



Classical Numerical Analysis, Chapter 25

Abner J. Salgado and Steven M. Wise

asalgad1@utk.edu swise1@utk.edu
The University of Tennessee



Chapter 25

Finite Element Methods for Elliptic Problems



The Finite Element Method Factsheet

- 1 Its formulation is based on integral — read weak — formulations of the problems to be discretized. As we have seen, weak formulations allow for more general problem data.
- 2 The idea of decomposing the domain into a finite number of pieces (the elements) and treating the properties and behavior of each piece as known and fixed not only has roots in the physical origins of many problems but also allows us to easily treat general geometries, unstructured meshes, and local mesh refinements without great complication.
- 3 The approximate solution to our problem is a piecewise polynomial function, i.e., a polynomial on each of the elements. Having an actual function instead of, say, a grid one is sometimes desirable, as this function can be easily evaluated, differentiated, etc., without any special considerations.
- 4 As opposed to a finite difference methodology, in finite elements one is not discretizing the partial differential operators, but rather searching for an approximation of the solution in a subspace of the solution space. For this reason, many of the properties of the discrete problem are automatically inherited from the continuous one.

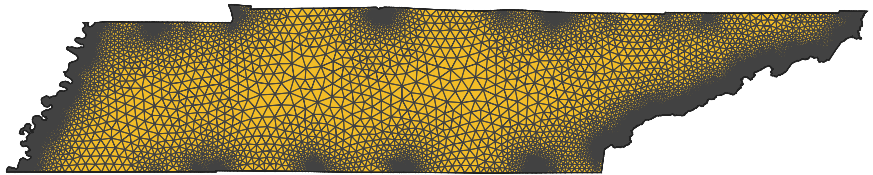


Figure: A sophisticated triangulation of the state of Tennessee (one of the 50 constituent states of the United States of America).



The Galerkin Method



Basic Assumptions

Let us assume that \mathcal{H} is a Hilbert space and $\mathcal{A}: \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ satisfies the following properties:

- ❶ Bilinear: In other words, linear in each argument.
- ❷ Bounded: There is a constant $M > 0$ such that, for every $v_1, v_2 \in \mathcal{H}$, we have

$$|\mathcal{A}(v_1, v_2)| \leq M \|v_1\|_{\mathcal{H}} \|v_2\|_{\mathcal{H}}.$$

- ③ Coercive: There is a constant $\alpha > 0$ such that, for every nonzero $v \in \mathcal{H}$,

$$\alpha \|v\|_{\mathcal{H}}^2 \leq \mathcal{A}(v, v).$$

Assume that $F \in \mathcal{H}'$, in other words, $F : \mathcal{H} \rightarrow \mathbb{R}$ is linear there is a constant, $C > 0$, such that

$$|F(v)| \leq C \|v\|_{\mathcal{H}}, \quad \forall v \in \mathcal{H}.$$



The Galerkin Approximation

In this chapter, we consider the following general problem: Find $u \in \mathcal{H}$ such that

$$\mathcal{A}(u, v) = F(v), \quad \forall v \in \mathcal{H}. \quad (1)$$

According to the Lax-Milgram Theorem, based on the assumptions about \mathcal{A} and F , it has a unique solution in \mathcal{H} .

We seek to approximate solutions to this problem.

Definition (Galerkin approximation)

Let $n \in \mathbb{N}$ and $\mathcal{H}_n \subseteq \mathcal{H}$, with $\dim \mathcal{H}_n = n$. We say that $u_n \in \mathcal{H}_n$ is a **Galerkin approximation** to u , solution of (1), if and only if

$$\mathcal{A}(u_n, v_n) = F(v_n), \quad \forall v_n \in \mathcal{H}_n.$$



Ritz Approximation

In Chapter 23 we showed that, if \mathcal{A} is symmetric, we can define the energy

$$E(v) = \frac{1}{2} \mathcal{A}(v, v) - F(v). \quad (2)$$

The Lax-Milgram Theorem (Chapter 23) showed that $u \in \mathcal{H}$ solves (1) if and only if it minimizes E over \mathcal{H} . The idea behind a Ritz approximation is to minimize this energy over a subspace.

Definition (Ritz approximation)

Let $n \in \mathbb{N}$ and $\mathcal{H}_n \leq \mathcal{H}$, with $\dim \mathcal{H}_n = n$. We say that $u_n \in \mathcal{H}_n$ is a **Ritz approximation** to u , solution of (1), if and only if

$$u_n = \underset{v_n \in \mathcal{H}_n}{\operatorname{argmin}} E(v_n).$$



Basis Representation

Before we begin to provide an analysis of these methods, we must realize what we have gained with a Galerkin, or Ritz, approach. Let us introduce a basis of \mathcal{H}_n ,

$$B_n = \{\phi_i\}_{i=1}^n, \quad \mathcal{H}_n = \text{span } B_n.$$

Then every $v_n \in \mathcal{H}_n$ has a unique representation of the form

$$v_n = \sum_{i=1}^n v_i \phi_i \quad \longleftrightarrow \quad \mathbf{v} = [v_1, \dots, v_n]^T \in \mathbb{R}^n.$$

Since this representation is dependent on the basis, we will often write

$$v_n \in \mathcal{H}_n \xleftrightarrow{B_n} \mathbf{v} \in \mathbb{R}^n$$

to emphasize this connection. The object on the right is called the coordinate vector.



The Stiffness Matrix

Now let us use such representation in definition of the Galerkin method and notice that, by linearity, it is sufficient to set $v_n = \phi_j \in B_n$. We obtain

$$\mathcal{A}(u_n, \phi_j) = \mathcal{A}\left(\sum_{i=1}^n u_i \phi_i, \phi_j\right) = \sum_{i=1}^n \mathcal{A}(\phi_i, \phi_j) u_i = F(\phi_j), \quad j = 1, \dots, n.$$

We have obtained a linear system of equations for $\mathbf{u} \overset{B_n}{\longleftrightarrow} u_n$, namely

$$A\mathbf{u} = \mathbf{f},$$

where $f_i = F(\phi_i)$, for $i = 1, \dots, n$. All the properties of this system are contained in the so-called *stiffness matrix*, A .



Lemma (properties of A)

Let \mathcal{H}_n be a finite-dimensional subspace of \mathcal{H} . The stiffness matrix $A = [a_{ij}] \in \mathbb{R}^{n \times n}$ relative to the basis $B_n = \{\phi_1, \dots, \phi_n\}$ is positive definite. In addition, if \mathcal{A} is symmetric, so is A .

