

Basic Concepts

The finite element method provides a formalism for generating discrete (finite) algorithms for approximating the solutions of differential equations. It should be thought of as a black box into which one puts the differential equation (boundary value problem) and out of which pops an algorithm for approximating the corresponding solutions. Such a task could conceivably be done automatically by a computer, but it necessitates an amount of mathematical skill that today still requires human involvement. The purpose of this book is to help people become adept at working the magic of this black box. The book does *not* focus on how to turn the resulting algorithms into computer codes, but this topic is being pursued by several groups. In particular, the FEniCS project (on the web at fenics.org) utilizes the mathematical structure of the finite element method to automate the generation of finite element codes.

In this chapter, we present a microcosm of a large fraction of the book, restricted to one-dimensional problems. We leave many loose ends, most of which will be tied up in the theory of Sobolev spaces to be presented in the subsequent chapter. These loose ends should provide motivation and guidance for the study of those spaces.

0.1 Weak Formulation of Boundary Value Problems

Consider the two-point boundary value problem

$$(0.1.1) \quad \begin{aligned} -\frac{d^2u}{dx^2} &= f \text{ in } (0, 1) \\ u(0) &= 0, \quad u'(1) = 0. \end{aligned}$$

If u is the solution and v is any (sufficiently regular) function such that $v(0) = 0$, then integration by parts yields

$$\begin{aligned}
 (0.1.2) \quad (f, v) &:= \int_0^1 f(x)v(x)dx = \int_0^1 -u''(x)v(x)dx \\
 &= \int_0^1 u'(x)v'(x)dx =: a(u, v).
 \end{aligned}$$

Let us define (formally, for the moment, since the notion of derivative to be used has not been made precise)

$$V = \{v \in L^2(0, 1): a(v, v) < \infty \text{ and } v(0) = 0\}.$$

Then we can say that the solution u to (0.1.1) is characterized by

$$(0.1.3) \quad u \in V \quad \text{such that} \quad a(u, v) = (f, v) \quad \forall v \in V,$$

which is called the *variational* or *weak* formulation of (0.1.1).

The relationship (0.1.3) is called “variational” because the function v is allowed to vary arbitrarily. It may seem somewhat unusual at first; later we will see that it has a natural interpretation in the setting of *Hilbert spaces*. (A Hilbert space is a vector space whose topology is defined using an inner-product.) One example of a Hilbert space is $L^2(0, 1)$ with inner-product (\cdot, \cdot) . Although it is by no means obvious, we will also see that the space V may be viewed as a Hilbert space with inner-product $a(\cdot, \cdot)$, which was defined in (0.1.2).

One critical question we have not yet dealt with is *what sort of derivative is to be used* in the definition of the bilinear form $a(\cdot, \cdot)$. Should this be the classical derivative

$$u'(x) = \lim_{h \rightarrow 0} \frac{u(x+h) - u(x)}{h} ?$$

Or should the “almost everywhere” definition valid for functions of bounded variation (BV) be used? We leave this point hanging for the moment and hope this sort of question motivates you to study the following chapter on Sobolev spaces. Of course, the central issue is that (0.1.3) still embodies the original problem (0.1.1). The following theorem verifies this under some simplifying assumptions.

(0.1.4) Theorem. *Suppose $f \in C^0([0, 1])$ and $u \in C^2([0, 1])$ satisfy (0.1.3). Then u solves (0.1.1).*

Proof. Let $v \in V \cap C^1([0, 1])$. Then integration by parts gives

$$(0.1.5) \quad (f, v) = a(u, v) = \int_0^1 (-u'')v dx + u'(1)v(1).$$

Thus, $(f - (-u''), v) = 0$ for all $v \in V \cap C^1([0, 1])$ such that $v(1) = 0$. Let $w = f + u'' \in C^0([0, 1])$. If $w \not\equiv 0$, then $w(x)$ is of one sign in some interval $[x_0, x_1] \subset [0, 1]$, with $x_0 < x_1$ (continuity). Choose $v(x) = (x - x_0)^2(x - x_1)^2$

in $[x_0, x_1]$ and $v \equiv 0$ outside $[x_0, x_1]$. But then $(w, v) \neq 0$, which is a contradiction. Thus, $-u'' = f$. Now apply (0.1.5) with $v(x) = x$ to find $u'(1) = 0$. Of course, $u \in V$ implies $u(0) = 0$, so u solves (0.1.1). \square

(0.1.6) Remark. The boundary condition $u(0) = 0$ is called *essential* as it appears in the variational formulation explicitly, i.e., in the definition of V . This type of boundary condition also frequently goes by the proper name “Dirichlet.” The boundary condition $u'(1) = 0$ is called *natural* because it is incorporated implicitly. This type of boundary condition is often referred to by the name “Neumann.” We summarize the different kinds of boundary conditions encountered so far, together with their various names in the following table:

Table 0.1. Naming conventions for two types of boundary conditions

Boundary Condition	Variational Name	Proper Name
$u(x) = 0$	essential	Dirichlet
$u'(x) = 0$	natural	Neumann

The assumptions $f \in C^0([0, 1])$ and $u \in C^2([0, 1])$ in the theorem allow (0.1.1) to be interpreted in the usual sense. However, we will see other ways in which to interpret (0.1.1), and indeed the theorem says that the formulation (0.1.3) is a way to interpret it that is valid with much less restrictive assumptions on f . For this reason, (0.1.3) is also called a *weak* formulation of (0.1.1).

0.2 Ritz-Galerkin Approximation

Let $S \subset V$ be any (finite dimensional) subspace. Let us consider (0.1.3) with V replaced by S , namely

$$(0.2.1) \quad u_S \in S \quad \text{such that} \quad a(u_S, v) = (f, v) \quad \forall v \in S.$$

It is remarkable that a discrete scheme for approximating (0.1.1) can be defined so easily. This is only one powerful aspect of the Ritz-Galerkin method. However, we first must see that (0.2.1) does indeed *define* an object. In the process we will indicate how (0.2.1) represents a (square, finite) system of equations for u_S . These will be done in the following theorem and its proof.

