

MATH 8445, University of Minnesota  
Numerical Analysis of Differential Equations

**Lecture notes on  
Numerical Analysis of  
Partial Differential Equations**

– version of 2012-12-15 –

Douglas N. Arnold



## Contents

Chapter 1. Introduction	1
1. Basic examples of PDEs	1
1.1. Heat flow and the heat equation	1
1.2. Elastic membranes	3
1.3. Elastic plates	3
2. Some motivations for studying the numerical analysis of PDE	4
Chapter 2. The finite difference method for the Laplacian	7
1. The 5-point difference operator	7
2. Analysis via a maximum principle	10
3. Consistency, stability, and convergence	11
4. Fourier analysis	13
5. Analysis via summation by parts	15
6. Extensions	17
6.1. Curved boundaries	17
6.2. More general PDEs	20
6.3. More general boundary conditions	21
6.4. Nonlinear problems	21
6.5. Three dimensions	21
Chapter 3. Linear algebraic solvers	23
1. Classical iterations	23
2. The conjugate gradient method	29
2.1. Line search methods and the method of steepest descents	29
2.2. The conjugate gradient method	31
2.3. Preconditioning	38
3. Multigrid methods	40
Chapter 4. Finite element methods for elliptic equations	49
1. Weak and variational formulations	49
2. Galerkin method and finite elements	50
3. Lagrange finite elements	51
4. Finite element assembly	54
5. Coercivity, inf-sup condition, and well-posedness	55
5.1. The symmetric coercive case	56
5.2. The coercive case	57
5.3. The inf-sup condition	58
6. Stability, consistency, and convergence	58

7.	Finite element approximation theory	59
8.	Error estimates for finite elements	65
8.1.	Estimate in $H^1$	65
8.2.	Estimate in $L^2$	66
9.	A posteriori error estimates and adaptivity	67
9.1.	The Clément interpolant	67
9.2.	The residual and the error	70
9.3.	Estimating the residual	71
9.4.	A posteriori error indicators and adaptivity	73
9.5.	Examples of adaptive finite element computations	73
9.6.	Nonlinear problems	75
Chapter 5.	Time-dependent problems	81
1.	Finite difference methods for the heat equation	81
1.1.	Forward differences in time	82
1.2.	Backward differences in time	84
1.3.	Fourier analysis	85
1.4.	Crank–Nicolson	85
2.	Finite element methods for the heat equation	86
2.1.	Analysis of the semidiscrete finite element method	87
2.2.	Analysis of a fully discrete finite element method	89

## CHAPTER 1

### Introduction

Galileo wrote that the great book of nature is written in the language of mathematics. The most precise and concise description of many physical systems is through partial differential equations.

#### 1. Basic examples of PDEs

**1.1. Heat flow and the heat equation.** We start with a typical physical application of partial differential equations, the modeling of heat flow. Suppose we have a solid body occupying a region  $\Omega \subset \mathbb{R}^3$ . The temperature distribution in the body can be given by a function  $u : \Omega \times J \rightarrow \mathbb{R}$  where  $J$  is an interval of time we are interested in and  $u(x, t)$  is the temperature at a point  $x \in \Omega$  at time  $t \in J$ . The heat content (the amount of thermal energy) in a subbody  $D \subset \Omega$  is given by

$$\text{heat content of } D = \int_D cu \, dx$$

where  $c$  is the product of the specific heat of the material and the density of the material. Since the temperature may vary with time, so can the heat content of  $D$ . The change of heat energy in  $D$  from a time  $t_1$  to a time  $t_2$  is given by

$$\begin{aligned} \text{change of heat in } D &= \int_D cu(x, t_2) \, dx - \int_D cu(x, t_1) \, dx \\ &= \int_{t_1}^{t_2} \frac{\partial}{\partial t} \int_D cu \, dx \, dt = \int_{t_1}^{t_2} \int_D \frac{\partial(cu)}{\partial t}(x, t) \, dx \, dt. \end{aligned}$$

Now, by conservation of energy, any change of heat in  $D$  must be accounted for by heat flowing in or out of  $D$  through its boundary or by heat entering from external sources (e.g., if the body were in a microwave oven). The heat flow is measured by a vector field  $\sigma(x, t)$  called the *heat flux*, which points in the direction in which heat is flowing with magnitude the rate energy flowing across a unit area per unit time. If we have a surface  $S$  embedded in  $D$  with normal  $n$ , then the heat flowing across  $S$  in the direction pointed to by  $n$  in unit time is  $\int_S \sigma \cdot n \, ds$ . Therefore the heat that flows out of  $D$ , i.e., across its boundary, in the time interval  $[t_1, t_2]$ , is given by

$$\text{heat flow out of } D = \int_{t_1}^{t_2} \int_{\partial D} \sigma \cdot n \, ds \, dt = \int_{t_1}^{t_2} \int_D \operatorname{div} \sigma \, dx \, dt,$$

where we have used the divergence theorem. We denote the heat entering from external sources by  $f(x, t)$ , given as energy per unit volume per unit time, so that  $\int_{t_1}^{t_2} \int_D f(x, t) \, dx \, dt$

gives amount external heat added to  $D$  during  $[t_1, t_2]$ , and so conservation of energy is expressed by the equation

$$(1.1) \quad \int_{t_1}^{t_2} \int_D \frac{\partial(cu)}{\partial t}(x, t) dx dt = - \int_{t_1}^{t_2} \int_D \operatorname{div} \sigma ds dt + \int_{t_1}^{t_2} \int_D f(x, t) dx dt,$$

for all subbodies  $D \subset \Omega$  and times  $t_1, t_2$ . Thus the quantity

$$\frac{\partial(cu)}{\partial t} + \operatorname{div} \sigma - f$$

must vanish identically, and so we have established the differential equation

$$\frac{\partial(cu)}{\partial t} = - \operatorname{div} \sigma + f, \quad x \in \Omega, \forall t.$$

To complete the description of heat flow, we need a *constitutive equation*, which tells us how the heat flux depends on the temperature. The simplest is Fourier's law of heat conduction, which says that heat flows in the direction opposite the temperature gradient with a rate proportional to the magnitude of the gradient:

$$\sigma = -\lambda \operatorname{grad} u,$$

where the positive quantity  $\lambda$  is called the conductivity of the material. (Usually  $\lambda$  is just a scalar, but if the material is thermally anisotropic, i.e., it has preferred directions of heat flow, as might be a fibrous or laminated material,  $\lambda$  can be a  $3 \times 3$  positive-definite matrix.) Therefore we have obtained the equation

$$\frac{\partial(cu)}{\partial t} = \operatorname{div}(\lambda \operatorname{grad} u) + f \text{ in } \Omega \times J.$$

The source function  $f$ , the material coefficients  $c$  and  $\lambda$  and the solution  $u$  can all be functions of  $x$  and  $t$ . If the material is homogeneous (the same everywhere) and not changing with time, then  $c$  and  $\lambda$  are constants and the equation simplifies to the *heat equation*,

$$\mu \frac{\partial u}{\partial t} = \Delta u + \tilde{f},$$

where  $\mu = c/\lambda$  and we have  $\tilde{f} = f/\lambda$ . If the material coefficients depend on the temperature  $u$ , as may well happen, we get a nonlinear PDE generalizing the heat equation.

The heat equation not only governs heat flow, but all sorts of diffusion processes where some quantity flows from regions of higher to lower concentration. The heat equation is the prototypical *parabolic* differential equation.

Now suppose our body reaches a steady state: the temperature is unchanging. Then the time derivative term drops and we get

$$(1.2) \quad - \operatorname{div}(\lambda \operatorname{grad} u) = f \text{ in } \Omega,$$

where now  $u$  and  $f$  are functions of  $f$  alone. For a homogeneous material, this becomes the Poisson equation

$$- \Delta u = \tilde{f},$$

the prototypical *elliptic* differential equation. For an inhomogeneous material we can leave the steady state heat equation in *divergence form* as in (1.2), or differentiate out to obtain

$$-\lambda \Delta u + \operatorname{grad} \lambda \cdot \operatorname{grad} u = f.$$

To determine the steady state temperature distribution in a body we need to know not only the sources and sinks within the body (given by  $f$ ), but also what is happening at the boundary  $\Gamma := \partial\Omega$ . For example a common situation is that the boundary is held at a given temperature

$$(1.3) \quad u = g \text{ on } \Gamma.$$

The PDE (1.2) together with the *Dirichlet boundary condition* (1.3) form an elliptic boundary value problem. Under a wide variety of circumstances this problem can be shown to have a unique solution. The following theorem is one example (although the smoothness requirements can be greatly relaxed).

**THEOREM 1.1.** *Let  $\Omega$  be a smoothly bounded domain in  $\mathbb{R}^n$ , and let  $\lambda : \bar{\Omega} \rightarrow \mathbb{R}_+$ ,  $f : \bar{\Omega} \rightarrow \mathbb{R}$ ,  $g : \Gamma \rightarrow \mathbb{R}$  be  $C^\infty$  functions. Then there exists a unique function  $u \in C^2(\bar{\Omega})$  satisfying the differential equation (1.2) and the boundary condition (1.3). Moreover  $u$  is  $C^\infty$ .*

Instead of the Dirichlet boundary condition of imposed temperature, we often see the Neumann boundary condition of imposed heat flux (flow across the boundary):

$$\frac{\partial u}{\partial n} = g \text{ on } \Gamma.$$

For example if  $g = 0$ , this says that the boundary is insulated. We may also have a Dirichlet condition on part of the boundary and a Neumann condition on another.

**1.2. Elastic membranes.** Consider a taut (homogeneous isotropic) elastic membrane affixed to a flat or nearly flat frame and possibly subject to a transverse force distribution, e.g., a drum head hit by a mallet. We model this with a bounded domain  $\Omega \subset \mathbb{R}^2$  which represents the undisturbed position of the membrane if the frame is flat and no force is applied. At any point  $x$  of the domain and any time  $t$ , the transverse displacement is given by  $u(x, t)$ . As long as the displacements are small, then  $u$  approximately satisfies the membrane equation

$$\rho \frac{\partial^2 u}{\partial t^2} = k \Delta u + f,$$

where  $\rho$  is the density of the membrane (mass per unit area),  $k$  is the tension (force per unit distance), and  $f$  is the imposed transverse force density (force per unit area). This is a second order hyperbolic equation, *the wave equation*. If the membrane is in steady state, the displacement satisfies the Poisson equation

$$-\Delta u = \tilde{f},$$

$$f = \tilde{f}/k.$$

**1.3. Elastic plates.** The derivation of the membrane equation depends upon the assumption that the membrane resists stretching (it is under tension), but does not resist bending. If we consider a *plate*, i.e., a thin elastic body made of a material which resists bending as well as stretching, we obtain instead the plate equation

$$\rho \frac{\partial^2 u}{\partial t^2} = -D \Delta^2 u + f,$$

where  $D$  is the bending modulus, a constant which takes into account the elasticity of the material and the thickness of the plate ( $D = Et^3/[12(1 - \nu^2)]$  where  $E$  is Young's modulus and  $\nu$  is Poisson's ratio). Now the steady state equation is the *biharmonic equation*

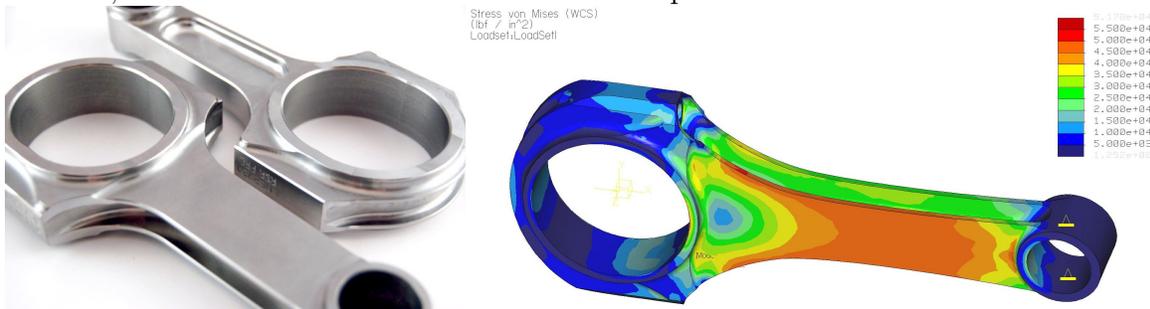
$$\Delta^2 u = \tilde{f}.$$

Later in this course we will study other partial differential equations, including the equations of elasticity, the Stokes and Navier–Stokes equations of fluid flow, and Maxwell's equations of electromagnetics.

## 2. Some motivations for studying the numerical analysis of PDE

In this course we will study algorithms for obtaining approximate solutions to PDE problems, for example, using the finite element method. Such algorithms are a hugely developed technology (we will, in fact, only skim the surface of what is known in this course), and there are thousands of computer codes implementing them. As an example of the sort of work that is done routinely, here is the result of a simulation using a finite element method to find a certain kind of force distribution, the so-called von Mises stress, engendered in a connecting rod of a Porsche race car when a certain load is applied. The von Mises stress predicts when and where the metal of the rod will deform, and was used to design the shape of the rod.

FIGURE 1.1. Connector rods designed by LN Engineering for Porsche race cars, and the stress distribution in a rod computed with finite elements.