Proof.

Let $v_n \in \mathcal{H}_n \xleftrightarrow{B_n} \mathbf{v} \in \mathbb{R}^n$ be arbitrary. Then, using the coercivity of $\mathcal{A}(\cdot, \cdot)$, we get

$$\mathbf{v}^\top \mathbf{A} \mathbf{v} = \sum_{i=1}^n \sum_{j=1}^n \mathcal{A}(\phi_j, \phi_i) v_i v_j = \mathcal{A}(\mathbf{v}_n, \mathbf{v}_n) \geq \alpha \|\mathbf{v}_n\|_{\mathcal{H}}^2.$$

The symmetry of A is obtained by observing that, since \mathcal{A} is symmetric,

$$a_{i,j} = \mathcal{A}(\phi_j, \phi_i) = \mathcal{A}(\phi_i, \phi_j) = a_{j,i}.$$





Theorem (uniform well-posedness)

Suppose that \mathcal{H}_n is a finite-dimensional subspace of \mathcal{H} , with basis $B_n = \{\phi_1, \dots, \phi_n\}$. There is a unique $u_n \in \mathcal{H}_n \xleftrightarrow{B_n} \mathbf{u} \in \mathbb{R}^n$ that is a Galerkin approximation of the solution to (1), which, moreover, satisfies

$$A\mathbf{u} = \mathbf{f}.$$

If, in addition, \mathcal{A} is symmetric, then $u_n \in \mathcal{H}_n$ is the Galerkin approximation to u if and only if it is the Ritz approximation to u . Finally, we have the estimate

$$\|u_n\|_{\mathcal{H}} \leq \frac{1}{\alpha} \|F\|_{\mathcal{H}'}.$$

Proof.

Exercise.





Theorem (Céa's Lemma)

Suppose that \mathcal{H}_n is a finite-dimensional subspace of \mathcal{H} , with basis $B_n = \{\phi_1, \dots, \phi_n\}$. Suppose that $u \in \mathcal{H}$ is the unique solution to (1) and $u_n \in \mathcal{H}_n \xrightarrow{B_n} \mathbf{u} \in \mathbb{R}^n$ is the unique Galerkin approximation. Then

$$\mathcal{A}(u - u_n, v_n) = 0, \quad \forall v_n \in \mathcal{H}_n, \quad (3)$$

which is called Galerkin orthogonality. Furthermore, we have the quasi-best approximation property

$$\|u - u_n\|_{\mathcal{H}} \leq \frac{M}{\alpha} \min_{v \in \mathcal{H}_n} \|u - v\|_{\mathcal{H}},$$

where M and α denote the boundedness and coercivity constants, respectively, of the bilinear form \mathcal{A} .

Proof.

Notice that, since $\mathcal{H}_n \subset \mathcal{H}$, we can set, in (1), the test function $v = v_n \in \mathcal{H}_n$. Subtracting this from the definition of Galerkin approximation, we then obtain the Galerkin orthogonality relation (3).



Remark (energy norm)

Notice that if \mathcal{A} is symmetric, then it defines the so-called energy norm

$$\|v\|_F = \mathcal{A}(v, v)^{1/2}, \quad \forall v \in \mathcal{H}.$$

The proof of the previous result shows that, in this case, we actually have the best approximation property

$$\|u - u_n\|_E = \min_{v \in \mathcal{H}_n} \|u - v\|_E.$$

Remark (best approximation)

Notice that the content of Céa's Lemma reduced the analysis of a Galerkin approximation to a question in approximation theory. In other words, after this result, it is no longer of relevance that we are trying to approximate the solution to (1). The only thing that matters is how well an object in \mathcal{H} can be approximated by elements in the subspace \mathcal{H}_n .



The Finite Element Method in One Dimension



The One-Dimensional Problem

Let $d = 1$, $\Omega = (0, 1)$, and consider the problem

$$-\frac{d}{dx} \left(a(x) \frac{du}{dx}(x) \right) = f(x), \quad x \in (0, 1), \quad \text{with} \quad u(0) = 0 = u(1).$$

where $a \in C(\bar{\Omega})$, such that, for some $\alpha, M > 0$,

$$0 < \alpha \leq a(x) \leq M, \quad \forall x \in [0, 1].$$

We will consider a weaken form of this problem: given $f \in L^2(\Omega)$, we seek a function $u \in H_0^1(\Omega)$ such that

$$\mathcal{A}(u, v) = \int_0^1 a u' v' \, dx = \int_0^1 f v \, dx, \quad \forall v \in H_0^1(\Omega), \quad (4)$$

where

$$H_0^1(\Omega) = \{v \in L^2(\Omega) \mid v' \in L^2(\Omega)\}.$$



Deriving the Weak Form

Strong Form: given $f : [0, 1] \rightarrow \mathbb{R}$, find $u : [0, 1] \rightarrow \mathbb{R}$ that solves

$$L[u](x) = -\frac{d}{dx} \left(a(x) \frac{du}{dx}(x) \right) = f(x), \quad 0 \leq x \leq 1,$$

with the boundary conditions $u(0) = 0 = u(1)$. Multiply the equation by a test function v and integrate to obtain

$$\int_0^1 a(x) u'(x) v'(x) dx - a(x) u'(x) v(x) \Big|_{x=0}^{x=1} = \int_0^1 f(x) v(x) dx,$$

where we used integration by parts. Now suppose that v has the same boundary conditions as u , that is, $v(0) = 0 = v(1)$. Then

$$\int_0^1 a(x) u'(x) v'(x) dx = \int_0^1 f(x) v(x) dx,$$

for all v with $v(0) = 0 = v(1)$. See Chapter 23 for more details.



Definition (mesh)

Let $\Omega = (0, 1)$, a **triangulation** or **mesh** of Ω is a partition of Ω into subintervals I_i , which we call **elements**:

$$\mathcal{T}_h = \{I_i\}_{i=0}^N, \quad I_i = (x_i, x_{i+1}), \quad h_i = x_{i+1} - x_i,$$

where the **nodes** are given by

$$0 = x_0 < x_1 < \cdots < x_{N+1} = 1.$$

Notice that, in the previous definition, we are not assuming that the nodes are equally spaced, i.e., h_i is not constant. We set

$$h = \max_{i=0, \dots, N} h_i,$$

which we call the mesh size. Subordinate to the mesh \mathcal{T}_h we define a finite element space.



