

Chapter 1

Weighted Residuals and Galerkin's Method for a Generic 1-D Problem

1.1 Introduction: Weighted Residual Methods

Let us start by considering a simple one-dimensional differential equation, written in abstract form

$$A(u) = f, \quad (1.1)$$

where, for example, $A(u) = \frac{d}{dx} \left(A_1 \frac{du}{dx} \right) + A_2 u$, where $A_1 = A_1(x)$ and $A_2 = A_2(x)$. Let us choose an approximate solution of the form

$$u^N = \sum_{i=1}^N a_i \phi_i(x), \quad (1.2)$$

where the ϕ_i 's are approximation functions, and the a_i 's are unknown constants that we will determine. Substituting the approximation leads to a "left over" amount called the residual:

$$r^N(x) = A(u^N) - f. \quad (1.3)$$

If we assume that the ϕ 's are given, we would like to choose the a_i 's to minimize the residual in an appropriate norm, denoted $\|r\|$. A primary question is which norm should be chosen to measure the solution and to determine its quality. Obviously, if the true solution is smooth enough to have a pointwise solution, and if we could take enough ϕ -functions, we could probably match the solution pointwise. However, as we shall see, this would be computationally prohibitively expensive to solve. Thus, we usually settle for a less stringent measure, for example a spatially averaged measure of solution quality. This is not a trivial point, and we will formalize the exact choice of the appropriate metric (a norm) momentarily. Let us make an obvious choice

$$\Pi(r^N) \stackrel{\text{def}}{=} ||r^N||^2 \stackrel{\text{def}}{=} \int_0^L (r^N(x))^2 dx. \quad (1.4)$$

Taking the derivative with respect to each a_i , and setting it to zero, we obtain for $i = 1, 2, \dots, N$

$$\frac{\partial \Pi}{\partial a_i} = \int_0^L 2r \frac{\partial r}{\partial a_i} dx = 0. \quad (1.5)$$

This leads to N equations and N unknowns. This method is called the “Method of Least Squares.” Another approach is to force the residual to be zero at a discrete number of locations, $i = 1, 2, \dots, N$

$$r^N(x_i) = 0, \quad (1.6)$$

which can also be written as

$$\int_0^L r^N(x) \delta(x - x_i) dx = 0, \quad (1.7)$$

where $\delta(x)$ is the Dirac Functional.¹ This approach is sometimes referred to as the “Collocation Method.” Notice that each method has the form

$$\int_0^L r^N(x) w(x) dx = 0, \quad (1.9)$$

where $w(x)$ is some “weight.” A general name for these methods is the “Method of Weighted Residuals.”

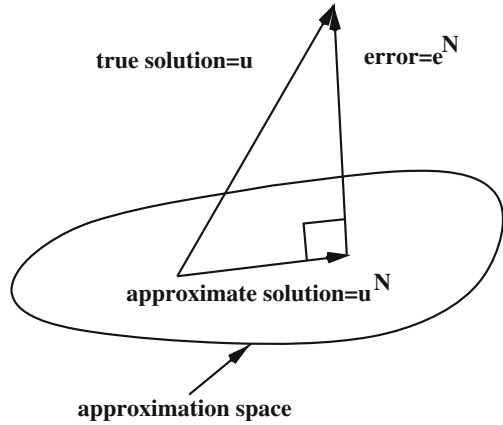
1.2 Galerkin's Method

Of all of the weighted residual methods used in the scientific community, one particular method, the Galerkin method, is by far the most widely used, and has been

¹ Recall, the Dirac Functional is defined via

$$\int_0^L \delta(x - x_i) f(x) dx = f(x_i). \quad (1.8)$$

Fig. 1.1 Orthogonality of the approximation error



shown to deliver the most accurate solutions on a wide variety of problems. We now explain the basic construction. Consider the true solution, approximate solution and the error, related through

$$u - u^N = e^N \Rightarrow u = u^N + e^N. \quad (1.10)$$

As a helpful mnemonic, now consider them as vectors (Fig. 1.1). Clearly, the error (e^N) is the smallest when e^N is orthogonal to u^N . The problem is that the error $e^N = u - u^N$ is unknown. However, the “next best thing,” the *residual*, is known.² This motivates Galerkin’s idea, namely to force u^N and r^N to be orthogonal. Mathematically, this can be expressed as

$$\int_0^L r^N(x) u^N(x) dx = \int_0^L r^N(x) \sum_{i=1}^N a_i \phi_i dx = 0. \quad (1.11)$$