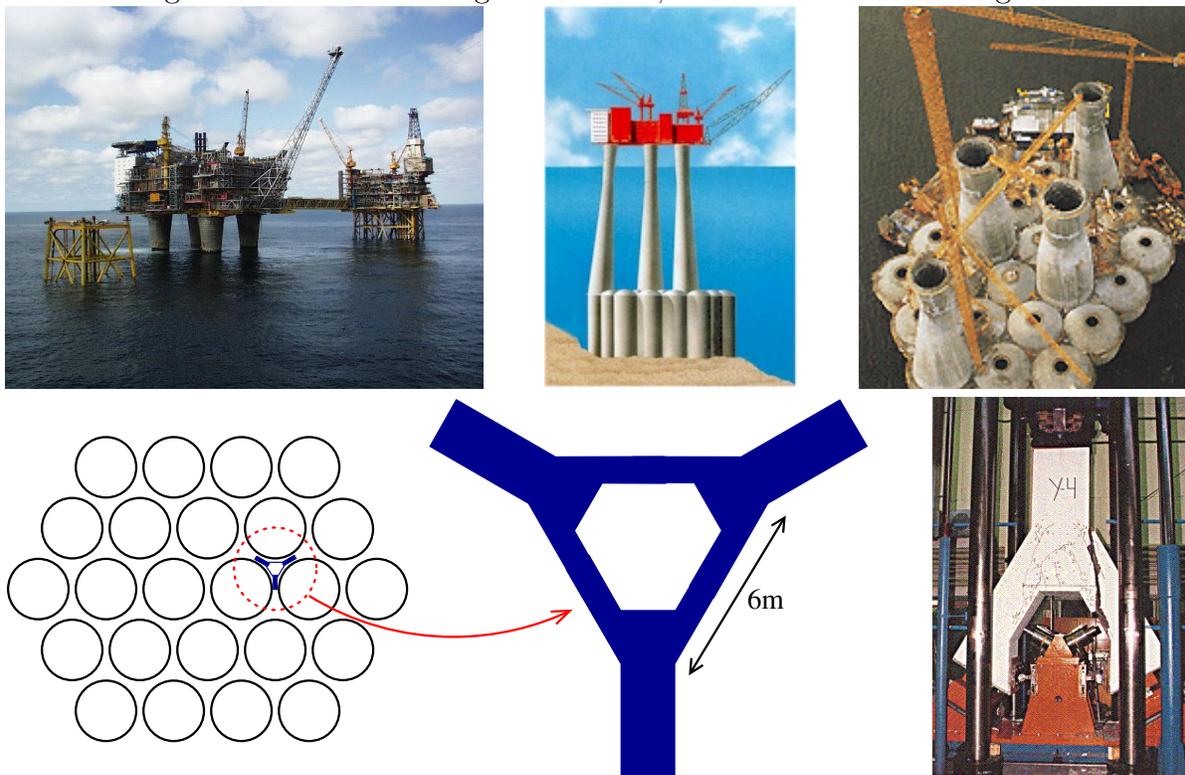
But one should not get the idea that it is straightforward to solve any reasonable PDE problem with finite elements. Not only do challenges constantly arise as practitioners seek to model new systems and solve new equations, but when used with insufficient knowledge and care, even advance numerical software can give disastrous results. A striking example is the sinking of the Sleipner A offshore oil platform in the North Sea in 1991. This occurred when the Norwegian oil company, Statoil, was slowly lowering to the sea floor an array of 24 massive concrete tanks, which would support the 57,000 ton platform (which was to accommodate about 200 people and 40,000 tons of drilling equipment). By flooding the tanks in a so-called controlled ballasting operation, they were lowered at the rate of about 5 cm per minute. When they reached a depth of about 65m the tanks imploded and crashed to the sea floor, leaving nothing but a pile of debris at 220 meters of depth. The crash did not result in loss of life, but did cause a seismic event registering 3.0 on the Richter scale, and an economic loss of about \$700 million.

An engineering research organization, SINTEF, was appointed to investigate the accident and released a sequence of 16 reports, which they summarized as follows:

*The conclusion of the investigation was that the loss was caused by a failure in a cell wall, resulting in a serious crack and a leakage that the pumps were not able to cope with. The wall failed as a result of a combination of a serious error in the finite element analysis and insufficient anchorage of the reinforcement in a critical zone.*

A better idea of what was involved can be obtained from this photo and sketch of the platform. The 24 cells and 4 shafts referred to above are shown to the left while at the sea surface. The cells are 12 meters in diameter. The cell wall failure was traced to a tricell, a triangular concrete frame placed where the cells meet, as indicated in the diagram below. To the right of the diagram is pictured a portion of tricell undergoing failure testing.

FIGURE 1.2. Top row: Offshore platform like the failed Sleipner design, diagram of structure, and concrete cells at sea surface. Bottom row: diagram showing the location and design of a tricell, and tricell failure testing.



The post accident investigation traced the error to inaccurate finite element approximation of one of the most basic PDEs used in engineering, the equations of linear elasticity, which were used to model the tricell (using the popular finite element program NASTRAN). The shear stresses were underestimated by 47%, leading to insufficient design. In particular, certain concrete walls were not thick enough. More careful finite element analysis, made after the accident, predicted that failure would occur with this design at a depth of 62m, which matches well with the actual occurrence at 65m.



## CHAPTER 2

### The finite difference method for the Laplacian

With the motivation of the previous section, let us consider the numerical solution of the elliptic boundary value problem

$$(2.1) \quad \Delta u = f \text{ in } \Omega, \quad u = g \text{ on } \Gamma.$$

For simplicity we will consider first a very simple domain  $\Omega = (0, 1) \times (0, 1)$ , the unit square in  $\mathbb{R}^2$ . Now this problem is so simplified that we can attack it analytically, e.g., by separation of variables, but it is a very useful *model problem* for studying numerical methods.

#### 1. The 5-point difference operator

Let  $N$  be a positive integer and set  $h = 1/N$ . Consider the *mesh* in  $\mathbb{R}^2$

$$\mathbb{R}_h^2 := \{ (mh, nh) : m, n \in \mathbb{Z} \}.$$

Note that each mesh point  $x \in \mathbb{R}_h^2$  has four *nearest neighbors* in  $\mathbb{R}_h^2$ , one each to the left, right, above, and below. We let  $\Omega_h = \Omega \cap \mathbb{R}_h^2$ , the set of interior mesh points, and we regard this a discretization of the domain  $\Omega$ . We also define  $\Gamma_h$  as the set of mesh points in  $\mathbb{R}_h^2$  which don't belong to  $\Omega_h$ , but which have a nearest neighbor in  $\Omega_h$ . We regard  $\Gamma_h$  as a discretization of  $\Gamma$ . We also let  $\bar{\Omega}_h := \Omega_h \cup \Gamma_h$

To discretize (2.1) we shall seek a function  $u_h : \bar{\Omega}_h \rightarrow \mathbb{R}$  satisfying

$$(2.2) \quad \Delta_h u_h = f \text{ on } \Omega_h, \quad u_h = g \text{ on } \Gamma_h.$$

Here  $\Delta_h$  is an operator, to be defined, which takes functions on  $\bar{\Omega}_h$  (*mesh functions*) to functions on  $\Omega_h$ . It should approximate the true Laplacian in the sense that if  $v$  is a smooth function on  $\bar{\Omega}$  and  $v_h = v|_{\bar{\Omega}_h}$  is the associated mesh function, then we want

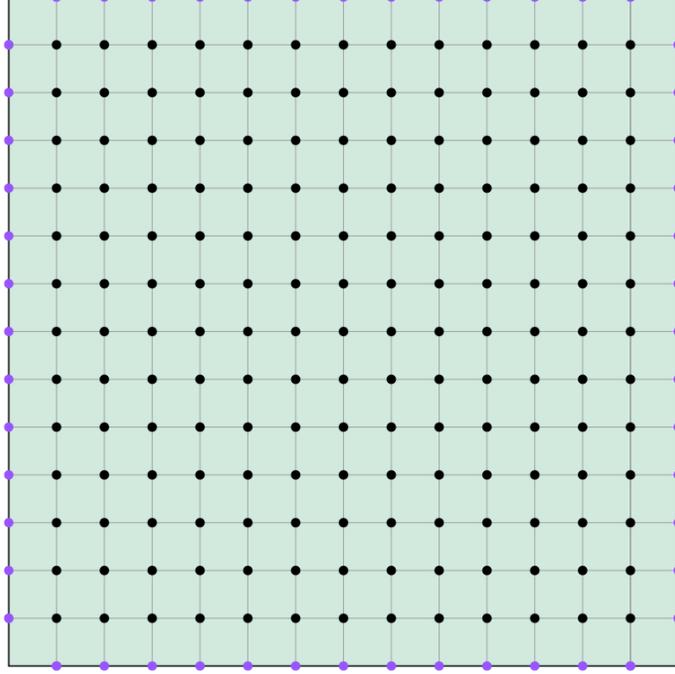
$$\Delta_h v_h \approx \Delta v|_{\Omega_h}$$

for  $h$  small.

Before defining  $\Delta_h$ , let us turn to the one-dimensional case. That is, given a function  $v_h$  defined at the mesh points  $nh, n \in \mathbb{Z}$ , we want to define a function  $D_h^2 v_h$  on the mesh points, so that  $D_h^2 v_h \approx v''|_{\mathbb{Z}h}$  if  $v_h = v|_{\mathbb{Z}h}$ . One natural procedure is to construct the quadratic polynomial  $p$  interpolating  $v_h$  at three consecutive mesh points  $(n-1)h, nh, (n+1)h$ , and let  $D_h^2 v_h(nh)$  be the constant value of  $p''$ . This gives the formula

$$D_h^2 v_h(nh) = 2v_h[(n-1)h, nh, (n+1)h] = \frac{v_h((n+1)h) - 2v_h(nh) + v_h((n-1)h)}{h^2}.$$

$D_h^2$  is known as the 3-point difference approximation to  $d^2/dx^2$ . We know that if  $v$  is  $C^2$  in a neighborhood of  $nh$ , then  $\lim_{h \rightarrow 0} v[x-h, x, x+h] = v''(x)/2$ . In fact, it is easy to show

FIGURE 2.1.  $\bar{\Omega}_h$  for  $h = 1/14$ : black: points in  $\Omega_h$ , purple: points in  $\Gamma_h$ .

by Taylor expansion (do it!), that

$$D_h^2 v(x) = v''(x) + \frac{h^2}{12} v^{(4)}(\xi), \text{ for some } \xi \in (x - h, x + h),$$

as long as  $v$  is  $C^4$  near  $x$ . Thus  $D_h^2$  is a second order approximation to  $d^2/dx^2$ .

Now returning to the definition of the  $\Delta_h \approx \Delta = \partial^2/\partial x^2 + \partial^2/\partial y^2$ , we simply use the 3-point approximation to  $\partial^2/\partial x^2$  and  $\partial^2/\partial y^2$ . Writing  $v_{mn}$  for  $v(mh, nh)$  we then have

$$\begin{aligned} \Delta_h v(mh, nh) &= \frac{v_{m+1,n} - 2v_{mn} + v_{m-1,n}}{h^2} + \frac{v_{m,n+1} - 2v_{mn} + v_{m,n-1}}{h^2} \\ &= \frac{v_{m+1,n} + v_{m-1,n} + v_{m,n+1} + v_{m,n-1} - 4v_{mn}}{h^2}. \end{aligned}$$

From the error estimate in the one-dimensional case we easily get that for  $v \in C^4(\bar{\Omega})$ ,

$$\Delta_h v(mh, nh) - \Delta v(mh, nh) = \frac{h^2}{12} \left[ \frac{\partial^4 v}{\partial x^4}(\xi, nh) + \frac{\partial^4 v}{\partial y^4}(mh, \eta) \right],$$

for some  $\xi, \eta$ . Thus:

**THEOREM 2.1.** *If  $v \in C^2(\bar{\Omega})$ , then*

$$\lim_{h \rightarrow 0} \|\Delta_h v - \Delta v\|_{L^\infty(\Omega_h)} = 0.$$

If  $v \in C^4(\bar{\Omega})$ , then

$$\|\Delta_h v - \Delta v\|_{L^\infty(\Omega_h)} \leq \frac{h^2}{6} M_4,$$

where  $M_4 = \max(\|\partial^4 v / \partial x^4\|_{L^\infty(\bar{\Omega})}, \|\partial^4 v / \partial y^4\|_{L^\infty(\bar{\Omega})})$ .

