

Lecture Notes: African Institute of Mathematics–Senegal, January 2016

Topic Title: A short introduction to numerical methods for elliptic PDEs

Authors and Lecturers: Gerard Awanou (University of Illinois-Chicago) and Johnny Guzmán
(Brown University)

1. Description of Course

In this course we will mostly focus on second-order elliptic Partial Differential Equations (PDEs). We start by developing finite difference methods (a five point stencil method) for Poisson problems in one and two dimensions. Then, we develop finite element methods for general second order problems. Finally, we develop finite element methods for Stokes problems.

Students are expected to write computer code to implement the numerical methods. Also, theoretical aspects such as stability and error analysis will be covered and students are expected to be able to follow all the main arguments.

Two Point Boundary Value Problem: finite difference and finite element methods

Before we jump into higher dimensional problems it is instructive to consider a two point boundary value problem

1. Two Point Boundary Value problem

$$(1.1a) \quad -u''(x) = f(x) \quad \text{for all } 0 < x < 1,$$

$$(1.1b) \quad u(0) = a,$$

$$(1.1c) \quad u(1) = b,$$

Here the function $f(x)$ is given. This is commonly known as the *source term*. The data a, b is also normally given and this is known as the Dirichlet boundary data. The unknown function is u . For now we are assuming that u belongs to $C^2([0, 1]) = \{u : u, u', \text{ and } u'' \text{ are continuous on } [0, 1]\}$.

We would like to prove a crucial property of these type of equations *weak maximum principle* for this problem. Before we do that let us recall basic properties of calculus.

PROPOSITION 1. *Suppose that $w \in C^2((0, 1))$ and suppose that w has a local maximum at $0 < z < 1$ then*

$$(1.2) \quad w'(z) = 0 \text{ and } w''(z) \leq 0$$

Similarly if w has a local minimum at z then

$$(1.3) \quad w'(z) = 0 \text{ and } w''(z) \geq 0.$$

Theorem 1. *Suppose u solves (1.1) and $f(x) \leq 0$ for all $x \in (0, 1)$ then*

$$(1.4) \quad \max_{x \in [0, 1]} u(x) \leq \max\{a, b\}.$$

On the other hand, if $f(x) \geq 0$ for all $x \in (0, 1)$ then

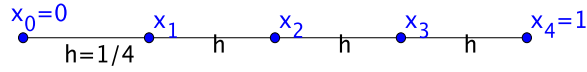
$$(1.5) \quad \min_{x \in [0, 1]} u(x) \geq \min\{a, b\}.$$

PROOF. We only prove (1.4) and leave the proof of (1.5) to the reader. We start by defining the function $\phi(x) = \frac{(x-1/2)^2}{2}$ and we let $w_\epsilon(x) = u(x) + \epsilon\phi(x)$ where $\epsilon > 0$. Noting that $\phi''(x) = 1$ for all x . Then we see that w_ϵ satisfies

$$(1.6a) \quad -w_\epsilon''(x) = f(x) - \epsilon \quad \text{for all } 0 < x < 1,$$

$$(1.6b) \quad w_\epsilon(0) = a + \frac{\epsilon}{8},$$

$$(1.6c) \quad w_\epsilon(1) = b + \frac{\epsilon}{8},$$

FIGURE 1. Example: $N = 3$, $h = 1/4$

Since $f(x) \leq 0$ for all $0 < x < 1$ we have that $w_\epsilon''(x) = -(f(x) - \epsilon) > 0$. We then see that w_ϵ cannot have a maximum at any point $0 < x < 1$. Indeed, if w_ϵ had a maximum at x then by (1.2) we have $w_\epsilon''(x) \leq 0$ contradicting that $w_\epsilon''(x) > 0$. Hence, the maximum must occur at $x = 0$ or $x = 1$. That is,

$$u(x) \leq u(x) + \epsilon\phi(x) = w_\epsilon(x) \leq \max\{w_\epsilon(0), w_\epsilon(1)\} = \max\{a + \epsilon/8, b + \epsilon/8\} \text{ for all } 0 \leq x \leq 1.$$

Therefore,

$$u(x) \leq \max\{a, b\} + \epsilon/8 \text{ for all } 0 \leq x \leq 1.$$

Now taking the limit at $\epsilon \rightarrow 0$ we have

$$u(x) \leq \max\{a, b\} \text{ for all } 0 \leq x \leq 1.$$

Taking the maximum we get (1.4). □

1.1. Finite Difference method. Now we develop a computational method to approximate (1.1). To do that, we let N be an integer and define $h = \frac{1}{N+1}$ and define $x_n = nh$ for $n = 0, 1, 2, \dots, N+1$ (see Figure 1). We will approximate $u(x_n)$ for $n = 0, 1, \dots, N+1$. Indeed, the finite difference method will find numbers v_0, v_1, \dots, v_{N+1} such that $v_n \approx u(x_n)$ for $n = 0, 1, \dots, N+1$.

In order to derive the method, we use Taylor expansions on the solution u

$$\begin{aligned} u''(x_i) &\approx (u'(x_{i+1}) - u'(x_i))/h \\ &\approx \left(\frac{u(x_{i+1}) - u(x_i)}{h} - \frac{u(x_i) - u(x_{i-1}))}{h} \right)/h \\ &= (u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))/h^2. \end{aligned}$$

In fact, we can show the following lemma (we leave the details to the reader).

LEMMA 1. *It holds*

$$(1.7) \quad u''(x_i) = (u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))/h^2 - \frac{h^2}{24}(u^{(4)}(\theta_1) + u^{(4)}(\theta_2))$$

where $x_{i-1} \leq \theta_1 \leq x_i$, $x_i \leq \theta_2 \leq x_{i+1}$.

Therefore, we define

Now, the finite difference method finds v_0, v_1, \dots, v_{N+1} such that:

$$(1.8a) \quad (-v_{n+1} + 2v_n - v_{n-1})/h^2 = f(x_n) \quad \text{for all } n = 1, 2, \dots, N,$$

$$(1.8b) \quad v_0 = a,$$

$$(1.8c) \quad v_{N+1} = b.$$

This will give rise to a system of N equations with N unknowns.

$$\begin{aligned} \frac{1}{h^2}(2v_1 - v_2) &= f(x_1) + a\frac{1}{h^2} \\ \frac{1}{h^2}(-v_1 + 2v_2 - v_3) &= f(x_2) \\ &\vdots \\ \frac{1}{h^2}(-v_{N-1} + 2v_N) &= f(x_N) + b\frac{1}{h^2} \end{aligned}$$

In matrix form we can write

$$\frac{1}{h^2} \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & & \\ 0 & \ddots & \ddots & \ddots & \\ \vdots & & & & -1 \\ 0 & \dots & & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_N \end{bmatrix} = \begin{bmatrix} f(x_1) \\ f(x_2) \\ f(x_3) \\ \vdots \\ f(x_N) \end{bmatrix} + \frac{1}{h^2} \begin{bmatrix} a \\ 0 \\ \vdots \\ 0 \\ b \end{bmatrix}$$

You can easily write a computer code to implement this method.

Next, we will analyze the finite difference method described above. To do this we use some notation. We let $S = \{x_0, x_1, \dots, x_{N+1}\}$ and we define all *discrete functions* as

$$P_h = \{v : v \text{ is a real valued function with domain } S\}.$$

That is, a function in P_h only has to be defined on the discrete points S . We use the notation $v_n = v(x_n)$.

We also denote the i -th derivative of a function as $D^i u(x) = u^{(i)}(x)$. In particular, we denote the second derivative $D^2 u(x) = u''(x)$. We then define the *discrete* second derivative as

$$D_h^2 v(x_n) \equiv \frac{1}{h^2}(v(x_{n+1}) - 2v(x_n) + v(x_{n-1})) \text{ for all } n = 1, \dots, N.$$

We also use the notation $D_h^2 v_n \equiv D_h^2 v(x_n)$. We can then rewrite (1.8) with the new notation. Find $v \in P_h$ satisfying

$$(1.9a) \quad -D_h^2 v_n = f(x_n) \quad \text{for all } n = 1, 2, \dots, N,$$

$$(1.9b) \quad v_0 = a,$$

$$(1.9c) \quad v_{N+1} = b.$$

1.1.1. *Discrete Maximum Principle.* We now state and prove the discrete analogue to Theorem 1.

Theorem 2. *Suppose v solves (1.9) and $f(x_n) \leq 0$ for all $n = 1, 2, \dots, N$ then*

$$(1.10) \quad \max_{n=0,1,\dots,N+1} v_n \leq \max\{a, b\}.$$

On the other hand, if $f(x_n) \geq 0$ for all $n = 1, 2, \dots, N$ then

$$(1.11) \quad \min_{n=0,1,\dots,N+1} v_n \geq \min\{a, b\}.$$

PROOF. We will prove (1.10) and leave the proof (1.11) to the reader. Let $M = \max_{n=0,1,\dots,N+1} v_n$. First suppose that $\max_{n=1,\dots,N} v_n < M$ then (1.10) holds.

On the other hand, suppose there that there exists an $1 \leq n \leq N$ such that $v_n = M$. Then from (1.9a) we have

$$M = v_n = \frac{1}{2}(v_{n+1} + v_{n-1}) + h^2 f(x_n) \leq \frac{1}{2}(v_{n+1} + v_{n-1}) \leq M$$

since $f(x_n) \leq 0$ and by our hypothesis $\frac{1}{2}(v_{n+1} + v_{n-1}) \leq M$. Hence, we must have

$$M = \frac{1}{2}(v_{n+1} + v_{n-1}),$$

or that

$$\frac{1}{2}(M - v_{n-1}) + \frac{1}{2}(M - v_{n+1}) = 0.$$

Since by our hypothesis $(M - v_{n-1}) \geq 0$ and $(M - v_{n+1}) \geq 0$ it must be $v_{n-1} = v_n = v_{n+1} = M$. We can continue this process to show that $v_i = M = \max_{n=0,1,\dots,N+1} v_n$ for all $i = 0, 1, \dots, N+1$. In other words, $v \equiv M$ is a constant discrete function. Therefore, (1.10) trivially holds. \square

We can use the discrete maximum principle to prove a stability result.