Definition (finite element spaces)

Given a mesh \mathcal{T}_h , we define the finite element spaces of **continuous piecewise linear functions**

$$\begin{aligned}\mathcal{S}^{1,0}(\mathcal{T}_h) &= \left\{ v_h \in C([0, 1]) \mid v_h|_{I_j} \in \mathbb{P}_1, \ 0 \leq j \leq N \right\}, \\ \mathcal{S}_0^{1,0}(\mathcal{T}_h) &= \mathcal{S}^{1,0}(\mathcal{T}_h) \cap H_0^1(\Omega).\end{aligned}$$

The **Lagrange nodal basis** of $\mathcal{S}^{1,0}(\mathcal{T}_h)$ is given by

$$\{\phi_i\}_{i=0}^{N+1} \subset \mathcal{S}^{1,0}(\mathcal{T}_h), \quad \phi_i(x_j) = \delta_{ij},$$

and the **Lagrange nodal basis** of $\mathcal{S}_0^{1,0}(\mathcal{T}_h)$ is $\{\phi_i\}_{i=1}^N$.

Remark

It is straightforward to show that $\tilde{B}_h = \{\phi_i\}_{i=0}^{N+1}$ is a basis of $\mathcal{S}^{1,0}(\mathcal{T}_h)$ and $B_h = \{\phi_i\}_{i=1}^N$ is a basis of $\mathcal{S}_0^{1,0}(\mathcal{T}_h)$. These facts show that $\dim(\mathcal{S}^{1,0}(\mathcal{T}_h)) = N + 2$ and $\dim(\mathcal{S}_0^{1,0}(\mathcal{T}_h)) = N$.

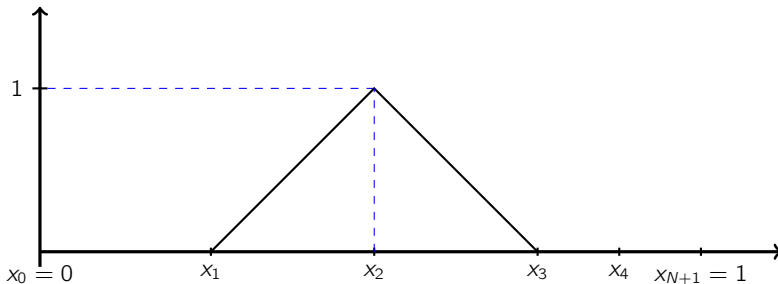


Figure: A one-dimensional hat function, i.e., a member of the Lagrange nodal basis for $\mathcal{P}^{1,0}(\mathcal{T}_h)$.

Remark (hat function)

A typical depiction of a function ϕ_i is given in the figure. This motivates us to call them hat functions.



Remark (implementation)

Recall that the support of a function is defined as

$$\text{supp } v = \overline{\{x \mid v(x) \neq 0\}}.$$

Notice that

$$\text{supp } \phi_i \cap \text{supp } \phi_j \neq \emptyset \iff |i - j| \leq 1.$$

This means that the stiffness matrix will be very sparse; in fact, tridiagonal in the present case. In addition, because of the way the basis functions are defined, when computing the entries of the stiffness matrix or load vector, we can subdivide the integral into elements and operate locally, where the basis functions are linear.



Remark (implementation, Cont.)

For instance,

$$[A]_{i,j} = a_{i,j} = \int_0^1 a(x) \phi_j' \phi_i' dx = \sum_{k=0}^N \int_{I_k} a(x) \phi_j' \phi_i' dx.$$

Now, because of how the hat functions are defined, we observe that

$$\phi_i'(x) = \begin{cases} \frac{1}{h_{i-1}}, & x_{i-1} < x < x_i, \\ -\frac{1}{h_i}, & x_i < x < x_{i+1}, \\ 0, & x \notin [x_{i-1}, x_{i+1}]. \end{cases}$$

This can be used to efficiently implement the finite element method.



Interpolation Operators

Let us now proceed to analyze the finite element method in one dimension. Recall that, from C ea's Lemma, we immediately obtain that

$$\|u - u_h\|_{H_0^1(0,1)} \leq C \inf_{v_h \in \mathcal{S}_0^{1,0}(\mathcal{T}_h)} \|u - v_h\|_{H_0^1(0,1)}.$$

It is now necessary to find the *best approximation error*, i.e., the right-hand side of the previous inequality. To do so, we will construct an *interpolation operator*

$$\Pi_h: H_0^1(0,1) \rightarrow \mathcal{S}_0^{1,0}(\mathcal{T}_h),$$

and we will show that it has suitable approximation properties. That the interpolation operator is defined for $H_0^1(0,1)$ functions is a consequence of the following observation.



Remark (one-dimensional embedding)

In one space dimension, and in one dimension only, we have

$$H^1(\Omega) \hookrightarrow C(\overline{\Omega}).$$

This means that $H^1(\Omega) \subset C(\overline{\Omega})$, and there exists a constant $C > 0$, independent of u , such that

$$\|u\|_{L^\infty(0,1)} \leq C \|u\|_{H^1(0,1)}.$$



Definition (Lagrange interpolant)

The **Lagrange nodal interpolation** operator

$$\Pi_h: C([0, 1]) \rightarrow \mathcal{S}^{1,0}(\mathcal{T}_h)$$

is (uniquely) defined by $\Pi_h v(x_i) = v(x_i)$ for all $i = 0, \dots, N + 1$.



Theorem (properties of Π_h)

Let Π_h be as in the last definition. There is a constant $C > 0$, independent of the mesh spacings, such that, for all $j = 0, \dots, N$, and all $v \in H^1(0, 1)$ such that $v'' \in L^2(0, 1)$,

$$\begin{aligned}\|v - \Pi_h v\|_{L^2(I_j)} &\leq Ch_j^2 \|v''\|_{L^2(I_j)}, \\ \| (v - \Pi_h v)' \|_{L^2(I_j)} &\leq Ch_j \|v''\|_{L^2(I_j)}.\end{aligned}$$

Proof.

To alleviate the notation, let us denote by $I = (x_l, x_r)$ a generic element of the mesh and $h = |I| = x_r - x_l$.