However, this only gives one equation. Thus, we enforce this for each of the individual approximation functions, which collectively form the space of approximations comprising u^N ,

$$\int_0^L r^N(x) a_i \phi_i(x) dx = a_i \int_0^L r^N(x) \phi_i(x) dx = 0 \Rightarrow \int_0^L r^N(x) \phi_i(x) dx = 0. \quad (1.12)$$

This leads to N equations and N unknowns, in order to solve for the a_i ’s. It is usual practice in Galerkin’s method to use approximation functions that are what

² Although the error and residual are not the same, we note that when the residual is zero, the error is zero.

is called “kinematically admissible,” which we define as functions that satisfy the (primal solution variable, u) displacement boundary condition a priori. Kinematically admissible functions do not have to satisfy boundary conditions that involve derivatives of the solution beforehand.

1.3 An Overall Framework

The basic “recipe” for the Galerkin process is as follows:

- Step 1: Compute the residual: $A(u^N) - f = r^N(x)$.
- Step 2: Force the residual to be orthogonal to each of the approximation functions: $\int_0^L r^N(x) \phi_i(x) dx = 0$.
- Step 3: Solve the set of coupled equations. The equations will be linear if the differential equation is linear, and nonlinear if the differential equation is nonlinear.

The primary problem with such a general framework is that it provides no systematic way of choosing the approximation functions, which is strongly dependent on issues of possible nonsmoothness of the true solution. The basic Finite Element Method has been designed to embellish and extend the fundamental Galerkin method by constructing ϕ_i 's in order to deal with such issues. In particular:

- It is based upon Galerkin's method;
- It is computationally systematic and efficient and;
- It is based on reformulations of the differential equation that remove the problems of restrictive differentiability requirements.

The approach that we will follow in this monograph is to introduce the basic concepts first in one dimension. We then present three-dimensional formulations, which extend naturally from one-dimensional formulations.

Remark: Two Appendices containing some essential background information on vector calculus, linear algebra, and basic mechanics, exemplified by linearized elasticity, are provided. Linearized elasticity will serve as our model problem in the chapters that follow.

Chapter 2

A Model Problem: 1-D Elastostatics

2.1 Introduction: A Model Problem

In most problems of mathematical physics the true solutions are nonsmooth, i.e., they are not continuously differentiable. *Thus, we cannot immediately apply a Galerkin approach.* For example, in the equation of static mechanical equilibrium¹

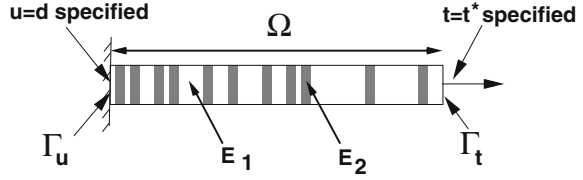
$$\nabla \cdot \sigma + f = 0, \quad (2.1)$$

there is an implicit requirement that the stress, σ , is differentiable in the classical sense. Virtually the same mathematical structure form holds for other partial differential equations of mathematical physics describing diffusion, heat conduction, etc. *In many applications, differentiability is too strong a requirement.* Therefore, when solving such problems we have two options: (1) enforcement of jump conditions at every interface or; (2) weak formulations (weakening the regularity requirements). Weak forms, which are designed to accommodate irregular data and solutions, are usually preferred. *Numerical techniques employing weak forms, such as the FEM, have been developed with the essential property that whenever a smooth classical solution exists, it is also a solution to the weak form problem.* Therefore, we lose nothing by reformulating a problem in a more general way, by weakening the a priori smoothness requirements of the solution.