Theorem 3. *Let $v \in P_h$ solve (1.9) then we have*

$$(1.12) \quad \max_{n=0,1,\dots,N+1} |v_n| \leq \max\{|a|, |b|\} + \frac{1}{8} \max_{n=1,\dots,N} |f(x_n)|.$$

PROOF. Let $Q = \max_{n=1,\dots,N} |f(x_n)|$ and define $\phi \in P_h$ as follows

$$\phi_n = \phi(x_n) \equiv \frac{Q}{2}(x_n - 1/2)^2 \quad \text{for all } n = 0, 1, 2, \dots, N+1.$$

Not difficult to see that $D_h^2 \phi_n = Q$ for $n = 1, \dots, N$. Define $w \in P_h$ as $w = v + \phi$. Then we see that w satisfies

$$(1.13a) \quad -D_h^2 w_n = f(x_n) - Q \quad \text{for all } n = 1, 2, \dots, N,$$

$$(1.13b) \quad v_0 = a + \frac{Q}{8},$$

$$(1.13c) \quad v_{N+1} = b + \frac{Q}{8}.$$

Now, note that $f(x_n) - Q \leq 0$ for all n and hence by (1.10) we have

$$\max_{n=0,1,\dots,N+1} w_n \leq \max\left\{a + \frac{Q}{8}, b + \frac{Q}{8}\right\}.$$

Since $\phi_n \geq 0$ we have,

$$\max_{n=0,1,\dots,N+1} v_n \leq \max_{n=0,1,\dots,N+1} w_n$$

Moreover, we have

$$\max\left\{a + \frac{Q}{8}, b + \frac{Q}{8}\right\} \leq \max\{|a|, |b|\} + \frac{Q}{8}.$$

Hence, combining the last two inequalities we get □

$$(1.14) \quad \max_{n=0,1,\dots,N+1} v_n \leq \max\{|a|, |b|\} + \frac{Q}{8}.$$

Noting that

$$\begin{aligned} -D_h^2(-v_n) &= -f(x_n) \quad \text{for all } n = 1, 2, \dots, N, \\ -v_0 &= -a, \\ -v_{N+1} &= -b. \end{aligned}$$

we can apply the previous argument verbatim to get

$$(1.15) \quad \max_{n=0,1,\dots,N+1} (-v_n) \leq \max\{|a|, |b|\} + \frac{Q}{8}.$$

Combining inequality (1.14) and (1.15) we get (1.12).

1.1.2. Error Estimate for Finite Difference Method. Now that we have established the stability result (1.12), we can prove an error estimate for the finite difference method. Let us define the max norm for functions in $w \in C([0, 1])$ as

$$\|w\|_{C([0,1])} = \max_{0 \leq x \leq 1} |w(x)|.$$

Theorem 4. *Suppose that u solves (1.1) and $u \in C^4([0, 1])$. If v solves (1.9) we have*

$$\max_{n=0,1,\dots,N+1} |u(x_n) - v_n| \leq \frac{h^2}{104} \|D^4 u\|_{C([0,1])}.$$

PROOF. Define $w \in P_h$ as $w_n = u(x_n) - v_n$ for $n = 0, 1, \dots, N + 1$. Then, we see that

$$\begin{aligned} -D_h^2(w_n) &= \tau_n \quad \text{for all } n = 1, 2, \dots, N, \\ w_0 &= 0, \\ w_{N+1} &= 0. \end{aligned}$$

where $\tau_n = -D_h^2 u(x_n) - f(x_n) = D^2 u(x_n) - D_h^2 u(x_n)$. By (1.7) we have

$$(1.16) \quad |\tau_n| \leq \frac{h^2}{12} \|D^4 u\|_{C([0,1])} \quad \text{for all } n = 1, 2, \dots, N.$$

Now applying (1.12) we have

$$\max_{n=0,1,\dots,N+1} |w_n| \leq \frac{1}{8} \max_{n=1,\dots,N} |\tau_n|.$$

Combining this with (1.16) proves the result. □

1.2. Finite element method. In this section we derive the finite element method for (1.1). The methodology is quite different from the finite difference method but at the end they give similar methods.

We first need to define the weak formulation of (1.1). Let us recall an integration formula

$$(1.17) \quad \int_0^1 v'(x)w(x)dx = - \int_0^1 v(x)w'(x)dx + v(1)w(1) - v(0)w(0).$$

Now let $v \in C^1([0, 1])$ with $v(0) = 0$ and $v(1) = 0$ then if we multiply (1.1a) by v and integrate we get

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

If we apply (1.17) we get

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

where we used that $v(0) = 0 = v(1)$. Hence, we have that if u solves (1.1) then it solves

$$(1.18) \quad \int_0^1 u'(x)v'(x)dx = \int_0^1 f(x)v(x)dx \quad \text{for all } v \in C^1([0, 1]), v(0) = 0 = v(1).$$

The finite element method is based on (1.18).

The next step is to build a finite dimensional space of functions.

We let

$$V_h = \{v \in C([0, 1]) : \text{for all } i = 0, 1, \dots, N, v|_{(x_i, x_{i+1})} \text{ is a linear function}\}.$$

We can define a basis easily for this space in the following way. For each $i = 0, 1, 2, \dots$ we define ψ_i (see Figures 1 and 4)

$$\psi_i(x_j) = \begin{cases} 1 & \text{if } j = i \\ 0 & \text{if } j \neq i \end{cases}$$

Then we see that if $v \in V_h$ then we have

$$v(x) = \sum_{i=0}^{N+1} v(x_i)\psi_i(x) \text{ for } 0 \leq x \leq 1.$$

We see that the dimension of V_h is exactly $N + 1$. We also need to define

$$V_h^0 = \{v \in V_h : v(0) = 0 = v(1)\}.$$

The finite element method is as follows:

Find $u_h \in V_h$ such that

$$(1.19) \quad \int_0^1 u_h'(x)v_h'(x)dx = \int_0^1 f(x)v_h(x)dx \quad \text{for all } v_h \in V_h^0,$$

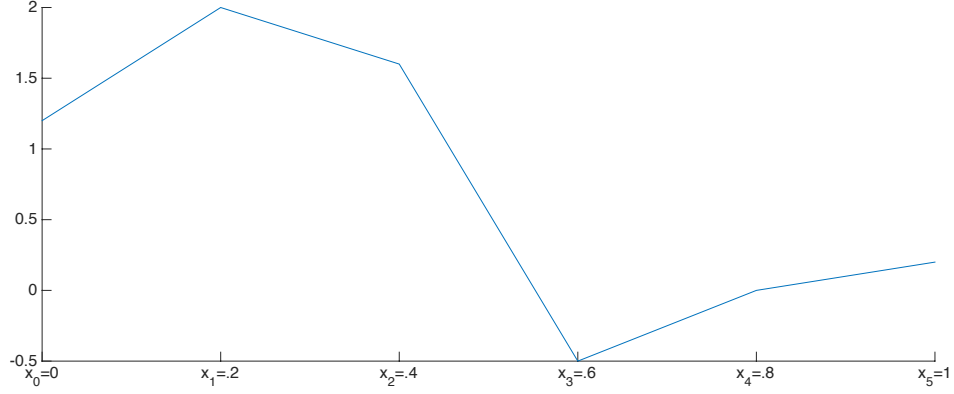
where $u_h(0) = u_h(x_0) = a$ and $u_h(1) = u_h(x_{N+1}) = b$.

Of course, since the space V_h^0 is finite dimensional then (1.21) is equivalent to:

$$(1.20) \quad \int_0^1 u_h'(x)\psi_j'(x)dx = \int_0^1 f(x)\psi_j(x)dx \quad \text{for all } j = 1, 2, \dots, N,$$

and $u_h(0) = u_h(x_0) = a$ and $u_h(1) = u_h(x_{N+1}) = b$. Writing

$$u_h'(x) = \sum_{i=0}^{N+1} u_h(x_i)\psi_i'(x) \text{ for } 0 \leq x \leq 1.$$

FIGURE 2. Example of a function in V_h , $h=.2$

we get

$$(1.21) \quad \int_0^1 u'_h(x) \psi'_j(x) dx = \int_0^1 \sum_{i=0}^{N+1} u_h(x_i) \psi'_i(x) \psi'_j(x) dx = \sum_{i=0}^{N+1} u_h(x_i) \int_0^1 \psi'_i(x) \psi'_j(x) dx.$$

Therefore, we have

$$\begin{aligned} \sum_{i=1}^N u_h(x_i) \int_0^1 \psi'_i(x) \psi'_j(x) dx &= \int_0^1 f(x) \psi_j(x) dx \\ &\quad - a \int_0^1 \psi'_0(x) \psi'_j(x) - b \int_0^1 \psi'_{N+1}(x) \psi'_j(x) \quad \text{for all } j = 1, 2, \dots, N, \end{aligned}$$

To write this method in matrix form we define

$$A = \begin{bmatrix} \int_0^1 \psi'_1(x) \psi'_1(x) dx & \int_0^1 \psi'_1(x) \psi'_2(x) dx & \dots & \int_0^1 \psi'_1(x) \psi'_N(x) dx \\ \int_0^1 \psi'_2(x) \psi'_1(x) dx & \int_0^1 \psi'_2(x) \psi'_2(x) dx & \dots & \\ \vdots & & \ddots & \\ \int_0^1 \psi'_N(x) \psi'_1(x) dx & \dots & & \int_0^1 \psi'_N(x) \psi'_N(x) dx \end{bmatrix}$$

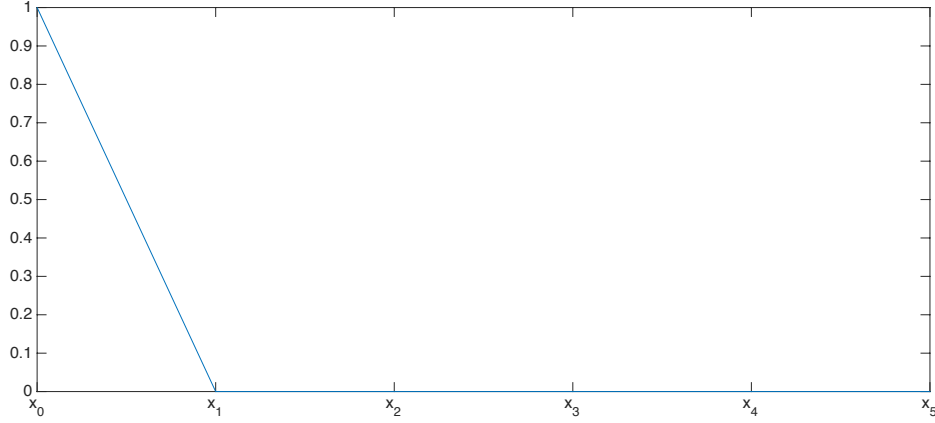


FIGURE 3. Graph of ψ_0

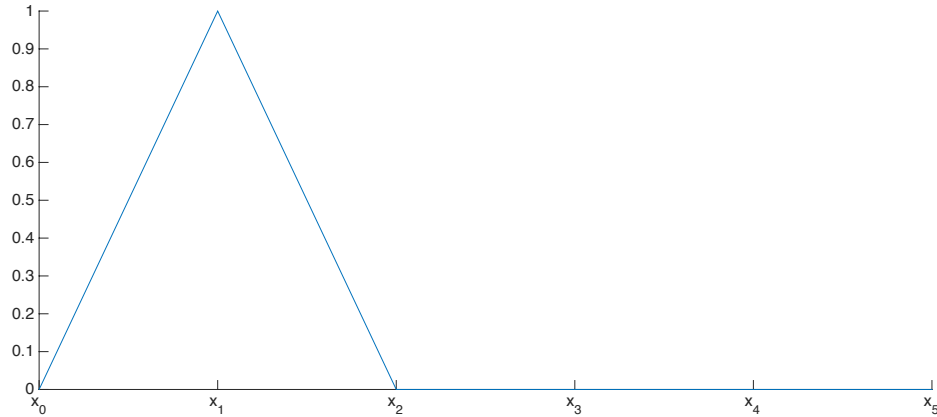
$$U = \begin{bmatrix} u_h(x_1) \\ u_h(x_2) \\ \vdots \\ u_h(x_N) \end{bmatrix}$$

$$F = \begin{bmatrix} \int_0^1 f(x)\psi_1(x)dx \\ \int_0^1 f(x)\psi_2(x)dx \\ \vdots \\ \int_0^1 f(x)\psi_N(x)dx \end{bmatrix}$$

$$G = -a \begin{bmatrix} \int_0^1 \psi'_0(x)\psi'_1(x)dx \\ \int_0^1 \psi'_0(x)\psi'_2(x)dx \\ \vdots \\ \int_0^1 \psi'_0(x)\psi'_N(x)dx \end{bmatrix} - b \begin{bmatrix} \int_0^1 \psi'_{N+1}(x)\psi'_1(x)dx \\ \int_0^1 \psi'_{N+1}(x)\psi'_2(x)dx \\ \vdots \\ \int_0^1 \psi'_{N+1}(x)\psi'_N(x)dx \end{bmatrix}$$