(0.2.2) Theorem. *Given $f \in L^2(0, 1)$, (0.2.1) has a unique solution.*

Proof. Let us write (0.2.1) in terms of a basis $\{\phi_i : 1 \leq i \leq n\}$ of S . Let $u_S = \sum_{j=1}^n U_j \phi_j$; let $K_{ij} = a(\phi_j, \phi_i)$, $F_i = (f, \phi_i)$ for $i, j = 1, \dots, n$. Set

$\mathbf{U} = (U_j)$, $\mathbf{K} = (K_{ij})$ and $\mathbf{F} = (F_i)$. Then (0.2.1) is equivalent to solving the (square) matrix equation

$$(0.2.3) \quad \mathbf{KU} = \mathbf{F}.$$

For a square system such as (0.2.3) we know that uniqueness is equivalent to existence, as this is a *finite dimensional* system. Nonuniqueness would imply that there is a nonzero \mathbf{V} such that $\mathbf{KV} = \mathbf{0}$. Write $v = \sum V_j \phi_j$ and note that the equivalence of (0.2.1) and (0.2.3) implies that $a(v, \phi_j) = 0$ for all j . Multiplying this by V_j and summing over j yields $0 = a(v, v) = \int_0^1 (v')^2(x) dx$, from which we conclude that $v' \equiv 0$. Thus, v is constant, and, since $v \in S \subset V$ implies $v(0) = 0$, we must have $v \equiv 0$. Since $\{\phi_i : 1 \leq i \leq n\}$ is a basis of S , this means that $\mathbf{V} = \mathbf{0}$. Thus, the solution to (0.2.3) must be unique (and hence must exist). Therefore, the solution u_S to (0.2.1) must also exist and be unique. \square

(0.2.4) Remark. Two subtle points are hidden in the “proof” of Theorem (0.2.2). Why is it that “thus v is constant”? And, moreover, why does $v \in V$ really imply $v(0) = 0$ (even though it is in the definition, i.e., why does the definition make sense)? The first question should worry those familiar with the Cantor function whose derivative is zero almost everywhere, but is certainly not constant (it also vanishes at the left of the interval in typical constructions). Thus, something about our definition of V must rule out such functions as members. V is an example of a *Sobolev* space, and we will see that such problems do not occur in these spaces. It is clear that functions such as the Cantor function should be ruled out (in a systematic way) as candidate solutions for differential equations since it would be a nontrivial solution to the o.d.e. $u' = 0$ with initial condition $u(0) = 0$.

(0.2.5) Remark. The matrix \mathbf{K} is often referred to as the *stiffness* matrix, a name coming from corresponding matrices in the context of structural problems. It is clearly symmetric, since the *energy* inner-product $a(\cdot, \cdot)$ is symmetric. It is also *positive definite*, since

$$\sum_{i,j=1}^n k_{ij} v_i v_j = a(v, v) \quad \text{where} \quad v = \sum_{j=1}^n v_j \phi_j.$$

Clearly, $a(v, v) \geq 0$ for all (v_j) and $a(v, v) = 0$ was already “shown” to imply $v \equiv 0$ in the proof of Theorem 0.2.3.

0.3 Error Estimates

Let us begin by observing the fundamental *orthogonality* relation between u and u_S . Subtracting (0.2.1) from (0.1.3) implies

$$(0.3.1) \quad a(u - u_S, w) = 0 \quad \forall w \in S.$$

Equation (0.3.1) and its subsequent variations are the key to the success of all Ritz-Galerkin/finite-element methods. Now define

$$\|v\|_E = \sqrt{a(v, v)}$$

for all $v \in V$, the energy *norm*. A critical relationship between the energy norm and inner-product is Schwarz' inequality:

$$(0.3.2) \quad |a(v, w)| \leq \|v\|_E \|w\|_E \quad \forall v, w \in V.$$

This inequality is a cornerstone of Hilbert space theory and will be discussed at length in Sect. 2.1. Then, for any $v \in S$,

$$\begin{aligned} \|u - u_S\|_E^2 &= a(u - u_S, u - u_S) \\ &= a(u - u_S, u - v) + a(u - u_S, v - u_S) \\ &= a(u - u_S, u - v) \quad (\text{from 0.3.1 with } w = v - u_S) \\ &\leq \|u - u_S\|_E \|u - v\|_E \quad (\text{from 0.3.2}). \end{aligned}$$

If $\|u - u_S\|_E \neq 0$, we can divide by it to obtain $\|u - u_S\|_E \leq \|u - v\|_E$, for any $v \in S$. If $\|u - u_S\|_E = 0$, this inequality is trivial. Taking the infimum over $v \in S$ yields

$$\|u - u_S\|_E \leq \inf\{\|u - v\|_E : v \in S\}.$$

Since $u_S \in S$, we have

$$\inf\{\|u - v\|_E : v \in S\} \leq \|u - u_S\|_E.$$

Therefore,

$$\|u - u_S\|_E = \inf\{\|u - v\|_E : v \in S\}.$$

Moreover, there is an element (u_S) for which the infimum is attained, and we indicate this by replacing ‘‘infimum’’ with ‘‘minimum.’’ Thus, we have proved the following.

(0.3.3) Theorem. $\|u - u_S\|_E = \min\{\|u - v\|_E : v \in S\}.$

This is the basic error estimate for the Ritz-Galerkin method, and it says that the error is optimal in the energy norm. We will use this later to derive more concrete estimates for the error based on constructing approximations to u in S for particular choices of S . Now we consider the error in another norm.

Define $\|v\| = (v, v)^{\frac{1}{2}} = (\int_0^1 v(x)^2 dx)^{\frac{1}{2}}$, the $L^2(0, 1)$ -norm. We wish to consider the size of the error $u - u_S$ in this norm. You might guess that the $L^2(0, 1)$ -norm is weaker than the energy norm, as the latter is the $L^2(0, 1)$ -norm of the *derivative* (this is the case, on V , although it is not completely obvious and makes use of the essential boundary condition incorporated in V). Thus, the error in the $L^2(0, 1)$ -norm will be at least comparable with the error measured in the energy norm. In fact, we will find it is considerably smaller.

