

# The finite element method

## 14.1 Introduction: the model problem

In Chapter 13 we explored finite difference methods for the numerical solution of two-point boundary value problems. The present chapter is devoted to the foundations of the theory of finite element methods. For the sake of simplicity the exposition will be, at least initially, confined to the second-order ordinary differential equation

$$-\frac{d}{dx} \left( p(x) \frac{du}{dx} \right) + r(x)u = f(x), \quad a < x < b, \quad (14.1)$$

where  $p \in C^1[a, b]$ ,  $r \in C[a, b]$ ,  $f \in L^2(a, b)$  and  $p(x) \geq c_0 > 0$ ,  $r(x) \geq 0$  for all  $x \in [a, b]$ , subject to the boundary conditions

$$u(a) = A, \quad u(b) = B. \quad (14.2)$$

Later on in the chapter, in Section 14.5, we shall also consider the ordinary differential equation

$$-\frac{d}{dx} \left( p(x) \frac{du}{dx} \right) + q(x) \frac{du}{dx} + r(x)u = f(x), \quad a < x < b, \quad (14.3)$$

subject to the boundary conditions (14.2). Indeed, much of the material discussed here can be extended to partial differential equations; for pointers to the relevant literature we refer to the Notes at the end of the chapter.

The finite element method was proposed in a paper by Richard Courant in the early 1940s,<sup>1</sup> although the historical roots of the method can be traced back to earlier work by Galerkin<sup>2</sup> in 1915; unfortunately, the relevance of Courant's article was not recognised at the time and the idea was forgotten. In the early 1950s the method was rediscovered by engineers, but its systematic mathematical analysis began only a decade later. Since then, the finite element method has been developed into one of the most general and powerful techniques for the numerical solution of differential equations which is widely used in engineering design and analysis.

Unlike finite difference schemes which seek to approximate the unknown analytical solution to a differential equation at a finite number of selected points, the grid points or mesh points in the computational domain, the finite element method supplies an approximation to the analytical solution in the form of a piecewise polynomial function, defined over the entire computational domain. For example, in the case of the boundary value problem (14.1), (14.2), the simplest finite element method uses a linear spline, defined over the interval  $[a, b]$ , to approximate the analytical solution  $u$ .

We shall consider two techniques for the construction of finite element approximations: the **Rayleigh–Ritz principle** and the **Galerkin principle**. In the case of the boundary value problem (14.1), (14.2) the approximations which stem from these two principles will be seen to coincide. We note, however, that since the Rayleigh–Ritz principle relies on the fact that the boundary value problem under consideration can be restated as a variational problem involving the minimisation of a certain quadratic functional over a function space, its use is restricted to *symmetric* boundary value problems, such as (14.1), (14.2) where (14.1) does not contain a first-derivative term; for example, the Rayleigh–Ritz principle is not applicable to (14.3), (14.2) unless  $q(x) \equiv 0$ . The precise sense in which the word *symmetric* is to be interpreted here will be clar-

<sup>1</sup> R. Courant, Variational methods for the solution of problems in equilibrium and vibrations, *Bull. Amer. Math. Soc.* **49**, 1–23, 1943; Richard Courant (8 January 1888, Lublinitz, Prussia, Germany (now Lubliniec, Poland) – 27 January 1972, New Rochelle, New York, USA). For an illuminating account of the lives of Richard Courant and David Hilbert, see the book of Constance Reid: *Hilbert–Courant*, Springer, New York, 1986.

<sup>2</sup> Boris Grigorievich Galerkin (4 March 1871, Polotsk, Russia (now in Belarus) – 12 June 1945, Moscow, USSR) studied mathematics and engineering at the St Petersburg Technological Institute. During his studies he supported himself by private tutoring and working as a designer. His ideas on the approximate solution of differential equations were published in 1915. From 1940 until his death, Galerkin was head of the Institute of Mechanics of the Soviet Academy of Sciences.

ified later in the chapter. On the other hand, as we shall see in Section 14.5, the Galerkin principle is more generally applicable and does not require symmetry of the boundary value problem.

To make these observations rigorous, we recall from Chapter 11 the concept of **Sobolev space**.

**Definition 14.1** For a positive integer  $k$ , we define the **Sobolev space**  $H^k(a, b)$  as the set of real-valued functions  $v$  defined on  $[a, b]$  such that  $v$  and all of its derivatives of order up to and including  $k-1$  are absolutely continuous on  $[a, b]$  and

$$v^{(k)} = \frac{d^k v}{dx^k} \in L^2(a, b).$$

Here  $L^2(a, b)$  denotes the set of all functions defined on  $(a, b)$  such that

$$\|v\|_2 = \|v\|_{L^2(a, b)} = \left( \int_a^b |v(x)|^2 dx \right)^{1/2}$$

is finite. We equip  $H^k(a, b)$  with the **Sobolev norm**

$$\|v\|_{H^k(a, b)} = \left( \sum_{m=0}^k \|v^{(m)}\|_{L^2(a, b)}^2 \right)^{1/2},$$

where  $v^{(0)} = v$ .

The Sobolev spaces  $H^1(a, b)$  and  $H^2(a, b)$  corresponding, respectively, to  $k = 1$  and  $k = 2$  will be particularly relevant in this chapter. The next definition introduces variants of the space  $H^1(a, b)$  required for the imposition of the boundary conditions (14.2).

**Definition 14.2** (i) Given that  $A$  and  $B$  are real numbers,  $H_E^1(a, b)$  will denote the set of all functions  $v \in H^1(a, b)$  such that  $v(a) = A$  and  $v(b) = B$ .

(ii)  $H_0^1(a, b)$  will signify the set of all functions  $v \in H^1(a, b)$  such that  $v(a) = 0$  and  $v(b) = 0$ .

In the next section we shall state, using Sobolev spaces, the Rayleigh–Ritz and Galerkin principles associated with the boundary value problem (14.1), (14.2), and explore their relationship.

## 14.2 Rayleigh–Ritz and Galerkin principles

The Rayleigh–Ritz principle relies on converting the boundary value problem (14.1), (14.2) into a variational problem involving the minimisation of a certain quadratic functional over a function space.

Let us define the quadratic functional  $\mathcal{J}: H_E^1(a, b) \rightarrow \mathbb{R}$  by

$$\mathcal{J}(w) = \frac{1}{2} \int_a^b [p(x)(w')^2 + r(x)w^2] dx - \int_a^b f(x)w(x) dx$$

where  $w \in H_E^1(a, b)$ , and consider the following *variational problem*:

$$(RR) \quad \text{find } u \in H_E^1(a, b) \text{ such that } \mathcal{J}(u) = \min_{w \in H_E^1(a, b)} \mathcal{J}(w),$$

which we shall henceforth refer to as the **Rayleigh–Ritz principle**. For the sake of notational simplicity we define

$$\mathcal{A}(w, v) = \int_a^b [p(x) w'(x) v'(x) + r(x) w(x) v(x)] dx$$

and recall from Chapter 9 the definition of inner product on  $L^2(a, b)$ :

$$\langle w, v \rangle = \int_a^b w(x) v(x) dx.$$

Using these, we can rewrite  $\mathcal{J}(w)$  as follows:

$$\mathcal{J}(w) = \frac{1}{2} \mathcal{A}(w, w) - \langle f, w \rangle, \quad w \in H_E^1(a, b). \quad (14.4)$$

The mapping  $\mathcal{A}: H^1(a, b) \times H^1(a, b) \rightarrow \mathbb{R}$  is a **bilinear functional** in the following sense:

- ❶  $\mathcal{A}(\lambda_1 w_1 + \lambda_2 w_2, v) = \lambda_1 \mathcal{A}(w_1, v) + \lambda_2 \mathcal{A}(w_2, v)$   
for all  $\lambda_1, \lambda_2 \in \mathbb{R}$  and all  $w_1, w_2, v \in H^1(a, b)$ ;
- ❷  $\mathcal{A}(w, \mu_1 v_1 + \mu_2 v_2) = \mu_1 \mathcal{A}(w, v_1) + \mu_2 \mathcal{A}(w, v_2)$   
for all  $\mu_1, \mu_2 \in \mathbb{R}$  and all  $w, v_1, v_2 \in H^1(a, b)$ .

We note, in addition, that the bilinear functional  $\mathcal{A}(\cdot, \cdot)$  is **symmetric**, in that

$$\mathcal{A}(w, v) = \mathcal{A}(v, w) \quad \forall w, v \in H^1(a, b). \quad (14.5)$$

Our next result provides an equivalent characterisation of the Rayleigh–Ritz principle; it relies on the fact that the bilinear functional  $\mathcal{A}(\cdot, \cdot)$  is symmetric in the sense of (14.5).

**Theorem 14.1** A function  $u$  in  $H_E^1(a, b)$  minimises  $\mathcal{J}(\cdot)$  over  $H_E^1(a, b)$  if, and only if,

$$(G) \quad \mathcal{A}(u, v) = \langle f, v \rangle \quad \forall v \in H_0^1(a, b). \quad (14.6)$$

This identity will be referred to as the **Galerkin principle**.

*Proof of theorem* Suppose that  $u \in H_E^1(a, b)$  minimises  $\mathcal{J}(\cdot)$  over  $H_E^1(a, b)$ ; that is,  $\mathcal{J}(u) \leq \mathcal{J}(w)$  for all  $w \in H_E^1(a, b)$ . Noting that  $w = u + \lambda v$  belongs to  $H_E^1(a, b)$  for all  $\lambda \in \mathbb{R}$  and all  $v \in H_0^1(a, b)$ , we deduce that

$$\begin{aligned} \mathcal{J}(u) &\leq \mathcal{J}(u + \lambda v) = \frac{1}{2} \mathcal{A}(u + \lambda v, u + \lambda v) - \langle f, u + \lambda v \rangle \\ &= \mathcal{J}(u) + \lambda [\mathcal{A}(u, v) - \langle f, v \rangle] + \frac{1}{2} \lambda^2 \mathcal{A}(v, v) \end{aligned} \quad (14.7)$$

for all  $v \in H_0^1(a, b)$  and all  $\lambda \in \mathbb{R}$ . Here, in the transition from the first line to the second we made use of the fact that  $\mathcal{A}(u, v) = \mathcal{A}(v, u)$  for all  $v$  in  $H_0^1(a, b)$ , which follows from (14.5). Now, (14.7) implies that

$$-\frac{1}{2} \lambda^2 \mathcal{A}(v, v) \leq \lambda [\mathcal{A}(u, v) - \langle f, v \rangle]$$

for all  $v \in H_0^1(a, b)$  and all  $\lambda \in \mathbb{R}$ . Let us suppose that  $\lambda > 0$ , divide both sides of the last inequality by  $\lambda$  and pass to the limit  $\lambda \rightarrow 0$  to deduce that

$$0 \leq \mathcal{A}(u, v) - \langle f, v \rangle \quad \forall v \in H_0^1(a, b). \quad (14.8)$$

On replacing  $v$  by  $-v$  in (14.8), we have that also

$$0 \geq \mathcal{A}(u, v) - \langle f, v \rangle \quad \forall v \in H_0^1(a, b). \quad (14.9)$$

We conclude from (14.8) and (14.9) that

$$\mathcal{A}(u, v) = \langle f, v \rangle \quad \forall v \in H_0^1(a, b), \quad (14.10)$$

as required.