Then, U solves

$$AU = F + G.$$

FIGURE 4. Graph of ψ_1

In order to be able to implement the method we need to find need to be able to evaluate the entries of the matrix A and the vectors F and G . Lets start with A . An easy calculation shows that (see Figure 5

$$\psi'_i(x) = \begin{cases} 0 & x \in [0, x_{i-1}) \\ 1/h & x \in [x_{i-1}, x_i) \\ -1/h & x \in [x_i, x_{i+1}) \\ 0 & x \in (x_{i+1}, 1]. \end{cases}$$

Since

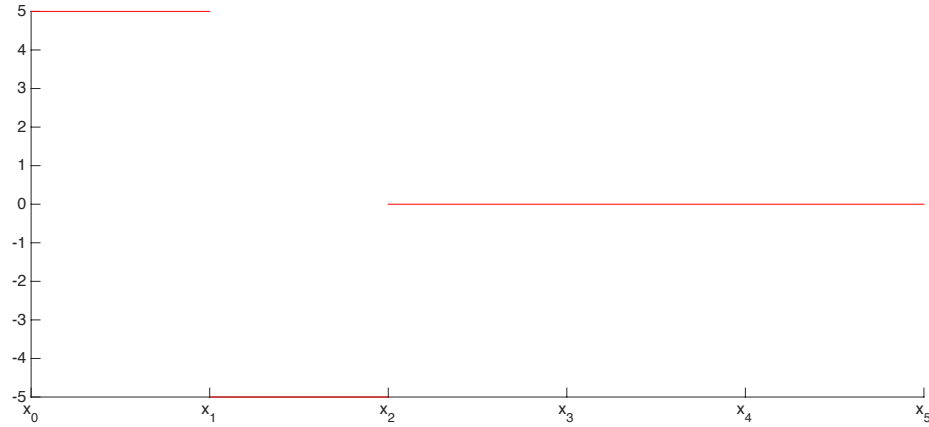
$$A_{ij} = \int_0^1 \psi'_i(x)\psi'_j(x)dx.$$

First note that since the support of ψ'_i is in (x_{i-1}, x_{i+1}) we see that

$$A_{ij} = 0 \text{ when } j = 1, \dots, i-2 \text{ and } j = i+2, \dots, N.$$

We thus only have to consider the case $j = i-1, i, i+1$. Let us start the $j = i$, In this case

$$A_{ii} = \int_{x_{i-1}}^{x_i} (1/h)^2 dx + \int_{x_i}^{x_{i+1}} (-1/h)^2 dx = 2/h.$$

FIGURE 5. Graph of ψ_1' , $h=0.2$

For $j = i + 1$ we have

$$A_{i,i+1} = \int_{x_i}^{x_{i+1}} (-1/h)(1/h)dx = -1/h.$$

Finally, for $j = i - 1$ we have

$$A_{i,i-1} = \int_{x_{i-1}}^{x_i} (1/h)(-1/h)dx = -1/h.$$

In other words, we get

$$A = \frac{1}{h} \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & & \\ 0 & \ddots & \ddots & \ddots & \\ \vdots & & & & -1 \\ 0 & \dots & & -1 & 2 \end{bmatrix}$$

We can easily calculate G to get

$$G = \frac{1}{h} \begin{bmatrix} a \\ 0 \\ \vdots \\ 0 \\ b \end{bmatrix}$$

Finally, in general we cannot compute exactly F . Instead we approximate F_i using the midpoint rule:

$$F_i = \int_0^1 f(x)\psi_i(x)dx = \int_{x_{i-1}}^{x_i} f(x)\psi_i(x)dx + \int_{x_i}^{x_{i+1}} f(x)\psi_i(x)dx \approx h \frac{f(\bar{x}_i) + f(\bar{x}_{i+1})}{2}.$$

where $\bar{x}_i = \frac{x_{i-1} + x_i}{2}$. Hence, computationally we can solve

$$\frac{1}{h^2} \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & & \\ 0 & \ddots & \ddots & \ddots & \\ \vdots & & & & -1 \\ 0 & \dots & & -1 & 2 \end{bmatrix} \begin{bmatrix} u_h(x_1) \\ u_h(x_2) \\ u_h(x_3) \\ \vdots \\ u_h(x_N) \end{bmatrix} = \begin{bmatrix} \frac{f(\bar{x}_1) + f(\bar{x}_2)}{2} \\ \frac{f(\bar{x}_2) + f(\bar{x}_3)}{2} \\ \frac{f(\bar{x}_3) + f(\bar{x}_4)}{2} \\ \vdots \\ \frac{f(\bar{x}_N) + f(\bar{x}_{N+1})}{2} \end{bmatrix} + \frac{1}{h^2} \begin{bmatrix} a \\ 0 \\ \vdots \\ 0 \\ b \end{bmatrix}$$

It is worth noting that the finite element method (where we approximate F with the midpoint rule) is very similar to the finite difference method. The only difference is the right hand side.

Finite difference method for poisson problem in two dimensions

In this chapter we consider a standard finite difference method for Poisson's problem. For simplicity, we will consider the domain $\Omega = (0,1)^2$. We let $\partial\Omega$ denote the boundary of Ω . Throughout we let $x = (x_1, x_2) \in \mathbb{R}^2$. We define the Laplacian Δ applied to a smooth function w :

$$\Delta w(x) = \partial_{x_1}^2 w(x) + \partial_{x_2}^2 w(x).$$

The **Poisson** problem is given by

$$(0.1a) \quad -\Delta u(x) = f(x) \quad \text{for } x \in \Omega$$

$$(0.1b) \quad u(x) = g(x) \quad \text{for } x \in \partial\Omega.$$

Here the functions f and g are given. The function f is known as the source term and g as the Dirichlet data. We will assume that f and g are smooth functions. The unknown is the function u for this section assume it belongs to $C^2(\Omega) \cap C(\bar{\Omega})$.

As we did in the one-dimensional case we can easily prove the weak maximum principle

Theorem 5. *Suppose u solves (0.1) and $f(x) \leq 0$ for $x \in \Omega$ then*

$$(0.2) \quad \max_{x \in \bar{\Omega}} u(x) \leq \max_{x \in \partial\Omega} g(x).$$

On the other hand, if $f(x) \geq 0$ for all $x \in \Omega$ then

$$(0.3) \quad \min_{x \in \bar{\Omega}} u(x) \geq \min_{x \in \partial\Omega} g(x).$$

PROOF. We only prove (0.2) and leave (0.3) to the reader. Define $\phi(x) = \frac{(x_1-1/2)^2 + (x_2-1/2)^2}{4}$ and let $w_\epsilon(x) = u(x) + \epsilon\phi(x)$ where $\epsilon > 0$. We see $\Delta\phi(x) = 1$ for all x and therefore we get

$$\begin{aligned} -\Delta w_\epsilon(x) &= f(x) - \epsilon & \text{for } x \in \Omega, \\ w_\epsilon(x) &= g(x) + \epsilon\phi(x) & \text{for } x \in \partial\Omega. \end{aligned}$$

Assume that w_ϵ attains a local maximum at $x \in \Omega$. Then by (1.2) we must have $\partial_{x_1}^2 w_\epsilon(x) \leq 0$ and $\partial_{x_2}^2 w_\epsilon(x) \leq 0$ which would imply that $\Delta w_\epsilon(x) \leq 0$. However, $\Delta w_\epsilon(x) = -(f(x) - \epsilon) > 0$ by our hypothesis which reaches a contradiction. Therefore, w_ϵ cannot attain a local maximum in Ω let alone its global maximum in Ω . Hence, w_ϵ must attain its global maximum on $\partial\Omega$. In other words,

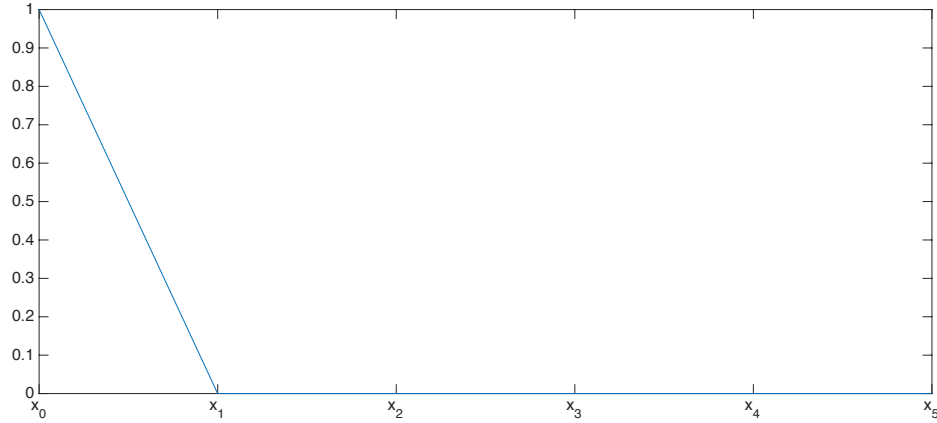
$$u(x) \leq u(x) + \epsilon\phi(x) = w_\epsilon(x) \leq \max_{y \in \partial\Omega} w_\epsilon(y) \quad \text{for all } x \in \bar{\Omega}.$$

On the other hand,

$$\max_{y \in \partial\Omega} w_\epsilon(y) \leq \max_{y \in \partial\Omega} g(y) + \epsilon \max_{y \in \partial\Omega} \phi(y) \leq \max_{x \in \partial\Omega} g(y) + \frac{\epsilon}{8}.$$

Where we used that $\max_{y \in \partial\Omega} \phi(y) \leq \frac{1}{8}$. Combining the last two inequalities and taking the limit as $\epsilon \rightarrow 0$ gives

$$u(x) \leq \max_{y \in \partial\Omega} g(y) \quad \text{for all } x \in \bar{\Omega}.$$

FIGURE 1. Graph of ψ_0

This proves (0.2). □

1. Finite difference method

In order to define a finite difference we need a grid defined on Ω . Let $N + 1$ be given and define $h = 1/(N + 1)$. Then, we define the grid as the collection of points

$$S = \{(nh, mh) : n, m = 0, 1, \dots, N + 1\}.$$

We denote the interior points as S_I

$$S_I = \{(nh, mh) : n, m = 1, \dots, N\}.$$

and the boundary points as $S_b = S \setminus S_I$. As we did in the one-dimensional case we define the space of discrete functions:

$$P_h = \{v : v \text{ is a real valued function with domain } S\}.$$

We denote the principle directions of \mathbb{R}^2 by $e_1 = (1, 0)$ and $e_2 = (0, 1)$.