To estimate $\|u - u_S\|$, we use what is known as a “duality” argument. Let w be the solution of

$$-w'' = u - u_S \quad \text{on } [0, 1] \quad \text{with} \quad w(0) = w'(1) = 0.$$

Integrating by parts, we find

$$\begin{aligned} \|u - u_S\|^2 &= (u - u_S, u - u_S) \\ &= (u - u_S, -w'') \\ &= a(u - u_S, w) \quad (\text{since } (u - u_S)(0) = w'(1) = 0) \\ &= a(u - u_S, w - v) \quad (\text{from 0.3.1}) \end{aligned}$$

for all $v \in S$. Thus, Schwarz' inequality (0.3.2) implies that

$$\begin{aligned} \|u - u_S\| &\leq \|u - u_S\|_E \|w - v\|_E / \|u - u_S\| \\ &= \|u - u_S\|_E \|w - v\|_E / \|w''\|. \end{aligned}$$

We may now take the infimum over $v \in S$ to get

$$\|u - u_S\| \leq \|u - u_S\|_E \inf_{v \in S} \|w - v\|_E / \|w''\|.$$

Thus, we see that the L^2 -norm of the error can be much smaller than the energy norm, provided that w can be approximated well by some function in S . It is reasonable to assume that we can take $v \in S$ close to w , which we formalize in the following *approximation assumption*:

$$(0.3.4) \quad \inf_{v \in S} \|w - v\|_E \leq \epsilon \|w''\|.$$

Of course, we envisage that this holds with ϵ being a small number. Applying (0.3.4) yields

$$\|u - u_S\| \leq \epsilon \|u - u_S\|_E,$$

and applying (0.3.4) again, with w replaced by u , and using Theorem 0.3.3 gives

$$\|u - u_S\|_E \leq \epsilon \|u''\|.$$

Combining these estimates, and recalling (0.1.1), yields

(0.3.5) Theorem. *Assumption (0.3.4) implies that*

$$\|u - u_S\| \leq \epsilon \|u - u_S\|_E \leq \epsilon^2 \|u''\| = \epsilon^2 \|f\|.$$

The point of course is that $\|u - u_S\|_E$ is of order ϵ whereas $\|u - u_S\|$ is of order ϵ^2 . We now consider a family of spaces S for which ϵ may be made arbitrarily small.

0.4 Piecewise Polynomial Spaces – The Finite Element Method

Let $0 = x_0 < x_1 < \dots < x_n = 1$ be a partition of $[0, 1]$, and let S be the linear space of functions v such that

- i) $v \in C^0([0, 1])$
- ii) $v|_{[x_{i-1}, x_i]}$ is a linear polynomial, $i = 1, \dots, n$, and
- iii) $v(0) = 0$.

We will see later that $S \subset V$. For each $i = 1, \dots, n$ define ϕ_i by the requirement that $\phi_i(x_j) = \delta_{ij}$ = the Kronecker delta, as shown in Fig. 0.1.

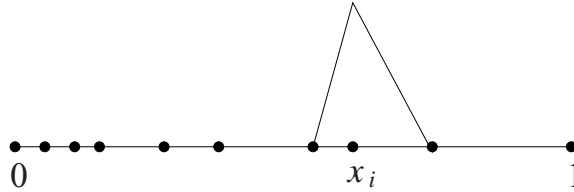


Fig. 0.1. piecewise linear basis function ϕ_i

(0.4.1) **Lemma.** $\{\phi_i : 1 \leq i \leq n\}$ is a basis for S .

(0.4.2) *Remark.* $\{\phi_i\}$ is called a **nodal** basis for S , and $\{v(x_i)\}$ are the **nodal values** of a function v . (The points $\{x_i\}$ are called the **nodes**.)

Proof. The set $\{\phi_i\}$ is linearly independent since $\sum_{i=1}^n c_i \phi_i(x_j) = 0$ implies $c_j = 0$. To see that it spans S , consider the following:

(0.4.3) **Definition.** Given $v \in C^0([0, 1])$, the **interpolant** $v_I \in S$ of v is determined by $v_I := \sum_{i=1}^n v(x_i) \phi_i$.

Clearly, the set $\{\phi_i\}$ spans S if the following is true.

(0.4.4) **Lemma.** $v \in S \Rightarrow v = v_I$.

Proof. $v - v_I$ is linear on each $[x_{i-1}, x_i]$ and zero at the endpoints, hence must be identically zero. \square

We will now prove the following approximation theorem for the interpolant.

(0.4.5) **Theorem.** Let $h = \max_{1 \leq i \leq n} (x_i - x_{i-1})$. Then

$$\|u - u_I\|_E \leq Ch \|u''\|$$

for all $u \in V$, where C is independent of h and u .

Proof. Recalling the definitions of the two norms, it is clearly sufficient to prove the estimate piecewise, i.e., that

$$\int_{x_{j-1}}^{x_j} (u - u_I)'(x)^2 dx \leq c(x_j - x_{j-1})^2 \int_{x_{j-1}}^{x_j} u''(x)^2 dx$$

as the stated result follows by summing over j , with $C = \sqrt{c}$. Let $e = u - u_I$ denote the error; since u_I is a linear polynomial on the interval $[x_{j-1}, x_j]$, the above is equivalent to

$$\int_{x_{j-1}}^{x_j} e'(x)^2 dx \leq c(x_j - x_{j-1})^2 \int_{x_{j-1}}^{x_j} e''(x)^2 dx.$$

Changing variables by an affine mapping of the interval $[x_{j-1}, x_j]$ to the interval $[0, 1]$, we see that this is equivalent to showing

$$\int_0^1 \tilde{e}'(\tilde{x})^2 d\tilde{x} \leq c \int_0^1 \tilde{e}''(\tilde{x})^2 d\tilde{x},$$

where $x = x_{j-1} + \tilde{x}(x_j - x_{j-1})$ and

$$\tilde{e}(\tilde{x}) = e(x_{j-1} + \tilde{x}(x_j - x_{j-1})).$$

Note that we have arrived at an equivalent estimate that does not involve the mesh size at all. The technique of reducing a mesh-length dependent estimate to a mesh-independent one in this way is called a *homogeneity argument* (or scaling argument) and will be used frequently in Chapter 4 and thereafter.