Conversely, if  $u \in H_E^1(a, b)$  is such that  $\mathcal{A}(u, v) = \langle f, v \rangle$  for all  $v$  in  $H_0^1(a, b)$ , then

$$\mathcal{J}(u + \lambda v) = \mathcal{J}(u) + \lambda [\mathcal{A}(u, v) - \langle f, v \rangle] + \frac{1}{2} \lambda^2 \mathcal{A}(v, v) \geq \mathcal{J}(u)$$

for all  $v \in H_0^1(a, b)$  and all  $\lambda \in \mathbb{R}$ ; therefore,  $u$  minimises  $\mathcal{J}(\cdot)$  over  $H_E^1(a, b)$ .  $\square$

Thus we have shown that, as long as  $\mathcal{A}(\cdot, \cdot)$  is a symmetric bilinear functional,  $u \in H_E^1(a, b)$  satisfies the Rayleigh–Ritz principle if, and only if, it satisfies the Galerkin principle.<sup>1</sup> Our next task is to explain the

<sup>1</sup> In the language of the calculus of variations, (G) is the Euler–Lagrange equation for the minimisation problem (RR).

relationship between (RR) and (G) on the one-hand and (14.1), (14.2) on the other. Since in the case of a symmetric bilinear functional  $\mathcal{A}(\cdot, \cdot)$  the principles (RR) and (G) are equivalent, it is sufficient to clarify the connection between (G), for example, and the boundary value problem (14.1), (14.2).

We begin with the following definition.

**Definition 14.3** *If a function  $u \in H_E^1(a, b)$  satisfies the Galerkin principle (14.6), it is called a **weak solution** to the boundary value problem (14.1), (14.2), and the Galerkin principle is referred to as the **weak formulation** of the boundary value problem (14.1), (14.2).*

Let us justify this terminology. Suppose that  $u \in H^2(a, b) \cap H_E^1(a, b)$  is a solution to the boundary value problem (14.1), (14.2). Then,

$$-\frac{d}{dx} \left( p(x) \frac{du}{dx} \right) + r(x)u = f(x), \quad (14.11)$$

for almost every  $x \in (a, b)$  (see the discussion prior to Example 11.1 for a definition of **almost every**). Multiplying this equality by an arbitrary function  $v \in H_0^1(a, b)$ , and integrating over  $(a, b)$ , we conclude that

$$-\int_a^b \frac{d}{dx} \left( p(x) \frac{du}{dx} \right) v \, dx + \int_a^b r(x)uv \, dx = \int_a^b f(x)v(x) \, dx.$$

On integration by parts in the first term on the left-hand side,

$$-\int_a^b \frac{d}{dx} \left( p(x) \frac{du}{dx} \right) v \, dx = \left[ p(x) \frac{du}{dx} v \right]_{x=a}^b + \int_a^b p(x) \frac{du}{dx} \frac{dv}{dx} \, dx.$$

Since, by hypothesis,  $v(a) = 0$  and  $v(b) = 0$ , it follows that

$$\int_a^b p(x) \frac{du}{dx} \frac{dv}{dx} \, dx + \int_a^b r(x)uv \, dx = \int_a^b f(x)v(x) \, dx$$

for all  $v \in H_0^1(a, b)$ . Thus, we have shown the following result.

**Theorem 14.2** *If  $u \in H^2(a, b) \cap H_E^1(a, b)$  is a solution to the boundary value problem (14.1), (14.2), then  $u$  is a weak solution to this problem; that is,*

$$\mathcal{A}(u, v) = \langle f, v \rangle \quad \forall v \in H_0^1(a, b). \quad (14.12)$$

The converse implication, namely that any weak solution  $u \in H_E^1(a, b)$  of (14.1), (14.2) belongs to  $H^2(a, b) \cap H_E^1(a, b)$  and solves (14.1), (14.2) in the usual (pointwise) sense, is not true in general, unless the weak

solution can be shown to be sufficiently smooth to belong to  $H^2(a, b)$ . It is for this reason that any function  $u \in H_E^1(a, b)$  satisfying (14.12) is called a *weak* solution of the original boundary value problem.

Thus, Theorem 14.1 shows that  $u \in H_E^1(a, b)$  is a weak solution to (14.1), (14.2) if, and only if, it minimises  $\mathcal{J}(\cdot)$  over  $H_E^1(a, b)$ . Next, we show that if a weak solution exists then it must be unique.

**Theorem 14.3** *The boundary value problem (14.1), (14.2) possesses at most one weak solution in  $H_E^1(a, b)$ .*

*Proof* The proof is by contradiction. Suppose that  $u \in H_E^1(a, b)$  and  $\tilde{u} \in H_E^1(a, b)$  are two weak solutions to (14.1), (14.2). Then,  $u - \tilde{u}$  belongs to  $H_0^1(a, b)$ , and

$$\mathcal{A}(u - \tilde{u}, v) = \mathcal{A}(u, v) - \mathcal{A}(\tilde{u}, v) = \langle f, v \rangle - \langle f, v \rangle = 0$$

for all  $v \in H_0^1(a, b)$ . In particular,

$$\mathcal{A}(u - \tilde{u}, u - \tilde{u}) = 0.$$

However, since  $p(x) \geq c_0 > 0$  and  $r(x) \geq 0$  for all  $x$  in  $[a, b]$ ,

$$\mathcal{A}(v, v) = \int_a^b [p(x)(v')^2 + r(x)v^2] dx \geq c_0 \int_a^b |v'|^2 dx.$$

On choosing  $v = u - \tilde{u}$ , this implies that

$$0 = \mathcal{A}(u - \tilde{u}, u - \tilde{u}) \geq c_0 \int_a^b |(u - \tilde{u})'|^2 dx.$$

Since the right-hand side in the last inequality is nonnegative, it follows that  $(u - \tilde{u})'(x) = 0$  for almost every  $x$  in  $(a, b)$ ; as  $u - \tilde{u}$  is absolutely continuous on  $[a, b]$  and  $(u - \tilde{u})(a) = (u - \tilde{u})(b) = 0$ , we conclude that  $u = \tilde{u}$ , and hence we get the desired uniqueness of a weak solution.  $\square$

It turns out that under the present hypotheses on  $p$ ,  $q$  and  $f$  the *existence* of a weak solution  $u \in H_E^1(a, b)$  is also ensured, although the proof of this is less simple and is omitted here; the interested reader is referred to the literature listed in the Notes at the end of the chapter.

### 14.3 Formulation of the finite element method

In the previous section we showed that the weak solution to the boundary value problem (14.1), (14.2) minimises  $\mathcal{J}(\cdot)$  over  $H_E^1(a, b)$ . The finite element method is based on constructing an approximate solution  $u^h$  to

the problem by minimising  $\mathcal{J}(\cdot)$  over a finite-dimensional subset  $S_E^h$  of  $H_E^1(a, b)$ , instead.

A simple way of constructing  $S_E^h$  is to choose any function  $\psi \in H_E^1(a, b)$ , for example,

$$\psi(x) = \frac{B-A}{b-a}(x-a) + A \quad (14.13)$$

and a finite set of linearly independent functions  $\varphi_j$ ,  $j = 1, \dots, n-1$ , in  $H_0^1(a, b)$  for  $n \geq 2$ , and then define

$$S_E^h = \{v^h \in H_E^1(a, b): v^h(x) = \psi(x) + \sum_{i=1}^{n-1} v_i \varphi_i(x), \\ \text{where } (v_1, \dots, v_{n-1})^T \in \mathbb{R}^{n-1}\}.$$

We consider the following approximation of problem (RR):

$$(RR)^h \quad \text{find } u^h \in S_E^h \text{ such that } \mathcal{J}(u^h) = \min_{w^h \in S_E^h} \mathcal{J}(w^h).$$

Our next result is a finite-dimensional analogue of Theorem 14.1.

**Theorem 14.4** *A function  $u^h \in S_E^h$  minimises  $\mathcal{J}(\cdot)$  over  $S_E^h$  if, and only if,*

$$(G)^h \quad \mathcal{A}(u^h, v^h) = \langle f, v^h \rangle \quad \forall v^h \in S_0^h. \quad (14.14)$$

Here,

$$S_0^h = \{v^h \in H_0^1(a, b): v^h(x) = \sum_{i=1}^{n-1} v_i \varphi_i(x), \\ \text{where } (v_1, \dots, v_{n-1})^T \in \mathbb{R}^{n-1}\}.$$

The problem  $(G)^h$  can be thought of as an approximation to the Galerkin principle (G), and is therefore referred to as the **Galerkin method**. For a similar reason,  $(RR)^h$  is called the Rayleigh–Ritz method, or just **Ritz method**. Thus, in complete analogy with the equivalence of (RR) and (G) formulated in Theorem 14.1, Theorem 14.4 now expresses the equivalence of  $(RR)^h$  and  $(G)^h$ , the approximations to (RR) and (G), respectively. Of course, as in the case of (RR) and (G), the equivalence of  $(RR)^h$  and  $(G)^h$  relies on the assumption that the bilinear functional  $\mathcal{A}(\cdot, \cdot)$  is symmetric. The proof is identical to that of Theorem 14.1, and is left as an exercise.

Theorem 14.4 provides no information about the existence and uniqueness of  $u^h$  that minimises  $\mathcal{J}(\cdot)$  over  $S_E^h$  (or, equivalently, of the existence



and uniqueness of  $u^h$  that satisfies (14.14)). This question is settled by our next result.

**Theorem 14.5** *There exists a unique function  $u^h \in S_E^h$  that minimises  $\mathcal{J}(\cdot)$  over  $S_E^h$ ; this  $u^h$  is called the **Ritz approximation** to  $u$ . Equivalently, there exists a unique function  $u^h \in S_E^h$  that satisfies (14.14); this  $u^h$  is called the **Galerkin approximation** to  $u$ . The Ritz and Galerkin approximations to  $u$  coincide.*

*Proof* We shall prove the second of these two equivalent statements: we shall show that there exists a unique  $u^h \in S_E^h$  that satisfies (14.14). The proof of uniqueness of  $u^h \in S_E^h$  is analogous to the proof of Theorem 14.3, with  $u$ ,  $\tilde{u}$ ,  $H_E^1(a, b)$  and  $H_0^1(a, b)$ , replaced by  $u^h$ ,  $\tilde{u}^h$ ,  $S_E^h$  and  $S_0^h$ , respectively. Since  $S_E^h$  is finite-dimensional, the uniqueness of  $u^h$  satisfying (14.14) implies its existence.  $\square$

Having shown the existence and uniqueness of  $u^h$  minimising  $\mathcal{J}(\cdot)$  over  $S_E^h$  (or, equivalently, satisfying (14.14)), we adopt the following definition.