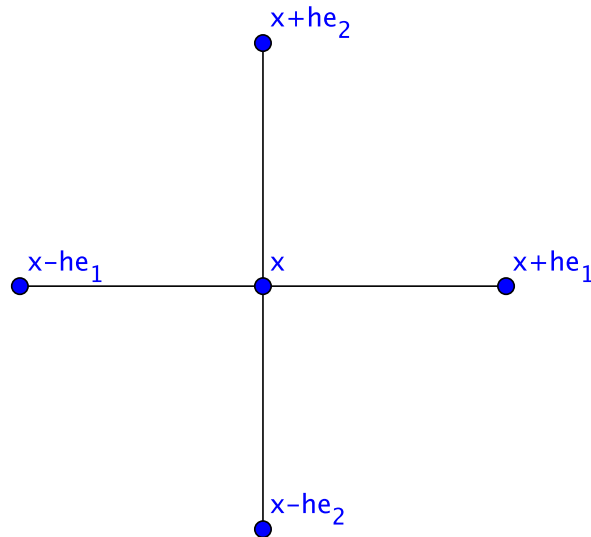


FIGURE 2. Five point stencil

We define the discrete Laplacian for every $x \in S_I$ as follows

$$\begin{aligned} \Delta_h v(x) &\equiv \frac{1}{h^2}(v(x + he_1) - 2v(x) + v(x - he_1)) + \frac{1}{h^2}(v(x + he_2) - 2v(x) + v(x - he_2)) \\ &= \frac{1}{h^2}(-4v(x) + v(x - he_1) + v(x - he_2) + v(x + he_1) + v(x + he_2)). \end{aligned}$$

In fact, we can prove the following result using Taylor's theorem

LEMMA 2. Let $u \in C^4(\bar{\Omega})$

$$(1.1) \quad \max_{x \in S_I} |\Delta u(x) - \Delta_h u(x)| \leq C h^2 \max\{\|\partial_{x_1}^4 u\|_{C(\Omega)}, \|\partial_{x_2}^4 u\|_{C(\Omega)}\}$$

where the constant C is independent of h and u . Here we define

$$\|w\|_{C(\Omega)} = \max_{x \in \bar{\Omega}} |w(x)|.$$

The standard finite difference method finds $v \in P_h$ such that

$$(1.2a) \quad -\Delta_h v(x) = f(x) \quad \text{for all } x \in S_I$$

$$(1.2b) \quad v(x) = g(x) \quad \text{for all } x \in S_b.$$

In order to be able to write a computer code we need to write the equivalent linear system that the finite difference gives rise to. To do that we should enumerate the vertices of S_I (note that there are N^2 vertices belonging to S_I). We define z_1, z_2, \dots, z_{N^2} as (see Figure 3)

$$z_{n+(m-1)N} = (nh, mh) \quad \text{for all } n, m = 1, \dots, N.$$

We then define $v(z_k) = v_k$ for $k = 1, \dots, N^2$.

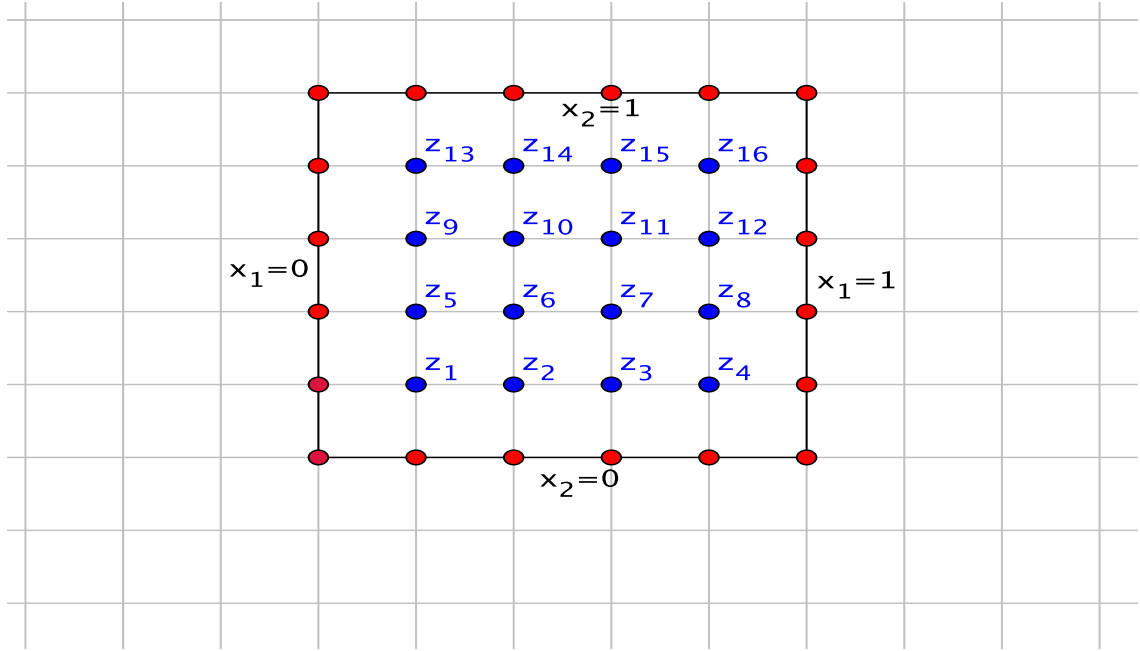


FIGURE 3. Grid, $N = 4$, grid points S_b are in red and S_I in blue

To write the linear system we need to distinguish between vertices that share an edge with a boundary vertex and those that don't. For vertices that don't share an edge with a boundary vertex we have the following linear equation

$$\frac{1}{h^2}(4v_k - v_{k+1} - v_{k+N} - v_{k-1} - v_{k-N}) = f(z_k)$$

for $k = n + (m - 1)N$ with $2 \leq n \leq N - 1$ and $2 \leq m \leq N - 1$.

Then for those interior vertices that share an edge with vertices belonging to line $x = 0$ we have the following linear system

$$\begin{aligned} \frac{1}{h^2}(4v_1 - v_2 - v_{N+1}) &= f(z_1) + \frac{1}{h^2}(g(z_1 - he_1) + g(z_1 - he_2)) \\ \frac{1}{h^2}(4v_2 - v_3 - v_{N+2} - v_1) &= f(z_2) + \frac{1}{h^2}(g(z_1 - he_2)) \\ &\vdots \\ \frac{1}{h^2}(4v_{N-1} - v_N - v_{2N-1} - v_{N-2}) &= f(z_{N-1}) + \frac{1}{h^2}(g(z_{N-1} - he_2)) \\ \frac{1}{h^2}(4v_N - v_{2N} - v_{N-1}) &= f(z_N) + \frac{1}{h^2}(g(z_N - he_2) + g(z_N + he_1)) \end{aligned}$$

Similarly, we can write the linear equations for the rest of vertices sharing an edge with a boundary vertex.

2. Error Analysis of Finite difference method

In order to obtain the error estimates of the finite difference method we need a stability result. And, before that we need a discrete maximum principle

2.1. Discrete Maximum Principle.

Theorem 6. *Suppose v solves (1.2) and $f(x) \leq 0$ for $x \in S_I$ then*

$$(2.1) \quad \max_{x \in S} v(x) \leq \max_{x \in S_b} g(x).$$

On the other hand, if $f(x) \geq 0$ for all $x \in S_I$ then

$$(2.2) \quad \min_{x \in S} v(x) \geq \min_{x \in S_b} g(x).$$

PROOF. We will prove (2.1) and leave the proof (2.2) to the reader. Let $M = \max_{x \in S} v(x)$. First suppose that then $\max_{x \in S_I} v(x) < M$ then (2.1) holds.

On the other hand, suppose there that there exists $x \in S_I$ such that $v(x) = M$. Then from (1.2a) we have

$$\begin{aligned} M &= v(x) \\ &= \frac{1}{4}(v(x + he_1) + v(x + he_2) + v(x - he_1) + v(x - he_2)) + h^2 f(x) \\ &\leq \frac{1}{4}(v(x + he_1) + v(x + he_2) + v(x - he_1) + v(x - he_2)) \leq M \end{aligned}$$

since $f(x) \leq 0$ and by our hypothesis $\frac{1}{4}(v(x + he_1) + v(x + he_2) + v(x - he_1) + v(x - he_2)) \leq M$. Hence, we must have

$$M = \frac{1}{4}(v(x + he_1) + v(x + he_2) + v(x - he_1) + v(x - he_2)),$$

or that

$$\frac{1}{4}(M - v(x + he_1)) + \frac{1}{4}(M - v(x + he_2)) + \frac{1}{4}(M - v(x - he_1)) + \frac{1}{4}(M - v(x - he_2)) = 0.$$

By our hypothesis each term on the left hand side is non-negative as so it must be $v(x) = v(x + he_1) = v(x + he_2) = v(x - he_1) = v(x - he_2) = M$. We can repeat the argument for the $v(x + he_1), v(x + he_2), v(x - he_1), v(x - he_2)$ and have all its neighbors also equal to M . If we repeat this process we will arrive at $v \equiv M$ is a constant discrete function. Therefore, (2.1) trivially holds. \square

2.2. Discrete Stability Result.

Theorem 7. *Suppose v solves (1.2) then*

$$(2.3) \quad \max_{x \in S} |v(x)| \leq \max_{x \in S_b} |g(x)| + \frac{1}{8} \max_{x \in S_I} |f(x)|.$$

PROOF. Let $\phi(x) = \frac{Q}{4}(x_1 - 1/2)^2 + (x_2 - 1/2)^2$ where $Q = \max_{x \in S_I} |f(x)|$. The it is not difficult to show that $\Delta_h \phi(x) = Q$ for all $x \in S_I$. Then we define $w \in P_h$ as follows $w(x) = v(x) + \phi(x)$. We see that w solves \square

$$(2.4a) \quad -\Delta_h w(x) = f(x) - Q \quad \text{for all } x \in S_I$$

$$(2.4b) \quad w(x) = g(x) + \phi(x) \quad \text{for all } x \in S_b.$$

Since $f(x) - Q \leq 0$ for all $x \in S_I$ we can apply (2.1) to get

$$v(x) \leq v(x) + \phi(x) \leq \max_{y \in S_b} (g(y) + \phi(y)) \quad \text{for all } x \in S.$$

However,

$$\max_{y \in S_b} (g(y) + \phi(y)) \leq \max_{y \in S_b} |g(y)| + \max_{y \in S_b} |\phi(y)| \leq \max_{y \in S_b} |g(y)| + Q/8.$$

Hence, we have

$$v(x) \leq \max_{y \in S_b} |g(y)| + Q/8. \quad \text{for all } x \in S.$$

Similarly, we can show that

$$-v(x) \leq \max_{y \in S_b} |g(y)| + Q/8. \quad \text{for all } x \in S.$$

Hence, we have

$$|v(x)| \leq \max_{y \in S_b} |g(y)| + Q/8. \quad \text{for all } x \in S,$$

which proves the result.