Consider now $v \in H^1(0, 1)$ such that $v'' \in L^2(0, 1)$. One can show that $v \in C^1([0, 1])$. Define the first-order Taylor polynomial of v about x_l ,

$$Q_1 v(x) = v(x_l) + v'(x_l)(x - x_l) \in \mathbb{P}_1.$$



Proof (Cont.)

Notice that:

- ① Π_h is *polynomial space preserving* in the following sense: $\Pi_h Q_1 v$, when restricted to I , equals $Q_1 v$.
- ② The operator Π_h is *max-norm stable*: i.e.,

$$\|\Pi_h v\|_{L^\infty(I)} = \max\{|v(x_l)|, |v(x_r)|\} \leq \|v\|_{L^\infty(I)}.$$

Now, to bound the interpolation error $v - \Pi_h v$, we proceed as follows:

$$\|v - \Pi_h v\|_{L^\infty(I)} \leq \|v - Q_1 v\|_{L^\infty(I)} + \|Q_1 v - \Pi_h v\|_{L^\infty(I)} \leq 2\|v - Q_1 v\|_{L^\infty(I)}.$$

By Taylor's Theorem with integral remainder,

$$(v - Q_1 v)(x) = \int_{x_l}^x (x - y)v''(y)dy,$$

which implies that

$$\|v - Q_1 v\|_{L^\infty(I)} \leq h \int_I |v''| dy.$$



Proof (Cont.)

To bound the derivatives, we must note that

$$(\Pi_h v)'(x) = \frac{1}{h} \int_I v'(y) dy,$$

so that

$$(\Pi_h v - v)'(x) = \frac{1}{h} \int_I (v'(y) - v'(x)) dy,$$

which implies that

$$\begin{aligned} \|(\Pi_h v - v)'\|_{L^2(I)}^2 &= \frac{1}{h^2} \int_I \left(\int_I (v'(y) - v'(x)) dy \right)^2 dx \\ &= \frac{1}{h^2} \int_I \left(\int_I \int_x^y v''(s) ds dy \right)^2 dx \\ &\leq \frac{1}{h^2} \int_I \left(\int_I |x - y|^{1/2} \left(\int_{\min\{x,y\}}^{\max\{x,y\}} |v''(s)|^2 ds \right)^{1/2} dy \right)^2 dx. \end{aligned}$$



Proof (Cont.)

Now, since $x, y \in I$, we can bound $|x - y| \leq h$ and

$$\int_{\min\{x,y\}}^{\max\{x,y\}} |v''(s)|^2 ds \leq \int_I |v''(s)|^2 ds.$$

Using these upper bounds, we obtain

$$\|(\Pi_h v - v)'\|_{L^2(I)}^2 \leq \frac{1}{h} \int_I |v''(s)|^2 ds \int_I \left(\int_I dy \right)^2 dx \leq h^2 \int_I |v''(s)|^2 ds,$$

as we needed to show. □



Remark (stable and space-preserving operator)

Notice that, in the course of the proof of the last theorem, the particular form of the Lagrange interpolation operator was, ultimately, inconsequential. All that was needed was that the operator was stable and that it preserved the polynomial space. In consequence, the last theorem also holds for any other operator that satisfies these two properties.



Corollary (convergence of finite element method)

Let $u \in H_0^1(0, 1)$ solve (4) and $u_h \in \mathcal{S}_0^{1,0}(\mathcal{T}_h)$ be its finite element approximation. If u is such that $u'' \in L^2(0, 1)$, then we have

$$\|u - u_h\|_{H_0^1(0,1)} \leq Ch \|u''\|_{L^2(0,1)},$$

where the constant $C > 0$ is independent of $h > 0$, u , and u_h .

Proof.

By C  a's Lemma, we have that

$$\|u - u_h\|_{H_0^1(0,1)} \leq C \inf_{v_h \in \mathcal{S}_h^1(\mathcal{T}_h)} \|u - v_h\|_{H_0^1(0,1)} \leq C \|u - \Pi_h u\|_{H_0^1(0,1)},$$

where we used the Lagrange interpolation operator Π_h . Notice that if $u \in H_0^1(0, 1)$, then $\Pi_h u(0) = \Pi_h u(1) = 0$, so that $\Pi_h u \in \mathcal{S}_0^{1,0}(\mathcal{T}_h)$.



Proof (Cont.)

By the *local* properties of the Lagrange interpolation operator, we see that

$$\|u - \Pi_h u\|_{H_0^1(0,1)}^2 = \sum_{j=0}^N \int_{I_j} |(u - \Pi_h u)'|^2 dx \leq C \sum_{j=0}^N h_j^2 \int_{I_j} |u''|^2 dx.$$

Using that $h = \max_j h_j$ implies the result.



Towards An L^2 Error Estimate



Notice that the previous result, in conjunction with the Poincaré inequality (See Chapter 23), implies that

$$\|u - u_h\|_{L^2(0,1)} \leq C_P \|u - u_h\|_{H_0^1(0,1)} \leq Ch \|u''\|_{L^2(0,1)}.$$

However, our previous theorem shows that the interpolation error is $\mathcal{O}(h^2)$. How can we regain the missing power? For that, we need to study the *dual problem*, i.e., given $g \in L^2(0, 1)$, we need to find $z_g \in H_0^1(0, 1)$ such that

$$\mathcal{A}(v, z_g) = \int_0^1 g v \, dx, \quad \forall v \in H_0^1(0, 1). \quad (5)$$

Notice that the order of the arguments in the bilinear form is switched. This is irrelevant if the bilinear form is symmetric, as in our present case. If it is not symmetric, the order is quite important.



Theorem ($L^2(\Omega)$ -estimate)

Assume that, for every $g \in L^2(0, 1)$, there is a unique solution to the dual problem (5); furthermore, $z_g'' \in L^2(0, 1)$, with the estimate

$$\|z_g''\|_{L^2(0,1)} \leq C \|g\|_{L^2(0,1)}$$

for some constant $C > 0$. In this case, if $u \in H_0^1(0, 1)$ solves (4), it is such that $u'' \in L^2(0, 1)$, and $u_h \in \mathcal{S}_0^{1,0}(\mathcal{T}_h)$ is its finite element approximation, then we have

$$\|u - u_h\|_{L^2(0,1)} \leq Ch^2 \|u''\|_{L^2(0,1)}.$$

Proof.