**Definition 14.4** *The functions  $\varphi_i$ ,  $i = 1, 2, \dots, n-1$ , appearing in the definitions of  $S_E^h$  and  $S_0^h$  are called the **Galerkin basis functions**.*

Since any function  $v^h \in S_0^h$  can be represented as a linear combination of the Galerkin basis functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , it is clear that (14.14) is equivalent to

$$\mathcal{A}(u^h, \varphi_i) = \langle f, \varphi_i \rangle, \quad 1 \leq i \leq n-1. \quad (14.15)$$

As  $u^h$  belongs to  $S_E^h$ , it can be expressed in terms of  $\psi$  and the Galerkin basis functions as

$$u^h(x) = \psi(x) + \sum_{j=1}^{n-1} u_j \varphi_j(x),$$

where  $u_j \in \mathbb{R}$ ,  $j = 1, \dots, n-1$ , are to be determined. On substituting this expansion of  $u^h$  into (14.15), we arrive at the following system of simultaneous linear equations:

$$\sum_{j=1}^{n-1} M_{ij} u_j = b_i, \quad 1 \leq i \leq n-1, \quad (14.16)$$

where

$$M_{ij} = \mathcal{A}(\varphi_j, \varphi_i), \quad b_i = \langle f, \varphi_i \rangle - \mathcal{A}(\psi, \varphi_i). \quad (14.17)$$

The coefficients  $u_j$ ,  $1 \leq j \leq n-1$ , in the representation of the approximate solution are thus obtained by solving the system of linear equations (14.16). The matrix  $M$  is, clearly, symmetric (since the bilinear form  $\mathcal{A}(\cdot, \cdot)$  is symmetric by hypothesis) and positive definite, because

$$\mathbf{v}^T M \mathbf{v} = \mathcal{A}(v, v) > 0,$$

where  $\mathbf{v} = (v_1, v_2, \dots, v_{n-1})^T \in \mathbb{R}^{n-1}$  is any nonzero vector and  $v = v_1\varphi_1 + \dots + v_{n-1}\varphi_{n-1} \in S_0^h$ .

The Ritz and Galerkin methods can be used to compute an approximation  $u^h$  to  $u$  as a linear combination of *any* finite set of linearly independent functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , in  $H_0^1(a, b)$ . We obtain the **Ritz finite element method** and the **Galerkin finite element method**, respectively, when we select the approximating subspaces  $S_E^h$  and  $S_0^h$  in the Ritz or the Galerkin method to be spaces of spline functions (see Chapter 11). Here we only consider the simplest case of linear splines, and choose the basis functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , to be the hat functions (11.4). We begin by fixing a set of points  $x_k$ ,  $k = 0, 1, \dots, n$ ,  $n \geq 2$ , in the interval  $[a, b]$  such that

$$a = x_0 < x_1 < \dots < x_n = b. \quad (14.18)$$

The intervals  $[x_{i-1}, x_i]$ ,  $1 \leq i \leq n$ , are referred to as **elements**; hence the name *finite element method*. In the theory of the finite element methods (14.18) is called a **subdivision** of the computational domain  $[a, b]$ , and the points  $x_k$  are called **mesh points**. The function  $\varphi_i$  is the piecewise linear function which takes the value 0 at all the mesh points except  $x_i$ , where it takes the value 1. Thus,

$$\varphi_i(x) = \begin{cases} (x - x_{i-1})/h_i & \text{if } x_{i-1} \leq x \leq x_i, \\ (x_{i+1} - x)/h_{i+1} & \text{if } x_i \leq x \leq x_{i+1}, \\ 0 & \text{otherwise,} \end{cases} \quad (14.19)$$

where  $h_i = x_i - x_{i-1}$ . The functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , are called the (piecewise linear) **finite element basis functions** and the associated Galerkin approximation  $u^h$  is referred to as the (piecewise linear) **finite element approximation** of  $u$ . The closure of the interval  $(x_{i-1}, x_{i+1})$  over which  $\varphi_i$  is nonzero is called the **support** of the function  $\varphi_i$ . The piecewise linear finite element basis function  $\varphi_i$ ,  $1 \leq i \leq n-1$ , with support  $[x_{i-1}, x_{i+1}]$ , is depicted in Figure 14.1.

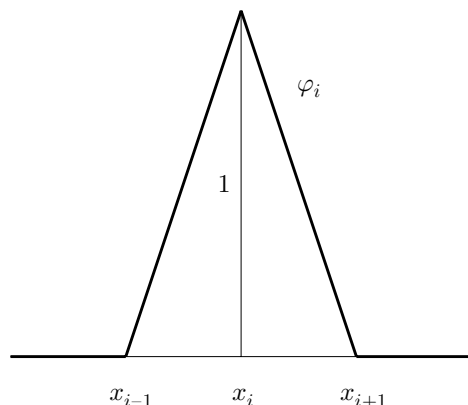


Fig. 14.1. A piecewise linear finite element basis function,  $\varphi_i$ ,  $1 \leq i \leq n-1$ .

For the finite element method the important property of the basis functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , is that they have *local* support, being nonzero only in one pair of adjacent intervals,  $(x_{i-1}, x_i]$  and  $[x_i, x_{i+1})$ . This means that, in the matrix  $M$ ,

$$M_{ij} = 0 \quad \text{if } |i - j| > 1.$$

The matrix  $M$  is, therefore, symmetric, positive definite and tridiagonal, and the associated system of linear equations can be solved very efficiently by the methods of Section 3.3, the most efficient algorithm being LU decomposition, without any use of symmetry. The fact that  $M$  is positive definite means that no interchanges are necessary.

The function  $\psi$  in (14.13), which is included in the definition of  $S_E^h$  to ensure that  $u^h$  satisfies the boundary conditions at  $x = a$  and  $x = b$ , is then given by

$$\psi(x) = A\varphi_0(x) + B\varphi_n(x),$$

which is also piecewise linear; clearly,  $\psi(a) = A$  and  $\psi(b) = B$ . Here,  $\varphi_0$  and  $\varphi_n$  are defined by setting, respectively,  $i = 0$  and  $i = n$  in (14.19) and restricting the resulting functions to the interval  $[a, b] = [x_0, x_n]$ . In (14.17) we see that the term  $\mathcal{A}(\psi, \varphi_i)$  is nonzero only for  $i = 1$  and  $i = n-1$ .

Before attempting to solve the system of linear equations we must, of course, first compute the elements of the matrix  $M$ , and the quantities on the right-hand side,  $b_i$ ,  $i = 1, \dots, n-1$ ; see (14.16) and (14.17). The

matrix elements are obtained from

$$M_{ij} = \mathcal{A}(\varphi_j, \varphi_i) = \int_a^b p(x) \varphi_j'(x) \varphi_i'(x) dx + \int_a^b r(x) \varphi_j(x) \varphi_i(x) dx,$$

with  $1 \leq i, j \leq n-1$ . We have written this as the sum of two terms, as the matrix  $M$  is often written in this way as the sum of two matrices which, for historical reasons, are often known as the **stiffness matrix** and the **mass matrix**, respectively. The terms  $M_{ij}$  are very simple; in fact in the first integral the derivatives  $\varphi_j'$  and  $\varphi_i'$  are piecewise constant functions over  $[a, b]$ .

It may be possible to compute these integrals analytically, but more generally some form of numerical quadrature will be necessary. It is then easy to show that if we use certain types of quadrature formulae we shall be led to the same system of equations as in the finite difference method of Section 13.5. Consider the particularly simple case where the mesh points are equally spaced, so that  $x_j = a + jh$ ,  $j = 0, 1, \dots, n$ ,  $h = (b-a)/n$ . If we then approximate the integrals involved in the stiffness matrix by the midpoint rule (see Chapter 10), we obtain

$$\begin{aligned} \int_{x_{i-1}}^{x_i} p(x) \varphi_{i-1}'(x) \varphi_i'(x) dx &= -(1/h^2) \int_{x_{i-1}}^{x_i} p(x) dx \\ &\approx -p_{i-1/2}/h, \end{aligned}$$

where  $p_{i-1/2} = p(x_{i-1/2})$ , and similarly for the other integrals involved. For the integrals in the mass matrix we use the trapezium rule, and then

$$\int_{x_{i-1}}^{x_i} r(x) \varphi_{i-1}(x) \varphi_i(x) dx \approx 0,$$

since  $\varphi_i$  is zero at  $x_{i-1}$  and  $\varphi_{i-1}$  is zero at  $x_i$ . In the same way

$$\int_{x_{i-1}}^{x_i} r(x) [\varphi_i(x)]^2 dx \approx \frac{1}{2} h r_i,$$

where  $r_i = r(x_i)$ , since  $\varphi_i$  is zero at one end of the interval and unity at the other. The other part of the integral is, similarly,

$$\int_{x_i}^{x_{i+1}} r(x) [\varphi_i(x)]^2 dx \approx \frac{1}{2} h r_i. \quad (14.20)$$

Assuming that  $f \in C[a, b]$ , approximating the integral on the right-hand side by the trapezium rule in the same way, and putting all the parts together, equation (14.14) now takes the approximate form

$$-\frac{p_{i-1/2}}{h} u_{i-1} + \frac{p_{i-1/2} + p_{i+1/2}}{h} u_i - \frac{p_{i+1/2}}{h} u_{i+1} + h r_i u_i = h f_i,$$

for  $i = 1, 2, \dots, n-1$ , with the notational convention that  $u_0 = A$  and  $u_n = B$ , and  $f_i = f(x_i)$ ; clearly, this is the same as the finite difference equation (13.19). Of course, had we used a different set of basis functions  $\varphi_i$ ,  $1 \leq i \leq n-1$ , or different numerical quadrature rules, the finite element and finite difference methods would have no longer been identical. Indeed, this example is just an illustration of the relation between the two methods; we should normally expect to compute the entries of the matrix  $M$  by using some more accurate quadrature method, such as a two-point Gauss formula.

In the next two sections we shall assess the accuracy of the finite element method. Our goal is to quantify the amount of reduction in the error  $u - u^h$  as the mesh spacing  $h$  is reduced.

#### 14.4 Error analysis of the finite element method

We begin with a fundamental result that underlies the error analysis of finite element methods.

**Theorem 14.6 (Céa's Lemma)** *Suppose that  $u$  is the function that minimises  $\mathcal{J}(u)$  over  $H_E^1(a, b)$  (or, equivalently, that  $u$  satisfies (14.6)), and that  $u^h$  is its Galerkin approximation obtained by minimising  $\mathcal{J}(\cdot)$  over  $S_E^h$  (or, equivalently, that  $u^h$  satisfies (14.14)). Then,*

$$\mathcal{A}(u - u^h, v^h) = 0 \quad \forall v^h \in S_0^h, \quad (14.21)$$

and

$$\mathcal{A}(u - u^h, u - u^h) = \min_{v^h \in S_E^h} \mathcal{A}(u - v^h, u - v^h). \quad (14.22)$$

The identity (14.21) is referred to as **Galerkin orthogonality**. The terminology stems from the fact that, since the bilinear functional  $\mathcal{A}(\cdot, \cdot)$  is symmetric and  $\mathcal{A}(v, v) > 0$  for all  $v \in H_0^1(a, b) \setminus \{0\}$ ,  $\mathcal{A}(\cdot, \cdot)$  is an inner product in the linear space  $H_0^1(a, b)$ . Therefore, by virtue of Definition 9.2, (14.21) means that  $u - u^h$  is orthogonal to  $S_0^h$  in  $H_0^1(a, b)$ . A geometrical illustration of Galerkin orthogonality is given in Figure 14.2. Given that  $\psi$  is a fixed element of  $H_E^1(a, b)$ , the mapping

$$R^h: u - \psi \in H_0^1(a, b) \mapsto u^h - \psi \in S_0^h$$

which assigns a  $u^h \in S_E^h$  to  $u \in H_E^1(a, b)$  (where  $u$  and  $u^h$  are as in Theorem 14.6) is called the **Ritz projector**.