The verification of the latter estimate is a simple calculus exercise. Let $w = \tilde{e}$ to simplify the notation, and write x for \tilde{x} . Note that w vanishes at both ends of the interval (the interpolation error is zero at all nodes). By Rolle's Theorem, $w'(\xi) = 0$ for some ξ satisfying $0 < \xi < 1$. Thus,

$$w'(y) = \int_{\xi}^y w''(x) dx.$$

By Schwarz' inequality,

$$\begin{aligned} |w'(y)| &= \left| \int_{\xi}^y w''(x) dx \right| \\ &= \left| \int_{\xi}^y 1 \cdot w''(x) dx \right| \\ (0.4.6) \quad &\leq \left| \int_{\xi}^y 1 dx \right|^{1/2} \cdot \left| \int_{\xi}^y w''(x)^2 dx \right|^{1/2} \\ &= |y - \xi|^{1/2} \left| \int_{\xi}^y w''(x)^2 dx \right|^{1/2} \\ &\leq |y - \xi|^{1/2} \left(\int_0^1 w''(x)^2 dx \right)^{1/2}. \end{aligned}$$

Squaring and integrating with respect to y completes the verification, with

$$c = \sup_{0 < \xi < 1} \int_0^1 |y - \xi| dy = \frac{1}{2}. \quad \square$$

(0.4.7) Corollary. $\|u - u_S\| + Ch \|u - u_S\|_E \leq 2(Ch)^2 \|u''\|$.

Proof. Theorem 0.4.5 implies that the approximation assumption (0.3.4) holds with $\epsilon = Ch$. \square

(0.4.8) Remark. The interpolant defines a linear operator $\mathcal{I}: C^0([0, 1]) \rightarrow S$ where $\mathcal{I}v = v_I$. Lemma 0.4.4 says that \mathcal{I} is a *projection* (i.e., $\mathcal{I}^2 = \mathcal{I}$). The estimate (0.4.6) for w' in the proof of (0.4.5) is an example of Sobolev's inequality, in which the pointwise values of a function can be estimated in terms of integrated quantities involving its derivatives. Estimates of this type will be considered at length in Chapter 1.

0.5 Relationship to Difference Methods

The stiffness matrix \mathbf{K} as defined in (0.2.3), using the basis $\{\phi_i\}$ described above, can be interpreted as a difference operator. Let $h_i = x_i - x_{i-1}$. Then the matrix entries $K_{ij} = a(\phi_i, \phi_j)$ can be easily calculated to be

$$(0.5.1) \quad K_{ii} = h_i^{-1} + h_{i+1}^{-1}, K_{i,i+1} = K_{i+1,i} = -h_{i+1}^{-1} \quad (i = 1, \dots, n-1)$$

and $K_{nn} = h_n^{-1}$ with the rest of the entries of \mathbf{K} being zero. Similarly, the entries of \mathbf{F} can be approximated if f is sufficiently smooth:

$$(0.5.2) \quad (f, \phi_i) = \frac{1}{2}(h_i + h_{i+1})(f(x_i) + \mathcal{O}(h))$$

where $h = \max h_i$. (This follows easily from Taylor's Theorem since the integral of ϕ_i is $(h_i + h_{i+1})/2$. Note that the error is *not* $\mathcal{O}(h^2)$ unless $1 - (h_i/h_{i+1}) = \mathcal{O}(h)$.) Thus, the i -th equation of $\mathbf{KU} = \mathbf{F}$ (for $1 \leq i \leq n-1$) can be written as

$$(0.5.3) \quad \frac{-2}{h_i + h_{i+1}} \left[\frac{U_{i+1} - U_i}{h_{i+1}} - \frac{U_i - U_{i-1}}{h_i} \right] = \frac{2(f, \phi_i)}{h_i + h_{i+1}} = f(x_i) + \mathcal{O}(h).$$

The difference operator on the left side of this equation can also be seen to be an $\mathcal{O}(h)$ accurate approximation to the differential operator $-d^2/dx^2$ (and *not* $\mathcal{O}(h^2)$ accurate in the usual sense unless $1 - h_i/h_{i+1} = \mathcal{O}(h)$.) For a uniform mesh, the equations reduce to the familiar difference equations

$$(0.5.4) \quad -\frac{U_{i+1} - 2U_i + U_{i-1}}{h^2} = f(x_i) + \mathcal{O}(h^2)$$

which are well known to be second-order accurate. However, for a general mesh (e.g., $h_i = h$ for i even and $h_i = h/2$ for i odd), we know from Corollary 0.4.7 that the answer is still second-order accurate (in $L^2(0, 1)$ at least, but it will also be proved to be so in the maximum norm in Sect. 0.7), even though the difference equations are formally only consistent to first order. This phenomenon has been studied in detail by Spijker (Spijker 1971), and related work has recently been done by (Kreiss, et.al. 1986). See exercises 0.x.11 through 0.x.15 for more details.

We will take this opportunity to philosophize about some powerful characteristics of the finite element formalism for generating discrete schemes for approximating the solutions to differential equations. Being based on the variational formulation of boundary value problems, it is quite systematic, handling different boundary conditions with ease; one simply replaces infinite dimensional spaces with finite dimensional subspaces. What results, as in (0.5.3), is the same as a finite difference equation, in keeping with the *dictum* that different numerical methods are usually more similar than they are distinct. However, we were able to derive very quickly the convergence properties of the finite element method. Finally, the notation for the discrete scheme is quite compact in the finite element formulation. This could be utilized to make coding the algorithm much more efficient if only the appropriate computer language and compiler were available. This latter characteristic of the finite element method is one that has not yet been exploited extensively, but an initial attempt has been made in the system `fec` (Bagheri, Scott & Zhang 1992). (One could also argue that finite element practitioners have already taken advantage of this by developing their own “languages” through extensive software libraries of their own, but this applies equally well to the finite-difference practitioners.)

0.6 Computer Implementation of Finite Element Methods

One key to the success of the finite element method, as developed in engineering practice, was the systematic way that computer codes could be implemented. One important step in this process is the *assembly* of the inner-product $a(u, v)$ by summing its constituent parts over each sub-interval, or *element*, which are computed separately. This is facilitated through the use of a numbering scheme called the *global-to-local* index. This index, $i(e, j)$, relates the local node number, j , on a particular element, e , to its position in the global data structure. In our one-dimensional example with piecewise linear functions, this index is particularly simple: the “elements” are based

on the intervals $I_e := [x_{e-1}, x_e]$ where e is an integer in the range $1, \dots, n$ and

$$i(e, j) := e + j - 1 \quad \text{for } e = 1, \dots, n \quad \text{and } j = 0, 1.$$

That is, for each element there are two nodal parameters of interest, one corresponding to the left end of the interval ($j = 0$) and one at the right ($j = 1$). Their relationship is represented by the mapping $i(e, j)$.