Define the error $e = u - u_h \in H_0^1(0, 1) \subset L^2(0, 1)$. Let z_e be the solution to the dual problem (5) with data $g = e$. If that is the case, then we have

$$\|e\|_{L^2(0,1)}^2 = \mathcal{A}(e, z_e) = \mathcal{A}(u - u_h, z_e) = \mathcal{A}(u - u_h, z_e - \Pi_h z_e),$$

where the last equality follows from Galerkin orthogonality.



Proof (Cont.)

Now, using the boundedness of the bilinear form, we obtain

$$\begin{aligned} \|e\|_{L^2(0,1)}^2 &= \mathcal{A}(u - u_h, z_e - \Pi_h z_e) \\ &\leq M \|u - u_h\|_{H_0^1(0,1)} \|z_e - \Pi_h z_e\|_{H_0^1(0,1)} \\ &\leq Ch^2 \|u''\|_{L^2(0,1)} \|z_e''\|_{L^2(0,1)}, \end{aligned}$$

where we used the convergence estimate of the previous corollary and, since $z_e'' \in L^2(0,1)$, the interpolation estimates. Using the estimate on the second derivatives of z_e then yields the result. □



Remark (higher order elements)

In our presentation, we chose the space $\mathcal{S}^{1,0}(\mathcal{T}_h)$ to make the discussion as transparent as possible. It is also possible to define finite element spaces of higher order. For $p \in \mathbb{N}$, we define

$$\begin{aligned}\mathcal{S}^{p,0}(\mathcal{T}_h) &= \left\{ v_h \in C([0, 1]) \mid v_h|_{I_j} \in \mathbb{P}_p, \ 0 \leq j \leq N \right\}, \\ \mathcal{S}_0^{p,0}(\mathcal{T}_h) &= \mathcal{S}^{p,0}(\mathcal{T}_h) \cap H_0^1(\Omega).\end{aligned}$$

The analysis of finite element methods with these spaces follows verbatim what we have done here. The only difference is in the way that the Lagrange interpolation operator is defined. In short, provided that $u^{(p+1)} \in L^2(0, 1)$, one can prove that

$$\|u - u_h\|_{H_0^1(0,1)} \leq Ch^p \|u^{(p+1)}\|_{L^2(0,1)}.$$



Remark (spaces of variable degree)

The spaces $\mathcal{S}^{p,0}(\mathcal{T}_h)$ can be even further generalized and allowed to have a different polynomial degree within each element. To define them, we let $\mathbf{p} \in \mathbb{N}^{N+1}$, which is called the degree vector. Then

$$\mathcal{S}^{p,0}(\mathcal{T}_h) = \left\{ v_h \in C([0, 1]) \mid v_h|_{I_j} \in \mathbb{P}_{p_{j+1}}, \ 0 \leq j \leq N \right\},$$

*with $\mathcal{S}_0^{p,0}(\mathcal{T}_h) = \mathcal{S}^{p,0}(\mathcal{T}_h) \cap H_0^1(0, 1)$. These spaces form the building block of what is known as *hp* finite element methods, where, to increase the accuracy of our numerical approximation, one is allowed to either reduce the local mesh size h_j or increase the polynomial degree p_j . The reader is referred to the textbook and references therein for more details.*



The Finite Element Method in Two Dimensions

Our Goal



As we saw in the one-dimensional case of the previous section, the finite element method is a particular version of the Galerkin method. More specifically, the finite element method gives a particular subspace where we seek the approximate solution, and a particular basis for it. In this section, we will present, mostly without proofs, the construction and analysis of finite element methods in two dimensions. We refer the reader to the textbook and references therein for full details, and further developments.

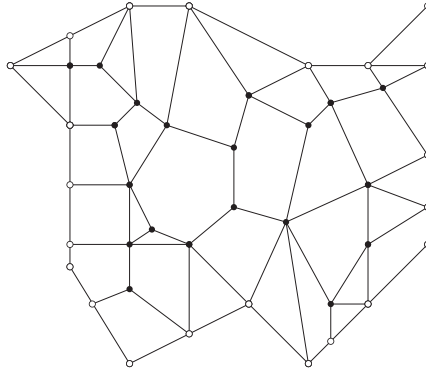


Figure: A conforming polygonal partition of a polygonal domain Ω . For every pair of elements, only one of the following possibilities holds: either $\overline{K}_i \cap \overline{K}_j = \emptyset$ or $\overline{K}_i \cap \overline{K}_j = e$, a complete edge of both K_i and K_j , or $\overline{K}_i \cap \overline{K}_j = \{x\}$, a shared vertex of K_i and K_j . The filled circles are the interior vertices and the unfilled circles are the boundary vertices.



Definition (polygonal partition)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, polygonal domain. Let

$$\mathcal{T}_h = \{K_i \mid i = 1, \dots, N_e\}$$

be a disjoint collection of open subsets of Ω such that

$$\overline{\Omega} = \bigcup_{i=1}^{N_e} \overline{K_i}.$$

The members $K_i \in \mathcal{T}_h$ are called **elements**. \mathcal{T}_h is called a **mesh** or **polygonal partition** of Ω if and only if each K_i is a convex polygon. \mathcal{T}_h is called a **triangulation** of Ω if and only if each element K_i is a triangle. Define the **element diameters** via

$$h_i = \text{diam}(K_i), \quad i = 1, \dots, N_e.$$

The value

$$h = \max_{1 \leq i \leq N_e} h_i$$

is called the **global mesh size**.



Definition (polygonal partition (Cont.))

By the set

$$\mathcal{N}_v = \{\mathbf{x}_j \mid j = 1, \dots, N_v\},$$

we denote the set of all **vertices** of \mathcal{T}_h , i.e., all the vertex points of the polygons K_i . By

$$\mathcal{N}_v^i = \mathcal{N}_v \cap \Omega = \left\{ \mathbf{x}_j \mid j = 1, \dots, N_v^i \right\},$$

we denote the set of all **interior vertices**. The set $\mathcal{N}_v \setminus \mathcal{N}_v^i$ is the set of **boundary vertices**.