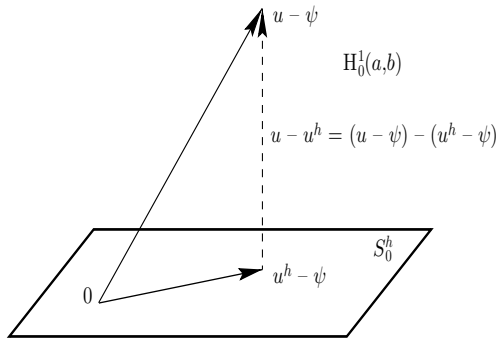


Fig. 14.2. Illustration of the Galerkin orthogonality property of the finite element method.  $\mathcal{A}((u - \psi) - (u^h - \psi), v^h) = \mathcal{A}(u - u^h, v^h) = 0$  for all  $v^h$  in  $S_0^h$ . Here,  $\psi(x) = A\varphi_0(x) + B\varphi_n(x)$ , so that  $u - \psi \in H_0^1(a, b)$  and  $u^h - \psi \in S_0^h$ . The 0 in the figure denotes the zero element of the linear space  $S_0^h$  (and, simultaneously, that of  $H_0^1(a, b)$ ), namely the function that is identically zero on the interval  $(a, b)$ .

*Proof of theorem* By the definition of the Galerkin method  $(G)^h$ ,

$$\mathcal{A}(u^h, v^h) = \langle f, v^h \rangle \quad \forall v^h \in S_0^h.$$

On the other hand, we deduce from  $(G)$  that

$$\mathcal{A}(u, v^h) = \langle f, v^h \rangle \quad \forall v^h \in S_0^h,$$

since  $v^h \in S_0^h \subset H_0^1(a, b)$ . The Galerkin orthogonality property (14.21) follows by subtraction.

Now suppose that  $v^h$  is any function in  $S_E^h$ ; then,

$$\begin{aligned} \mathcal{A}(u - v^h, u - v^h) &= \mathcal{A}(u - u^h + u^h - v^h, u - u^h + u^h - v^h) \\ &= \mathcal{A}(u - u^h, u - u^h) + \mathcal{A}(u^h - v^h, u^h - v^h) \\ &\quad + 2\mathcal{A}(u - u^h, u^h - v^h) \\ &= \mathcal{A}(u - u^h, u - u^h) + \mathcal{A}(u^h - v^h, u^h - v^h), \end{aligned}$$

by Galerkin orthogonality, given that  $u^h - v^h \in S_0^h$ . In the transition from the first line to the second, we made use of the fact that the bilinear functional  $\mathcal{A}$  is symmetric. As the term  $\mathcal{A}(u^h - v^h, u^h - v^h)$  is nonnegative, we deduce that

$$\mathcal{A}(u - u^h, u - u^h) \leq \mathcal{A}(u - v^h, u - v^h) \quad \forall v^h \in S_0^h,$$

with equality when  $v^h = u^h$ ; hence (14.22).  $\square$

Motivated by the minimisation property (14.22), we define the **energy norm**  $\|\cdot\|_{\mathcal{A}}$  on  $H_0^1(a, b)$  via

$$\|v\|_{\mathcal{A}} = [\mathcal{A}(v, v)]^{1/2}. \quad (14.23)$$

Under our hypotheses on  $p$  and  $q$ , it is easy to see that  $\|\cdot\|_{\mathcal{A}}$  satisfies all axioms of norm (see Chapter 2). The result we have just proved shows that  $u^h$  is the *best approximation* from  $S_E^h$  to the true solution  $u \in H_E^1(a, b)$  of our problem, when we measure the error of the approximation in the energy norm:

$$\|u - u^h\|_{\mathcal{A}} = \min_{v^h \in S_E^h} \|u - v^h\|_{\mathcal{A}}. \quad (14.24)$$

A particularly relevant question is how the error  $u - u^h$  depends on the spacing  $h$  of the subdivision of the computational domain  $[a, b]$ . We can obtain a bound on the error  $u - u^h$ , measured in the energy norm, by choosing a particular function  $v^h \in S_E^h$  in (14.24) whose closeness to  $u$  is easy to assess. For this purpose, we introduce the **finite element interpolant**  $\mathcal{I}^h u \in S_E^h$  of  $u \in H_E^1(a, b)$  by

$$\mathcal{I}^h u(x) = \psi(x) + \sum_{i=1}^{n-1} u(x_i) \varphi_i(x), \quad x \in [a, b].$$

Clearly,

$$\mathcal{I}^h u(x_j) = u(x_j), \quad j = 0, 1, \dots, n,$$

which justifies our use of the word *interpolant*.

We then deduce from (14.24) that

$$\|u - u^h\|_{\mathcal{A}} \leq \|u - \mathcal{I}^h u\|_{\mathcal{A}}; \quad (14.25)$$

hence, in order to quantify  $\|u - u^h\|_{\mathcal{A}}$ , we only need to estimate the size of  $\|u - \mathcal{I}^h u\|_{\mathcal{A}}$ . This leads us to the next theorem.

**Theorem 14.7** *Suppose that  $u \in H^2(a, b) \cap H_E^1(a, b)$  and let  $\mathcal{I}^h u$  be the finite element interpolant of  $u$  from  $S_E^h$  defined above; then, the following error bounds hold:*

$$\begin{aligned} \|u - \mathcal{I}^h u\|_{L^2(x_{i-1}, x_i)} &\leq \left(\frac{h_i}{\pi}\right)^2 \|u''\|_{L^2(x_{i-1}, x_i)}, \\ \|u' - (\mathcal{I}^h u)'\|_{L^2(x_{i-1}, x_i)} &\leq \frac{h_i}{\pi} \|u''\|_{L^2(x_{i-1}, x_i)}, \end{aligned}$$

for  $i = 1, 2, \dots, n$ , where  $h_i = x_i - x_{i-1}$ .

*Proof* Consider an element  $[x_{i-1}, x_i]$ ,  $1 \leq i \leq n$ , and define  $\zeta(x) = u(x) - \mathcal{I}^h u(x)$  for  $x \in [x_{i-1}, x_i]$ . Then,  $\zeta \in H^2(x_{i-1}, x_i)$  and  $\zeta(x_{i-1}) = \zeta(x_i) = 0$ . Therefore  $\zeta$  can be expanded into a convergent Fourier sine-series,

$$\zeta(x) = \sum_{k=1}^{\infty} a_k \sin \frac{k\pi(x - x_{i-1})}{h_i}, \quad x \in [x_{i-1}, x_i].$$

Here, convergence is to be understood in the norm  $\|\cdot\|_{L^2(x_{i-1}, x_i)}$ . Hence,

$$\begin{aligned} \int_{x_{i-1}}^{x_i} [\zeta(x)]^2 dx &= \int_{x_{i-1}}^{x_i} \zeta(x) \zeta(x) dx \\ &= \sum_{k, \ell=1}^{\infty} a_k a_{\ell} \int_{x_{i-1}}^{x_i} \sin \frac{k\pi(x - x_{i-1})}{h_i} \sin \frac{\ell\pi(x - x_{i-1})}{h_i} dx \\ &= h_i \sum_{k, \ell=1}^{\infty} a_k a_{\ell} \int_0^1 \sin k\pi t \sin \ell\pi t dt \\ &= \frac{h_i}{2} \sum_{k, \ell=1}^{\infty} a_k a_{\ell} \delta_{k\ell} \\ &= \frac{h_i}{2} \sum_{k=1}^{\infty} |a_k|^2, \end{aligned}$$

where  $\delta_{k\ell}$  is the Kronecker delta. Differentiating the Fourier sine series of  $\zeta$  twice, we find that the Fourier coefficients of  $\zeta'$  are  $(k\pi/h_i)a_k$ , while those of  $\zeta''$  are  $-(k\pi/h_i)^2 a_k$ . Thus, proceeding in the same way as above,

$$\begin{aligned} \int_{x_{i-1}}^{x_i} [\zeta'(x)]^2 dx &= \frac{h_i}{2} \sum_{k=1}^{\infty} \left( \frac{k\pi}{h_i} \right)^2 |a_k|^2, \\ \int_{x_{i-1}}^{x_i} [\zeta''(x)]^2 dx &= \frac{h_i}{2} \sum_{k=1}^{\infty} \left( \frac{k\pi}{h_i} \right)^4 |a_k|^2. \end{aligned}$$

Because  $k^4 \geq k^2 \geq 1$ , it follows that

$$\begin{aligned} \int_{x_{i-1}}^{x_i} [\zeta(x)]^2 dx &\leq \left( \frac{h_i}{\pi} \right)^4 \int_{x_{i-1}}^{x_i} [\zeta''(x)]^2 dx, \\ \int_{x_{i-1}}^{x_i} [\zeta'(x)]^2 dx &\leq \left( \frac{h_i}{\pi} \right)^2 \int_{x_{i-1}}^{x_i} [\zeta''(x)]^2 dx. \end{aligned}$$

However,  $\zeta''(x) = u''(x) - (\mathcal{I}^h u)''(x) = u''(x)$  for  $x \in (x_{i-1}, x_i)$ , and hence the desired bounds on the interpolation error.  $\square$



Now, substituting the bounds from Theorem 14.7 into the definition of the norm  $\|u - \mathcal{I}^h u\|_{\mathcal{A}}$ , we arrive at the following estimate of the interpolation error in the energy norm.

**Corollary 14.1** *Suppose that  $u \in H^2(a, b) \cap H_E^1(a, b)$ . Then,*

$$\|u - \mathcal{I}^h u\|_{\mathcal{A}}^2 \leq \sum_{i=1}^n \left\{ \left( \frac{h_i}{\pi} \right)^2 P_i + \left( \frac{h_i}{\pi} \right)^4 R_i \right\} \|u''\|_{L^2(x_{i-1}, x_i)}^2,$$

where  $P_i = \max_{x \in [x_{i-1}, x_i]} p(x)$  and  $R_i = \max_{x \in [x_{i-1}, x_i]} r(x)$ .

*Proof* Let us observe that

$$\begin{aligned} \|v\|_{\mathcal{A}}^2 &= \mathcal{A}(v, v) \\ &= \int_a^b [p(x)|v'(x)|^2 + r(x)|v(x)|^2] dx \\ &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} [p(x)|v'(x)|^2 + r(x)|v(x)|^2] dx \\ &\leq \sum_{i=1}^n \left\{ P_i \|v'\|_{L^2(x_{i-1}, x_i)}^2 + R_i \|v\|_{L^2(x_{i-1}, x_i)}^2 \right\}. \end{aligned}$$

On letting  $v = u - \mathcal{I}^h u$  and applying the preceding theorem on the right-hand side of the last inequality, with  $v'$  and  $v$  replaced by  $u' - (\mathcal{I}^h u)'$  and  $u - \mathcal{I}^h u$ , respectively, the result follows.  $\square$

Inserting this estimate into (14.25) leads to the desired bound on the error between the analytical solution  $u$  and its finite element approximation  $u^h$  in the energy norm.