2.3. Error Estimate. We can prove the desired error estimate.

Theorem 8. Suppose v solves (1.2) and suppose that $u \in C^4(\bar{\Omega})$ solves (0.1) then we have

$$\max_{x \in S} |u(x) - v(x)| \leq C h^2 \max\{\|\partial_{x_1}^4 u\|_{C(\Omega)}, \|\partial_{x_2}^4 u\|_{C(\Omega)}\}$$

where C is independent of h and u

PROOF. Let $w \in P_h$ and let $w(x) = u(x) - v(x)$ for all $x \in S$.

Then, w solves

$$\begin{aligned} -\Delta_h w(x) &= -\Delta_h u(x) + f(x) & \text{for all } x \in S_I \\ w(x) &= 0 & \text{for all } x \in S_b. \end{aligned}$$

Therefore, we get by (2.3) we get

$$\max_{x \in S} |w(x)| \leq \frac{1}{8} \max_{x \in S_I} |-\Delta_h u(x) + f(x)| = \frac{1}{8} \max_{x \in S_I} |\Delta u(x) - \Delta_h u(x)|.$$

The result now follows from (1.1). □

Finite Element Methods for second order elliptic Problems

In this chapter we study finite element methods for second order elliptic problems.

1. Function Spaces

We first need to study some function spaces. We let $\Omega \subset \mathbb{R}^2$ be a bounded domain with Lipschitz boundary. We recall the space of m -th order continuous functions:

$C^m(\Omega) = \{v : \text{all the partial derivatives of order less than or equal to } m \text{ of } v \text{ are continuous in } \Omega\}$.

$$C_c^m(\Omega) = \{v \in C^m(\Omega) : v \text{ vanishes outside of } \Omega_0 \subset\subset \Omega\}.$$

Here Ω_0 is a compact domain and $\Omega_0 \subset\subset \Omega$ means $\partial\Omega_0 \cap \partial\Omega = \emptyset$.

We will need to make sense of derivatives of functions that are not differentiable everywhere. For example, the function $u(x) = |x|$. To do that we need to define the *weak derivative* of a function if it exists. This is done via integration by parts.

Let us recall integration by parts formula for smooth functions. Let $\mathbf{n} = (\mathbf{n}_1, \mathbf{n}_2)$ be the unit outward pointing normal to Ω . Suppose that $u, v \in C^\infty(\mathbb{R}^2)$

$$\int_{\Omega} \partial_{x_i} u(x)v(x)dx = - \int_{\Omega} u(x)\partial_{x_i} v(x)dx - \int_{\partial\Omega} u(x)v(x)n_i ds \quad \text{for } i = 1, 2.$$

Let u be an integrable function on Ω then we say that $\partial_{x_i} u$ exists if there is an integrable function ϕ such that

$$\int_{\Omega} \phi(x)v(x)dx = - \int_{\Omega} u(x)\partial_{x_i} v(x)dx \text{ for all } v \in C_c^\infty(\Omega).$$

in which case we define $\partial_{x_i} u(x) = \phi(x)$. In other words, if $\partial_{x_i} u$ then

$$(1.1) \quad \int_{\Omega} \partial_{x_i} u(x)v(x)dx = - \int_{\Omega} u(x)\partial_{x_i} v(x)dx \text{ for all } v \in C_c^\infty(\Omega).$$

Or course $\partial_{x_i} u(x)$ are uniquely defined up to a set of measure zero.

We will also need to define the following function spaces

$$L^2(\Omega) = \{v : v \text{ is integrable on } \Omega \text{ and } \int_{\Omega} v^2(x)dx < \infty\}.$$

We define the Sobolev space H^m as follows

$$H^m(\Omega) = \{v \in L^2(\Omega) : \partial_{x_1}^{\alpha_1} \partial_{x_2}^{\alpha_2} u \in L^2(\Omega), \text{ for all } \alpha_1 + \alpha_2 \leq m\}.$$

The associated norm is given as follows

$$\|u\|_{H^m(\Omega)}^2 = \sum_{j=0}^m |u|_{H^j(\Omega)}^2 \text{ where } |u|_{H^j(\Omega)}^2 = \sum_{\alpha_1 + \alpha_2 = j} \|\partial_{x_1}^{\alpha_1} \partial_{x_2}^{\alpha_2} u\|_{L^2(\Omega)}^2.$$

2. Linear second-order elliptic problems

Recall that $\operatorname{div} F(x) = \partial_{x_1} F_1(x) + \partial_{x_2} F_2(x)$ where $F(x) = (F_1(x), F_2(x))$. We will study second order elliptic problems of the form

$$(2.1a) \quad -\operatorname{div} (A(x)\nabla u(x)) = f(x) \quad \text{for } x \in \Omega$$

$$(2.1b) \quad u(x) = g(x) \quad \text{for } x \in \partial\Omega.$$

Here we assume that $A(x) \in \mathbb{R}^{2 \times 2}$ is a symmetric matrix for every $x \in \bar{\Omega}$ and is smooth on $\bar{\Omega}$. Notice that if $A(x)$ is the identity matrix for all x then $\operatorname{div} (A(x)\nabla u(x)) = -\Delta u(x)$. We can write out $\operatorname{div} (A(x)\nabla u(x))$

$$\operatorname{div} (A(x)\nabla u(x)) = \partial_{x_1} (A_{11}(x)\partial_{x_1} u(x) + A_{12}(x)\partial_{x_2} u(x)) + \partial_{x_2} (A_{21}(x)\partial_{x_1} u(x) + A_{22}(x)\partial_{x_2} u(x)).$$

In particular, if A is a constant matrix independent of x we have

$$\operatorname{div} (A(x)\nabla u(x)) = A_{11}\partial_{x_1}^2 u(x) + 2A_{12}\partial_{x_1}\partial_{x_2} u(x) + A_{22}\partial_{x_2}^2 u(x).$$

We assume that A is uniformly elliptic. That, is there exists a constant $\theta > 0$ such that

$$(2.2) \quad y^t A(x)y \geq \theta y^t y \quad \text{for all } y \in \mathbb{R}^2 \text{ and } x \in \Omega.$$

One can easily show that this is equivalent to $A(x)$ having uniformly positive eigenvalues on Ω . Since A is smooth we also have that there exists a constant $\gamma > 0$ such that

$$(2.3) \quad y^t A(x)z \leq \gamma |y||z| \quad \text{for all } y, z \in \mathbb{R}^2 \text{ and } x \in \Omega.$$

2.1. Weak Formulation of second-order problems. We define $H_0^1(\Omega) = \{v \in H^1(\Omega) : v \text{ vanishes on } \partial\Omega\}$. If we use the integration by parts formula (1.1) we can derive Gauss' integral formula

$$\int_{\Omega} \operatorname{div} F(x)v(x)dx = - \int_{\Omega} F(x) \cdot \nabla v(x)dx + \int_{\partial\Omega} F(x) \cdot \mathbf{nv}(\mathbf{x})d\mathbf{s}.$$

If we apply this to (2.1a) we have

$$\int_{\Omega} A(x)\nabla u(x) \cdot \nabla v(x)dx = \int_{\Omega} f(x)v(x)dx \quad \text{for all } v \in H_0^1(\Omega).$$

Hence, we can state the weak formulation of (2.1) as follows Find $u \in H^1(\Omega)$ with $u = g$ on $\partial\Omega$ such that

$$(2.4) \quad \int_{\Omega} A(x)\nabla u(x) \cdot \nabla v(x)dx = \int_{\Omega} f(x)v(x)dx \quad \text{for all } v \in H_0^1(\Omega).$$

We use the notation

$$(2.5) \quad a(u, v) = \int_{\Omega} A(x)\nabla u(x) \cdot \nabla v(x)dx.$$

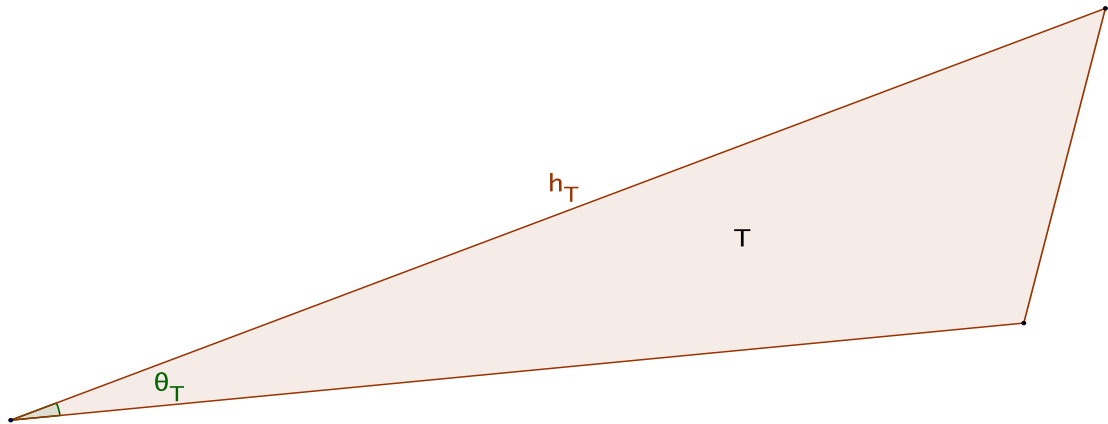
and

$$(f, u) = \int_{\Omega} f(x)v(x)dx.$$

Hence, we can write the weak form as follows Find $u \in H^1(\Omega)$ with $u = g$ on $\partial\Omega$ such that

$$(2.6) \quad a(u, v) = (f, v) \quad \text{for all } v \in H_0^1(\Omega).$$

The nice property of the weak formulation is that solutions could exist even though they do not have second derivatives. They only require that $u \in H^1(\Omega)$. Therefore, the weak formulation is a generalization of (2.1).