In the following few chapters, we shall initially consider a one-dimensional structure that occupies an open bounded domain in $\Omega \in \mathbb{R}$, with boundary $\partial\Omega$. The boundary consists of Γ_u on which the displacements (u), or any other primal variable (temperature in heat conduction applications, concentration in diffusion applications, etc. (see Appendix B) are prescribed and a part Γ_t on which tractions ($t \stackrel{\text{def}}{=} \sigma n$, n being the outward normal) are prescribed ($t = t^*$, Fig. 2.1). We now focus on weak forms of a one-dimensional version of Eq. 2.1

¹ Here f are the body forces.

Fig. 2.1 A one-dimensional body



$$\frac{d\sigma}{dx} + f = 0, \quad (\sigma = E \frac{du}{dx}), \quad (2.2)$$

where $E = E(x)$ is a spatially varying coefficient (Fig. 2.1). Thereafter, we will discuss three-dimensional problems.

2.2 Weak Formulations in One Dimension

To derive a direct weak formulation for a body, we take Eq. 2.2 (denoting the strong form), form a product with an arbitrary smooth scalar valued function ν , and integrate over the body

$$\int_{\Omega} \left(\frac{d\sigma}{dx} + f \right) \nu dx = \int_{\Omega} r \nu dx = 0, \quad (2.3)$$

where r is called the residual. We call ν a “test” function. If we were to add a condition that we do this for all ($\stackrel{\text{def}}{=} \forall$) possible “test” functions then

$$\int_{\Omega} \left(\frac{d\sigma}{dx} + f \right) \nu dx = \int_{\Omega} r \nu dx = 0, \quad (2.4)$$

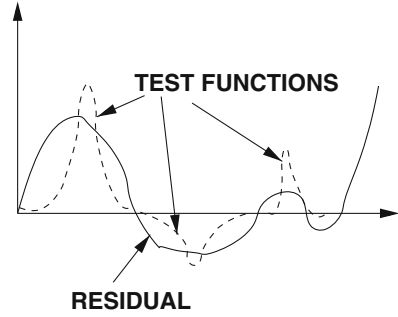
$\forall \nu$, implies $r(x) = 0$. Therefore, if every possible test function were considered, then $r = \frac{d\sigma}{dx} + f = 0$ on any finite region in (Ω) . Consequently, the weak and strong statements would be equivalent provided the true solution is smooth enough to have a strong solution. Clearly, r can never be nonzero over any finite region in the body, because the test function will “find” them (Fig. 2.2). Using the product rule of differentiation

$$\frac{d}{dx}(\sigma \nu) = \left(\frac{d\sigma}{dx} \right) \nu + \sigma \frac{d\nu}{dx} \quad (2.5)$$

leads to, $\forall \nu$

$$\int_{\Omega} \left(\frac{d}{dx}(\sigma \nu) - \sigma \frac{d\nu}{dx} \right) dx + \int_{\Omega} f \nu dx = 0, \quad (2.6)$$

Fig. 2.2 Test function actions on residuals



where we choose the ν from an admissible set, to be discussed momentarily. This leads to, $\forall \nu$

$$\int_{\Omega} \frac{d\nu}{dx} \sigma dx = \int_{\Omega} f \nu dx + \sigma \nu|_{\partial\Omega}, \quad (2.7)$$

which leads to

$$\int_{\Omega} \frac{d\nu}{dx} \sigma dx = \int_{\Omega} f \nu dx + t \nu|_{\partial\Omega}. \quad (2.8)$$

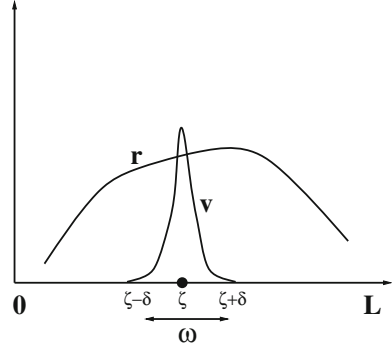
If we decide to restrict our choices of ν 's to those such that $\nu|_{\Gamma_u} = 0$ and $u|_{\Gamma_u} = d$, where d is the applied boundary displacement, on a displacement part of the boundary, Γ_u , we have