We may write the interpolant of a continuous function for the space of all piecewise linear functions (no boundary conditions imposed) via

$$(0.6.1) \quad f_I := \sum_e \sum_{j=0}^1 f(x_{i(e,j)}) \phi_j^e$$

where $\{\phi_j^e : j = 0, 1\}$ denotes the set of basis functions for linear functions on the single interval $I_e = [x_{e-1}, x_e]$:

$$\phi_j^e(x) = \phi_j((x - x_{e-1})/(x_e - x_{e-1}))$$

where

$$\phi_0(x) := \begin{cases} 1 - x & x \in [0, 1] \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \phi_1(x) := \begin{cases} x & x \in [0, 1] \\ 0 & \text{otherwise.} \end{cases}$$

Note that we have related all of the “local” basis functions ϕ_j^e to a fixed set of basis functions on a “reference” element, $[0, 1]$, via an affine mapping of $[0, 1]$ to $[x_{e-1}, x_e]$. (By definition, the local basis functions, ϕ_j^e , are extended by zero outside the interval I_e .)

The expression (0.6.1) for the interpolant shows (cf. Lemma 0.4.4) that any piecewise linear function f (no boundary conditions imposed) can be written in the form

$$(0.6.2) \quad f := \sum_e \sum_{j=0}^1 f_{i(e,j)} \phi_j^e$$

where $f_i = f(x_i)$ for all i . In particular, the cardinality of the image of the index mapping $i(e, j)$ is the dimension of the space of piecewise linear functions. Note that the expression (0.6.2) represents f incorrectly at the nodal points, but this has no effect on the evaluation of multilinear forms involving integrals of f .

The bilinear forms defined in (0.1.2) can be easily evaluated (assembled) using this representation as well. For example,

$$a(v, w) = \sum_e a_e(v, w)$$

where the “local” bilinear form is defined (and evaluated) via

$$\begin{aligned}
a_e(v, w) &:= \int_{I_e} v' w' dx \\
&= (x_e - x_{e-1})^{-1} \int_0^1 (\sum_j v_{i(e,j)} \phi_j)' (\sum_j w_{i(e,j)} \phi_j)' dx \\
&= (x_e - x_{e-1})^{-1} \begin{pmatrix} v_{i(e,0)} \\ v_{i(e,1)} \end{pmatrix}^t \mathbf{K} \begin{pmatrix} w_{i(e,0)} \\ w_{i(e,1)} \end{pmatrix}.
\end{aligned}$$

Here, the *local stiffness matrix*, \mathbf{K} , is given by

$$K_{i,j} := \int_0^1 \phi'_{i-1} \phi'_{j-1} dx \quad \text{for } i, j = 1, 2.$$

Note that we have identified the space of piecewise linear functions, v , with the vector space of values, (v_i) , at the nodes. The subspace, S , of piecewise linear functions that vanish at $x = 0$, defined in Sect. 0.4, can be identified with the subspace $\{(v_i) : v_0 = 0\}$. Including v_0 in the data structure (with a value of zero) makes the assembly of bilinear forms equally easy in the presence of boundary conditions.

0.7 Local Estimates

We wish to derive estimates for the error, $u - u_S$, in the pointwise sense. As in the case for the L^2 -norm, we begin by writing the error that we wish to bound in terms of the energy bilinear form applied to $u - u_S$ and some other function. In this case, this other function is the so-called Green's function for the problem (0.1.1), which in this case is simply

$$g_x(t) := \begin{cases} t & t < x \\ x & \text{otherwise} \end{cases}$$

where x is any point in $[0, 1]$. Integration by parts shows that

$$v(x) = a(v, g_x) \quad \forall v \in V$$

since g_x'' is identically zero on either side of x . Therefore,

$$\begin{aligned}
(u - u_S)(x) &= a(u - u_S, g_x) \\
&= a(u - u_S, g_x - v) \quad \forall v \in S.
\end{aligned}$$

One conclusion is that, if S is the space of piecewise linear functions defined on a partition $\{x_i : i = 1, \dots, n\}$ as in Sect. 0.4, then

$$(u - u_S)(x_i) = 0 \quad \forall i = 1, \dots, n$$

since $g_{x_i} \in S$ in this case. Thus, we conclude that $u_S = u_I$, and a variant of Theorem 0.4.5 yields

$$(0.7.1) \quad \|u - u_I\|_{\max} \leq Ch^2 \|u''\|_{\max}.$$

(Recall that $\|f\|_{\max} = \max_{0 \leq x \leq 1} |f(x)|$.) Combining the above estimates, we have proved the following.

(0.7.2) Theorem. *Let u_S be determined by (0.2.1) using the space of piecewise linear functions defined in Sect. 0.4. Then*

$$\|u - u_S\|_{\max} \leq Ch^2 \|u''\|_{\max}.$$

Local estimates for higher-dimensional problems are much more difficult to derive, but the use of the Green's function is similar. However, the local character of the singularity of the one-dimensional Green's function disappears, and the distributed nature of the higher-dimensional Green's function requires techniques that are illustrated in the next section.

0.8 Adaptive Approximation

In many cases, the solution to a differential equation is rapidly varying only in restricted regions. For such problems, it makes sense to adapt the mesh to match the variation in the solution. The difference in approximation power between a mesh chosen to solve general problems versus one adapted to a particular one can be substantial. We present a particularly simple approximation problem here to illustrate this effect. For more complex results, see (DeVore, Howard & Micchelli 1989).

Let us consider the problem of approximating functions of one variable whose derivatives are integrable. This is an even weaker condition than what we used in section 0.3, and we wish to consider approximation in a stronger norm, the maximum norm. We consider approximation by the space S_Δ of piecewise constant functions on a partition

$$(0.8.1) \quad \Delta = \{x_0, x_1, \dots, x_n : 0 = x_0 < x_1 < \dots < x_n = 1\}.$$

In this case, we will say that $\text{size}(\Delta) = n$. It is not hard to see that the best result of the form

$$(0.8.2) \quad \inf_{v \in S_\Delta} \|u - v\|_{\max} \leq Cn^{-p} \int_0^1 |u'(x)| dx$$

to hold for *all* u (with a fixed mesh) is to have $p = 0$. Indeed, whatever the mesh, we can let u go from zero at x_0 to one at x_1 (and stay at one the rest of the interval). This particular u has $\int_0^1 |u'(x)| dx = 1$ and yet

$$\inf_{v \in S_\Delta} \|u - v\|_{\max} = \frac{1}{2}.$$

Of course, writing u as the integral of u' (cf. (0.4.6)) allows us to prove (0.8.2) with $C = 1$ and $p = 0$, simply by taking $v \equiv 0$.