The discrete PDE  $\Delta_h u_h = f$  on  $\Omega_h$  is a system of  $M = (N - 1)^2$  linear equations in the unknown values of  $u_h$  at the mesh points. Since the values of  $u_h$  are given on the boundary mesh points, we may regard (2.2) as a system of  $M^2$  linear equations in  $M$  unknowns. For example, in the case  $N = 4$ ,  $M = 9$ , the system is

$$\begin{pmatrix} -4 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & -4 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & -4 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & -4 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & -4 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & -4 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & -4 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & -4 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & -4 \end{pmatrix} \begin{pmatrix} u_{1,1} \\ u_{2,1} \\ u_{3,1} \\ u_{1,2} \\ u_{2,2} \\ u_{3,2} \\ u_{1,3} \\ u_{2,3} \\ u_{3,3} \end{pmatrix} = \begin{pmatrix} h^2 f_{1,1} - u_{1,0} - u_{0,1} \\ h^2 f_{2,1} - u_{2,0} \\ h^2 f_{3,1} - u_{3,0} - u_{4,1} \\ h^2 f_{1,2} - u_{0,2} \\ h^2 f_{2,2} \\ h^2 f_{3,2} - u_{4,2} \\ h^2 f_{1,3} - u_{0,3} - u_{1,4} \\ h^2 f_{2,3} - u_{2,4} \\ h^2 f_{3,3} - u_{4,3} - u_{3,4} \end{pmatrix}$$

The matrix may be rewritten as

$$\begin{pmatrix} A & I & O \\ I & A & I \\ O & I & A \end{pmatrix}$$

where  $I$  is the  $3 \times 3$  identity matrix,  $O$  is the  $3 \times 3$  zero matrix, and

$$A = \begin{pmatrix} -4 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & -4 \end{pmatrix}.$$

For general  $N$  the matrix can be partitioned into  $(N - 1) \times (N - 1)$  blocks, each in  $\mathbb{R}^{(N-1) \times (N-1)}$ :

$$\begin{pmatrix} A & I & O & \cdots & O & O \\ I & A & I & \cdots & O & O \\ O & I & A & \cdots & O & O \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ O & O & O & \cdots & I & A \end{pmatrix},$$

where  $I$  and  $O$  are the identity and zero matrix in  $\mathbb{R}^{(N-1) \times (N-1)}$ , respectively, and  $A \in \mathbb{R}^{(N-1) \times (N-1)}$  is the tridiagonal matrix with  $-4$  on the diagonal and 1 above and below the diagonal. This assumes the unknowns are ordered

$$u_{1,1}, u_{2,1}, \dots, u_{N-1,1}, u_{1,2}, \dots, u_{N-1,N-1},$$

and the equations are ordered similarly.

The matrix can be created as in Matlab with the following code.

```

I = speye(n-1);
e = ones(n-1,1);
A = spdiags([e,-4*e,e],[-1,0,1],n-1,n-1);
J = spdiags([e,e],[-1,1],n-1,n-1);
Lh = kron(I,A) + kron(J,I)

```

Notice that the matrix has many special properties:

- it is sparse with at most 5 elements per row nonzero
- it is block tridiagonal, with tridiagonal and diagonal blocks
- it is symmetric
- it is diagonally dominant
- its diagonal elements are negative, all others nonnegative
- it is negative definite

## 2. Analysis via a maximum principle

We will now prove that the problem (2.2) has a unique solution and prove an error estimate. The key will be a discrete maximum principle.

**THEOREM 2.2** (Discrete Maximum Principle). *Let  $v$  be a function on  $\bar{\Omega}_h$  satisfying*

$$\Delta_h v \geq 0 \text{ on } \Omega_h.$$

*Then  $\max_{\Omega_h} v \leq \max_{\Gamma_h} v$ . Equality holds if and only if  $v$  is constant.*

**PROOF.** Suppose  $\max_{\Omega_h} v \geq \max_{\Gamma_h} v$ . Take  $x_0 \in \Omega_h$  where the maximum is achieved. Let  $x_1, x_2, x_3,$  and  $x_4$  be the nearest neighbors. Then

$$4v(x_0) = \sum_{i=1}^4 v(x_i) - h^2 \Delta_h v(x_0) \leq \sum_{i=1}^4 v(x_i) \leq 4v(x_0),$$

since  $v(x_i) \leq v(x_0)$ . Thus equality holds throughout and  $v$  achieves its maximum at all the nearest neighbors of  $x_0$  as well. Applying the same argument to the neighbors in the interior, and then to their neighbors, etc., we conclude that  $v$  is constant.  $\square$

**REMARKS.** 1. The analogous discrete minimum principle, obtained by reversing the inequalities and replacing max by min, holds. 2. This is a discrete analogue of the maximum principle for the Laplace operator.

**THEOREM 2.3.** *There is a unique solution to the discrete boundary value problem (2.2).*

**PROOF.** Since we are dealing with a square linear system, it suffices to show nonsingularity, i.e., that if  $\Delta_h u_h = 0$  on  $\Omega_h$  and  $u_h = 0$  on  $\Gamma_h$ , then  $u_h \equiv 0$ . Using the discrete maximum and the discrete minimum principles, we see that in this case  $u_h$  is everywhere 0.  $\square$

The next result is a statement of maximum norm stability.

**THEOREM 2.4.** *The solution  $u_h$  to (2.2) satisfies*

$$(2.3) \quad \|u_h\|_{L^\infty(\bar{\Omega}_h)} \leq \frac{1}{8} \|f\|_{L^\infty(\Omega_h)} + \|g\|_{L^\infty(\Gamma_h)}.$$

This is a stability result in the sense that it states that the mapping  $(f, g) \mapsto u_h$  is bounded uniformly with respect to  $h$ .

PROOF. We introduce the comparison function  $\phi(x) = [(x_1 - 1/2)^2 + (x_2 - 1/2)^2]/4$ , which satisfies  $\Delta_h \phi = 1$  on  $\Omega_h$ , and  $0 \leq \phi \leq 1/8$  on  $\bar{\Omega}_h$ . Set  $M = \|f\|_{L^\infty(\Omega_h)}$ . Then

$$\Delta_h(u_h + M\phi) = \Delta_h u_h + M \geq 0,$$

so

$$\max_{\Omega_h} u_h \leq \max_{\Omega_h} (u_h + M\phi) \leq \max_{\Gamma_h} (u_h + M\phi) \leq \max_{\Gamma_h} g + \frac{1}{8}M.$$

Thus  $u_h$  is bounded above by the right-hand side of (2.3). A similar argument applies to  $-u_h$  giving the theorem.  $\square$

By applying the stability result to the error  $u - u_h$  we can bound the error in terms of the consistency error  $\Delta_h u - \Delta u$ .

THEOREM 2.5. *Let  $u$  be the solution of the Dirichlet problem (1.2) and  $u_h$  the solution of the discrete problem (2.2). Then*

$$\|u - u_h\|_{L^\infty(\bar{\Omega}_h)} \leq \frac{1}{8} \|\Delta u - \Delta_h u\|_{L^\infty(\bar{\Omega}_h)}.$$

PROOF. Since  $\Delta_h u_h = f = \Delta u$  on  $\Omega_h$ ,  $\Delta_h(u - u_h) = \Delta_h u - \Delta u$ . Also,  $u - u_h = 0$  on  $\Gamma_h$ . Apply Theorem 2.4 (with  $u_h$  replaced by  $u - u_h$ ), we obtain the theorem.  $\square$

Combining with Theorem 2.1, we obtain error estimates.

COROLLARY 2.6. *If  $u \in C^2(\bar{\Omega})$ , then*

$$\lim_{h \rightarrow 0} \|u - u_h\|_{L^\infty(\bar{\Omega}_h)} = 0.$$

*If  $u \in C^4(\bar{\Omega})$ , then*

$$\|u - u_h\|_{L^\infty(\bar{\Omega}_h)} \leq \frac{h^2}{48} M_4,$$

*where  $M_4 = \max(\|\partial^4 u / \partial x_1^4\|_{L^\infty(\bar{\Omega})}, \|\partial^4 u / \partial x_2^4\|_{L^\infty(\bar{\Omega})})$ .*

### 3. Consistency, stability, and convergence

Now we introduce an abstract framework in which to understand the preceding analysis. It is general enough that it applies, or can be adapted to, a huge variety of numerical methods for PDE. We will keep in mind, as a basic example, the 5-point difference discretization of the Poisson equation with homogeneous boundary conditions, so the PDE problem to be solved is

$$\Delta u = f \text{ in } \Omega, \quad u = 0 \text{ on } \Gamma,$$

and the numerical method is

$$\Delta_h u_h = f_h \text{ in } \Omega_h, \quad u_h = 0 \text{ on } \Gamma_h.$$

Let  $X$  and  $Y$  be vector spaces and  $L : X \rightarrow Y$  a linear operator. Given  $f \in Y$ , we seek  $u \in X$  such that  $Lu = f$ . This is the problem we are trying to solve. So, for the homogeneous Dirichlet BVP for Poisson's equation, we could take  $X$  to be the space of  $C^2$  functions on  $\bar{\Omega}$  which vanish on  $\Gamma$ ,  $Y = C(\bar{\Omega})$ , and  $L = \Delta$ . (Actually, slightly more sophisticated spaces

should be taken if we wanted to get a good theory for the Poisson equation, but that won't concern us now.) We shall assume that there is a solution  $u$  of the original problem.

Now let  $X_h$  and  $Y_h$  be finite dimensional normed vector spaces and  $L_h : X_h \rightarrow Y_h$  a linear operator. Our numerical method, or discretization, is:

$$\text{Given } f_h \in Y_h \text{ find } u_h \in X_h \text{ such that } L_h u_h = f_h.$$

Of course, this is a very minimalistic framework so far. Without some more hypotheses, we do not know if this finite dimensional problem has a solution, or if the solution is unique. And we certainly don't know that  $u_h$  is in any sense an approximation of  $u$ .

In fact, up until now, there is no way to compare  $u$  to  $u_h$ , since they belong to different spaces. For this reason, we introduce a *representative* of  $u$ ,  $r_h u \in X_h$ . We can then talk about the *error*  $r_h u - u_h$  and its norm  $\|r_h u - u_h\|_{X_h}$ . If this error norm is small, that means that  $u_h$  is close to  $u$ , or at least close to our representative  $r_h u$  of  $u$ , in the sense of the norm.

In short, we would like the error to be small in norm. To make this precise we do what is always done in numerical analysis: we consider not a single discretization, but a sequence of discretizations. To keep the notation simple, we will now think of  $h > 0$  as a parameter tending to 0, and suppose that we have the normed spaces  $X_h$  and  $Y_h$  and the linear operator  $L_h : X_h \rightarrow Y_h$  and the element  $f_h \in Y_h$  for each  $h$ . This family of discretizations is called *convergent* if the norm  $\|r_h u - u_h\|_{X_h}$  tends to 0 as  $h \rightarrow 0$ .

In our example, we take  $X_h$  to be the grid functions in  $L^\infty(\bar{\Omega}_h)$  which vanish on  $\Gamma_h$ , and  $Y_h$  to be the grid functions in  $L^\infty(\Omega)$ , and equip both with the maximum norm. We also simply define  $r_h u = u|_{\Omega_h}$ . Thus a small error means that  $u_h$  is close to the true solution  $u$  at all the grid points, which is a desirable result.

Up until this point there is not enough substance to our abstract framework for us to be able to prove a convergence result, because the only connection between the original problem  $Lu = f$  and the discrete problems  $L_h u_h = f_h$  is that the notations are similar. We surely need some hypotheses. The first of two key hypotheses is *consistency*, which says that, in some sense, the discrete problem is reasonable, in that the solution of the original problem almost satisfies the discrete problem. More precisely, we define the *consistency error* as  $L_h r_h u - f_h \in Y_h$ , a quantity which we can measure using our norm in  $Y_h$ . The family of discretizations is called *consistent* if the norm  $\|L_h r_h u - f_h\|_{Y_h}$  tends to 0 as  $h \rightarrow 0$ .

Not every consistent family of discretizations is convergent (as you can easily convince yourself, since consistency involves the norm in  $Y_h$  but not the norm in  $X_h$  and for convergence it is the opposite). There is a second key hypothesis, uniform well-posedness of the discrete problems. More precisely, we assume that each discrete problem is uniquely solvable (nonsingular): for every  $g_h \in Y_h$  there is a unique  $v_h \in X_h$  with  $L_h v_h = g_h$ . Thus the operator  $L_h^{-1} : Y_h \rightarrow X_h$  is defined and we call its norm  $c_h = \|L_h^{-1}\|_{\mathcal{L}(Y_h, X_h)}$  the *stability constant* of the discretization. The family of discretizations is called *stable* if the stability constants are bounded uniformly in  $h$ :  $\sup_h c_h < \infty$ . Note that stability is a property of the discrete problems and depends on the particular choice of norms, but it does not depend on the true solution  $u$  in any way.

With these definitions we get a theorem which is trivial to prove, but which captures the underlying structure of many convergence results in numerical PDE.

**THEOREM 2.7.** *Let there be given normed vector spaces  $X_h$  and  $Y_h$ , an invertible linear operator  $L_h : X_h \rightarrow Y_h$ , an element  $f_h \in Y_h$ , and a representative  $r_h u \in X_h$ . Define  $u_h \in X_h$*

by  $L_h u_h = f_h$ . Then the norm of the error is bounded by the stability constant times the norm of the consistency error. If a family of such discretizations is consistent and stable, then it is convergent.