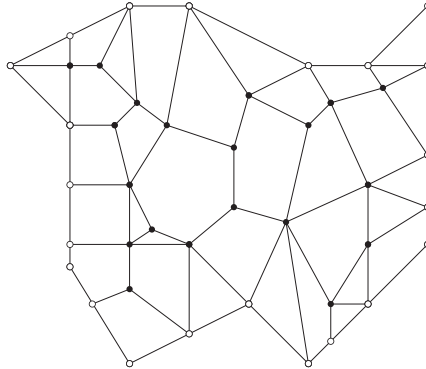


Figure: A conforming polygonal partition of a polygonal domain Ω . For every pair of elements, only one of the following possibilities holds: either $\overline{K}_i \cap \overline{K}_j = \emptyset$ or $\overline{K}_i \cap \overline{K}_j = e$, a complete edge of both K_i and K_j , or $\overline{K}_i \cap \overline{K}_j = \{x\}$, a shared vertex of K_i and K_j . The filled circles are the interior vertices and the unfilled circles are the boundary vertices.



Definition (conforming partition)

A polygonal partition \mathcal{T}_h is called **conforming** if and only if for every pair of elements, only one of the following possibilities holds:

- ❶ $\overline{K}_i \cap \overline{K}_j = \emptyset$,
- ❷ $\overline{K}_i \cap \overline{K}_j = e$, a complete edge of both K_i and K_j , or
- ❸ $\overline{K}_i \cap \overline{K}_j = \{\mathbf{x}\}$, a shared vertex of K_i and K_j .

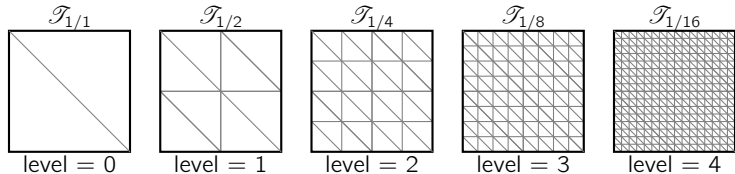


Figure: A family of nested uniform triangulations of a square. Starting at level = 0, the level = 1 triangulation is obtained by connecting the three midpoints of each triangle to form four congruent sub-triangles. This process can continue indefinitely. The global mesh size decreases by two as the level increases by one.

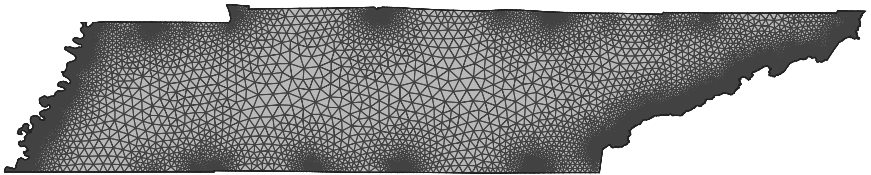


Figure: A sophisticated triangulation of the state of Tennessee (one of the 50 constituent states of the United States of America).



Definition (finite element space)

Suppose that $p \in \mathbb{N}$ and $\Omega \subset \mathbb{R}^2$ is an open polygonal domain. Let $\mathcal{T}_h = \{K_i\}$ be a conforming triangulation of Ω . Define

$$\mathcal{S}^{p,0}(\mathcal{T}_h) = \{v \in C(\overline{\Omega}) \mid v|_K \in \mathbb{P}_p, \forall K \in \mathcal{T}_h\}.$$

The set $\mathcal{S}^{p,0}(\mathcal{T}_h)$ is called the **piecewise polynomial (of degree p) finite element space**. By the set

$$\mathcal{S}_0^{p,0}(\mathcal{T}_h) = \{v \in \mathcal{S}^{p,0}(\mathcal{T}_h) \mid v|_{\partial\Omega} = 0\},$$

we denote the subspace of functions that vanish on the boundary.

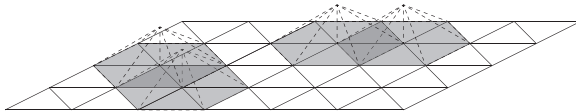


Figure: Some of the hat basis functions from $\mathcal{S}_0^{1,0}(\mathcal{T}_h)$ subordinate to a uniform conforming triangulation of a rectangle. The supports of the basis functions are the grey shaded triangles in the mesh. The dark grey triangles indicate the regions where the supports intersect.



Theorem (embedding)

Suppose that $p \in \mathbb{N}$ and $\Omega \subset \mathbb{R}^2$ is an open polygonal domain. Let $\mathcal{T}_h = \{K_i\}$ be a conforming triangulation of Ω . Then $\mathcal{S}^{p,0}(\mathcal{T}_h)$ is a subspace of $H^1(\Omega)$ and $\mathcal{S}_0^{p,0}(\mathcal{T}_h)$ is a subspace of $H_0^1(\Omega)$.

Proof.

See the books by Braess (2007) and Brenner and Scott (2007).





Theorem (basis of $\mathcal{S}^{1,0}(\mathcal{T}_h)$)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open polygonal domain and $\mathcal{T}_h = \{K_i\}$ is a conforming triangulation of Ω . Then

$$\dim(\mathcal{S}^{1,0}(\mathcal{T}_h)) = N_v, \quad \dim(\mathcal{S}_0^{1,0}(\mathcal{T}_h)) = N_v^i.$$

In particular, defining $\phi_k \in \mathcal{S}_0^{1,0}(\mathcal{T}_h)$ via

$$\phi_k(\zeta_j) = \delta_{j,k}, \quad \forall \zeta_j \in \mathcal{N}_v^i,$$

we see that $\{\phi_1, \dots, \phi_{N_v^i}\}$ is a basis for $\mathcal{S}_0^{1,0}(\mathcal{T}_h)$. The basis for $\mathcal{S}^{1,0}(\mathcal{T}_h)$ is constructed similarly.

Proof.

Exercise. □



Definition (midpoints)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open polygonal domain and $\mathcal{T}_h = \{K_i\}$ is a conforming triangulation of Ω . By

$$\mathcal{N}_m = \{\xi_j \mid j = 1, \dots, N_m\},$$

we denote the set of midpoints of all edges in the triangulation, i.e., the **midpoints set**. By

$$\mathcal{N}_m^i = \mathcal{N}_m \cap \Omega = \{\xi_j \mid j = 1, \dots, N_m^i\},$$

we denote the set of all interior edge midpoints, i.e., the **interior midpoints set**. The set $\mathcal{N}_m \setminus \mathcal{N}_m^i$ is the **boundary midpoints set**.