Find $u, u|_{\Gamma_u} = d$, such that $\forall \nu, \nu|_{\Gamma_u} = 0$

$$\underbrace{\int_{\Omega} \frac{d\nu}{dx} E \frac{du}{dx} dx}_{\stackrel{\text{def}}{=} \mathcal{B}(u, \nu)} = \underbrace{\int_{\Omega} f \nu dx + t \nu|_{\Gamma_t}}_{\stackrel{\text{def}}{=} \mathcal{F}(\nu)}. \quad (2.9)$$

This is called a *weak* form because it does not require the differentiability of σ . In other words, the differentiability requirements have been *weakened*. It is clear that we are able to consider problems with quite irregular solutions. We observe that if we test the solution with all possible test functions of sufficient smoothness, then the weak solution is equivalent to the strong solution. *We emphasize that provided the true solution is smooth enough, the weak and strong forms are equivalent, which can be seen by the above constructive derivation.* To explain the point more clearly, we consider a simple example.

Fig. 2.3 A residual function and a test function



2.3 An Example

Let us define a one-dimensional continuous function $r \in C^0(\Omega)$, on a one-dimensional domain, $\Omega = (\Omega)$. Our claim that

$$\int_{\Omega} r \nu \, dx = 0, \quad (2.10)$$

$\forall \nu \in C^0(\Omega)$, implies $r = 0$. This can be easily proven by contradiction. Suppose $r \neq 0$ at some point $\zeta \in \Omega$. Since $r \in C^0(\Omega)$, there must exist a subdomain (subinterval), $\omega \in \Omega$, defined through δ , $\omega \stackrel{\text{def}}{=} \zeta \pm \delta$ such that r has the same sign as at point ζ . Since ν is arbitrary, we may choose ν to be zero outside of this interval, and positive inside (Fig. 2.3). This would imply that

$$0 < \int_{\Omega} r \nu \, dx = \int_{\omega} r \nu \, dx = 0, \quad (2.11)$$

which is a contradiction. Now select

$$r = \frac{d\sigma}{dx} + f \in C^0(\Omega) \Rightarrow \frac{d}{dx} \left(E \frac{du}{dx} \right) + f \in C^0(\Omega) \Rightarrow u \in C^2(\Omega). \quad (2.12)$$

Therefore, for this model problem, the equivalence of weak and strong forms occurs if $u \in C^2(\Omega)$.

2.4 Some Restrictions

A key question is the selection of the sets of functions in the weak form. Somewhat naively, the answer is simple, the integrals must remain finite. Therefore, the following restrictions hold ($\forall \nu$), $\int_{\Omega} f \nu \, dx < \infty$, $\int_{\partial\Omega} t \nu \, dx < \infty$ and $\int_{\Omega} \frac{d\nu}{dx} \sigma \, dx < \infty$,

and govern the selection of the approximation spaces. These relations simply mean that the functions must be square integrable. In order to make precise statements one must have a method of “book keeping.” Such a system is to employ so-called Hilbertian Sobolev spaces. We recall that a norm has three main characteristics for any vectors u and v such that $\|u\| < \infty$ and $\|v\| < \infty$ are (1) $\|u\| > 0$, $\|u\| = 0$ if and only if $u = 0$ (“positivity”), (2) $\|u + v\| \leq \|u\| + \|v\|$ (Triangle Inequality) and (3) $\|\alpha u\| = |\alpha| \|u\|$, where α is a scalar constant (“scalability”). Certain types of norms, so-called Hilbert space norms, are frequently used in mathematical physics. Following standard notation, we denote $H^1(\Omega)$ as the usual space of scalar functions with generalized partial derivatives of order ≤ 1 in $L^2(\Omega)$, i.e., it is square integrable. In other words, $u \in H^1(\Omega)$ if

$$\|u\|_{H^1(\Omega)}^2 \stackrel{\text{def}}{=} \int_{\Omega} \frac{\partial u}{\partial x} \frac{\partial u}{\partial x} dx + \int_{\Omega} uu dx < \infty. \quad (2.13)$$