PROOF. Since  $L_h u_h = f_h$ ,

$$L_h(r_h u - u_h) = L_h r_h u - f_h.$$

Applying  $L_h^{-1}$  we obtain

$$r_h u - u_h = L_h^{-1}(L_h r_h u - f_h),$$

and taking norms we get

$$\|r_h u - u_h\|_{X_h} = \|L_h^{-1}\|_{\mathcal{L}(Y_h, X_h)} \|L_h r_h u - f_h\|_{Y_h},$$

which is the desired result.  $\square$

REMARK. We emphasize that the concepts of convergence, consistency, and stability depend on the choice of norms in  $X_h$ ,  $Y_h$ , and both, respectively. The norm in  $X_h$  should be chosen so that the convergence result gives information that is desired. Choosing a weak norm may make the hypotheses easier to verify, but the result of less interest. Similarly,  $f_h$  must be chosen in a practical way. We need  $f_h$  to compute  $u_h$ , so it should be something we know before we solve the problem, typically something easily computed from  $f$ . Similarly as well,  $r_h u$  should be chosen in a reasonable way. For example, choosing  $r_h u = L_h^{-1} f_h$  would give  $r_h u = u_h$  so we definitely have a convergent method, but this is cheating: convergence is of no interest with this choice. The one other choice we have at our disposal is the norm on  $Y_h$ . This we are free to choose in order to make the hypotheses of consistency and stability possible to verify. Note that weakening the norm on  $Y_h$  makes it easier to prove consistency, while strengthening it makes it easier to prove stability.

Returning to our example, we see that the first statement of Theorem 2.1 is just the statement that the method is consistent for any solution  $u \in C^2(\bar{\Omega})$ , and the second statement says that the consistency error is  $O(h^2)$  if  $u \in C^4(\bar{\Omega})$ . On the other hand, if we apply Theorem 2.4 with  $g = 0$ , it states that the stability constant  $c_h \leq 1/8$  for all  $h$ , and so the method is stable. We then obtain the convergence result in Corollary 2.6 by the basic result of Theorem 2.7.

#### 4. Fourier analysis

Define  $L(\Omega_h)$  to be the set of functions  $\Omega_h \rightarrow \mathbb{R}$ , which is isomorphic to  $\mathbb{R}^M$ ,  $M = (N-1)^2$ . Sometimes we think of these as functions on  $\bar{\Omega}_h$  extended by zero to  $\Gamma_h$ . The discrete Laplacian then defines an isomorphism of  $L(\Omega_h)$  onto itself. As we just saw, the  $L^\infty$  stability constant,  $\|\Delta_h^{-1}\|_{\mathcal{L}(L^\infty, L^\infty)} \leq 1/8$ . In this section we use Fourier analysis to establish a similar  $L^2$  stability result.

First consider the one-dimensional case. With  $h = 1/N$  let  $I_h = \{h, 2h, \dots, (N-1)h\}$ , and let  $L(I_h)$  be the space of functions on  $I_h$ , which is an  $N-1$  dimensional vectorspace. On  $L(I_h)$  we define the inner product

$$\langle u, v \rangle_h = h \sum_{k=1}^{N-1} u(kh)v(kh),$$

with the corresponding norm  $\|v\|_h$ .

The space  $L(I_h)$  is a discrete analogue of  $L^2(I)$  where  $I$  is the unit interval. On this latter space the functions  $\sin \pi m x$ ,  $m = 1, 2, \dots$ , form an orthogonal basis consisting of eigenfunctions of the operator  $-d^2/dx^2$ . The corresponding eigenvalues are  $\pi^2, 4\pi^2, 9\pi^2, \dots$ . We now establish the discrete analogue of this result.

Define  $\phi_m \in L(I_h)$  by  $\phi_m(x) = \sin \pi m x$ ,  $x \in I_h$ . It turns out that these mesh functions are precisely the eigenvectors of the operator  $D_h^2$ . Indeed

$$D_h^2 \phi_m(x) = \frac{\sin \pi m(x+h) - 2 \sin \pi m x + \sin \pi m(x-h)}{h^2} = \frac{2}{h^2} (\cos \pi m h - 1) \sin \pi m x.$$

Thus

$$D_h^2 \phi_m = -\lambda_m \phi_m, \quad \lambda_m = \frac{2}{h^2} (1 - \cos \pi m h) = \frac{4}{h^2} \sin^2 \frac{\pi m h}{2}.$$

Note that

$$0 < \lambda_1 < \lambda_2 < \dots < \lambda_{N-1} < \frac{4}{h^2}.$$

Note also that for small  $m \ll N$ ,  $\lambda_m \approx \pi^2 m^2$ . In particular  $\lambda_1 \approx \pi^2$ . To get a strict lower bound we note that  $\lambda_1 = 8$  for  $N = 2$  and  $\lambda_1$  increases with  $N$ .

Since the operator  $D_h^2$  is symmetric with respect to the inner product on  $L(I_h)$ , and the eigenvalues  $\lambda_m$  are distinct, it follows that the eigenvectors  $\phi_m$  are mutually orthogonal. (This can also be obtained using trigonometric identities, or by expressing the sin functions in terms of complex exponentials and using the discrete Fourier transform.) Since there are  $N - 1$  of them, they form a basis of  $L(I_h)$ .

**THEOREM 2.8.** *The functions  $\phi_m$ ,  $m = 1, 2, \dots, N - 1$  form an orthogonal basis of  $L(I_h)$ . Consequently, any function  $v \in L(I_h)$  can be expanded as  $v = \sum_{m=1}^{N-1} a_m \phi_m$  with  $a_m = \langle v, \phi_m \rangle_h / \|\phi_m\|_h^2$ , and  $\|v\|_h^2 = \sum_{m=1}^{N-1} a_m^2 \|\phi_m\|_h^2$ .*

From this we obtain immediately a stability result for the one-dimensional Laplacian. If  $v \in L(I_h)$  and  $D_h^2 v = f$ , we expand  $v$  in terms of the  $\phi_m$ :

$$v = \sum_{m=1}^{N-1} a_m \phi_m, \quad \|v\|_h^2 = \sum_{m=1}^{N-1} a_m^2 \|\phi_m\|_h^2.$$

Then

$$f = - \sum_{m=1}^{N-1} \lambda_m a_m \phi_m, \quad \|f\|_h^2 = \sum_{m=1}^{N-1} \lambda_m^2 a_m^2 \|\phi_m\|_h^2 \geq 8^2 \|v\|_h^2.$$

Thus  $\|v\|_h \leq \|f\|_h / 8$ .

The extension to the two-dimensional case is straightforward. We use the basis  $\phi_{mn} = \phi_m \otimes \phi_n$ , i.e.,

$$\phi_{mn}(x, y) := \phi_m(x) \phi_n(y), \quad m, n = 1, \dots, N - 1,$$

for  $L(\Omega_h)$ . It is easy to see that these  $(N - 1)^2$  functions form an orthogonal basis for  $L(\Omega_h)$  equipped with the inner product

$$\langle u, v \rangle_h = h^2 \sum_{m=1}^{N-1} \sum_{n=1}^{N-1} u(mh, nh) v(mh, nh)$$

and corresponding norm  $\|\cdot\|_h$ . Moreover  $\phi_{mn}$  is an eigenvector of  $-\Delta_h$  with eigenvalue  $\lambda_{mn} = \lambda_m + \lambda_n \geq 16$ . The next theorem follows immediately.

**THEOREM 2.9.** *The operator  $\Delta_h$  defines an isomorphism from  $L(\Omega_h)$  to itself. Moreover  $\|\Delta_h^{-1}\| \leq 1/16$  where the operator norm is with respect to the norm  $\|\cdot\|_h$  on  $L(\Omega_h)$ .*

Since the  $\|v\|_h \leq \|v\|_{L^\infty(\Omega_h)}$  we also have consistency with respect to the discrete 2-norm. We leave it to the reader to complete the analysis with a convergence result.

### 5. Analysis via summation by parts

Fourier analysis is not the only approach to get an  $L^2$  stability result. Another uses *summation by parts*, the discrete analogue of integration by parts.

Let  $v$  be a mesh function. Define the backward difference operator

$$\partial_x v(mh, nh) = \frac{v(mh, nh) - v((m-1)h, nh)}{h}, \quad 1 \leq m \leq N, \quad 0 \leq n \leq N.$$

In this section we denote

$$\langle v, w \rangle_h = h^2 \sum_{m=1}^N \sum_{n=1}^N v(mh, nh) w(mh, nh),$$

with the corresponding norm  $\|\cdot\|_h$  (this agrees with the notation in the last section for mesh functions which vanish on  $\Gamma_h$ ).

**LEMMA 2.10.** *If  $v \in L(\Omega_h)$  (the set of mesh functions vanishing on  $\Gamma_h$ ), then*

$$\|v\|_h \leq \frac{1}{2} (\|\partial_x v\|_h + \|\partial_y v\|_h).$$

**PROOF.** It is enough to show that  $\|v\|_h \leq \|\partial_x v\|_h$ . The same will similarly hold for  $\partial_y$  as well, and we can average the two results.

For  $1 \leq m \leq N$ ,  $0 \leq n \leq N$ ,

$$\begin{aligned} |v(mh, nh)|^2 &\leq \left( \sum_{i=1}^N |v(ih, nh) - v((i-1)h, nh)| \right)^2 \\ &= \left( h \sum_{i=1}^N |\partial_x v(ih, nh)| \right)^2 \\ &\leq \left( h \sum_{i=1}^N |\partial_x v(ih, nh)|^2 \right) \left( h \sum_{i=1}^N 1^2 \right) \\ &= h \sum_{i=1}^N |\partial_x v(ih, nh)|^2. \end{aligned}$$

Therefore

$$h \sum_{m=1}^N |v(mh, nh)|^2 \leq h \sum_{i=1}^N |\partial_x v(ih, nh)|^2$$

and

$$h^2 \sum_{m=1}^N \sum_{n=1}^N |v(mh, nh)|^2 \leq h^2 \sum_{i=1}^N \sum_{n=1}^N |\partial_x v(ih, nh)|^2,$$

i.e.,  $\|v\|_h^2 \leq \|\partial_x v\|_h^2$ , as desired.  $\square$

This result is a discrete analogue of Poincaré's inequality, which bounds a function in terms of its gradient as long as the function vanishes on a portion of the boundary. The constant of  $1/2$  in the bound can be improved. The next result is a discrete analogue of Green's Theorem (essentially, integration by parts).

LEMMA 2.11. *If  $v, w \in L(\Omega_h)$ , then*

$$-\langle \Delta_h v, w \rangle_h = \langle \partial_x v, \partial_x w \rangle_h + \langle \partial_y v, \partial_y w \rangle_h.$$

PROOF. Let  $v_0, v_1, \dots, v_N, w_0, w_1, \dots, w_N \in \mathbb{R}$  with  $w_0 = w_N = 0$ . Then

$$\begin{aligned} \sum_{i=1}^N (v_i - v_{i-1})(w_i - w_{i-1}) &= \sum_{i=1}^N v_i w_i + \sum_{i=1}^N v_{i-1} w_{i-1} - \sum_{i=1}^N v_{i-1} w_i - \sum_{i=1}^N v_i w_{i-1} \\ &= 2 \sum_{i=1}^{N-1} v_i w_i - \sum_{i=1}^{N-1} v_{i-1} w_i - \sum_{i=1}^{N-1} v_{i+1} w_i \\ &= - \sum_{i=1}^{N-1} (v_{i+1} - 2v_i + v_{i-1}) w_i. \end{aligned}$$

Hence,

$$\begin{aligned} -h \sum_{i=1}^{N-1} \frac{v((i+1)h, nh) - 2v(ih, nh) + v((i-1)h, nh)}{h^2} w(ih, nh) \\ = h \sum_{i=1}^N \partial_x v(ih, nh) \partial_x w(ih, nh), \end{aligned}$$

and thus

$$-\langle D_x^2 v, w \rangle_h = \langle \partial_x v, \partial_x w \rangle_h.$$

Similarly,  $-\langle D_y^2 v, w \rangle_h = \langle \partial_y v, \partial_y w \rangle_h$ , so the lemma follows.  $\square$

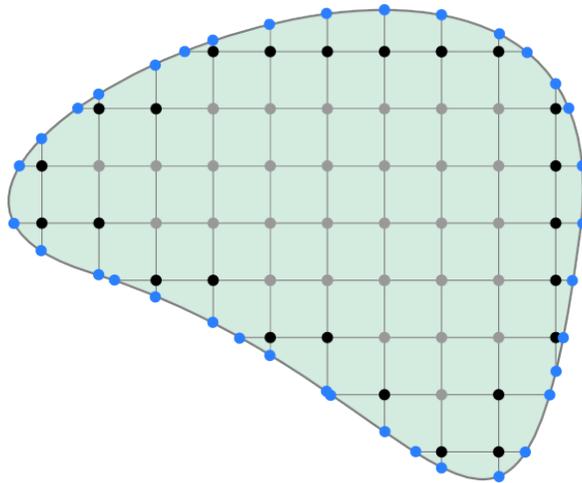
Combining the discrete Poincaré inequality with the discrete Green's theorem, we immediately get a stability result. If  $v \in L(\Omega_h)$ , then

$$\|v\|_h^2 \leq \frac{1}{2} (\|\partial_x v\|_h^2 + \|\partial_y v\|_h^2) = -\frac{1}{2} \langle \Delta_h v, v \rangle_h \leq \frac{1}{2} \|\Delta_h v\|_h \|v\|_h.$$

Thus

$$\|v\|_h \leq \|\Delta_h v\|_h, \quad v \in L(\Omega_h),$$

which is a stability result.