**Corollary 14.2** *Suppose that  $u \in H^2(a, b) \cap H_E^1(a, b)$ . Then,*

$$\|u - u^h\|_{\mathcal{A}}^2 \leq \sum_{i=1}^n \left\{ \left( \frac{h_i}{\pi} \right)^2 P_i + \left( \frac{h_i}{\pi} \right)^4 R_i \right\} \|u''\|_{L^2(x_{i-1}, x_i)}^2,$$

where  $P_i = \max_{x \in [x_{i-1}, x_i]} p(x)$  and  $R_i = \max_{x \in [x_{i-1}, x_i]} r(x)$ . Further,

$$\|u - u^h\|_{\mathcal{A}} \leq \frac{h}{\pi} \left\{ P + \left( \frac{h}{\pi} \right)^2 R \right\}^{1/2} \|u''\|_{L^2(a, b)}, \quad (14.26)$$

where  $P = \max_{x \in [a, b]} p(x)$ ,  $R = \max_{x \in [a, b]} r(x)$ , and  $h = \max_{1 \leq i \leq n} h_i$ .

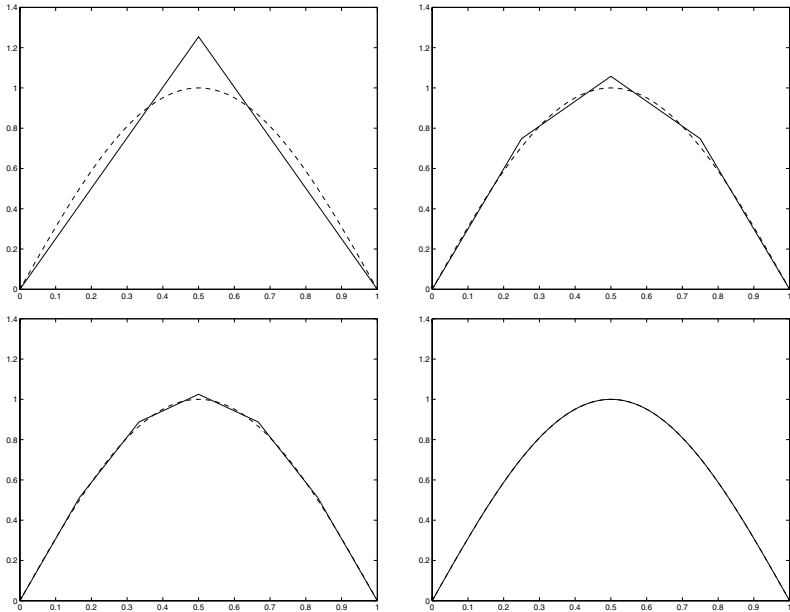


Fig. 14.3. Graph of the finite element approximation  $u^h$  to the analytical solution  $u$  of the boundary value problem (14.27) on a uniform subdivision of  $[0, 1]$  of spacing  $h = 1/n$ , with  $n = 2$  (top left),  $n = 4$  (top right),  $n = 6$  (bottom left), and  $n = 100$  (bottom right). In each of the four subfigures, the dashed curve is the graph of the analytical solution  $u(x) = \sin(\pi x)$ . In the last figure the approximation error is so small that  $u$  and  $u^h$  are indistinguishable.

In order to illustrate the performance of the finite element method, we consider the following example:

$$-u'' + r(x)u = f(x), \quad x \in (0, 1), \quad u(0) = 0, \quad u(1) = 0. \quad (14.27)$$

If  $r(x) \equiv 1$  and  $f(x) = (1 + \pi^2)\sin(\pi x)$ , the unique solution to this problem is  $u(x) = \sin(\pi x)$ . Let us pretend that we do not know the analytical solution  $u$ , and solve the boundary value problem numerically, using the finite element method on a subdivision of  $[0, 1]$  of uniform spacing  $h = 1/n$ , for various values of  $n$ . The integrals  $\langle f, \varphi_i \rangle$  involved in the definition of  $b_i$  in (14.17) have been approximated, on each of the elements  $[x_{i-1}, x_i]$ ,  $1 \leq i \leq n$ , by means of the trapezium rule. The resulting approximations  $u^h$ , for  $n = 2, 4, 6, 100$ , are shown in Figure 14.3.

We see from Figure 14.3 that, as the spacing  $h$  of the subdivision is reduced, the finite element solution  $u^h$  approximates the analytical solution  $u(x) = \sin(\pi x)$  with increasing accuracy. Indeed, the results corresponding to  $n = 2$  and  $n = 4$  in Figure 14.3 indicate that as the number of intervals in the subdivision is doubled (*i.e.*,  $h$  is halved), the maximum error between  $u(x)$  and  $u^h(x)$  is reduced by a factor of about 4. This reduction in the error cannot be explained by Corollary 14.2 which merely implies that halving  $h$  should lead to a reduction in  $\|u - u^h\|_A$  by a factor no less than 2. If you would like to learn more about the source of the observed enhancement of accuracy, consult Exercise 5 at the end of the chapter.

### 14.5 A *a posteriori* error analysis by duality

The bound on the error between the analytical solution  $u$  and its finite element approximation  $u^h$  formulated in Corollary 14.2 shows that, in the limit of  $h \rightarrow 0$ , the error  $\|u - u^h\|_A$  will tend to zero as  $\mathcal{O}(h)$ . This is a useful result from the theoretical point of view: it reassures us that the unknown analytical solution may be approximated arbitrarily well by making  $h$  sufficiently small. On the other hand, asymptotic error bounds of this kind are not particularly helpful for the purpose of precisely quantifying the size of the error between  $u$  and  $u^h$  for a given, *fixed*, mesh size  $h > 0$ : as  $u$  is unknown, it is difficult to tell just how large the right-hand side of (14.26) really is.

The aim of the present section is, therefore, to derive a computable bound on the error, and to demonstrate how such a bound may be implemented into an adaptive mesh-refinement algorithm, capable of reducing the error  $u - u^h$  below a certain prescribed tolerance in an automated manner, without human intervention. The approach is based on seeking a bound on  $u - u^h$  in terms of the computed solution  $u^h$  rather than in terms of norms of the unknown analytical solution  $u$ . A bound on the error in terms of  $u^h$  is referred to as an ***a posteriori*** error bound, due to the fact that it becomes *computable* only *after* the numerical solution  $u^h$  has been obtained.

In order to illuminate the key ideas while avoiding technical difficulties, we shall consider the two-point boundary value problem

$$-(p(x)u')' + q(x)u' + r(x)u = f(x), \quad a < x < b, \quad (14.28)$$

$$u(a) = A, \quad u(b) = B, \quad (14.29)$$

where  $p, q \in C^1[a, b]$ ,  $r \in C[a, b]$  and  $f \in L^2(a, b)$ . We shall assume, as

at the beginning of the chapter, that  $p(x) \geq c_0 > 0$ ,  $x \in [a, b]$ ; however, instead of supposing that  $r(x) \geq 0$ , we shall now demand that

$$r(x) - \frac{1}{2}q'(x) \geq c_1, \quad x \in [a, b], \quad (14.30)$$

where  $c_1$  is assumed to be a positive constant.<sup>1</sup>

Letting

$$\mathcal{A}(w, v) = \int_a^b [p(x)w'(x)v'(x) + q(x)w'(x)v(x) + r(x)w(x)v(x)]dx,$$

the weak formulation of (14.28), (14.29) is as follows:

$$\text{find } u \in H_E^1(a, b) \text{ such that } \mathcal{A}(u, v) = \langle f, v \rangle \quad \forall v \in H_0^1(a, b). \quad (14.31)$$

Here, the bilinear functional  $\mathcal{A}(\cdot, \cdot)$  is *not* symmetric, unless  $q(x) \equiv 0$ : indeed,  $\mathcal{A}(w, v) = \mathcal{A}(v, w)$  for all  $v, w \in H^1(a, b)$  if, and only if,  $q \equiv 0$ . Hence, in general, the boundary value problem (14.28), (14.29) cannot be assigned a Ritz principle. On the other hand, the Galerkin principle (weak formulation) (14.31) is perfectly meaningful for any choice of  $q$ .

The Galerkin finite element approximation of (14.31) is constructed by introducing a (possibly nonuniform) subdivision of the interval  $[a, b]$  defined by the points

$$a = x_0 < x_1 < \cdots < x_{n-1} < x_n = b$$

and considering the finite element space  $S_E^h \subset H_E^1(a, b)$  consisting of all continuous piecewise linear functions  $v^h$  on this subdivision that satisfy the boundary conditions  $v^h(a) = A$  and  $v^h(b) = B$ . The Galerkin finite element approximation of the boundary value problem is

$$\text{find } u^h \in S_E^h \text{ such that } \mathcal{A}(u^h, v^h) = \langle f, v^h \rangle \quad \forall v^h \in S_0^h. \quad (14.32)$$

We let  $h_i = x_i - x_{i-1}$ ,  $i = 1, \dots, n$ , and put  $h = \max_i h_i$ .

We wish to derive an *a posteriori* bound on the error in the  $\|\cdot\|_{L^2(a,b)}$  norm; that is, our aim is to quantify the size of  $\|u - u^h\|_{L^2(a,b)}$  in terms of the mesh parameter  $h$  and the computed solution  $u^h$  (rather than in terms of the analytical solution  $u$  as was the case in the *a priori* error analysis developed in the previous section). For this purpose, we

<sup>1</sup> At the expense of slight technical complications in the subsequent discussion, the requirement that  $c_1 > 0$  can be relaxed to  $c_1 > -\lambda_1$ , where  $\lambda_1$  is the smallest (positive) eigenvalue for the Sturm–Liouville eigenvalue problem  $-(p(x)w')' = \lambda w$  for  $x \in (a, b)$ ,  $w(a) = 0$ ,  $w(b) = 0$ .

consider the auxiliary boundary value problem

$$-(p(x)z')' - (q(x)z)' + r(x)z = (u - u^h)(x), \quad a < x < b, \quad (14.33)$$

$$z(a) = 0, \quad z(b) = 0, \quad (14.34)$$

called the **dual problem** (or adjoint problem).

