{\partial N_k}{\partial \eta} x_k, \quad \frac{\partial y}{\partial \eta} = \sum_{k=1}^4 \frac{\partial N_k}{\partial \eta} y_k$$

But analytical evaluations are not really required in this case. We can recognize that the integrand in (2.42) is quadratic in  $(x, y)$  and that a 4-point Gauss-Legendre quadrature leads to an exact evaluation.