FIGURE 2.2. Gray points:  $\mathring{\Omega}_h$ . Black points:  $\Omega_h^\partial$ . Blue points:  $\Gamma_h$ .

## 6. Extensions

**6.1. Curved boundaries.** Thus far we have studied as a model problem the discretization of Poisson's problem on the square. In this subsection we consider a variant which can be used to discretize Poisson's problem on a fairly general domain.

Let  $\Omega$  be a smoothly bounded open set in  $\mathbb{R}^2$  with boundary  $\Gamma$ . We again consider the Dirichlet problem for Poisson's equation, (2.1), and again set  $\Omega_h = \Omega \cap \mathbb{R}_h^2$ . If  $(x, y) \in \Omega_h$  and the segment  $(x + sh, y)$ ,  $0 \leq s \leq 1$  belongs to  $\Gamma$ , then the point  $(x + h, y)$ , which belongs to  $\Omega_h$ , is a *neighbor* of  $(x, y)$  to the right. If this segment doesn't belong to  $\Omega$  we define another sort of neighbor to the right, which belongs to  $\Gamma$ . Namely we define the neighbor to be the point  $(x + sh, y)$  where  $0 < s \leq 1$  is the largest value for which  $(x + th, y) \in \Omega$  for all  $0 \leq t < s$ . The points of  $\Gamma$  so constructed (as neighbors to the right or left or above or below points in  $\Omega_h$ ) constitute  $\Gamma_h$ . Thus every point in  $\Omega_h$  has four nearest neighbors all of which belong to  $\bar{\Omega}_h := \Omega_h \cup \Gamma_h$ . We also define  $\mathring{\Omega}_h$  as those points in  $\Omega_h$  all four of whose neighbor belong to  $\bar{\Omega}_h$  and  $\Omega_h^\partial$  as those points in  $\Omega_h$  with at least one neighbor in  $\Gamma_h$ . See Figure 2.2.

In order to discretize the Poisson equation we need to construct a discrete analogue of the Laplacian  $\Delta_h v$  for mesh functions  $v$  on  $\bar{\Omega}_h$ . Of course on  $\mathring{\Omega}_h$ ,  $\Delta_h v$  is defined as the usual 5-point Laplacian. For  $(x, y) \in \Omega_h^\partial$ , let  $(x + h_E, y)$ ,  $(x, y + h_N)$ ,  $(x - h_W, y)$ , and  $(x, y - h_S)$  be the nearest neighbors (with  $0 < h_E, h_N, h_W, h_S \leq h$ ), and let  $v_E, v_N, v_W$ , and  $v_S$  denote the value of  $v$  at these four points. Setting  $v_0 = v(x, y)$  as well, we will define  $\Delta_h v(x, y)$  as a linear combination of these five values of  $v$ . In order to derive the formula, we first consider approximating  $d^2v/dx^2(0)$  by a linear combination of  $v(-h_-)$ ,  $v(0)$ , and  $v(h_+)$ , for a function  $v$  of one variable. By Taylor's theorem

$$\begin{aligned} \alpha_- v(-h_-) + \alpha_0 v(0) + \alpha_+ v(h_+) &= (\alpha_- + \alpha_0 + \alpha_+)v(0) + (\alpha_+ h_+ - \alpha_- h_-)v'(0) \\ &\quad + \frac{1}{2}(\alpha_+ h_+^2 + \alpha_- h_-^2)v''(0) + \frac{1}{6}(\alpha_+ h_+^3 - \alpha_- h_-^3)v'''(0) + \dots \end{aligned}$$

Thus, to obtain a consistent approximation we must have

$$\alpha_- + \alpha_0 + \alpha_+ = 0, \quad \alpha_+ h_+ - \alpha_- h_- = 0, \quad \frac{1}{2}(\alpha_+ h_+^2 + \alpha_- h_-^2) = 1,$$

which give

$$\alpha_- = \frac{2}{h_-(h_- + h_+)}, \quad \alpha_+ = \frac{2}{h_+(h_- + h_+)}, \quad \alpha_0 = \frac{-2}{h_- h_+}.$$

Note that we have simply recovered the usual divided difference approximation to  $d^2v/dx^2$ :

$$\alpha_- v(-h_-) + \alpha_0 v(0) + \alpha_+ v(h_+) = \frac{[v(h_+) - v(0)]/h_+ - [v(0) - v(-h_-)]/h_-}{(h_+ + h_-)/2} = 2v[-h_-, 0, h_+].$$

Returning to the 2-dimensional case, and applying the above considerations to both  $\partial^2 v/\partial x^2$  and  $\partial^2 v/\partial y^2$  we arrive at the *Shortley–Weller* formula for  $\Delta_h v$ :

$$\begin{aligned} \Delta_h v(x, y) &= \frac{2}{h_E(h_E + h_W)} v_E + \frac{2}{h_N(h_N + h_S)} v_N \\ &\quad + \frac{2}{h_W(h_E + h_W)} v_W + \frac{2}{h_S(h_N + h_S)} v_S - \left( \frac{2}{h_E h_W} + \frac{2}{h_N h_S} \right) v_0. \end{aligned}$$

Using Taylor's theorem with remainder we easily calculate that for  $v \in C^3(\bar{\Omega})$ ,

$$\|\Delta v - \Delta_h v\|_{L^\infty(\Omega_h)} \leq \frac{2M_3}{3}h,$$

where  $M_3$  is the maximum of the  $L^\infty$  norms of the third derivatives of  $v$ . Of course at the mesh points in  $\bar{\Omega}_h$ , the consistency error is bounded by  $M_4 h^2/6 = O(h^2)$ , as before, but for mesh points neighboring the boundary, it is reduced to  $O(h)$ .

The approximate solution to (2.1) is  $u_h : \bar{\Omega}_h \rightarrow \mathbb{R}$  determined again by 2.2. This is a system of linear equations with one unknown for each point of  $\Omega_h$ . In general the matrix won't be symmetric, but it maintains other good properties from the case of the square:

- it is sparse, with at most five elements per row
- it has negative diagonal elements and non-negative off-diagonal elements
- it is diagonally dominant.

Using these properties we can obtain the discrete maximum principle with virtually the same proof as for Theorem 2.2, and then a stability result as in Theorem 2.4 follows as before. (By contrast, the Fourier analysis approach to stability does not apply when the mesh spacing is not uniform.) In this way we can easily obtain an  $O(h)$  convergence result.

We now show how, with a more carefully analysis, we can improve this convergence result. We shall show that, *even though the consistency error is only  $O(h)$  at some points, the error is  $O(h^2)$  at all mesh points.*

Let  $X_h$  denote the space of mesh functions defined on  $\bar{\Omega}_h$  and which vanish on the mesh points in  $\Gamma_h$ . On this space we continue to use the maximum norm. Let  $Y_h$  denote the space of mesh functions defined on the interior mesh points only, i.e., on  $\Omega_h$ . On this space we shall use a different norm, namely,

$$(2.4) \quad \|f\|_{Y_h} := \max \left\{ \max_{x \in \bar{\Omega}_h} |f(x)|, h \max_{x \in \Omega_h^\circ} |f(x)| \right\}.$$

Thus we use the maximum norm except with a weight which decreases the emphasis on the points with a neighbor on the boundary. This norm is smaller than the maximum norm, and, measured in this norm, the consistency error is not just  $O(h)$  but rather  $O(h^2)$ :

$$\|\Delta_h u - \Delta u\|_{Y_h} \leq \max\left(\frac{M_4}{6}h^2, h\frac{2M_3}{3}h\right) = O(h^2).$$

Now, with respect to a smaller norm in  $Y_h$  the condition of stability is more stringent. So the key point is to show that *the Shortley-Weller discrete Laplacian is stable from  $X_h$  to  $Y_h$*  with this new norm. For the argument we will use the maximum principle with a slightly more sophisticated comparison function.

Before we used as a comparison function  $\phi : \bar{\Omega}_h \rightarrow \mathbb{R}$  defined by  $\phi(x_1, x_2) = [(x_1 - 1/2)^2 + (x_2 - 1/2)^2]/4$ , where  $(1/2, 1/2)$  was chosen as the vertex because it was the center of the square (making  $\|\phi\|_{L^\infty}$  as small as possible while satisfying  $\Delta_h \phi \equiv 1$ ). Now, suppose that  $\Omega$  is contained in the disk of some radius  $r$  about some point  $p$ . Then we define

$$(2.5) \quad \phi(x) = \begin{cases} [(x_1 - p_1)^2 + (x_2 - p_2)^2]/4, & x \in \Omega_h, \\ [(x_1 - p_1)^2 + (x_2 - p_2)^2]/4 + h, & x \in \Gamma_h \end{cases}$$

Thus we perturb the quadratic comparison function by adding  $h$  on the boundary. Then  $\phi$  is bounded independent of  $h$  ( $\|\phi\|_{L^\infty} \leq r^2/4 + h \leq r^2/4 + 2r$ ). Moreover  $\Delta_h \phi(x) = 1$ , if  $x \in \mathring{\Omega}_h$ , since then  $\phi$  is just the simple quadratic at  $x$  and all its neighbors. However, if  $x \in \Omega_h^\partial$ , then there is an additional term in  $\Delta_h \phi(x)$  for each neighbor of  $x$  on the boundary (typically one or two). For example, if  $(x_1 - h_W, x_2) \in \Gamma_h$  is a neighbor of  $x$  and the other neighbors are in  $\mathring{\Omega}_h$ , then

$$\Delta_h \phi(x) = 1 + \frac{2}{h_W(h_W + h)}h \geq h^{-1},$$

since  $h_W \leq h$ . Thus we have

$$(2.6) \quad \Delta_h \phi(x) \geq \begin{cases} 1, & x \in \mathring{\Omega}_h, \\ h^{-1}, & x \in \Omega_h^\partial. \end{cases}$$

Now let  $v : \bar{\Omega}_h \rightarrow \mathbb{R}$  be a mesh function, and set  $M = \|\Delta_h v\|_{Y_h}$  (weighted max norm of the Shortley-Weller discrete Laplacian of  $v$ ). If  $x \in \mathring{\Omega}_h$ , then  $M \geq |\Delta_h v(x)|$  and  $\Delta_h \phi(x) = 1$ , so

$$\Delta_h(M\phi)(x) \geq |\Delta_h v(x)|.$$

If  $x \in \Omega_h^\partial$ , then  $M \geq h|\Delta_h v(x)|$  and  $\Delta_h \phi(x) \geq h^{-1}$ , so again

$$\Delta_h(M\phi)(x) \geq |\Delta_h v(x)|.$$

Therefore

$$\Delta_h(v + M\phi) \geq 0 \text{ on } \Omega_h.$$

We can then apply the maximum principle (which easily extends to the Shortley-Weller discrete Laplacian), to get

$$\max_{\Omega_h} v \leq \max_{\Omega_h} (v + M\phi) \leq \max_{\Gamma_h} (v + M\phi) \leq \max_{\Gamma_h} v + c\|\Delta_h v\|_{Y_h},$$

where  $c = \|\phi\|_{L^\infty}$ . Of course, we have a similar result for  $-v$ , so

$$\|v\|_{L^\infty(\bar{\Omega}_h)} \leq \|v\|_{L^\infty(\Gamma_h)} + c\|\Delta_h v\|_{Y_h}.$$

In particular, if  $v$  vanishes on  $\Gamma_h$ , then

$$\|v\|_{L^\infty(\bar{\Omega}_h)} \leq c \|\Delta_h v\|_{Y_h}, \quad v \in X_h,$$

which is the desired stability result. As usual, we apply the stability estimate to  $v = u - u_h$ , and so get the error estimate

$$\|u - u_h\|_{L^\infty(\bar{\Omega}_h)} \leq c \|\Delta_h u - \Delta u\|_{Y_h} = O(h^2).$$

**REMARK.** The perturbation of  $h$  on the boundary in the definition (2.5) of the comparison function  $\phi$ , allowed us to place a factor of  $h$  in front of the  $\Omega_h^\partial$  terms in the  $Y_h$  norm (2.4) and still obtain stability. For this we needed (2.6) and the fact that the perturbed comparison function  $\phi$  remained bounded independent of  $h$ . In fact, we could take a larger perturbation by replacing  $h$  with 1 in (2.5). This would lead to a strengthening of (2.6), namely we could replace  $h^{-1}$  by  $h^{-2}$ , and still have  $\phi$  bounded independently of  $h$ . In this way we can prove stability with the same  $L^\infty$  norm for  $X_h$  and an even weaker norm for  $Y_h$ :

$$\|f\|_{Y_h} := \max \left\{ \max_{x \in \bar{\Omega}_h} |f(x)|, h^2 \max_{x \in \Omega_h^\partial} |f(x)| \right\}.$$

We thus get an even stronger error bound, with  $T_h = \Delta_h u - \Delta u$  denoting the consistency error, we get

$$\|u - u_h\|_{L^\infty(\bar{\Omega}_h)} \leq c \max \left\{ \|T_h\|_{L^\infty(\bar{\Omega}_h)}, h^2 \|T_h\|_{L^\infty(\Omega_h^\partial)} \right\} \leq c \max \{M_4 h^2, M_3 h^3\} = O(h^2).$$

This estimate shows that the points with neighbors on the boundary, despite having the largest consistency error ( $O(h)$  rather than  $O(h^2)$  for the other grid points), contribute only a small portion of the error ( $O(h^3)$  rather than  $O(h^2)$ ).