We begin our error analysis by noting that the definition of the dual problem and straightforward integration by parts yield (recalling that  $(u - u^h)(a) = 0$ ,  $(u - u^h)(b) = 0$ )

$$\begin{aligned} \|u - u^h\|_{L^2(a,b)}^2 &= \langle u - u^h, u - u^h \rangle \\ &= \langle u - u^h, -(pz')' - (qz)' + rz \rangle \\ &= \mathcal{A}(u - u^h, z). \end{aligned}$$

On the other hand, (14.31) and (14.32) imply the Galerkin orthogonality property

$$\mathcal{A}(u - u^h, z^h) = 0 \quad \forall z^h \in S_0^h.$$

In particular, by choosing

$$z^h = \mathcal{I}^h z \in S_0^h,$$

the continuous piecewise linear interpolant of the function  $z \in H_0^1(a, b)$ , associated with the subdivision  $a = x_0 < x_1 < \cdots < x_{n-1} < x_n = b$ , we have that

$$\mathcal{A}(u - u^h, \mathcal{I}^h z) = 0.$$

Thus,

$$\begin{aligned} \|u - u^h\|_{L^2(a,b)}^2 &= \mathcal{A}(u - u^h, z - \mathcal{I}^h z) \\ &= \mathcal{A}(u, z - \mathcal{I}^h z) - \mathcal{A}(u^h, z - \mathcal{I}^h z) \\ &= \langle f, z - \mathcal{I}^h z \rangle - \mathcal{A}(u^h, z - \mathcal{I}^h z), \end{aligned} \quad (14.35)$$

where the last transition follows from (14.31) with  $v = z - \mathcal{I}^h z$ .

We observe that the right-hand side no longer involves the unknown analytical solution  $u$ . Furthermore,

$$\begin{aligned} \mathcal{A}(u^h, z - \mathcal{I}^h z) &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} p(x) (u^h)'(x) (z - \mathcal{I}^h z)'(x) \, dx \\ &\quad + \sum_{i=1}^n \int_{x_{i-1}}^{x_i} q(x) (u^h)'(x) (z - \mathcal{I}^h z)(x) \, dx \\ &\quad + \sum_{i=1}^n \int_{x_{i-1}}^{x_i} r(x) u^h(x) (z - \mathcal{I}^h z)(x) \, dx. \end{aligned}$$

Integrating by parts in each of the  $n$  integrals in the first sum on the right-hand side, noting that  $(z - \mathcal{I}^h z)(x_i) = 0$ ,  $i = 0, \dots, n$ , we deduce that

$$\begin{aligned} \mathcal{A}(u^h, z - \mathcal{I}^h z) \\ = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} [-(p(x)(u^h)')' + q(x)(u^h)' + r(x)u^h] (z - \mathcal{I}^h z)(x) \, dx. \end{aligned}$$

Furthermore,

$$\langle f, z - \mathcal{I}^h z \rangle = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} f(x) (z - \mathcal{I}^h z)(x) \, dx.$$

Substituting these two identities into (14.35), we deduce that

$$\|u - u^h\|_{L^2(a,b)}^2 = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} R(u^h)(x) (z - \mathcal{I}^h z)(x) \, dx, \quad (14.36)$$

where, for  $1 \leq i \leq n$ , and  $x \in (x_{i-1}, x_i)$ ,

$$R(u^h)(x) = f(x) - [-(p(x)(u_h)')' + q(x)(u^h)' + r(x)u^h].$$

The function  $R(u^h)$  is called **the finite element residual**; it measures the extent to which  $u^h$  fails to satisfy the differential equation

$$-(p(x)u')' + q(x)u' + r(x)u = f(x)$$

on the union of the intervals  $(x_{i-1}, x_i)$ ,  $i = 1, \dots, n$ . Now, applying the Cauchy–Schwarz inequality on the right-hand side of (14.36) yields

$$\|u - u^h\|_{L^2(a,b)}^2 \leq \sum_{i=1}^n \|R(u^h)\|_{L^2(x_{i-1}, x_i)} \|z - \mathcal{I}^h z\|_{L^2(x_{i-1}, x_i)}.$$

Recalling from Theorem 14.7 that

$$\|z - \mathcal{I}^h z\|_{L^2(x_{i-1}, x_i)} \leq \left(\frac{h_i}{\pi}\right)^2 \|z''\|_{L^2(x_{i-1}, x_i)}, \quad i = 1, 2, \dots, n,$$

we deduce that

$$\|u - u^h\|_{L^2(a,b)}^2 \leq \frac{1}{\pi^2} \sum_{i=1}^n h_i^2 \|R(u^h)\|_{L^2(x_{i-1}, x_i)} \|z''\|_{L^2(x_{i-1}, x_i)},$$

and consequently, using the Cauchy–Schwarz inequality for finite sums,

$$\sum_{i=1}^n a_i b_i \leq \left(\sum_{i=1}^n |a_i|^2\right)^{1/2} \left(\sum_{i=1}^n |b_i|^2\right)^{1/2}$$

with

$$a_i = h_i^2 \|R(u^h)\|_{L^2(x_{i-1}, x_i)} \quad \text{and} \quad b_i = \|z''\|_{L^2(x_{i-1}, x_i)},$$

we find that

$$\|u - u^h\|_{L^2(a,b)}^2 \leq \frac{1}{\pi^2} \left( \sum_{i=1}^n h_i^4 \|R(u^h)\|_{L^2(x_{i-1}, x_i)}^2 \right)^{1/2} \|z''\|_{L^2(0,1)}. \quad (14.37)$$

The rest of the discussion is aimed at eliminating  $\|z''\|_{L^2(a,b)}$  from the right-hand side of (14.37). The desired *a posteriori* bound on the error  $\|u - u^h\|_{L^2(a,b)}$  in terms of  $R(u^h)$  will then follow.

**Lemma 14.1** *Suppose that  $z$  is the solution of the dual problem (14.33), (14.34). Then, there exists a positive constant  $K$ , dependent only on  $p$ ,  $q$  and  $r$ , such that*

$$\|z''\|_{L^2(a,b)} \leq K \|u - u^h\|_{L^2(a,b)}.$$

*Proof* As

$$-pz'' - p'z' - qz' - q'z + rz = u - u^h,$$

it follows that

$$pz'' = u^h - u - (p' + q)z' + (r - q')z,$$

and therefore, recalling that  $p(x) \geq c_0 > 0$  for  $x \in [a, b]$ ,

$$\begin{aligned} c_0 \|z''\|_{L^2(a,b)} &\leq \|u - u^h\|_{L^2(a,b)} + \|p' + q\|_{\infty} \|z'\|_{L^2(a,b)} \\ &\quad + \|r - q'\|_{\infty} \|z\|_{L^2(a,b)}, \end{aligned} \quad (14.38)$$

where we used the notation  $\|w\|_{\infty} = \max_{x \in [a,b]} |w(x)|$ .

We shall show that both  $\|z'\|_{L^2(a,b)}$  and  $\|z\|_{L^2(a,b)}$  can be bounded in terms of  $\|u - u^h\|_{L^2(a,b)}$  and then, using (14.38), we shall deduce that the same is true of  $\|z''\|_{L^2(a,b)}$ . Let us observe that, by (14.33),

$$\langle -(pz')' - (qz)' + rz, z \rangle = \langle u - u^h, z \rangle. \quad (14.39)$$

Integrating by parts in the terms involving  $p$  and  $q$  and noting that  $z(0) = 0$  and  $z(1) = 0$  yields

$$\begin{aligned} \langle -(pz')' - (qz)' + rz, z \rangle &= \langle pz', z' \rangle + \langle qz, z' \rangle + \langle rz, z \rangle \\ &\geq c_0 \|z'\|_{L^2(a,b)}^2 + \frac{1}{2} \int_a^b q(x) [z^2(x)]' dx + \int_a^b r(x) [z(x)]^2 dx. \end{aligned}$$

Integrating by parts, again, in the second term on the right gives

$$\begin{aligned} \langle -(pz')' - (qz)' + rz, z \rangle &\geq c_0 \|z'\|_{L^2(a,b)}^2 - \frac{1}{2} \int_a^b q'(x) [z^2(x)] dx \\ &\quad + \int_a^b r(x) [z(x)]^2 dx. \end{aligned}$$

Hence, from (14.39),

$$c_0 \|z'\|_{L^2(a,b)}^2 + \int_a^b \left( r(x) - \frac{1}{2} q'(x) \right) [z(x)]^2 dx \leq \langle u - u^h, z \rangle,$$

and thereby, noting (14.30) and using the Cauchy–Schwarz inequality on the right-hand side,

$$\begin{aligned} \min\{c_0, c_1\} \left( \|z'\|_{L^2(a,b)}^2 + \|z\|_{L^2(a,b)}^2 \right) &\leq \langle u - u^h, z \rangle \\ &\leq \|u - u^h\|_{L^2(a,b)} \|z\|_{L^2(a,b)}. \end{aligned} \quad (14.40)$$

Therefore, also

$$\min\{c_0, c_1\} \|z\|_{H^1(a,b)}^2 \leq \|u - u^h\|_{L^2(a,b)} \|z\|_{H^1(a,b)},$$

which means that

$$\begin{aligned} \left( \|z'\|_{L^2(a,b)}^2 + \|z\|_{L^2(a,b)}^2 \right)^{1/2} &= \|z\|_{H^1(a,b)} \\ &\leq \frac{1}{\min\{c_0, c_1\}} \|u - u^h\|_{L^2(a,b)}. \end{aligned} \quad (14.41)$$

Now we substitute (14.41) into (14.38) to deduce that

$$\|z''\|_{L^2(a,b)} \leq K \|u - u^h\|_{L^2(a,b)}, \quad (14.42)$$

where

$$K = \frac{1}{c_0} \left( 1 + \frac{1}{\min\{c_0, c_1\}} (\|p' + q\|_\infty^2 + \|r - q'\|_\infty^2)^{1/2} \right).$$

□

It is important to observe here that  $K$  involves only known quantities: the coefficients in the differential equation under consideration. Therefore  $K$  can be computed, or at least bounded above, without difficulties. On inserting (14.42) into (14.37), we arrive at our final result, the computable *a posteriori* error bound,

$$\|u - u^h\|_{L^2(a,b)} \leq K_0 \left( \sum_{i=1}^n h_i^4 \|R(u^h)\|_{L^2(x_{i-1}, x_i)}^2 \right)^{1/2}, \quad (14.43)$$



where  $K_0 = K/\pi^2$ .

Next we shall describe the construction of an adaptive mesh refinement algorithm based on the *a posteriori* error bound (14.43).

Suppose that **TOL** is a prescribed tolerance and that our aim is to compute a finite element approximation  $u^h$  to the unknown solution  $u$  so that

$$\|u - u^h\|_{L^2(a,b)} \leq \text{TOL}. \quad (14.44)$$

We shall use the *a posteriori* error bound (14.43) to achieve this goal by systematically refining the subdivision, and computing a succession of numerical solutions  $u^h$  on this sequence of subdivisions, until the inequality

$$K_0 \left( \sum_{i=1}^n h_i^4 \|R(u^h)\|_{L^2(x_{i-1}, x_i)}^2 \right)^{1/2} \leq \text{TOL} \quad (14.45)$$

is satisfied. Clearly, if  $u^h$  satisfies (14.45), then, by virtue of (14.43), it also satisfies (14.44).

In order for the inequality (14.45) to hold it is sufficient to ensure that, on each interval  $[x_{i-1}, x_i]$ ,  $i = 1, 2, \dots, n$ , we have

$$h_i^4 \|R(u^h)\|_{L^2(x_{i-1}, x_i)}^2 \leq \frac{1}{n} \left( \frac{\text{TOL}}{K_0} \right)^2. \quad (14.46)$$

Thus, a sufficient condition for (14.44) is that (14.46) holds for all  $i = 1, 2, \dots, n$ .