On the other hand, suppose that we fix a particular u and ask that (0.8.2) hold for *some* partition Δ as in (0.8.1). That is, what if we are allowed to choose Δ based on properties of u ? To be more precise, we are making the distinction between a statement that $\forall u \exists \Delta$ such that (0.8.2) holds versus our earlier statement that, given Δ , (0.8.2) holds $\forall u$.

To see that there is a better estimate possible with an adaptively chosen mesh, suppose that we have a u such that $\int_0^1 |u'(x)| dx = 1$. The function

$$(0.8.3) \quad \phi(x) = \int_0^x |u'(t)| dt$$

vanishes at $x = 0$ and is a non-decreasing function. Moreover, $\phi(1) = 1$, so there must be a points x_i where $\phi(x_i) = i/n$ and such that $x_i < x_{i+1}$ for all i . If by chance we have $x_n < 1$ in this process, we set $x_n = 1$. One property of this partition is that

$$(0.8.4) \quad \int_{x_{i-1}}^{x_i} |u'(t)| dt = \phi(x_i) - \phi(x_{i-1}) = \frac{1}{n}$$

for all $i = 1, \dots, n$.

To approximate u on the interval $[x_{i-1}, x_i]$ we use the constant $c_i = u(x_{i-1})$. Then for $x \in [x_{i-1}, x_i]$

$$(0.8.3) \quad |u(x) - c_i| = \left| \int_{x_{i-1}}^x u'(t) dt \right| \leq \int_{x_{i-1}}^{x_i} |u'(t)| dt = \frac{1}{n}$$

proving that (0.8.2) holds for all n with $p = 1$ and $C = 1$, at least when $\int_0^1 |u'(x)| dx = 1$. In the general case, simply divide everything in (0.8.2) by $\int_0^1 |u'(x)| dx$.

Again to get the quantifiers right, let us define the approximation quotient

$$(0.8.5) \quad Q(u, \Delta) = \inf_{v \in S_\Delta} \|u - v\|_{\max} / \int_0^1 |u'(x)| dx$$

for a given u such that $0 < \int_0^1 |u'(x)| dx < \infty$ and a given partition Δ . Then the first result we proved is that

$$(0.8.6) \quad \forall \Delta \exists u \text{ such that } Q(u, \Delta) \geq \frac{1}{2}$$

and yet in the second result we constructed a Δ to prove that

$$(0.8.7) \quad \forall u \exists \Delta \text{ with } \text{size}(\Delta) = n \text{ such that } Q(u, \Delta) \leq \frac{1}{n}.$$

These results indicate what a dramatic difference in approximation power there can be in using a fixed mesh versus a mesh adapted to a particular function.

0.9 Weighted Norm Estimates

Suppose $h(x)$ is a function that measures the local mesh size near the point x . In particular, we will assume that h is a piecewise linear function satisfying

$$h(x_j) = h_j + h_{j+1}$$

where $h_j = x_j - x_{j-1}$ (and we set $h_{n+1} = h_n$ and $h_0 = h_1$). Note that for all $j = 1, \dots, n$

$$(0.9.1) \quad h(x) \geq h_j \quad \forall x \in [x_{j-1}, x_j],$$

since this holds at each endpoint of the interval and h is linear between them.

We begin by deriving a basic estimate analogous to (0.4.5). From its proof and (0.9.1), we have

$$\begin{aligned} \|u - u_I\|_E^2 &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} (u - u_I)'(x)^2 dx \\ &\leq \frac{1}{2} \sum_{i=1}^n h_i^2 \int_{x_{i-1}}^{x_i} u''(x)^2 dx \\ &\leq \frac{1}{2} \sum_{i=1}^n \int_{x_{i-1}}^{x_i} h(x)^2 u''(x)^2 dx \\ &= \frac{1}{2} \|hu''\|^2. \end{aligned}$$

Therefore,

$$(0.9.2) \quad \|u - u_S\|_E \leq \frac{1}{\sqrt{2}} \|hu''\|.$$

We next derive an L^2 estimate analogous to the first inequality in Theorem 0.3.5. Choosing w as was done in the proof of that result, we find

$$\|u - u_S\|^2 = a(u - u_S, w)$$

where w solves the boundary value problem (0.1.1) with $u - u_S$ as right-hand-side. For simplicity of notation, let $e := u - u_S$. Using the orthogonality relation (0.3.1) and Schwarz' inequality, we find

$$\begin{aligned}
a(e, w) &= a(e, w - w_I) \\
&= \int_0^1 h(u - u_S)'(w - w_I)' / h \, dx \\
&\leq \left(\int_0^1 (h(u - u_S)')^2 \, dx \right)^{1/2} \left(\int_0^1 ((w - w_I)' / h)^2 \, dx \right)^{1/2}.
\end{aligned}$$

From the results of Sect. 0.4 we have

$$\begin{aligned}
\int_0^1 ((w - w_I)'(x)/h(x))^2 \, dx &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} ((w - w_I)'(x)/h(x))^2 \, dx \\
&\leq \sum_{i=1}^n h_i^{-2} \int_{x_{i-1}}^{x_i} (w - w_I)'(x)^2 \, dx \\
&\leq \sum_{i=1}^n \frac{1}{2} \int_{x_{i-1}}^{x_i} w''(x)^2 \, dx.
\end{aligned}$$

Combining the previous inequalities, we have

$$(0.9.3) \quad a(e, w) \leq \|he'\| \left(\sum_{i=1}^n \frac{1}{2} \int_{x_{i-1}}^{x_i} w''(x)^2 \, dx \right)^{1/2}.$$

Recalling that $-w'' = e$, we find

$$\begin{aligned}
\|e\|^2 &= a(e, w) \\
&\leq \frac{1}{\sqrt{2}} \|he'\| \left(\sum_{i=1}^n \int_{x_{i-1}}^{x_i} e(x)^2 \, dx \right)^{1/2} \\
&= \frac{1}{\sqrt{2}} \|he'\| \|e\|.
\end{aligned}$$

Dividing by $\|e\|$ and recalling that $e = u - u_S$, we have proved

$$(0.9.4) \quad \|u - u_S\| \leq \frac{1}{\sqrt{2}} \left(\int_0^1 (h(u - u_S)')^2 \, dx \right)^{1/2}.$$

This says that the L^2 error can always be estimated in terms of a weighted integral of the squared derivative error, where the weight is given by the mesh function (0.9.1). Now we proceed to estimate the “weighted energy” norm on the right hand side of (0.9.4).