This example should be another warning to placing too much trust in a naive analysis of a numerical method by just using Taylor's theorem to expand the consistency error. Not only can a method perform worse than this might suggest, because of instability, it can also perform better, because of additional stability properties, as in this example.

**6.2. More general PDEs.** It is not difficult to extend the method and analysis to more general PDEs. For example, instead of the Poisson equation, we may take

$$\Delta u - a \frac{\partial u}{\partial x_1} - b \frac{\partial u}{\partial x_2} - cu = f,$$

where  $a$ ,  $b$ , and  $c$  are continuous coefficient functions on the square  $\bar{\Omega}$ . The difference method takes the obvious form:

$$\begin{aligned} \Delta_h u(x) - a(x) \frac{u(x_1 + h, x_2) - u(x_1 - h, x_2)}{h} - b(x) \frac{u(x_1, x_2 + h) - u(x_1, x_2 - h)}{h} \\ - c(x)u(x) = f(x), \quad x \in \Omega_h. \end{aligned}$$

It is easy to show that the consistency error is  $O(h^2)$ . As long as the coefficient  $c \geq 0$ , a version of the discrete maximum principle holds, and one thus obtains stability and convergence.

**6.3. More general boundary conditions.** It is also fairly easy to extend the method to more general boundary conditions, e.g., the Neumann condition  $\partial u/\partial n = g$  on all or part of the boundary, although some cleverness is needed to obtain a stable method with consistency error  $O(h^2)$  especially on a domain with curved boundary. We will not go into this topic here, but will treat Neumann problems when we consider finite elements.

**6.4. Nonlinear problems.** Consider, for example, the quasilinear equation

$$\Delta u = F(u, \partial u/\partial x_1, \partial u/\partial x_2),$$

with Dirichlet boundary conditions on the square. Whether this problem has a solution, and whether that solution is unique, or at least locally unique, depends on the nature of the nonlinearity  $F$ , and is beyond the scope of these notes. Supposing the problem does have a (locally) unique solution, we may try to compute it with finite differences. A simple scheme is

$$\Delta_h u_h = F(u_h, \partial_{x_1} u_h, \partial_{x_2} u_h), \quad x \in \Omega_h,$$

where we use, e.g., centered differences like

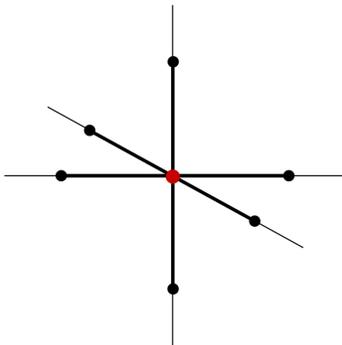
$$\partial_{x_1} u_h(x) = \frac{u(x_1 + h, x_2) - u(x_1 - h, x_2)}{2h}, \quad x \in \Omega_h.$$

Viewing the values of  $u_h$  at the  $M$  interior mesh points as unknowns, this is a system of  $M$  equations in  $M$  unknowns. The equations are not linear, but they have the same sparsity pattern as the linear systems we considered earlier: the equation associated to a certain grid point involves at most 5 unknowns, those associated to the grid point and its nearest neighbors.

The nonlinear system is typically solved by an iterative method, very often Newton's method or a variant of it. Issues like solvability, consistency, stability, and convergence can be studied for a variety of particular nonlinear problems. As for nonlinear PDE themselves, many issues arise which vary with the problem under consideration.

**6.5. Three dimensions.** The 5-point Laplacian on a square grid extends in a straightforward way to a 7-point Laplacian on a cubic grid. Figure 2.3 shows a grid point and its 6 nearest neighbors. The matrix for the 7-point Laplacian is roughly  $N^3 \times N^3$  where  $N = 1/h$ , so much larger for the same  $h$ , and solving the matrix equation which arises can be very challenging even for large fast computers.

FIGURE 2.3. A grid point and its six nearest neighbors on a 3-D grid.





## CHAPTER 3

### Linear algebraic solvers

The finite difference method reduces a boundary value problem for a PDE to a linear algebraic system  $Ax = f$ , with  $A \in \mathbb{R}^{n \times n}$  and  $f \in \mathbb{R}^n$ . The solution of this system dominates the computation time. (For the 5-point Laplacian on a square with  $h = 1/N$ , then  $n = (N - 1)^2$ .) The simplest way to solve this is through some variation of Gaussian elimination. Since the matrix  $A$  is symmetric positive definite (for the 5-point Laplacian on a square, for instance), we can use the Cholesky decomposition. Cholesky usually requires  $O(n^3) = O(N^6)$  floating point additions and multiplications (more precisely  $n^3/6 + O(n^2)$ ), but this is reduced in this case, because of the sparsity of the matrix. Gaussian elimination is not able to exploit the full sparsity of  $A$  (since when we factor  $A$  as  $LL^T$  with  $L$  lower triangular,  $L$  will be much less sparse than  $A$ ), but it is able to exploit the fact that  $A$  is *banded*: in the natural ordering all the nonzero entries are on the main diagonal or on one of the first  $N - 1$  sub- or super-diagonals. As a result, the storage count is reduced from  $O(n^2) = O(N^4)$  to  $O(nN) = O(N^3)$  and the operation count is reduced from  $O(N^6)$  to  $O(nN^2) = O(N^4)$ .

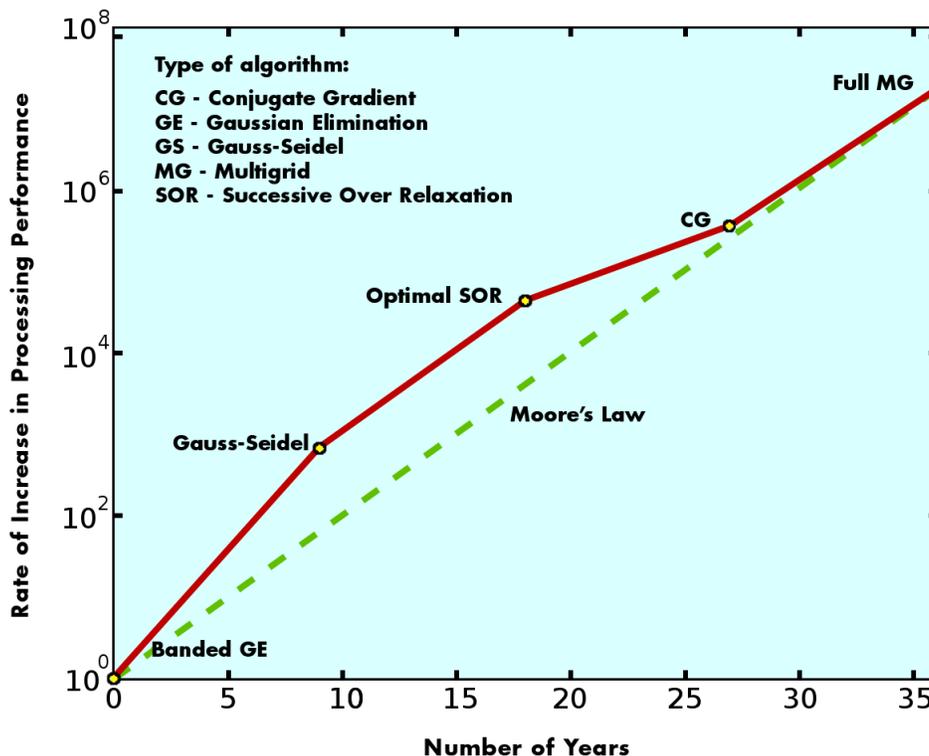
For the 3-dimensional case 7-point Laplacian on a cube, the matrix is  $n \times n$  with  $n = (N - 1)^3$ , with the bandwidth  $(N - 1)^2$ . In this case, elimination (e.g., Cholesky) would require storage  $O(nN^2) = O(N^5)$  and an operation count of  $O(nN^4) = O(N^7)$ .

Fortunately, far more efficient ways to solve the equations have been devised. In fact, algorithm improvements from the early 1960s to the late 1990s are estimate to account for a speed-up of about  $10^7$  when solving the 7-point Laplacian or similar problems on a  $64 \times 64 \times 64$  grid. This is summarized in the following table, taken from Figure 5, page 53 of *Computational Science: Ensuring America's Competitiveness*, a 2005 report to the President of the United States from the President's Information Technology Advisory Committee (PITAC). See also Figure 13, page 32 of the DOE Office of Science report *A science-based case for large-scale simulation*, 2003.

#### 1. Classical iterations

Gaussian elimination and its variants are called *direct methods*, meaning that they produce the exact solution of the linear system in finite number of steps. (This ignores the effects of round-off error, which is, in fact, a significant issue for some problems.) More efficient methods are *iterative methods*, which start from an initial guess  $u_0$  of the solution of the linear system, and produce a sequence  $u_1, u_2, \dots$ , of iterates which—hopefully—converge to the solution of the linear system. Stopping after a finite number of iterates, gives us an approximate solution to the linear system. This is very reasonable. Since the solution of the linear system is only an approximation for the solution of the PDE problem, there is little point in computing it exactly or nearly exactly. If the numerical discretization provides

FIGURE 3.1. Algorithmic speedup from early 1960s through late 1990s for solving the discrete Laplacian on a cubic mesh of size  $64 \times 64 \times 64$ . The comparison line labelled “Moore’s Law” is based on a speedup by a factor of two every 18 months.



about 4 significant digits, we would be happy if the linear solver provides 4 or maybe 5 digits. Further accuracy in the linear solver serves no purpose.

For an iterative method the goal is, of course, to design an iteration for which

- (1) the iteration is efficient, i.e., the amount of work to compute an iteration should not be too large: typically we want it to be proportional to the number  $n$  of unknowns;
- (2) the rate of convergence of the iterative method is fast, so that not too many iterations are needed.

First we consider some classical iterative methods to solve  $Au = f$ . One way to motivate such methods is to note that if  $u_0$  is some approximate solution, then the exact solution  $u$  may be written  $u = u_0 + e$  and the error  $e = u - u_0$  is related to the residual  $r = f - Au_0$  by the equation  $Ae = r$ . That is, we can express  $u$  as a *residual correction* to  $u_0$ :  $u = u_0 + A^{-1}(f - Au_0)$ . Of course, this merely rephrases the problem, since computing  $e = A^{-1}(f - Au_0)$  means solving  $Ae = r$  for  $e$ , which is as difficult as the original problem of solving  $Au = f$  for  $u$ . But suppose we can find some matrix  $B$  which approximates  $A^{-1}$  but is less costly to apply. We are then led to the iteration  $u_1 = u_0 + B(f - Au_0)$ , which can be repeated to give

$$(3.1) \quad u_{i+1} = u_i + B(f - Au_i), \quad i = 0, 1, 2, \dots$$

Of course the effectiveness of such a method will depend on the choice of  $B$ . For speed of convergence, we want  $B$  to be close to  $A^{-1}$ . For efficiency, we want  $B$  to be easy to apply. Some typical choices of  $B$  are:

- $B = \omega I$  for some  $\omega > 0$ . As we shall see, this method will converge for symmetric positive definite  $A$  if  $\omega$  is a sufficiently small positive number. This iteration is often called Richardson iteration.
- $B = D^{-1}$  where  $D$  is the diagonal matrix with the same diagonal elements as  $A$ . This is called the Jacobi method.
- $B = E^{-1}$  where  $E$  is the lower triangular matrix with the same diagonal and sub-diagonal elements of  $A$ . This is the Gauss–Seidel method.

Another way to derive the classical iterative methods, instead of residual correction, is to give a splitting of  $A$  as  $P + Q$  for two matrices  $P$  and  $Q$  where  $P$  is in some sense close to  $A$  but much easier to invert. We then write the equations as  $Pu = f - Qu$ , which suggests the iteration

$$u_{i+1} = P^{-1}(f - Qu_i).$$

Since  $Q = A - P$ , this iteration may also be written

$$u_{i+1} = u_i + P^{-1}(f - Au_i).$$

Thus this iteration coincides with (3.1) when  $B = P^{-1}$ .

Sometimes the iteration (3.1) is modified to

$$u_{i+1} = (1 - \alpha)u_i + \alpha[u_i + B(f - Au_i)], \quad i = 0, 1, 2, \dots,$$

for a real parameter  $\alpha$ . If  $\alpha = 1$ , this is the unmodified iteration. For  $0 < \alpha < 1$  the iteration has been *damped*, while for  $\alpha > 1$  the iteration is *amplified*. The damped Jacobi method will come up below when we study multigrid. The amplified Gauss–Seidel method is known as SOR (successive over-relaxation). This terminology is explained in the next two paragraphs.