Using these definitions, a complete boundary value problem can be written as follows. The input data loading are assumed to be such that $f \in L^2(\Omega)$ and $t \in L^2(\Gamma_t)$, but less smooth data can be considered without complications. In summary we assume that our solutions obey these restrictions, leading to the following weak form

$$\begin{aligned} &\text{Find } u \in H^1(\Omega), u|_{\Gamma_u} = d, \text{ such that } \forall v \in H^1(\Omega), v|_{\Gamma_u} = 0 \\ &\int_{\Omega} \frac{dv}{dx} E \frac{du}{dx} dx = \int_{\Omega} f v dx + t v|_{\Gamma_t}. \end{aligned} \quad (2.14)$$

We note that if the data in (2.14) are smooth and if (2.14) possesses a solution u that is sufficiently regular, then u is the solution of the classical problem in strong form

$$\begin{aligned} &\frac{d}{dx} \left(E \frac{du}{dx} \right) + f = 0, \quad \forall \mathbf{x} \in \Omega, \\ &u = d, \quad \forall \mathbf{x} \in \Gamma_u, \\ &\left(E \frac{du}{dx} \right) n = t, \quad \forall \mathbf{x} \in \Gamma_t. \end{aligned} \quad (2.15)$$

2.5 Remarks on Nonlinear Problems

The treatment of nonlinear problems is outside the scope of this introductory monograph. However, a few comments are in order. The literature of solving nonlinear problems with the FEM is vast. This is a complex topic, that is best illustrated for students with an extremely simple one-dimensional example with material nonlinearities. Starting with

$$\frac{d}{dx} \underbrace{\left(E \underbrace{\left(\frac{du}{dx} \right)^p}_{\substack{\text{def} \\ \equiv \epsilon}} \right)}_{\substack{\text{def} \\ \equiv \sigma}} + f = 0. \quad (2.16)$$

The weak form reads

$$\int_0^L \frac{d\nu}{dx} \sigma dx = \int_0^L f \nu dx + t\nu|_{\Gamma_l}. \quad (2.17)$$

Using a Taylor series expansion of $\sigma(\epsilon(u))$ about a trial solution $u^{(k)}$ yields (k will be used as an iteration counter)

$$\begin{aligned} \sigma(u^{(k+1)}) &= E\epsilon^p(u^{(k+1)}) \\ &\approx E \left(\epsilon^p(u^{(k)}) + p\epsilon^{p-1}(u^{(k)}) \times (\epsilon(u^{(k+1)}) - \epsilon(u^{(k)})) + \mathcal{O}(\|u^{(k+1)} - u^{(k)}\|^2) \right) \\ &\approx 0. \end{aligned} \quad (2.18)$$

and substituting this into the weak form yields

$$\begin{aligned} \int_0^L \frac{d\nu}{dx} \left(E p \epsilon^{p-1}(u^{(k)}) \right) \epsilon(u^{(k+1)}) dx &= \int_0^L f \nu dx + t\nu|_{\Gamma_l} \\ &\quad - \int_0^L \frac{d\nu}{dx} E \left(\epsilon^p(u^{(k)}) + p(\epsilon^p(u^{(k)})) \right) dx. \end{aligned} \quad (2.19)$$

One then iterates $k = 1, 2, \dots$, until $\|u^{(k+1)} - u^{(k)}\| \leq TOL$. Convergence of such a Newton-type formulation is of concern. We refer the reader to the seminal book of Oden [1], which developed and pioneered nonlinear formulations and convergence analysis. For example, consider a general abstract nonlinear equation of the form

$$\Pi(u^{(k)}) = 0, \quad (2.20)$$

and the expansion

$$\Pi(u^{(k+1)}) = \Pi(u^{(k)}) + \nabla_u \Pi(u^{(k)}) \cdot (u^{(k+1)} - u^{(k)}) + \mathcal{O}(\|u^{(k+1)} - u^{(k)}\|^2) \approx 0. \quad (2.21)$$