Theorem (basis of $\mathcal{S}^{2,0}(\mathcal{T}_h)$)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open polygonal domain and $\mathcal{T}_h = \{K_i\}$ is a conforming triangulation of Ω . Then

$$\dim(\mathcal{S}^{2,0}(\mathcal{T}_h)) = N_v + N_m, \quad \dim(\mathcal{S}_0^{2,0}(\mathcal{T}_h)) = N_v^i + N_m^i.$$

In particular, defining $\phi_k \in \mathcal{S}_0^{2,0}(\mathcal{T}_h)$ via

$$\phi_k(\zeta_j) = \delta_{j,k}, \quad \forall \zeta_j \in \mathcal{N}_v^i \cup \mathcal{N}_m^i,$$

we see that $\{\phi_1, \dots, \phi_{N_v^i + N_m^i}\}$ is a basis for $\mathcal{S}_0^{2,0}(\mathcal{T}_h)$. The basis for $\mathcal{S}^{2,0}(\mathcal{T}_h)$ is constructed similarly.

Proof.

The figure on the next page gives an illustration of this construction. □

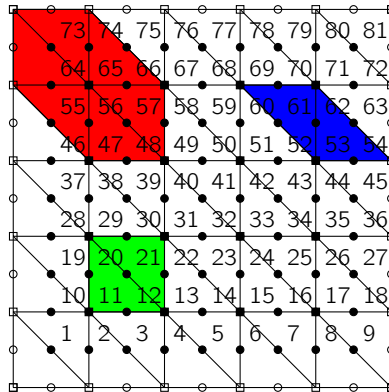


Figure: A uniform triangulation of a square domain Ω showing the nodes for the piecewise quadratic finite element space $\mathcal{S}_0^{2,0}(\mathcal{T}_h)$. The interior edge midpoint nodes are the filled circles; the unfilled circles are the boundary edge midpoint nodes. The interior vertex nodes are the filled squares; the unfilled squares are the boundary vertex nodes. The supports of the Lagrange nodal basis elements ϕ_{21} , ϕ_{62} , and ϕ_{65} are shown as the shaded regions. There are exactly 81 nodes in the mesh and $\dim(\mathcal{S}_0^{2,0}(\mathcal{T}_h)) = 81$. Observe that $N_v^i = 16$ and $N_m^i = 65$.



Remark (nodal basis)

The bases that we constructed in the previous theorems are called Lagrange nodal bases. These have the property that the basis elements satisfy $\phi_k(\xi_j) = \delta_{j,k}$, for $1 \leq j, k \leq N$, where N is the number of elements in the basis, and $\{\xi_j\}$ is the Lagrange nodal set. We have only defined this basis for $\mathcal{S}_0^{p,0}(\mathcal{T}_h)$ with $p = 1, 2$, where Ω is a polygonal set. But we can construct this type of basis for any $p \in \mathbb{N}$.



The Basic Problem in Two Dimensions

Having defined finite element spaces, we can use them to approximate weak solutions to elliptic problems. Let us illustrate this in the case of the Poisson problem. Thus, let $\Omega \subset \mathbb{R}^2$ be a bounded polygonal domain, $f \in L^2(\Omega)$, and we seek $u \in H_0^1(\Omega)$ such that

$$\mathcal{A}(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, d\mathbf{x} = \int_{\Omega} f v \, d\mathbf{x} = F(v), \quad \forall v \in H_0^1(\Omega). \quad (6)$$

The finite element method is then a Galerkin method, where we use as a subspace $\mathcal{S}_0^{p,0}(\mathcal{T}_h)$ for some $p \in \mathbb{N}$. Thus, we will seek $u_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)$ such that

$$\mathcal{A}(u_h, v_h) = F(v_h), \quad \forall v_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h). \quad (7)$$

Existence and uniqueness of discrete solutions, as well as a quasi-best approximation result

$$\|u - u_h\|_{H_0^1(\Omega)} \leq C \inf_{v_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)} \|u - v_h\|_{H_0^1(\Omega)},$$

follow from the general theory described in earlier. It remains then to provide error estimates.



Remark (sparsity)

The supports of the Lagrange nodal basis functions have minimal or no overlap. Consequently, the stiffness matrix A constructed in the abstract Galerkin framework will be sparse, i.e., having only a few nonzero elements.



Definition (Lagrange interpolant)

Let $\Omega \subset \mathbb{R}^2$ be an open, bounded, and polygonal domain. Assume that $\mathcal{T}_h = \{K_i\}$ is a conforming triangulation of Ω . Let $\{\phi_1, \dots, \phi_N\}$ be the Lagrange nodal basis for the space $\mathcal{S}_0^{p,0}(\mathcal{T}_h)$, with respect to the Lagrange nodal set $\{\zeta_j\}_{j=1}^N$, so that

$$\phi_j(\zeta_k) = \delta_{j,k}, \quad 1 \leq k, j \leq N.$$

The **Lagrange nodal interpolant**

$$\Pi_h: \{v \in C(\bar{\Omega}) \mid v|_{\partial\Omega} = 0\} \rightarrow \mathcal{S}_0^{p,0}(\mathcal{T}_h)$$

is defined as

$$\Pi_h v = \sum_{j=1}^N v(\zeta_j) \phi_j, \quad v \in C(\bar{\Omega}).$$



Proposition (projection)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, and polygonal domain; $\mathcal{T}_h = \{K_i\}$ is a conforming triangulation of Ω ; and $\{\phi_1, \dots, \phi_N\}$ is the Lagrange nodal basis for the space $\mathcal{S}_0^{p,0}(\mathcal{T}_h)$, with respect to the Lagrange nodal set $\{\zeta_j\}_{j=1}^N$. Then the Lagrange nodal interpolant Π_h is a projection operator, i.e., $\Pi_h^2 = \Pi_h$.

Proof.

An exercise. □

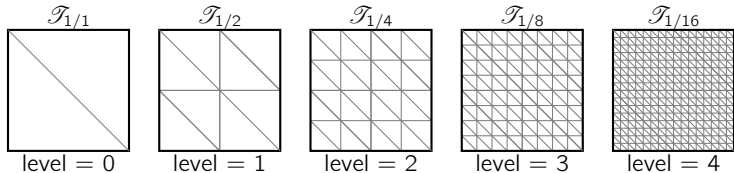


Definition (shape regularity)

Let $\Omega \subset \mathbb{R}^2$ be an open, bounded, and polygonal domain. Let $\{\mathcal{T}_h\}_{h>0}$ be a family of (not necessarily nested) conforming triangulations of Ω , parameterized by $h > 0$. The family $\{\mathcal{T}_h\}_{h>0}$ is called **shape regular** if and only if there is a constant $C > 0$ such that, for all $h > 0$ and all $K \in \mathcal{T}_h$,

$$1 \leq \frac{\rho_{\text{ext}}(K)}{\rho_{\text{int}}(K)} \leq C,$$

where $\rho_{\text{ext}}(K)$ is the radius of the smallest circle that circumscribes the triangle K and $\rho_{\text{int}}(K)$ is the radius of the largest circle that is inscribed in K .