Before investigating their convergence, let us particularize the classical iterations to the discrete Laplacian  $-\Delta_h^2$  in one or two dimensions. In one dimension, the equations are

$$\frac{-u^{m+1} + 2u^m - u^{m-1}}{h^2} = f^m, \quad m = 1, \dots, N-1,$$

where  $h = 1/N$  and  $u^0 = u^N = 0$ . The Jacobi iteration is then simply

$$u_{i+1}^m = \frac{u_i^{m-1} + u_i^{m+1}}{2} + \frac{h^2}{2}f^m, \quad m = 1, \dots, N-1,$$

The error satisfies

$$e_{i+1}^m = \frac{e_i^{m-1} + e_i^{m+1}}{2},$$

so at each iteration the error at a point is set equal to the average of the errors at the neighboring points at the previous iteration. The same holds true for the 5-point Laplacian in two dimensions, except that now there are four neighboring points. In an old terminology, updating the value at a point based on the values at the neighboring points is called *relaxing* the value at the point.

For the Gauss–Seidel method, the corresponding equations are

$$u_{i+1}^m = \frac{u_{i+1}^{m-1} + u_i^{m+1}}{2} + \frac{h^2}{2}f^m, \quad m = 1, \dots, N-1,$$

and

$$e_{i+1}^m = \frac{e_{i+1}^{m-1} + e_i^{m+1}}{2}, \quad m = 1, \dots, N-1.$$

We can think of the Jacobi method as updating the value of  $u$  at all the mesh points simultaneously based on the old values, while the Gauss–Seidel method updates the values of one point after another always using the previously updated values. For this reason the Jacobi method is sometimes referred to as *simultaneous relaxation* and the Gauss–Seidel method as *successive relaxation* (and amplified Gauss–Seidel as successive overrelaxation). Note that the Gauss–Seidel iteration gives different results if the unknowns are reordered. (In fact, from the point of view of convergence of Gauss–Seidel, there are better orderings than just the naive orderings we have taken so far.) By contrast, the Jacobi iteration is unaffected by reordering of the unknowns. The Jacobi iteration is very naturally a *parallel* algorithm: if we have many processors, each can independently update one or several variables.

Our next goal is to investigate the convergence of (3.1). Before doing so we make some preliminary definition and observations. First we recall that a sequence of vectors or matrices  $X_i$  *converges linearly* to a vector or matrix  $X$  if there exists a positive number  $r < 1$  and a number  $C$  such that

$$(3.2) \quad \|X - X_i\| \leq Cr^i, \quad i = 1, 2, \dots$$

In particular this holds (with  $C = \|X - X_0\|$ ) if  $\|X - X_{i+1}\| \leq r\|X - X_i\|$   $i = 0, 1, \dots$ . For a linearly convergent sequence, the *rate of linear convergence* is the infimum of all  $r$  for which there exists a  $C$  such that (3.2) holds. In a finite dimensional vector space, both the notion of linear convergence and the rate of linear convergence are independent of a choice of norm. In investigating iterative methods applied to problems with a mesh size parameter  $h$ , we will typically find that the rate of linear convergence depends on  $h$ . Typical is an estimate like  $\|X_i\| \leq Cr^i$  where all we can say about  $r$  is  $r \leq 1 - ch^p$  for some positive constants  $c$  and  $p$ . In order to interpret this, suppose that we want the error to be less than some tolerance  $\epsilon > 0$ . Thus we need to take  $m$  iterations with  $Cr^m \leq \epsilon$ , or  $r^m \leq C^{-1}\epsilon$ , or  $m \geq |\log(C^{-1}\epsilon)|/|\log r|$  (note that  $\log r < 0$  and  $\log(C^{-1}\epsilon) < 0$  unless already  $\|X - X_0\| \leq \epsilon$ ). Now, for  $r = 1 - ch^p$ ,  $|\log r| \approx |ch^p|$ , so the number of iterations needed will be about  $m = Kh^{-p}$ , with  $K = c^{-1}|\log(C^{-1}\epsilon)|$ . In short, linear convergence with rate  $r = 1 - ch^p$  means that the number of iterations required to reduce the error to a given tolerance will be  $O(h^{-p})$ .

Next we recall that the *spectrum*  $\sigma(G)$  of a matrix  $G \in \mathbb{R}^{n \times n}$  is its set of eigenvalues, a set of at most  $n$  complex numbers. The *spectral radius*  $\rho(G) = \max_{\lambda \in \sigma(G)} |\lambda|$ . Now consider the  $L^2$ -matrix norm  $\|G\|_2$  corresponding to the Euclidean norm on  $\mathbb{R}^n$ . Then

$$\|G\|_2^2 = \sup_{0 \neq x \in \mathbb{R}^n} \frac{(Gx)^T Gx}{x^T x} = \sup_{0 \neq x \in \mathbb{R}^n} \frac{x^T (G^T G)x}{x^T x} = \rho(G^T G),$$

( $G^T G$  is a symmetric positive semidefinite matrix and its spectral radius is the maximum of its Rayleigh quotient). That is,  $\|G\|_2 = \sqrt{\rho(G^T G)}$ . If  $G$  is symmetric, then  $G^T G = G^2$ , so its eigenvalues are just the squares of the eigenvalues of  $G$ , and  $\rho(G^T G) = \rho(G^2)$ , so  $\|G\|_2 = \rho(G)$ . Independently of whether  $G$  is symmetric or not, for any choice of norm on  $\mathbb{R}^n$ , the corresponding matrix norm certainly satisfies  $\|G\| \geq \rho(G)$ . The next theorem shows that we nearly have equality for some choice of norm.

**THEOREM 3.1.** *Let  $G \in \mathbb{R}^{n \times n}$  and  $\epsilon > 0$ . Then there exists a norm on  $\mathbb{R}^n$  such that the corresponding matrix norm satisfies  $\|G\| \leq \rho(G) + \epsilon$ .*

**PROOF.** We may use the Jordan canonical form to write  $SGS^{-1} = J$  where  $S$  is an invertible matrix and  $J$  has the eigenvalues of  $G$  on the diagonal, 0's and  $\epsilon$ 's on the first superdiagonal, and 0's everywhere else. (The usual Jordan canonical form is the case  $\epsilon = 1$ , but if we conjugate a Jordan block by the matrix  $\text{diag}(1, \epsilon, \epsilon^2, \dots)$  the 1's above the diagonal are changed to  $\epsilon$ .) We select as the vector norm  $\|x\| := \|Sx\|_\infty$ . This leads to  $\|G\| = \|SGS^{-1}\|_\infty = \|J\|_\infty \leq \rho(A) + \epsilon$  (the infinity matrix norm, is the maximum of the row sums).  $\square$

An important corollary of this result is a criterion for when the powers of a matrix tend to zero.

**THEOREM 3.2.** *For  $G \in \mathbb{R}^{n \times n}$ ,  $\lim_{i \rightarrow \infty} G^i = 0$  if and only if  $\rho(G) < 1$ , and in this case the convergence is linear with rate  $\rho(G)$ .*

**PROOF.** For any choice of vector norm  $\|G^n\| \geq \rho(G^n) = \rho(G)^n$ , so if  $\rho(G) \geq 1$ , then  $G^n$  does not converge to 0.

Conversely, if  $\rho(G) < 1$ , then for any  $\bar{\rho} \in (\rho(G), 1)$  we can find an operator norm so that  $\|G\| \leq \bar{\rho}$ , and then  $\|G^n\| \leq \|G\|^n = \bar{\rho}^n \rightarrow 0$ .  $\square$

We now apply this result to the question of convergence of the iteration (3.1), which we write as

$$u_{i+1} = (I - BA)u_i + Bf = Gu_i + Bf,$$

where the *iteration matrix*  $G = I - BA$ . The equation  $u = Gu + Bf$  is certainly satisfied (where  $u$  is the exact solution), and so we have another way to view a classical iteration: it is a one-point iteration for this fixed point equation. The error then satisfies  $e_{i+1} = Ge_i$ , and the method converges for all starting values  $e_0 = u - u_0$  if and only if  $\lim_{i \rightarrow \infty} G^i = 0$ , which, as we have just seen, holds if and only if  $\rho(G) < 1$ , in which case the convergence is linear with rate of linear convergence  $\rho(G)$ . Now the condition that the  $\rho(G) < 1$  means that all the eigenvalues of  $G = I - BA$  lie strictly inside the unit circle in the complex plane, or equivalently that all the eigenvalues of  $BA$  lie strictly inside the circle of radius 1 in the complex plane centered at the point 1. If  $BA$  has real eigenvalues, then the condition becomes that all the eigenvalues of  $BA$  belong to the interval  $(0, 2)$ . Note that, if  $A$  is symmetric positive definite (SPD) and  $B$  is symmetric, then  $BA$  is symmetric with respect to the inner product  $\langle u, v \rangle_A = u^T Av$ , so  $BA$  does indeed have real eigenvalues in that case.

As a first example, we consider the convergence of the Richardson method for an SPD matrix  $A$ . Since the matrix is SPD, it has a basis of eigenvectors with positive real eigenvalues

$$0 < \lambda_{\min}(A) = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = \lambda_{\max}(A) = \rho(A).$$

The eigenvalues of  $BA = \omega A$  are then  $\omega \lambda_i$ ,  $i = 1, \dots, n$ , and the iteration converges if and only if  $0 < \omega < 2/\lambda_{\max}$ .

**THEOREM 3.3.** *Let  $A$  be an SPD matrix. Then the Richardson iteration  $u_{m+1} = u_m + \omega(f - Au_m)$  is convergent for all choices of  $u_0$  if and only if  $0 < \omega < 2/\lambda_{\max}(A)$ . In this case the rate of convergence is*

$$\max(|1 - \omega \lambda_{\max}(A)|, |1 - \omega \lambda_{\min}(A)|).$$

Note that the optimal choice is given by  $\omega\lambda_{\max}(A) - 1 = 1 - \omega\lambda_{\min}(A)$ , i.e.,  $\omega_{\text{opt}} = 2/[\lambda_{\max}(A) + \lambda_{\min}(A)]$ , and, with this choice of  $\omega$ , the rate of convergence is

$$\frac{\lambda_{\max}(A) - \lambda_{\min}(A)}{\lambda_{\max}(A) + \lambda_{\min}(A)} = \frac{\kappa - 1}{\kappa + 1},$$

where  $\kappa = \lambda_{\max}(A)/\lambda_{\min}(A) = \|A\|_2\|A^{-1}\|_2$  is the *spectral condition number* of  $A$ . Of course, in practice we do not know the eigenvalues, so we cannot make the optimal choice. But even if we could, we would find very slow convergence when  $\kappa$  is large, as it typically is for discretizations of PDE.

For example, if we consider  $A = -D_h^2$ , then  $\lambda_{\min} \approx \pi^2$ ,  $\lambda_{\max} \approx 4/h^2$ , so  $\kappa = O(h^{-2})$ , and the rate of convergence is like  $1 - ch^2$  for some  $c$ . Thus the converge is indeed very slow (we will need  $O(h^{-2})$  iterations).

Note that for  $A = -D_h^2$  the Jacobi method coincides with the Richardson method with  $\omega = h^2/2$ . Since  $\lambda_{\max}(A) < 4/h^2$ , we have  $\omega < 2/\lambda_{\max}(A)$  and the Jacobi method is convergent. But again convergence is very slow, with a rate of  $1 - O(h^2)$ . In fact for any  $0 < \alpha \leq 1$ , the damped Jacobi method is convergent, since it coincides with the Richardson method with  $\omega = \alpha h^2/2$ .

For the Richardson, Jacobi, and damped Jacobi iterations, the approximate inverse  $B$  is symmetric, but this is not the case for Gauss–Seidel, in which  $B$  is the inverse of the lower triangle of  $A$ . Of course we get a similar method if we use  $B^T$ , the inverse of the upper triangle of  $A$ . If we take two steps of Gauss–Seidel, one with the lower triangle and one with the upper triangle, the iteration matrix is

$$(I - B^T A)(I - BA) = I - (B^T + B - B^T AB)A,$$

so this double iteration is itself a classical iteration with the approximate inverse

$$(3.3) \quad \bar{B} := B^T + B - B^T AB.$$

This iteration is called *symmetric Gauss–Seidel*. Now, from the definition of  $\bar{B}$ , we get the identity

$$(3.4) \quad \|v\|_A^2 - \|(I - BA)v\|_A^2 = \langle \bar{B}Av, v \rangle_A.$$

It follows that  $\langle \bar{B}Av, v \rangle_A \leq \|v\|_A^2$ , and hence that  $\lambda_{\max}(\bar{B}A) \leq 1$ . Thus the symmetrized Gauss–Seidel iteration is convergent if and only if  $\lambda_{\min}(\bar{B}A) > 0$ , i.e., if and only if  $\bar{B}A$  is SPD with respect to the  $A$  inner product. This is easily checked to be equivalent to  $\bar{B}$  being SPD with respect to the usual inner product. When this is the case (3.4) implies that  $\|(I - BA)v\|_A < \|v\|_A$  for all nonzero  $v$ , and hence the original iteration is convergent as well.

In fact the above argument didn't use any properties of the original approximate inverse  $B$ . So what we have really proved this more general theorem.