FIGURE 1. Illustration of θ_T, h_T

3. Finite Element Formulation

In this section we develop finite element method for (2.6). To do that we need to define some concepts. For simplicity we will assume that Ω is a polygonal domain. We will assume that we have a family of triangulations of Ω of $\{\mathcal{T}_h\}$. For every h we assume that

$$\bar{\Omega} = \cup_{T \in \mathcal{T}_h} \bar{T}.$$

If $T, K \in \mathcal{T}_h$ and $T \neq K$ then $\bar{T} \cap \bar{K}$ is either empty, a vertex or an edge of the triangulation.

We assume that the mesh is *shape regular*: There exists a constant $\kappa > 0$ such that

$$(3.1) \quad \theta_T \geq \kappa \quad \text{for all } T \in \{\mathcal{T}_h\},$$

where θ_T denotes the smallest angle of the triangle T .

We also also define

$$h_T = \text{diam}(T),$$

and

$$h = \sup_{T \in \mathcal{T}_h} h_T.$$

We now define finite element spaces:

$$V_h = \{v \in C(\Omega) : v|_T \in P^1(T) \text{ for all } T \in \mathcal{T}_h\},$$

where $P^1(T)$ is the space of affine function defined on T . We also define

$$V_h^0 = \{v \in V_h, v \equiv 0 \text{ on } \partial\Omega\}.$$

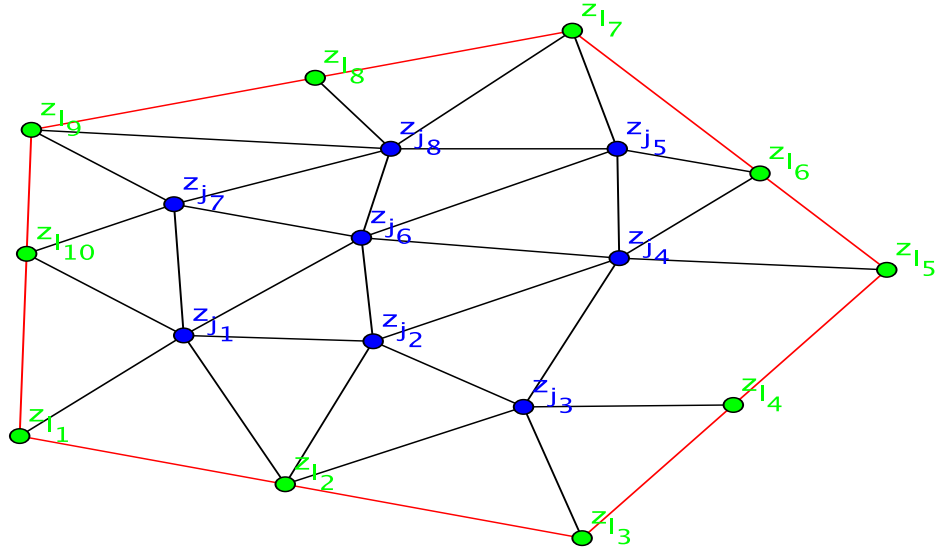


FIGURE 2. Example of a triangulation of a polygonal domain the vertices S_b are in green, and S_I in blue

We can also define a basis for the spaces V_h . Let us denote by

$$\begin{aligned} S &= \{\text{set of all vertices of } \mathcal{T}_h\} \\ S_I &= \{x \in S : x \notin \partial\Omega\} \\ S_b &= S \setminus S_I. \end{aligned}$$

Let us enumerate the vertices in S , S_I , S_b :

$$\begin{aligned} S &= \{z_1, z_2, \dots, z_M\} \\ S_I &= \{z_{j_1}, z_{j_2}, \dots, z_{j_Q}\} \\ S_b &= \{z_{\ell_1}, z_{\ell_2}, \dots, z_{\ell_R}\} \end{aligned}$$

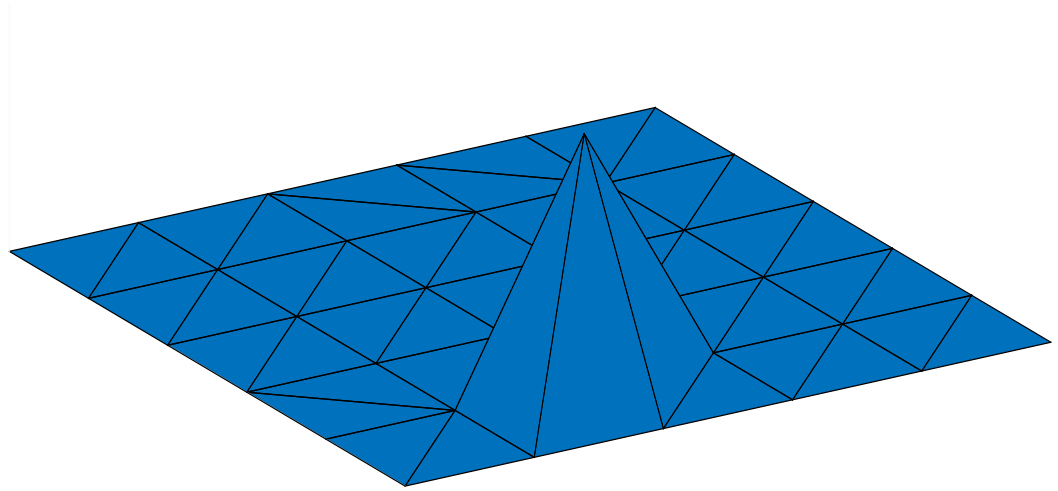
here $1 \leq j_1 \leq \dots \leq j_Q \leq M$, $1 \leq \ell_1 \leq \dots \leq \ell_R \leq M$, and $Q + R = M$.

For each $i = 1, \dots, M$ we define $\psi_i \in V_h$

$$\psi_i(z_j) = \begin{cases} 1 & \text{if } j = i \\ 0 & \text{if } j \neq i \end{cases}$$

Then we see that if $v \in V_h$ we have

$$v(x) = \sum_{1 \leq i \leq M} v(z_i) \psi_i(x) \quad \text{for all } x \in \bar{\Omega}.$$

FIGURE 3. An example of an ψ_i basis function

The finite element method solves: Find $u_h \in V_h$ with $u(z) = g(z)$ for all $z \in S_b$ such that

$$(3.2) \quad \int_{\Omega} A(x) \nabla u_h(x) \cdot \nabla v_h(x) dx = \int_{\Omega} f(x) v_h(x) dx \quad \text{for all } v_h \in V_h^0.$$

or using the $a(\cdot, \cdot)$ notation Find $u_h \in V_h$ with $u(z) = g(z)$ for all $z \in S_b$ such that

$$(3.3) \quad a(u_h, v_h) = (f, v) \quad \text{for all } v_h \in V_h^0.$$

We can easily show, using that $a(u_h, \cdot)$, (f, \cdot) are linear forms that

$$(3.4) \quad a(u_h, \psi_{j_k}) = (f, \psi_{j_k}) \quad 1 \leq k \leq Q.$$

If we write

$$u_h(x) = \sum_{1 \leq i \leq M} u_h(z_i) \psi_i(x) \quad \text{for all } x \in \bar{\Omega}.$$

we have

$$(3.5) \quad \sum_{1 \leq i \leq M} u_h(z_i) a(\psi_i, \psi_{j_k}) = (f, \psi_{j_k}) \quad 1 \leq k \leq Q.$$

or using that $u_h(z) = g(z)$ for $z \in S_b$ we get:

$$(3.6) \quad \sum_{1 \leq s \leq Q} u_h(z_{j_s}) a(\psi_{j_s}, \psi_{j_k}) = (f, \psi_{j_k}) - \sum_{1 \leq s \leq R} g(z_{\ell_s}) a(\psi_{\ell_s}, \psi_{j_k}) \quad 1 \leq k \leq Q.$$

We then have the matrix system

$$\begin{bmatrix} a(\psi_{j_1}, \psi_{j_1}) & a(\psi_{j_1}, \psi_{j_2}) & \dots & a(\psi_{j_1}, \psi_{j_Q}) \\ a(\psi_{j_2}, \psi_{j_1}) & a(\psi_{j_2}, \psi_{j_2}) & \dots & a(\psi_{j_2}, \psi_{j_Q}) \\ \vdots & & \ddots & \vdots \\ a(\psi_{j_Q}, \psi_{j_1}) & \dots & & a(\psi_{j_Q}, \psi_{j_Q}) \end{bmatrix} \begin{bmatrix} u_h(z_{j_1}) \\ u_h(z_{j_2}) \\ \vdots \\ u_h(z_{j_Q}) \end{bmatrix} = \begin{bmatrix} (f, \psi_{j_1}) \\ (f, \psi_{j_2}) \\ \vdots \\ (f, \psi_{j_Q}) \end{bmatrix} + G,$$

where

$$G = \begin{bmatrix} a(\psi_{j_1}, \psi_{\ell_1}) & a(\psi_{j_1}, \psi_{\ell_2}) & \dots & a(\psi_{j_1}, \psi_{\ell_R}) \\ a(\psi_{j_2}, \psi_{\ell_1}) & a(\psi_{j_2}, \psi_{\ell_2}) & \dots & a(\psi_{j_2}, \psi_{\ell_R}) \\ \vdots & & \ddots & \vdots \\ \vdots & & & \vdots \\ \vdots & & & \vdots \\ a(\psi_{j_Q}, \psi_{\ell_1}) & \dots & & a(\psi_{j_Q}, \psi_{\ell_R}) \end{bmatrix} \begin{bmatrix} g(z_{\ell_1}) \\ g(z_{\ell_2}) \\ \vdots \\ g(z_{\ell_R}) \end{bmatrix}.$$

3.1. Implementation of FEM.

3.2. Error Analysis of FEM. In this section we perform an error analysis of the finite element method. We will see that error estimates follows easily from, coercivity and boundedness of bilinear form and from Galerkin Orthogonality of the FEM.

LEMMA 3. **Coercivity:** *It holds*

$$(3.7) \quad \theta \|\nabla v\|_{L^2(\Omega)}^2 \leq a(v, v) \quad \text{for all } v \in H^1(\Omega).$$

PROOF. Recall that

$$a(v, v) = \int_{\Omega} A(x) \nabla v(x) \cdot \nabla v(x) dx = \int_{\Omega} (\nabla v(x))^t A(x) \nabla v(x) dx$$

then the result follows from (2.2). \square

LEMMA 4. **Boundedness:** *It holds,*

$$(3.8) \quad a(v, w) \leq \gamma \|\nabla v\|_{L^2(\Omega)} \|\nabla w\|_{L^2(\Omega)} \leq \quad \text{for all } v, w \in H^1(\Omega).$$

PROOF. Similar to the previous proof the result follows from (2.3). \square

LEMMA 5. **Galerkin Orthogonality** *Let u solve (2.5) and let u_h solve (3.3) then we have*

$$(3.9) \quad a(u - u_h, v_h) = 0 \quad \text{for all } v_h \in V_h^0.$$

PROOF. Since $V_h^0 \subset H_0^1(\Omega)$ we have by (2.5)

$$a(u, v_h) = (f, v_h) \quad \text{for all } v_h \in V_h^0.$$

However, (3.3) states that

$$a(u_h, v_h) = (f, v_h) \quad \text{for all } v_h \in V_h^0.$$

Therefore, we have

$$a(u, v_h) - a(u_h, v_h) = 0 \quad \text{for all } v_h \in V_h^0.$$

Using that $a(\cdot, v_h)$ is linear we get

$$a(u, v_h) - a(u_h, v_h) = a(u - u_h, v_h) \quad \text{for all } v_h \in V_h^0.$$

The result now follows. \square

Now we can easily prove the following result.