Example

The nested family of triangulations shown in the figure above is shape regular. All triangles in every triangulation are right and isosceles.



Theorem (interpolation error estimate)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, polygonal domain and $\{\mathcal{T}_h\}_{h>0}$ is a family of (not necessarily nested) conforming, shape regular triangulations of Ω , parameterized by $h > 0$, where $h = \max_{K \in \mathcal{T}_h} h_K$. Then there is a constant $C_1 > 0$, independent of h but may depend on p , such that, for any $v \in H^{p+1}(\Omega)$ and all $0 \leq m \leq p$,

$$\|v - \Pi_h v\|_{H^m(\Omega)} \leq C_1 h^{p+1-m} |v|_{H^{p+1}(\Omega)}.$$

Proof.

See the book by Brenner and Scott (2007). □



Theorem (error estimate)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, polygonal domain; $f \in L^2(\Omega)$; and $\{\mathcal{T}_h\}_{h>0}$ is a family of (not necessarily nested) conforming, shape regular triangulations of Ω , parameterized by $h > 0$. Suppose that $u \in H_0^1(\Omega) \cap H^{p+1}(\Omega)$ is a weak solution to the Poisson problem (6). Suppose that $u_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)$ is the finite element approximation defined by (7). Then

$$\|u - u_h\|_{H^1(\Omega)} \leq \frac{M}{\alpha} C_1 h^p |u|_{H^{p+1}(\Omega)},$$

where $C_1 > 0$ is the constant from the theorem on interpolation error.

Proof.

One needs to use Céa's Lemma, and the interpolation error estimate from a previous theorem. The details are left to the reader as an exercise. □



Theorem (duality)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, polygonal, and convex domain; $f \in L^2(\Omega)$; and $\{\mathcal{T}_h\}_{h>0}$ is a family of (not necessarily nested) conforming, shape regular triangulations of Ω , parameterized by $h > 0$. Suppose that $u \in H_0^1(\Omega)$ is the weak solution of the Poisson problem (6) and $u_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)$ is the finite element approximation defined by (7). Then there is a constant $C_2 > 0$, independent of h and u , such that

$$\|u - u_h\|_{L^2(\Omega)} \leq C_2 h \|u - u_h\|_{H_0^1(\Omega)}.$$

Proof.

Set $e = u - u_h \in H_0^1(\Omega)$. Let $z_e \in H_0^1(\Omega)$ be the unique solution of the *dual problem*

$$\mathcal{A}(v, z_e) = \int_{\Omega} e v \, d\mathbf{x}, \quad \forall v \in H_0^1(\Omega).$$

Notice that, since \mathcal{A} is symmetric, the dual problem is equivalent to the original problem.



Proof (Cont.)

Since Ω is assumed to be convex, by the elliptic regularity result in Chapter 23, we have that $z_e \in H^2(\Omega) \cap H_0^1(\Omega)$ with

$$|z_e|_{H^2(\Omega)} \leq C_R \|e\|_{L^2(\Omega)}.$$

Now suppose that $v_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)$ is arbitrary and set $v = e$ in the dual problem. Using Galerkin orthogonality, and the boundedness of \mathcal{A} , we have that

$$\|e\|_{L^2(\Omega)}^2 = \int_{\Omega} e^2 \, d\mathbf{x} = \mathcal{A}(e, z_e) = \mathcal{A}(e, z_e - v_h) \leq M \|e\|_{H_0^1(\Omega)} \|z_e - v_h\|_{H_0^1(\Omega)}.$$

Let us choose $v_h = \Pi_h z$, where Π_h is the Lagrange interpolation operator into $\mathcal{S}_0^{1,0}(\mathcal{T}_h)$, the piecewise linear finite element space. Observe that, for any $p \in \mathbb{N}$, we have $\mathcal{S}_0^{1,0}(\mathcal{T}_h) \subseteq \mathcal{S}_0^{p,0}(\mathcal{T}_h)$.



Proof (Cont.)

Then, by the interpolation error theorem,

$$\begin{aligned}\|e\|_{L^2(\Omega)}^2 &\leq M \|e\|_{H_0^1(\Omega)} \|z_e - \Pi_h z_e\|_{H_0^1(\Omega)} \\ &\leq Ch^{2-1} \|e\|_{H_0^1(\Omega)} |z_e|_{H^2(\Omega)} \\ &\leq C_2 h \|e\|_{H_0^1(\Omega)} \|e\|_{L^2(\Omega)} .\end{aligned}$$

Therefore,

$$\|e\|_{L^2(\Omega)} \leq C_2 h \|e\|_{H_0^1(\Omega)} ,$$

and the result follows. □



Corollary ($L^2(\Omega)$ estimate)

Suppose that $\Omega \subset \mathbb{R}^2$ is an open, bounded, convex, polygonal domain; $f \in L^2(\Omega)$; and $\{\mathcal{T}_h\}_{h>0}$ is a family of (not necessarily nested) conforming, shape regular triangulations of Ω , parameterized by $h > 0$. Suppose that, for some $p \in \mathbb{N}$, $u \in H_0^1(\Omega) \cap H^{p+1}(\Omega)$ is a weak solution of the Poisson problem (6). Suppose also that $u_h \in \mathcal{S}_0^{p,0}(\mathcal{T}_h)$ is the finite element approximation defined by (7). Then there is a constant $C > 0$, independent of h , such that

$$\|u - u_h\|_{L^2(\Omega)} \leq Ch^{p+1} |u|_{H^{p+1}(\Omega)}.$$



Remark (linear finite elements)

If we set $p = 1$ in the corollary, i.e., if we suppose that $u \in H_0^1(\Omega) \cap H^2(\Omega)$ is a weak solution of the Poisson problem, then

$$\|u - u_h\|_{L^2(\Omega)} \leq Ch^2 |u|_{H^2(\Omega)} ,$$

provided that $\Omega \subset \mathbb{R}^2$ is a convex, bounded, polygonal domain. In other words, we get the expected second-order convergence, as for the finite difference approximation.