The mesh adaptation algorithm, therefore, proceeds as follows:

**Step 1.** Choose an initial subdivision

$$\mathcal{T}_0: \quad a = x_0^{(0)} < x_1^{(0)} < \dots < x_{n_0-1}^{(0)} < x_{n_0}^{(0)} = b$$

of the interval  $[a, b]$ , with  $h_i^{(0)} = x_i^{(0)} - x_{i-1}^{(0)}$ , for  $i = 1, 2, \dots, n_0$ ; let  $h^{(0)} = \max_i h_i^{(0)}$ , and consider the associated finite element space  $S_E^{h^{(0)}}$  (of dimension  $n_0 - 1$ );

**Step 2.** Compute the corresponding solution  $u^{h^{(0)}} \in S_E^{h^{(0)}}$ ;

**Step 3.** Given a computed solution  $u^{h^{(m)}} \in S_E^{h^{(m)}}$  for some  $m \geq 0$ , defined on a subdivision  $\mathcal{T}_m$ , **STOP** if

$$K_0 \left( \sum_{i=1}^{n_m} \left( h_i^{(m)} \right)^4 \|R(u^{h^{(m)}})\|_{L^2(x_{i-1}^{(m)}, x_i^{(m)})}^2 \right)^{1/2} \leq \text{TOL}; \quad (14.47)$$

**Step 4.** If not, then halve those elements  $[x_{i-1}^{(m)}, x_i^{(m)}]$  in  $\mathcal{T}_m$ , with  $i$  in the set  $\{1, 2, \dots, n_m\}$ , for which

$$\left(h_i^{(m)}\right)^4 \|R(u^{h^{(m)}})\|_{L^2(x_{i-1}^{(m)}, x_i^{(m)})}^2 > \frac{1}{n_m} \left(\frac{\text{TOL}}{K_0}\right)^2, \quad (14.48)$$

denote by  $\mathcal{T}_{m+1}$  the resulting subdivision of  $[a, b]$  with  $n_{m+1}$  elements  $[x_{i-1}^{(m+1)}, x_i^{(m+1)}]$  of respective lengths

$$h_i^{(m+1)} = x_i^{(m+1)} - x_{i-1}^{(m+1)}, \quad i = 1, \dots, n_{m+1},$$

and consider the associated finite element space  $S_E^{h^{(m+1)}}$  of dimension  $n_{m+1} - 1$ ;

**Step 5.** Compute the finite element approximation  $u^{h^{(m+1)}} \in S_E^{h^{(m+1)}}$ , increase  $m$  by 1 and return to **Step 3**.

The inequality (14.47) is called the **stopping criterion** for the mesh adaptation algorithm, and (14.48) is referred to as the **refinement criterion**. According to the *a posteriori* error bound (14.43), when the adaptive algorithm terminates, the error  $\|u - u^h\|_{L^2(a,b)}$  is guaranteed not to exceed the prescribed tolerance TOL.

We conclude the body of this chapter with a numerical experiment which illustrates the performance of the adaptive algorithm.

**Example 14.1** *Let us consider the second-order ordinary differential equation*

$$-(p(x)u')' + q(x)u' + r(x)u = f(x), \quad x \in (0, 1), \quad (14.49)$$

*subject to the boundary conditions*

$$u(0) = 0, \quad u(1) = 0. \quad (14.50)$$

*Suppose, for example, that*

$$p(x) \equiv 1, \quad q(x) \equiv 20, \quad r(x) \equiv 10 \quad \text{and} \quad f(x) \equiv 1.$$

In this case, the analytical solution,  $u$ , can be expressed in closed form:

$$u(x) = C_1 e^{\lambda_1 x} + C_2 e^{\lambda_2 x} + \frac{1}{10},$$

where  $\lambda_1$  and  $\lambda_2$  are the two roots of the characteristic polynomial of the differential equation,  $-\lambda^2 + 20\lambda + 10 = 0$ , *i.e.*,

$$\lambda_1 = 10 + \sqrt{110}, \quad \lambda_2 = 10 - \sqrt{110},$$

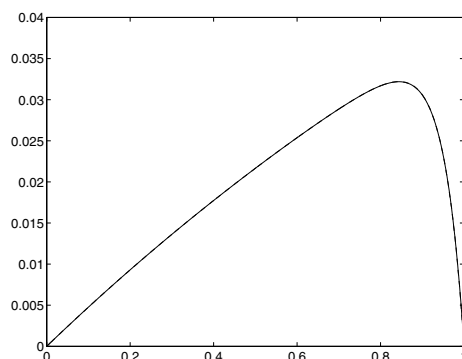


Fig. 14.4. Analytical solution of the boundary value problem (14.49), (14.50), with  $p(x) \equiv 1$ ,  $q(x) \equiv 20$ ,  $r(x) \equiv 10$  and  $f(x) \equiv 1$ .

and  $C_1$  and  $C_2$  are constants chosen so as to ensure that  $u(0) = 0$  and  $u(1) = 0$ ; hence,

$$C_1 = \frac{e^{\lambda_2} - 1}{10(e^{\lambda_1} - e^{\lambda_2})}, \quad C_2 = \frac{1 - e^{\lambda_1}}{10(e^{\lambda_1} - e^{\lambda_2})}.$$

The function  $u$  is shown in Figure 14.4.

Now, let us imagine for a moment that  $u$  is unknown, and let us compute a numerical approximation  $u^h$  to  $u$ , using the adaptive finite element algorithm described above, so that  $\|u - u^h\|_{L^2(0,1)} \leq \text{TOL}$ , where  $\text{TOL} = 10^{-4}$ . The computation begins on a coarse subdivision of the interval  $[0, 1]$  containing only 10 elements. This is then successively refined using the refinement criterion (14.48) until the stopping criterion (14.47) is satisfied; the resulting subdivisions are shown in Figure 14.5. In this example, the constant  $K_0$  appearing in (14.43) and (14.45)–(14.48) is  $(1 + \sqrt{500})/\pi^2 (\approx 2.367)$ .

Since we are in the fortunate (but highly idealised) position that, in addition to the numerical solution  $u^h$ , the analytical solution  $u$  is also available, we can assess the sharpness of our *a posteriori* error bound (14.43) by comparing the error  $\|u - u^h\|_{L^2(0,1)}$  appearing on the left-hand side of (14.43) with the computable *a posteriori* error bound on the right-hand side of (14.43). Figure 14.6 shows that the *a posteriori* bound consistently overestimates the error  $\|u - u^h\|_{L^2(0,1)}$  by about two orders of magnitude. By comparing the slopes of the two curves in Figure 14.6, we also see that the error and the *a posteriori* error bound decay at approximately the same rate as the number of mesh points increases in the course of mesh adaptation.  $\diamond$

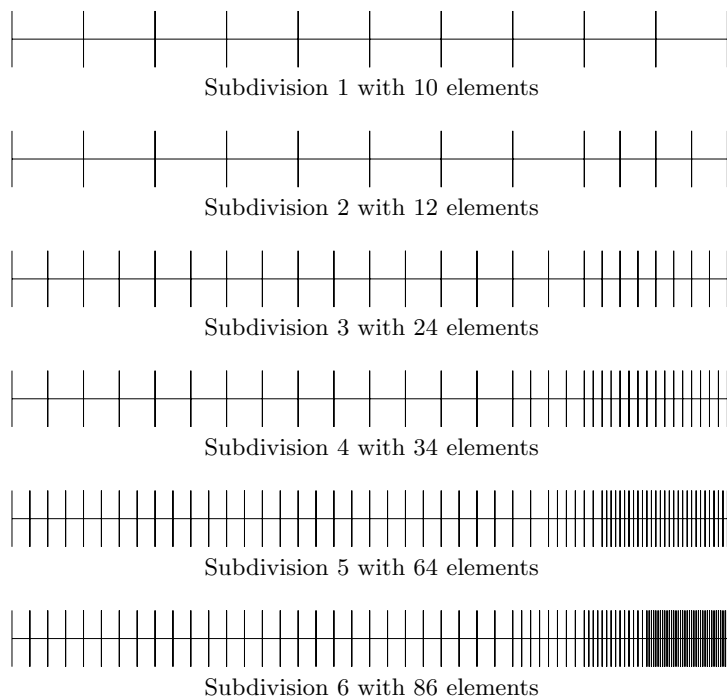


Fig. 14.5. Sequence of subdivisions of the interval  $[0, 1]$  designed by the adaptive algorithm with  $\text{TOL} = 10^{-4}$ .

## 14.6 Notes

For further details concerning the mathematical theory and the implementation of the finite element method we refer to the following books.

- ◆ D. BRAESS, *Finite Elements*, Cambridge University Press, Cambridge, 2001.
- ◆ S. BRENNER AND L.R. SCOTT, *The Mathematical Theory of Finite Element Methods*, Second Edition, Springer, New York, 2002.
- ◆ C. JOHNSON, *Numerical Solution of Partial Differential Equations by the Finite Element Method*, Cambridge University Press, Cambridge, 1996.

For recent results on the theory of *a posteriori* error estimation for finite element approximations of differential equations, based on duality arguments, the interested reader may wish to consult the following review articles.

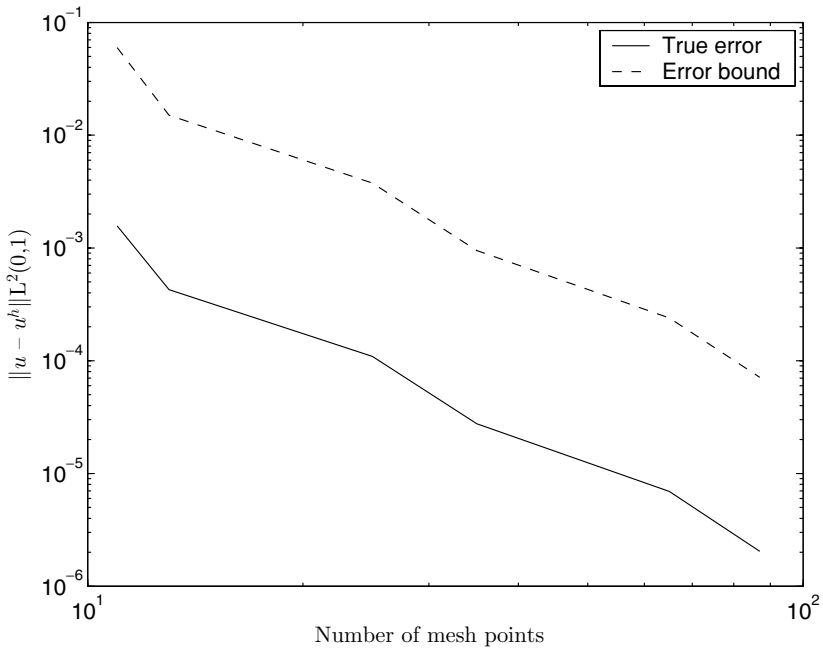


Fig. 14.6. Comparison of the true error  $\|u - u^h\|_{L^2(0,1)}$  (solid curve) with the *a posteriori* error bound delivered by the adaptive algorithm (dashed curve) with  $\text{TOL} = 10^{-4}$ .