Theorem 9. *Let u solve (2.5) and let u_h solve (3.3) then we have*

$$(3.10) \quad \|\nabla(u - u_h)\|_{L^2(\Omega)} \leq \frac{\gamma}{\theta} \min_{v \in V_h^g} \|\nabla(u - v_h)\|_{L^2(\Omega)}.$$

where $V_h^g = \{v_h \in V_h : v_h(z) = g(z) \text{ for all } z \in S_b\}$.

PROOF. By Coercivity (3.7) we obtain

$$\theta \|\nabla(u - u_h)\|_{L^2(\Omega)}^2 \leq a(u - u_h, u - u_h).$$

Now by linearity of $a(u - u_h, \cdot)$ we have

$$a(u - u_h, u - u_h) = a(u - u_h, u - v_h) - a(u - u_h, u_h - v_h) \quad \text{for all } v_h \in V_h^g.$$

Since $u_h, v_h \in V_h^g$ we have $u_h - v_h \in V_h^0$ and therefore, we get

By Galerkin orthogonality (3.9)

$$a(u - u_h, u_h - v_h) = 0 \quad \text{for all } v_h \in V_h^g.$$

Combining the above results we get

$$\theta \|\nabla(u - u_h)\|_{L^2(\Omega)}^2 \leq a(u - u_h, u - v_h) \quad \text{for all } v_h \in V_h^g.$$

Using (3.8) we get

$$\theta \|\nabla(u - u_h)\|_{L^2(\Omega)}^2 \leq \gamma \|\nabla(u - u_h)\|_{L^2(\Omega)} \|\nabla(u - v_h)\|_{L^2(\Omega)} \quad \text{for all } v_h \in V_h^g.$$

Dividing by $\|\nabla(u - u_h)\|_{L^2(\Omega)}$ we obtain

$$\theta \|\nabla(u - u_h)\|_{L^2(\Omega)} \leq \frac{\gamma}{\theta} \|\nabla(u - v_h)\|_{L^2(\Omega)} \quad \text{for all } v_h \in V_h^g.$$

The result follows by taking the minimum over $v_h \in V_h^g$. \square

The result (3.12) says that the FEM approximation u_h is almost the best approximation to u in the space V_h when measured in the $\|\nabla \cdot\|_{L^2(\Omega)}$ norm. Next, we give an approximation result in terms of the mesh size h . The proof of this result which is quite technical will be proved in the next section.

PROPOSITION 2. *If $u \in H^2(\Omega)$ then*

$$(3.11) \quad \min_{v_h \in V_h^g} \|\nabla(u - v_h)\|_{L^2(\Omega)} \leq Ch|u|_{H^2(\Omega)}.$$

Then the following corollary follows from (3.11) and (3.12).

COROLLARY 1. *Let u solve (2.5) and let u_h solve (3.3) and assume that $u \in H^2(\Omega)$, then*

$$(3.12) \quad \|\nabla(u - u_h)\|_{L^2(\Omega)} \leq Ch|u|_{H^2(\Omega)}.$$

where the constant C is independent of u and h .

We now want to study the error estimates $\|u - u_h\|_{L^2(\Omega)}$ instead. For simplicity we assume that $g \equiv 0$. To do this we need a result from PDE theory. The following is a standard H^2 regularity result.

PROPOSITION 3. *Let u solve (2.5) with $g \equiv 0$ and suppose that Ω is a convex polygon then*

$$(3.13) \quad \|u\|_{H^2(\Omega)} \leq C \|f\|_{L^2(\Omega)}.$$

The result seems very reasonable since formally $\Delta u = -f$. So, we think that some combination of second derivatives are controlled by f . The result says that all the second derivatives of u are controlled by f .

We can now state the desired estimate.

Theorem 10. *Let u solve (2.5) with $g \equiv 0$ and let u_h solve (3.3) and assume that Ω is a convex polygon, then*

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

Please note that this result says that $\|u - u_h\|_{L^2(\Omega)}$ converges to zero with one order higher than $\|\nabla(u - u_h)\|_{L^2(\Omega)}$.

PROOF. Let $\phi \in H_0^1(\Omega)$ solve the following problem

$$a(\phi, v) = (u - u_h, v) \quad \text{for all } v \in H_0^1(\Omega).$$

Then, setting $v = u - u_h$ we get

$$\|u - u_h\|_{L^2(\Omega)}^2 = a(\phi, u - u_h) = a(u - u_h, \phi).$$

If we use Galerkin Orthogonality (3.9) we get

$$\|u - u_h\|_{L^2(\Omega)}^2 = a(u - u_h, \phi - v_h) \quad \text{for all } v_h \in V_h^0.$$

By (3.8) we get

$$\|u - u_h\|_{L^2(\Omega)}^2 \leq \gamma \|\nabla(u - u_h)\|_{L^2(\Omega)} \|\nabla(\phi - v_h)\|_{L^2(\Omega)} \quad \text{for all } v_h \in V_h^0.$$

Using (3.11) and (3.13) we obtain

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

After dividing by $\|u - u_h\|_{L^2(\Omega)}$, using (3.12) and (3.13) we obtain the result. \square

3.3. Approximation of Piecewise linear functions. We are left to prove (3.11). To be precise we will use standard results that are used in the analysis of elliptic PDEs. These results are useful in many other context and not only important for the approximation theory we will give here.

We will fix a reference triangle: \hat{T} which is the triangle with vertices $(0,0)$, $(1, 0)$, $(0,1)$. We then state a few results.

The first is Poincare's inequality.

PROPOSITION 4. Poincare's inequality: *Suppose that $w \in H^1(\hat{T})$ and let $\bar{w} = \frac{1}{|\hat{T}|} \int_{\hat{T}} w(x) dx$ then,*

$$(3.14) \quad \|w - \bar{w}\|_{L^2(\hat{T})} \leq C |w|_{H^1(\hat{T})},$$

where C is independent of w .

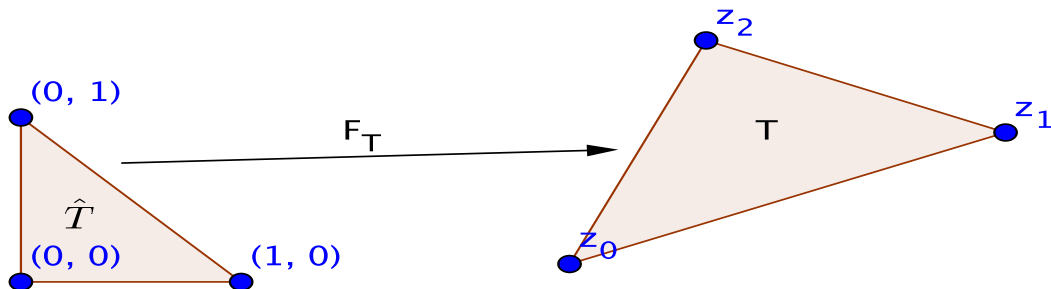
The next result is one type Sobolev inequality. There are many Sobolev embedding and inequalities.

PROPOSITION 5. A Sobolev inequality: *If $v \in H^2(\hat{T})$ then $v \in C(\overline{\hat{T}})$ and the following inequality holds*

$$(3.15) \quad \|v\|_{L^\infty(\hat{T})} \leq C \|v\|_{H^2(\hat{T})}.$$

We will be going to transfer estimates from $T \in \mathcal{T}_h$ to \hat{T} . The important point here is that \hat{T} is fixed. To do this, we need to define an affine transformation from $F_T : \hat{T} \rightarrow T$ define as

$$F_T(\hat{x}) = B_T \hat{x} + b_T.$$

FIGURE 4. An example of an ψ_i basis function

Not that if T has vertices z_0, z_1, z_2 then the first column of B_T is $z_1 - z_0$ while the second column is $z_2 - z_0$ and finally $b_T = z_0$. That is,

$$B_T = [z_1 - z_0 \quad z_2 - z_0], \quad b_T = z_0.$$

We will state a proof of the size of entries of B_T and B_T^{-1} . Here we use strongly the shape regularity, (3.1)

LEMMA 6. *There exists a constant C independent of h and $T \in \mathcal{T}_h$ such that*

$$(3.16) \quad |(B_T)_{ij}| \leq Ch_T \quad |(B_T^{-1})_{ij}| \leq Ch_T^{-1} \quad \text{for all } T \in \mathcal{T}_h, i, j = 1, 2.$$

We also use the following notation. Let w be defined as a function on T then define the function \hat{w} as function on \hat{T} given by

$$\hat{w}(\hat{x}) = w(F_T(\hat{x})).$$

Using a change of variable formula we have

$$(3.17) \quad \int_T w(x) dx = \int_{\hat{T}} w(F_T(\hat{x})) \det(B_T) d\hat{x}.$$

or

$$(3.18) \quad \int_T w(x) dx = \int_{\hat{T}} \hat{w}(\hat{x}) \det(B_T) d\hat{x}.$$

Now notice that using the chain rule we have

$$(3.19) \quad \nabla_x w(x) = \nabla_x \hat{w}(F_T^{-1}(x)) = B_T^{-t} \nabla_{\hat{x}} \hat{w}(F_T^{-1}(x)) = B_T^{-t} \nabla_{\hat{x}} \hat{w}(\hat{x}),$$

where $x = F_T(\hat{x})$. Note that we write ∇_x to denote the gradient with respect to the variable x . Whereas, $\nabla_{\hat{x}}$ is gradient with respect to \hat{x} .