A Newton update can be written in the following form

$$\mathbf{u}^{(k+1)} = \mathbf{u}^{(k)} - \left(\boldsymbol{\Pi}^{TAN}(\mathbf{u}^{(k)}) \right)^{-1} \cdot \boldsymbol{\Pi}(\mathbf{u}^{(k)}), \quad (2.22)$$

where $\boldsymbol{\Pi}^{TAN}(\mathbf{u}) \stackrel{\text{def}}{=} \nabla_{\mathbf{u}} \boldsymbol{\Pi}(\mathbf{u})$ is the so-called “tangent operator.” One immediately sees a potential difficulty, because of the possibility of a zero, or near zero, tangent when employing Newton’s method in a system that may have a nonmonotonic response, for example those involving material laws with softening. Specialized techniques can be developed for such problems, and we refer the reader to the state-of-the-art found in the book by Wriggers [2].

References

1. Oden, J. T. (1972). *Finite elements of non-linear continua*. New York: McGraw-Hill.
2. Wriggers, P. (2008). *Nonlinear finite element analysis*. Heidelberg: Springer.

Chapter 3

A Finite Element Implementation in One Dimension

3.1 Introduction

Classical techniques construct approximations from globally kinematically admissible functions, which we define as functions that satisfy the displacement boundary condition beforehand. Two main obstacles arise: (1) it may be very difficult to find a kinematically admissible function over the entire domain and (2) if such functions are found they lead to large, strongly coupled, and complicated systems of equations. These problems have been overcome by the fact that local approximations (posed over very small partitions of the entire domain) can deliver high quality solutions, and simultaneously lead to systems of equations which have an advantageous mathematical structure amenable to large-scale computation by high-speed computers. These piece-wise or “element-wise” approximations were recognized at least 60 years ago by Courant [1] as being quite advantageous. There have been a variety of such approximation methods to solve equations of mathematical physics. The most popular method of this class is the Finite Element Method. The central feature of the method is to partition the domain in a systematic manner into an assembly of discrete subdomains or “elements,” and then to approximate the solution of each of these pieces in a manner that couples them to form a global solution valid over the whole domain. The process is designed to keep the resulting algebraic systems as computationally manageable, and memory efficient, as possible.

3.2 Finite Element Method Implementation

Consider the following general form

$$\begin{aligned} &\text{Find } u \in H^1(\Omega) \text{ such that } \forall v \in H^1(\Omega), v|_{\Gamma_u} = 0 \\ &\int_{\Omega} \frac{dv}{dx} E \frac{du}{dx} dx = \int_{\Omega} f v dx + t v|_{\Gamma_t}. \end{aligned} \tag{3.1}$$

3.3 FEM Approximation

It is convenient to write the bilinear form in the following manner

$$\int_{\Omega} \frac{dv}{dx} E \frac{du}{dx} dx = \int_{\Omega} f v dx + t v|_{\Gamma_1}. \quad (3.2)$$

We approximate u by

$$u(x) = \sum_{j=1}^N a_j \phi_j(x). \quad (3.3)$$

If we choose v with the same approximation functions, but a different linear combination

$$v(x) = \sum_{i=1}^N b_i \phi_i(x), \quad (3.4)$$

then we may write

$$\begin{aligned} & \underbrace{\int_{\Omega} \frac{d}{dx} \left(\sum_{i=1}^N b_i \phi_i(x) \right) E \frac{d}{dx} \left(\sum_{j=1}^N a_j \phi_j(x) \right) dx}_{\stackrel{\text{def}}{=} \text{stiffness}} \\ &= \underbrace{\int_{\Omega} \left(\sum_{i=1}^N b_i \phi_i(x) \right) f dx}_{\stackrel{\text{def}}{=} \text{body load}} + \underbrace{\left(\left(\sum_{i=1}^N b_i \phi_i(x) \right) t \right) |_{\Gamma_1}}_{\stackrel{\text{def}}{=} \text{traction load}}. \end{aligned} \quad (3.5)$$