Let us write $e := u - u_S$ for simplicity. Then first observe that

$$\int_0^1 (h(u - u_S)')^2 \, dx = \|he'\|^2 = a(e, h^2e) - \int_0^1 2hh'ee' \, dx$$

simply by expanding the expression $a(e, h^2e)$. We will begin to make the assumption that h' is small, i.e., that the mesh does not change rapidly (for a uniform mesh, $h' \equiv 0$). This will allow us to neglect the term

$$\int_0^1 2hh'ee' dx$$

in comparison with the other terms in the preceding equation. To do so, we will make frequent use of the *arithmetic-geometric mean inequality*, which is nothing more than the simple observation that, for any real numbers a and b ,

$$ab \leq \frac{1}{2} (a^2 + b^2)$$

(just observe that $0 \leq (a-b)^2 = -2ab + a^2 + b^2$). A slightly more complicated version of the inequality comes by writing

$$ab = (\epsilon a)(b/\epsilon) \leq \frac{1}{2} ((\epsilon a)^2 + (b/\epsilon)^2).$$

Writing δ in place of ϵ^2 , we find

$$(0.9.5) \quad ab \leq \frac{\delta}{2} a^2 + \frac{1}{2\delta} b^2$$

for any $\delta > 0$.

Let $M := \|h'\|_{\max}$. Then Schwarz' inequality and the arithmetic-geometric mean inequality imply

$$\begin{aligned} \left| \int_0^1 2hh'ee' dx \right| &\leq 2M \int_0^1 |hee'| dx \\ &\leq 2M \|he'\| \|e\| \\ &\leq M (\|he'\|^2 + \|e\|^2). \end{aligned}$$

Therefore,

$$\|he'\|^2 \leq a(e, h^2e) + M (\|he'\|^2 + \|e\|^2)$$

and hence,

$$(1 - M)\|he'\|^2 \leq a(e, h^2e) + M\|e\|^2.$$

We now estimate the term $a(e, h^2e)$. Let $w := h^2e$. From (0.9.3) and the arithmetic-geometric mean inequality,

$$\begin{aligned} a(e, h^2e) &= a(e, w) \\ &\leq \frac{1}{\sqrt{2}} \|he'\| \left(\sum_{i=1}^n \int_{x_{i-1}}^{x_i} (w'')^2 dx \right)^{1/2} \\ &\leq \frac{1-M}{2} \|he'\|^2 + \frac{1}{4(1-M)} \sum_{i=1}^n \int_{x_{i-1}}^{x_i} (w'')^2 dx \end{aligned}$$

which, combined with the previous estimate, implies that

$$\frac{1-M}{2} \|he'\|^2 \leq \frac{1}{4(1-M)} \sum_{i=1}^n \int_{x_{i-1}}^{x_i} (w'')^2 dx + M \|e\|^2.$$

Expanding, we have (on each interval (x_{i-1}, x_i) separately)

$$w'' = h^2 e'' + 4hh'e' + 2(h')^2 e$$

since $h'' \equiv 0$. Expanding again, and using the arithmetic-geometric mean inequality, we find

$$\begin{aligned} (w'')^2 &\leq h^4 (e'')^2 + 16M^2 (he')^2 + 4M^4 e^2 \\ &\quad + 8Mh^3 |e''||e'| + 4M^2 h^2 |e''||e| + 16M^3 h |e'| |e| \\ &\leq 7h^4 (e'')^2 + 28M^2 (he')^2 + 14M^4 e^2. \end{aligned}$$

Integrating, we find

$$\begin{aligned} \frac{1-M}{2} \|he'\|^2 &\leq \frac{1}{4(1-M)} \sum_{i=1}^n \int_{x_{i-1}}^{x_i} 7h^4 (e'')^2 dx \\ &\quad + \frac{7M^2}{1-M} \|he'\|^2 + \left(M + \frac{7M^4}{2(1-M)} \right) \|e\|^2, \end{aligned}$$

which implies

$$\begin{aligned} \left(\frac{1-M}{2} - \frac{7M^2}{1-M} \right) \|he'\|^2 &\leq \frac{1}{4(1-M)} \sum_{i=1}^n \int_{x_{i-1}}^{x_i} 7h^4 (e'')^2 dx \\ &\quad + \left(M + \frac{7M^4}{2(1-M)} \right) \|e\|^2. \end{aligned}$$

Letting $c_1 = \left(\frac{1-M}{2} - \frac{7M^2}{1-M} \right)^{-1}$ and recalling that $e'' = u''$, we have

$$\|he'\|^2 \leq \frac{7c_1}{4(1-M)} \|h^2 u''\|^2 + c_1 \left(M + \frac{7M^4}{2(1-M)} \right) \|e\|^2$$

provided that

$$M < \frac{1}{1 + \sqrt{14}}.$$

Combining with estimate (0.9.4), we find that

$$\left(2 - c_1 \left(M + \frac{7M^4}{2(1-M)} \right) \right) \|u - u_S\|^2 \leq \frac{7c_1}{4(1-M)} \|h^2 u''\|^2.$$

Finally, we assume that M is sufficiently small so that

$$2 - c_1 \left(M + \frac{7M^4}{2(1-M)} \right) > 0$$

(observe that $c_1 \rightarrow 2$ as $M \rightarrow 0$), and we conclude that

$$(0.9.6) \quad \|u - u_S\|^2 \leq C(M) \|h^2 u''\|^2$$

where $C(M) \rightarrow 7/4$ as $M \rightarrow 0$.

We summarize the above results in the following theorem.

(0.9.7) Theorem. *Without any restrictions on the mesh, we have*

$$\|u - u_S\|_E \leq \frac{1}{\sqrt{2}} \|hu''\|$$

and

$$\|u - u_S\| \leq \frac{1}{\sqrt{2}} \|h(u - u_S)'\|.$$

Provided that the mesh-size variation, $M := \|h'\|_{\max}$, is sufficiently small, there is a constant, C , depending on M but otherwise independent of the mesh, such that

$$\|u - u_S\| \leq C \|h^2 u''\|.$$

The condition that the derivative of h be small is easy to interpret. From its definition,

$$h'|_{(x_{i-1}, x_i)} = \frac{h_{i+1} - h_{i-1}}{h_i} = r_{i+1} - \frac{1}{r_i},$$

where r_i is the ratio of lengths of adjacent mesh intervals, $r_i = h_i/h_{i-1}$. Thus, $|h'|$ is small whenever these ratios are sufficiently close to one. However, this does not preclude strong mesh gradings, e.g., a geometrically graded mesh, $x_i = e^{\delta(i-n)}$ for δ sufficiently small.