Hence using (3.17) and (3.19) we obtain

$$(3.20) \quad \int_T \nabla_x w(x) \cdot \nabla_x v(x) dx = \int_{\hat{T}} (B_T^{-t} \nabla_{\hat{x}} \hat{w}(\hat{x})) \cdot (B_T^{-t} \nabla_{\hat{x}} \hat{v}(\hat{x})) \det(B_T) d\hat{x}.$$

We can then prove the following results

LEMMA 7. *It holds,*

$$(3.21) \quad \frac{c_1}{h_T} \|w\|_{L^2(T)} \leq \|\hat{w}\|_{L^2(\hat{T})} \leq \frac{c_2}{h_T} \|w\|_{L^2(T)},$$

$$(3.22) \quad c_1 |w|_{H^1(T)} \leq |\hat{w}|_{H^1(\hat{T})} \leq c_2 |w|_{H^1(T)},$$

$$(3.23) \quad c_1 h_T |w|_{H^2(T)} \leq |\hat{w}|_{H^2(\hat{T})} \leq c_2 h_T \|w\|_{L^2(T)}.$$

PROOF. We only prove (3.22). The proof of the other two proofs follow a similar line of argument. By (3.20) we have

$$|w|_{H^1(T)}^2 = \|\nabla w\|_{L^2(T)}^2 = \det B_T \|B_T^{-t} \nabla \hat{w}\|_{L^2(\hat{T})}^2.$$

By (3.16) we have

$$\det B_T \leq C h_T^2,$$

and

$$\|B_T^{-t} \nabla \hat{w}\|_{L^2(\hat{T})}^2 \leq C h_T^{-2} \|\nabla \hat{w}\|_{L^2(\hat{T})}^2.$$

The result now follows. \square

We can now state a local approximation result known as the Bramble-Hilbert lemma. The result will follow

LEMMA 8. **Bramble–Hilbert:** *Let $w \in H^2(T)$ then there exists a $v \in P^1(T)$ such that*

$$\|w - v\|_{L^2(T)} + h_T \|\nabla(w - v)\|_{L^2(T)} \leq C h_T^2 |w|_{H^2(T)}.$$

PROOF. Let $v(x) = c_0 + c_1 x_1 + c_2 x_2$ where $x = (x_1, x_2)$. We define first c_1, c_2

$$c_i = \frac{1}{|T|} \int_T \partial_{x_i} w(x) dx \quad \text{for all } i = 1, 2.$$

Note that since $c_i = \partial_{x_i} v(x)$ we have that

$$(3.24) \quad \frac{1}{|T|} \int_T \partial_{x_i} v(x) dx = \frac{1}{|T|} \int_T \partial_{x_i} w(x) dx \quad \text{for all } i = 1, 2.$$

Now that we have define c_1, c_2 we define c_0 so that

$$(3.25) \quad \frac{1}{|T|} \int_T v(x) dx = \frac{1}{|T|} \int_T w(x) dx.$$

That is,

$$c_0 = \frac{1}{|T|} \int_T (w(x) - (c_1 x_1 + c_2 x_2)) dx.$$

If we define $q(x) = w(x) - v(x)$ and consider $\hat{q}(\hat{x})$. Note that by (3.24)

$$\int_T \partial_{x_i} q(x) = 0 \quad \text{for } i = 1, 2$$

Again, using (3.17) and (3.19) we get

$$\det B_T ((B_T^{-t})_{i1} \int_{\hat{T}} \partial_{\hat{x}_1} \hat{q}(\hat{x}) d\hat{x} + (B_T^{-t})_{i2} \int_{\hat{T}} \partial_{\hat{x}_2} \hat{q}(\hat{x}) d\hat{x}) = \int_T \partial_{x_i} q(x) dx,$$

for $i = 1, 2$. Hence, we get

$$(\det B_T) B_T^{-1} \begin{bmatrix} \int_T \partial_{\hat{x}_1} \hat{q}(\hat{x}) d\hat{x} \\ \int_T \partial_{\hat{x}_2} \hat{q}(\hat{x}) d\hat{x} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Hence, we get

$$\int_{\hat{T}} \partial_{\hat{x}_i} \hat{q}(\hat{x}) d\hat{x} = 0 \quad \text{for } i = 1, 2$$

Similarly, we have

$$\int_T \hat{q}(\hat{x}) d\hat{x} = 0.$$

Thus, $\overline{\partial_{\hat{x}_i} \hat{q}} \equiv 0$ and $\bar{q} \equiv 0$ and by Poincare's inequality we get

$$|\hat{q}|_{H^1(\hat{T})}^2 = \|\partial_{\hat{x}_1} \hat{q}\|_{L^2(\hat{T})}^2 + \|\partial_{\hat{x}_2} \hat{q}\|_{L^2(\hat{T})}^2 \leq C \sum_{\alpha_1 + \alpha_2 = 2} \|\partial_{\hat{x}_1}^{\alpha_1} \partial_{\hat{x}_2}^{\alpha_2} \hat{q}\|_{L^2(\hat{T})}^2 = C |\hat{q}|_{H^2(\hat{T})}^2.$$

Using this result with (3.22) gives

$$|q|_{H^1(T)} \leq C |\hat{q}|_{H^1(\hat{T})} \leq C |\hat{q}|_{H^2(\hat{T})}^2 \leq h_T |q|_{H^2(T)}.$$

Also, using (3.21), Poincare's inequality and the previous inequality we get

$$\|q\|_{L^2(T)} \leq Ch_T \|\hat{q}\|_{L^2(\hat{T})} \leq Ch_T |\hat{q}|_{H^1(\hat{T})} \leq Ch_T \|q\|_{L^2(T)} \leq Ch_T^2 |q|_{H^2(T)}.$$

□

We will also need an inverse estimate. We can prove this result with more elementary results. We give a proof that is generalizable to other finite element spaces.

LEMMA 9. Inverse inequality : *Let $T \in \mathcal{T}_h$ then for every $v \in P^1(T)$*

$$|v|_{H^1(T)} \leq \frac{C}{h_T} \|v\|_{L^2(T)}$$

where the constant C is independent of T and v .

PROOF. By (3.22)

$$(3.26) \quad |v|_{H^1(T)} \leq C |\hat{v}|_{H^1(\hat{T})}$$

Note that $\hat{v} \in P^1(\hat{T})$. Since $P^1(\hat{T})$ is a finite dimensional space we know that norms are equivalent. Hence,

$$|\hat{v}|_{H^1(\hat{T})} \leq \hat{C} \|\hat{v}\|_{L^2(\hat{T})}.$$

where the constant \hat{C} is independent \hat{v} . Note that we are using that \hat{T} is fixed. Using (3.21) we get

$$\|\hat{v}\|_{L^2(\hat{T})} \leq \frac{C}{h_T} \|v\|_{L^2(T)}.$$

Combining the above inequalities gives the result. □

We can state a simple corollary.

COROLLARY 2. *Let $T \in \mathcal{T}_h$ with vertices z_i ($i = 0, 1, 2$) then for every $v \in P^1(T)$*

$$(3.27) \quad |v|_{H^1(T)} \leq C \max\{v(z_0), v(z_1), v(z_2)\}.$$

PROOF. We easily see that

$$\|v\|_{L^2(T)}^2 = \int_T v^2(x) dx \leq |T| \|v\|_{L^\infty(T)}^2.$$

Hence,

$$\|v\|_{L^2(T)} \leq \sqrt{|T|} \|v\|_{L^\infty(T)}.$$

Since the maximum of $|v|$ occurs at the vertices and so we have

$$\|v\|_{L^2(T)} \leq \sqrt{|T|} \max\{v(z_0), v(z_1), v(z_2)\}.$$

Using that the mesh is shape regular (3.1) we easily have that $\sqrt{|T|} \leq Ch_T$. Hence, combining the above inequality with (3.26). \square

In order to prove our main result (3.11) we need to define the interpolant operator.

DEFINITION 3.1. For $w \in C(\bar{\Omega})$ we define $I_h w = V_h$ by

$$I_h w(z) = w(z) \quad \text{for all } z \in S.$$

3.4. Proof of (3.11). First we see that $I_h u \in V_h^g$ and so

$$\inf_{v \in V_h^g} \|\nabla(u - v)\|_{L^2(\Omega)} \leq \|\nabla(u - I_h u)\|_{L^2(\Omega)}.$$

Let us write

$$(3.28) \quad \|\nabla(u - I_h u)\|_{L^2(\Omega)}^2 = \sum_{T \in \mathcal{T}_h} \|\nabla(u - I_h u)\|_{L^2(T)}^2.$$

From the Bramble-Hilbert Lemma we have there exists $q \in P^1(T)$ so that

$$(3.29) \quad \|u - q\|_{L^2(T)} + h_T \|\nabla(u - q)\|_{L^2(T)} \leq C h_T^2 |u|_{H^2(T)}.$$

Using the triangle inequality we get

$$(3.30) \quad \|\nabla(u - I_h u)\|_{L^2(T)} \leq C(\|\nabla(u - q)\|_{L^2(T)} + \|\nabla(q - I_h u)\|_{L^2(T)}).$$

We let $w = (q - I_h u)|_T$ and use (3.27) to get

$$\|\nabla(q - I_h u)\|_{L^2(T)} = \|\nabla w\|_{L^2(T)} \leq C \max\{w(z_0), w(z_1), w(z_2)\},$$

where z_i ($i = 0, 1, 2$) are the vertices of T . However, $w(z_i) = q(z_i) - I_h u(z_i) = q(z_i) - u(z_i)$ by the definition of I_h . Hence,

$$\max\{w(z_0), w(z_1), w(z_2)\} \leq \|u - q\|_{L^\infty(T)}.$$

Therefore, we have

$$(3.31) \quad \|\nabla(q - I_h u)\|_{L^2(T)} \leq C \|u - q\|_{L^\infty(T)}.$$

Clearly we have

$$\|u - q\|_{L^\infty(T)} \leq \|\hat{u} - \hat{q}\|_{L^\infty(\hat{T})}.$$

By (3.15) and definition of H^2 -norm

$$\|\hat{u} - \hat{q}\|_{L^\infty(\hat{T})} \leq C \|\hat{u} - \hat{q}\|_{H^2(\hat{T})} \leq C(|\hat{u} - \hat{q}|_{H^2(\hat{T})} + |\hat{u} - \hat{q}|_{H^1(\hat{T})} + \|\hat{u} - \hat{q}\|_{L^2(\hat{T})}).$$

Then, combining the above inequalities and using (3.21), (3.22), (3.23) we have

$$\|u - q\|_{L^\infty(T)} \leq C(h_T |u - q|_{H^2(T)} + |u - q|_{H^1(T)} + \frac{1}{h_T} \|u - q\|_{L^2(T)}).$$

If we combine this inequality with this inequality with (3.31) and (3.30) we get

$$\|\nabla(u - I_h u)\|_{L^2(T)} \leq C(h_T |u - q|_{H^2(T)} + |u - q|_{H^1(T)} + \frac{1}{h_T} \|u - q\|_{L^2(T)}).$$

Using (3.29) we get

$$\|\nabla(u - I_h u)\|_{L^2(T)} \leq C h_T |u|_{H^2(T)},$$

where we also used that the second derivatives of $q \in P^1(T)$ are zero since. Hence, using (3.28) we get

$$\|\nabla(u - I_h u)\|_{L^2(\Omega)}^2 \leq C h^2 \sum_{T \in \mathcal{T}_h} |u|_{H^2(T)}^2 = C h^2 |u|_{H^2(\Omega)}^2.$$

The result now follows after taking the square root of both sides.

Bibliography

- [1] C. Johnson, *Numerical solution of partial differential equations by the finite element method*, Cambridge University Press, Cambridge, 1987.
- [2] D. Gilbarg and N. Trudinger, *Elliptic Partial Differential Equations of Second Order*, 2001, Springer.