Since the b_i are arbitrary, $\forall v \Rightarrow \forall b_i$, therefore

$$\begin{aligned} \sum_{i=1}^N b_i \left(\sum_{j=1}^N K_{ij} a_j - R_i \right) &= 0 \Rightarrow [K] \{a\} = \{R\}, \\ K_{ij} &\stackrel{\text{def}}{=} \int_{\Omega} \frac{d\phi_i}{dx} E \frac{d\phi_j}{dx} dx \\ R_i &\stackrel{\text{def}}{=} \int_{\Omega} \phi_i f dx + \phi_i t|_{\Gamma_1}, \end{aligned} \quad (3.6)$$

where $[K]$ is an $N \times N$ (“stiffness”) matrix with components K_{ij} and $\{R\}$ is an $N \times 1$ (“load”) vector with components R_i . This is the system of equations that is to

be solved. Thus, large N implies large systems of equations, and more computational effort. However, with increasing N , we obtain more accurate approximate solutions. Note that large N does not seem like much of a concern for one-dimensional problems, but it is of immense concern for three-dimensional problems.

3.4 Construction of FEM Basis Functions

As mentioned, a primary problem with Galerkin's method is that it provides no systematic way of constructing approximation functions. The difficulties that arise include (1) ill-conditioned systems due to poor choice of approximation functions and (2) domains with irregular geometries. To circumvent these problems, the FEM defines basis (approximation) functions in a piecewise manner over a subdomain, "the finite elements," of the entire domain. The basis functions are usually simple polynomials of low degree. The following three criteria are important:

- The basis functions are smooth enough to be members of $H^1(\Omega)$.
- The basis functions are simple piecewise polynomials, defined element by element.
- The basis functions form a simple nodal basis where $\phi_i(x_j) = 0$ ($i \neq j$) and $\phi_i(x_i) = 1$, furthermore, $\sum_{i=1}^N \phi_i(x) = 1$ for all x .

A set of candidate functions are defined by

$$\phi(x) = \frac{x - x_{i-1}}{h_i} \quad \text{for } x_{i-1} \leq x \leq x_i, \quad (3.7)$$

where $h_i = x_i - x_{i-1}$ and

$$\phi(x) = 1 - \frac{x - x_i}{h_{i+1}} \quad \text{for } x_i \leq x \leq x_{i+1}, \quad (3.8)$$

and $\phi(x) = 0$ otherwise. The derivative of the function is

$$\frac{d\phi}{dx} = \frac{1}{h_i} \quad \text{for } x_{i-1} \leq x \leq x_i, \quad (3.9)$$

and

$$\frac{d\phi}{dx} = -\frac{1}{h_{i+1}} \quad \text{for } x_i \leq x \leq x_{i+1}. \quad (3.10)$$

The functions are arranged so that the "apex" of the i th function coincides with the i th node (Fig. 3.1). This framework provides many advantages, for example, simple numerical integration (Fig. 3.2).

Fig. 3.1 A one-dimensional finite element basis. At the *top*, is a uniform mesh example and at the *bottom*, a nonuniform mesh

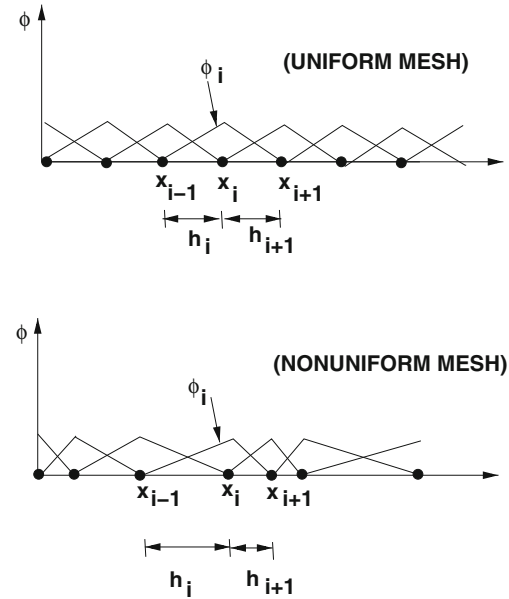


Fig. 3.2 Integration using Gaussian quadrature