- ◆ K. ERIKSON, D. ESTEP, P. HANSBO, AND C. JOHNSON, Introduction to adaptive methods for differential equations, in *Acta Numerica* **4** (A. Iserles, ed.), Cambridge University Press, Cambridge, 105–158, 1995.
- ◆ R. BECKER AND R. RANNACHER, An optimal control approach to a-posteriori error estimation in finite element methods, in *Acta Numerica* **10** (A. Iserles, ed.), Cambridge University Press, Cambridge, 1–102, 2001.
- ◆ M.B. GILES AND E. SÜLI, Adjoint methods for PDEs: superconvergence and adaptivity by duality, in *Acta Numerica* **11** (A. Iserles, ed.), Cambridge University Press, Cambridge, 145–236, 2002.

A detailed and general survey of the subject of *a posteriori* error estimation can be found in

- ◆ M. AINSWORTH AND J.T. ODEN, *A posteriori Error Estimation in Finite Element Analysis*, John Wiley & Sons, New York, 2000.

In this chapter we were concerned with the *a priori* error analysis of the piecewise linear finite element method in the energy norm, and its *a posteriori* error analysis in the  $L^2$  norm. Using similar techniques, one can establish an *a priori* error bound in the  $L^2$  norm and an *a posteriori* error bound in the energy norm. For extensions of the theory considered here to higher-order piecewise polynomial finite element approximations and generalisations to partial differential equations, the reader is referred to the books listed above.

### Exercises

14.1 Given that  $(a, b)$  is an open interval of the real line, let

$$H_{E_0}^1(a, b) = \{v \in H^1(a, b): v(a) = 0\}.$$

(i) By writing

$$v(x) = \int_a^x v'(\xi) d\xi,$$

for  $v \in H_{E_0}^1(a, b)$  and  $x \in [a, b]$ , show the following (**Poincaré–Friedrichs**) inequality:

$$\|v\|_{L^2(a, b)}^2 \leq \frac{1}{2}(b-a)^2 \|v'\|_{L^2(a, b)}^2 \quad \forall v \in H_{E_0}^1(a, b).$$

(ii) By writing

$$[v(x)]^2 = \int_a^x \frac{d}{d\xi} [v(\xi)]^2 d\xi = 2 \int_a^x v(\xi) v'(\xi) d\xi$$

for  $v \in H_{E_0}^1(a, b)$  and  $x \in [a, b]$ , show the following (**Agmon's**) inequality:

$$\max_{x \in [a, b]} |v(x)|^2 \leq 2 \|v\|_{L^2(a, b)} \|v'\|_{L^2(a, b)} \quad \forall v \in H_{E_0}^1(a, b).$$

14.2 Given that  $f \in L^2(0, 1)$ , state the weak formulation of each of the following boundary value problems on the interval  $(0, 1)$ :

- (a)  $-u'' + u = f(x)$ ,  $u(0) = 0$ ,  $u(1) = 0$ ;
- (b)  $-u'' + u = f(x)$ ,  $u(0) = 0$ ,  $u'(1) = 1$ ;
- (c)  $-u'' + u = f(x)$ ,  $u(0) = 0$ ,  $u(1) + u'(1) = 2$ .

In each case, show that there exists at most one weak solution.

14.3 Give a proof of Theorem 14.4.

14.4 Prove Corollary 14.2.

14.5 Consider the boundary value problem

$$-p_0 u'' + r_0 u = f(x), \quad u(0) = 0, \quad u(1) = 0,$$

on the interval  $[0, 1]$ , where  $p_0$  and  $r_0$  are positive constants and  $f \in C^4[0, 1]$ . Using equally spaced points

$$x_i = ih, \quad i = 0, 1, \dots, n, \quad \text{with } h = 1/n, \quad n \geq 2,$$

and the standard piecewise linear finite element basis functions (hat functions)  $\varphi_i$ ,  $i = 1, 2, \dots, n-1$ , show that the finite element equations for  $u_i = u^h(x_i)$  become

$$-p_0(u_{i-1} - 2u_i + u_{i+1})/h^2 + r_0(u_{i-1} + 4u_i + u_{i+1})/6 = \frac{1}{h} \langle f, \varphi_i \rangle$$

for  $i = 1, 2, \dots, n-1$ , with  $u_0 = 0$  and  $u_n = 0$ . By expanding in Taylor series, show that

$$\frac{1}{h} \langle f, \varphi_i \rangle = f(x_i) + \frac{1}{12} h^2 f''(x_i) + \mathcal{O}(h^4).$$

Interpreting this set of difference equations as a finite difference approximation to the boundary value problem, as in Chapter 13, show that the corresponding truncation error  $T_i$  satisfies

$$T_i = \frac{1}{12} h^2 r_0 u''(x_i) + \mathcal{O}(h^4), \quad i = 1, \dots, n-1,$$

and use the method of Exercise 13.2 to show that

$$\max_{0 \leq i \leq n} |u(x_i) - u^h(x_i)| \leq Mh^2,$$

where  $M$  is a positive constant.

14.6 In the notation of Exercise 5 suppose that all the integrals involved in the calculation are approximated by the trapezium rule. Show that the system of equations becomes identical to that obtained from the central difference approximation in Chapter 13, and deduce that

$$\max_{0 \leq i \leq n} |u(x_i) - u^h(x_i)| \leq Mh^2,$$

where  $M$  is a positive constant.

14.7 Consider the differential equation

$$-(p(x)u')' + r(x)u = f(x), \quad a < x < b,$$

with  $p$ ,  $r$  and  $f$  as at the beginning of the chapter, subject to the boundary conditions

$$-p(a)u'(a) + \alpha u(a) = A, \quad p(b)u'(b) + \beta u(b) = B,$$

where  $\alpha$  and  $\beta$  are positive real numbers, and  $A$  and  $B$  are real numbers. Show that the weak formulation of the boundary value problem is

find  $u \in H^1(a, b)$  such that  $\mathcal{A}(u, v) = \ell(v)$  for all  $v \in H^1(a, b)$ , where

$$\begin{aligned}\mathcal{A}(u, v) &= \int_a^b [p(x)u'(x)v'(x) + r(x)u(x)v(x)]dx \\ &\quad + \alpha u(a)v(a) + \beta u(b)v(b),\end{aligned}$$

and

$$\ell(v) = \langle f, v \rangle + Av(a) + Bv(b).$$

Construct a finite element approximation of the boundary value problem based on this weak formulation using piecewise linear finite element basis functions on the subdivision

$$a = x_0 < x_1 < \cdots < x_{n-1} < x_n = b$$

of the interval  $[a, b]$ . Show that the finite element method gives rise to a set of  $n + 1$  simultaneous linear equations with  $n + 1$  unknowns  $u_i = u^h(x_i)$ ,  $i = 0, 1, \dots, n$ . Show that this linear system has a unique solution.

Comment on the structure of the matrix  $M \in \mathbb{R}^{(n+1) \times (n+1)}$  of the linear system: (a) Is  $M$  symmetric? (b) Is  $M$  positive definite? (c) Is  $M$  tridiagonal?

- 14.8 Given that  $\alpha$  is a nonnegative real number, consider the differential equation

$$-u'' + u = f(x) \quad \text{for } x \in (0, 1),$$

subject to the boundary conditions

$$u(0) = 0, \quad \alpha u(1) + u'(1) = 0.$$

State the weak formulation of the problem. Using continuous piecewise linear basis functions on a uniform subdivision of  $[0, 1]$  into elements of size  $h = 1/n$ ,  $n \geq 2$ , write down the finite element approximation to this problem and show that this has a unique solution  $u^h$ . Expand  $u^h$  in terms of the standard piecewise linear finite element basis functions (hat functions)  $\varphi_i$ ,



$i = 1, 2, \dots, n$ , by writing

$$u^h(x) = \sum_{i=1}^n U_i \varphi_i(x)$$

to obtain a system of linear equations for the vector of unknowns  $(U_1, \dots, U_n)^T$ .

Suppose that  $\alpha = 0$ ,  $f(x) \equiv 1$  and  $h = 1/3$ . Solve the resulting system of linear equations and compare the corresponding numerical solution  $u^h(x)$  with the exact solution  $u(x)$  of the boundary value problem.

14.9 Consider the differential equation

$$-(p(x)u')' + r(x)u = f(x), \quad x \in (0, 1),$$

subject to the boundary conditions  $u(0) = 0$ ,  $u(1) = 0$ , where  $p(x) \geq c_0 > 0$ ,  $r(x) \geq 0$  for all  $x$  in the closed interval  $[0, 1]$ , with  $p \in C^1[0, 1]$ ,  $r \in C[0, 1]$  and  $f \in L^2(0, 1)$ . Given that  $u^h$  denotes the continuous piecewise linear finite element approximation to  $u$  on a uniform subdivision of  $[0, 1]$  into elements of size  $h = 1/n$ ,  $n \geq 2$ , show that

$$\|u - u^h\|_{H^1(0,1)} \leq C_1 h \|u''\|_{L^2(0,1)},$$

where  $C_1$  is a positive constant that you should specify. Show further that there exists a positive constant  $C$  such that

$$\|u - u^h\|_{H^1(0,1)} \leq Ch \|f\|_{L^2(0,1)}.$$

Calculate the right-hand sides in these inequalities in the case when

$$p(x) \equiv 1, \quad r(x) \equiv 0, \quad f(x) \equiv 1,$$

for  $x \in [0, 1]$ , and  $h = 10^{-3}$ .

14.10 Consider the two-point boundary value problem

$$-u'' + u = f(x), \quad x \in (0, 1), \quad u(0) = 0, \quad u(1) = 0,$$

with  $f \in C^2[0, 1]$ . State the piecewise linear finite element approximation to this problem on a nonuniform subdivision

$$0 = x_0 < x_1 < \dots < x_n = 1, \quad n \geq 2,$$

with  $h_i = x_i - x_{i-1}$ , assuming that, for a continuous piecewise

linear function  $v^h$ ,

$$\int_0^1 f(x)v^h(x)dx$$

has been approximated by applying the trapezium rule on each element  $[x_{i-1}, x_i]$ .

Verify that the following *a posteriori* bound holds for the error between  $u$  and its finite element approximation  $u^h$ :

$$\begin{aligned} \|u - u^h\|_{L^2(0,1)} &\leq K_0 \left( \sum_{i=1}^n h_i^4 \|R(u^h)\|_{L^2(x_{i-1}, x_i)}^2 \right)^{1/2} \\ &\quad + K_1 \max_{1 \leq i \leq n} h_i^2 \left( \max_{x \in [x_{i-1}, x_i]} |f''(x)|^2 + 4 \max_{x \in [x_{i-1}, x_i]} |f'(x)| \right)^{1/2}, \end{aligned}$$

where  $R(u^h) = f(x) - (-(u^h)''(x) + u^h(x))$  for  $x \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, n$ , and  $K_0, K_1$  are constants which you should specify.

How would you use this bound to compute  $u$  to within a specified tolerance TOL?