0.x Exercises

0.x.1 Verify the expressions (0.5.1) for the “stiffness” matrix \mathbf{K} for piecewise linear functions. If f is piecewise linear, i.e.,

$$f(x) = \sum_{i=1}^n f_i \phi_i(x)$$

determine the matrix \mathbf{M} (called the “mass” matrix) such that

$$\mathbf{K}\mathbf{U} = \mathbf{M}\mathbf{F}.$$

0.x.2 Give weak formulations of modifications of the two-point boundary-value problem (0.1.1) where

- a) the o. d. e. is $-u'' + u = f$ instead of $-u'' = f$ and/or
 b) the boundary conditions are $u(0) = u(1) = 0$.

0.x.3 Explain what is wrong in both the variational setting and the classical setting for the problem

$$-u'' = f \quad \text{with } u'(0) = u'(1) = 0.$$

That is, explain in both contexts why this problem is not well-posed.

0.x.4 Show that piecewise *quadratics* have a nodal basis consisting of values at the nodes x_i together with the midpoints $\frac{1}{2}(x_i + x_{i+1})$. Calculate the stiffness matrix for these elements.

0.x.5 Verify (0.5.2).

0.x.6 Under the same assumptions as in Theorem 0.4.5, prove that

$$\|u - u_I\| \leq Ch^2 \|u''\|.$$

(Hint: use a homogeneity argument as in the proof of Theorem 0.4.5. Using the notation of that proof, show further that

$$\int_0^1 w(x)^2 dx \leq \tilde{c} \int_0^1 w'(x)^2 dx,$$

by utilizing the fact that $w(0) = 0$. How small can you make \tilde{c} if you use both $w(0) = 0$ and $w(1) = 0$?)

0.x.7 Using only Theorems 0.3.5 and 0.4.5, prove that

$$\inf_{v \in S} \|u - v\| \leq Ch^2 \|u''\|.$$

Exercise 0.x.6 also would imply this result independently. Compare the different constants, C , derived with the different approaches.

0.x.8 Prove that (0.1.1) has a solution $u \in C^2([0, 1])$ provided $f \in C^0([0, 1])$. (Hint: write

$$u(x) = \int_0^x \left(\int_s^1 f(t) dt \right) ds$$

and verify the equations.)

0.x.9 Let V denote the space, and $a(\cdot, \cdot)$ the bilinear form, defined in Sect. 0.1. Prove the following *coercivity* result

$$\|v\|^2 + \|v'\|^2 \leq Ca(v, v) \quad \forall v \in V.$$

Give a value for C . (Hint: see the hint in exercise 0.x.6. For simplicity, restrict the result to $v \in V \cap C^1(0, 1)$.)

0.x.10 Let V denote the space, and $a(\cdot, \cdot)$ the bilinear form, defined in Sect. 0.1. Prove the following version of Sobolev's inequality:

$$\|v\|_{\max}^2 \leq Ca(v, v) \quad \forall v \in V.$$

Give a value for C . (Hint: see the hint in exercise 0.x.6. For simplicity, restrict the result to $v \in V \cap C^1(0, 1)$.)

0.x.11 Consider the difference method represented by (0.5.3), namely

$$\frac{-2}{h_i + h_{i+1}} \left(\frac{U_{i+1} - U_i}{h_{i+1}} - \frac{U_i - U_{i-1}}{h_i} \right) = f(x_i).$$

Prove $\tilde{u}_S := \sum U_i \phi_i$ satisfies the following modification to (0.2.1):

$$a(\tilde{u}_S, v) = Q(fv) \quad \forall v \in S$$

where $a(\cdot, \cdot)$ is the bilinear form defined in Sect. 0.1, S consists of piecewise linears as defined in Sect. 0.4 and Q denotes the quadrature approximation based on the trapezoidal rule

$$Q(w) := \sum_{i=0}^n \frac{h_i + h_{i+1}}{2} w(x_i).$$

Here ϕ_i , x_i and h_i are as defined in Sect. 0.4; we further define $h_0 = h_{n+1} = 0$ for simplicity of notation.

0.x.12 Let Q be defined as in exercise 0.x.11. Prove that

$$\left| Q(w) - \int_0^1 w(x) dx \right| \leq Ch^2 \sum_{i=1}^n \int_{x_{i-1}}^{x_i} |w''(x)| dx.$$

(Hint: observe that the trapezoidal rule is exact for piecewise linears and refer to the hint in exercise 0.x.6.)

0.x.13 Let u_S solve (0.2.1) where S consists of piecewise linears as defined in Sect. 0.4 and let \tilde{u}_S be as in exercise 0.x.11. Prove that

$$|a(u_S - \tilde{u}_S, v)| \leq Ch^2 (\|f'\| + \|f''\|) (\|v\| + \|v'\|) \quad \forall v \in S.$$

(Hint: apply exercise 0.x.12 and Schwarz' inequality.)

0.x.14 Let u_S and \tilde{u}_S be as in exercise 0.x.13. Prove that

$$\|u_S - \tilde{u}_S\|_E \leq Ch^2 (\|f'\| + \|f''\|).$$

(Hint: apply exercise 0.x.13, pick $v = u_S - \tilde{u}_S$ and apply exercise 0.x.9.)

0.x.15 Let \tilde{u}_S be as in exercise 0.x.11 and let u solve (0.1.1). Prove that

$$\|u - \tilde{u}_S\|_{\max} \leq Ch^2 (\|f\|_{\max} + \|f'\| + \|f''\|).$$

(Hint: apply exercise 0.x.14 and Theorem 0.7.2.)

- 0.x.16 Give weak formulation of modifications of the two-point boundary-value problem (0.1.1) where the boundary conditions are $u(0) = 0$ and $u'(1) = \lambda$. (Hint: show that $a(u, v) = F(v)$ where F is the linear functional $F(v) = \lambda v(1)$.)