**THEOREM 3.4.** *Let  $u_{i+1} = u_i + B(f - Au_i)$  be an iterative method in residual correction form, and consider the symmetrized iteration, i.e.,  $u_{i+1} = u_i + \bar{B}(f - Au_i)$  with  $\bar{B}$  given by (3.3). Then the symmetrized iteration is convergent if and only if  $\bar{B}$  is SPD, and, in that case, the original iteration is convergent as well.*

Returning to Gauss–Seidel, we write  $A = L + D + L^T$  where  $D$  is diagonal and  $L$  strictly lower diagonal, so  $B = (L + D)^{-1}$  and

$$\begin{aligned}\bar{B} &= B^T + B - B^T A B = B^T (B^{-1} + B^{-T} - A) B \\ &= B^T [(L + D) + (L^T + D) - (L + D + L^T)] B = B^T D B,\end{aligned}$$

which is clearly SPD whenever  $A$  is. Thus we have proven:

**THEOREM 3.5.** *The Gauss–Seidel and symmetric Gauss–Seidel iterations are convergent for any symmetric positive definite linear system.*

It is worth remarking that the same result is *not* true for the Jacobi iteration: although convergence can be proven for many of the SPD matrices that arise from discretizations of PDE, it is easy to construct an SPD matrix for which Jacobi iteration does not converge. As to the speed of convergence, for Gauss–Seidel applied to the discrete Laplacian, the analysis is much trickier than for Jacobi, but it can again be proven (or convincingly demonstrated via simple numerical experiments) that for  $A = -D_h^2$  the rate of convergence is again is about  $1 - ch^2$ , as for Jacobi, although the constant  $c$  is about twice as big for Gauss–Seidel as for Jacobi.

For both of these iterations, applied to the 5-point Laplacian, the cost of an iteration is  $O(n) = O(N^2)$ , and we need  $O(h^{-2}) = O(N^2)$  iterations to achieve a given decrease in the error. Thus the total cost will be  $O(N^4)$  operations to achieve a given reduction factor, the same order as for banded Cholesky. In 3 dimensions, the situation is more favorable for the iterative methods. In this case, the cost of an iteration is  $O(n) = O(N^3)$ , and we will again need  $O(N^2)$  iterations, for a total cost of  $O(N^5)$ , compared to  $O(N^7)$  for banded Cholesky.

For SOR, the analysis is more complicated, but can be carried out in a similar way. A careful analysis for  $\Delta_h$ , which can be found in many texts, shows that there is an optimal value of the relaxation parameter  $\alpha$ , and for that value, the spectral radius behaves like  $1 - ch$  rather than  $1 - ch^2$ . This is significantly more efficient, giving  $O(N)$  rather than  $O(N^2)$  operations. However, in practice it can be difficult or impossible to find the optimal relaxation parameter, and the convergence is quite sensitive to the choice of parameter.

## 2. The conjugate gradient method

**2.1. Line search methods and the method of steepest descents.** We now restrict to the case where  $A$  is SPD. In this case the solution  $u$  of  $Au = f$  is also the unique minimizer of the function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$ ,

$$F(v) = \frac{1}{2} v^T A v - v^T f$$

This is a quadratic functional with a unique minimum, which can be found by solving the equation  $\nabla F(u) = 0$ , i.e.,  $Au = f$ . Now, for any  $v, w \in \mathbb{R}^n$ , we can write

$$\frac{1}{2} v^T A v = \frac{1}{2} [w + (v - w)]^T A [w + (v - w)] = \frac{1}{2} w^T A w + \frac{1}{2} (v - w)^T A (v - w) + (v - w)^T A w,$$

so

$$F(v) = F(w) + \frac{1}{2} (v - w)^T A (v - w) + (v - w)^T (A w - f).$$

If we take  $w = u$  the last term vanishes, giving

$$F(v) = F(u) + \frac{1}{2}(v - u)^T A(v - u),$$

which again shows that  $u$  is the unique minimizer of  $F$ , and helps us to visualize  $F(u)$ . Its graph is an upward opening paraboloid with vertex at  $v = u$  and height  $F(v) = -v^T Av/2$ .

Now one way to try to find a point in a vector space is through a *line search* method:

```

choose initial iterate  $u_0$ 
for  $i = 0, 1, \dots$ 
  choose  $s_i \in \mathbb{R}^n$ 
  choose  $\lambda_i \in \mathbb{R}$ 
  set  $u_{i+1} = u_i + \lambda_i s_i$ 
end

```

At each step the *search direction*  $s_i$  and *step length*  $\lambda_i$  are chosen to, hopefully, get us nearer to the desired solution vector. If the goal is to minimize a function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  (quadratic or not), a reasonable choice (but certainly not the only reasonable choice) of search direction is the direction of steepest descent of  $F$  at  $u_i$ , i.e.,  $s_i = -\nabla F(u_i)$ . In our quadratic case, the steepest descent direction is  $s_i = f - Au_i = r_i$ , the residual. Thus the Richardson iteration can be viewed as a line search method with steepest descent as search direction, and a fixed step size.

Also for a general minimization problem, for any choice of search direction, there is an obvious choice of stepsize, namely we can do an *exact line search* by minimizing the function of one variable  $\lambda \mapsto F(u_i + \lambda s_i)$ . Thus we must solve  $s_i^T \nabla F(u_i + \lambda s_i) = 0$ , which, in the quadratic case, gives

$$(3.5) \quad \lambda_i = \frac{s_i^T r_i}{s_i^T A s_i}.$$

If we choose the steepest descent direction with exact line search, we get  $s_i = r_i$ ,  $\lambda_i = r_i^T r_i / r_i^T A r_i$ , giving the *method of steepest descents*:

```

choose initial iterate  $u_0$ 
for  $i = 0, 1, \dots$ 
  set  $r_i = f - Au_i$ 
  set  $u_{i+1} = u_i + \frac{r_i^T r_i}{r_i^T A r_i} r_i$ 
end

```

Thus the method of steepest descents is a variant of the Richardson iteration  $u_{i+1} = u_i + \omega(f - Au_i)$  in which the parameter  $\omega$  depends on  $i$ . It does not fit in the category of simple iterations  $u_{i+1} = Gu_i + Bf$  with a fixed iteration matrix  $G$  which we analyzed in the previous section, so we shall need to analyze it by other means.

Let us consider the work per iteration of the method of steepest descents. As written above, it appears to require two matrix-vector multiplications per iteration, one to compute

$Ar_i$  used in defining the step length, and one to compute  $Au_i$  used to compute the residual, and one to compute  $Ar_i$  used in defining the step length. However, once we have computed  $p_i := Ar_i$  and the step length  $\lambda_i$  we can compute the next residual without an additional matrix-vector multiplication, since  $u_{i+1} = u_i + \lambda_i r_i$  implies that  $r_{i+1} = r_i - \lambda_i p_i$ . Thus we can write the algorithm as

```

choose  $u_0$ 
set  $r_0 = f - Au_0$ 
for  $i = 0, 1, \dots$ 
  set  $p_i = Ar_i$ 
  set  $\lambda_i = \frac{r_i^T r_i}{r_i^T p_i}$ 
  set  $u_{i+1} = u_i + \lambda_i r_i$ 
  set  $r_{i+1} = r_i - \lambda_i p_i$ 
end

```

Thus, for each iteration we need to compute one matrix-vector multiplication, two Euclidean inner products, and two operations which consist of a scalar-vector multiplication and a vector-vector additions (referred to as a SAXPY operation). The matrix-vector multiplication involves roughly one addition and multiplication for each nonzero in the matrix, while the inner products and SAXPY operations each involve  $n$  multiplications and additions. If  $A$  is sparse with  $O(n)$  nonzero elements, the entire per iteration cost is  $O(n)$  operations.

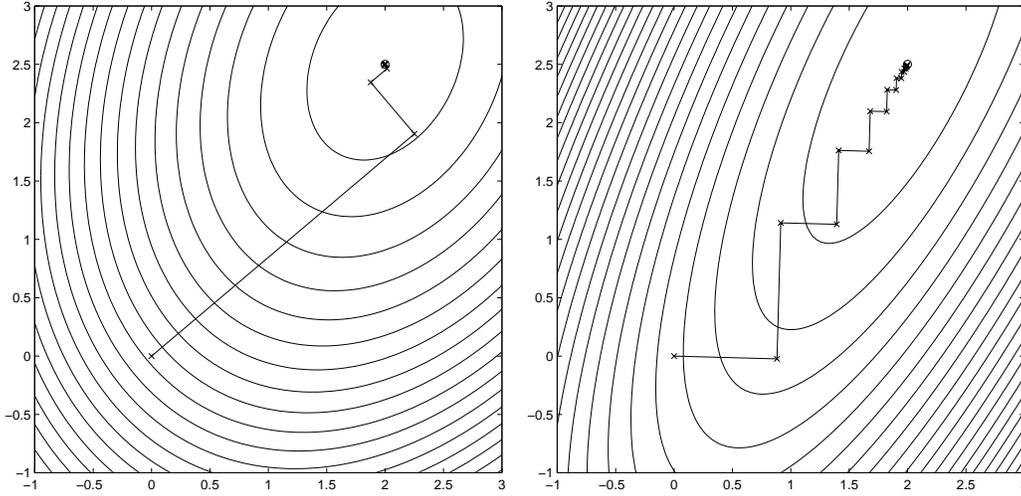
We shall show below that if the matrix  $A$  is SPD, the method of steepest descents converges to the solution of  $Au = f$  for any initial iterate  $u_0$ , and that the convergence is linear with the same rate of convergence as we found for Richardson extrapolation with the optimal parameter, namely  $(\kappa - 1)/(\kappa + 1)$  where  $\kappa$  is the spectral condition number of  $A$ . This means, again, that the convergence is slow if the condition number is large. This is quite easy to visualize already for  $2 \times 2$  matrices. See Figure 3.2.

**2.2. The conjugate gradient method.** The slow convergence of the method of steepest descents motivates a far superior line search method, the *conjugate gradient method*. CG also uses exact line search to choose the step length, but uses a more sophisticated choice of search direction than steepest descents.

For any line search method with exact line search,  $u_1 = u_0 + \lambda_0 s_0$  minimizes  $F$  over the 1-dimensional affine space  $u_0 + \text{span}[s_0]$ , and then  $u_2 = u_0 + \lambda_0 s_0 + \lambda_1 s_1$  minimizes  $F$  over the 1-dimensional affine space  $u_0 + \lambda_0 s_0 + \text{span}[s_1]$ . However  $u_2$  does not minimize  $F$  over the 2-dimensional affine space  $u_0 + \text{span}[s_0, s_1]$ . If that were the case, then for 2-dimensional problems we would have  $u_2 = u$  and we saw that that was far from the case for steepest descents.

However, it turns out that there is a simple condition on the search directions  $s_i$  that ensures that  $u_2$  is the minimizer of  $F$  over  $u_0 + \text{span}[s_0, s_1]$ , and more generally that  $u_i$  is the minimizer of  $F$  over  $u_0 + \text{span}[s_0, \dots, s_{i-1}]$ . Such a choice of search directions is very favorable. While we only need do 1-dimensional minimizations, after  $k$  steps we end up finding the minimizer in an  $k$ -dimensional space. In particular, as long as the search directions are linearly independent, this implies that  $u_n = u$ .

FIGURE 3.2. Convergence of steepest descents with a quadratic cost function. Left: condition number 2; right: condition number: 10.



THEOREM 3.6. Suppose that  $u_i$  are defined by exact line search using search directions which are  $A$ -orthogonal:  $s_i^T A s_j = 0$  for  $i \neq j$ . Then

$$F(u_i) = \min \{ F(u) \mid u \in u_0 + \text{span}[s_0, \dots, s_{i-1}] \}.$$

PROOF. Write  $W_i$  for  $\text{span}[s_0, \dots, s_{i-1}]$ , so  $u_i \in u_0 + W_i$  and we wish to prove that  $u_i$  minimizes  $F$  over  $u_0 + W_i$ . This is at least true for  $i = 1$ , since we use exact line search. The proof is by induction on  $i$ , so we assume that it is true and must prove that  $u_{i+1}$  minimizes  $F$  over  $u_0 + W_{i+1}$ . Now  $u_0 + W_0 = \{ y + \lambda s_i \mid y \in u_0 + W_i, \lambda \in \mathbb{R} \}$ , so we need to show that

$$F(u_{i+1}) = \min_{\substack{y \in u_0 + W_i \\ \lambda \in \mathbb{R}}} F(y + \lambda s_i).$$

The key point is that the function  $(y, \lambda) \mapsto F(y + \lambda s_i)$  decouples into the sum of a function of  $y$  which does not depend on  $\lambda$  plus a function of  $\lambda$  which does not depend on  $y$ . This is because  $u_i \in u_0 + W_i$ , so  $s_i^T A u_i = s_i^T A u_0 = s_i^T A y$  for any  $y \in u_0 + W_i$ , thanks to the  $A$ -orthogonality of the search directions. Thus

$$\begin{aligned} F(y + \lambda s_i) &= \frac{1}{2} y^T A y + \lambda s_i^T A y + \frac{\lambda^2}{2} s_i^T A s_i - y^T f - \lambda s_i^T f \\ &= F(y) + \left[ \frac{\lambda^2}{2} s_i^T A s_i - \lambda s_i^T (f - A u_i) \right]. \end{aligned}$$

Thus the minimum is obtained when  $y \in u_0 + W_i$  minimizes  $F(y)$ , which by the inductive hypothesis occurs when  $y = u_i$ , and when  $\lambda \in \mathbb{R}$  minimizes the term in brackets, which just gives us  $\lambda = s_i^T (f - A u_i) / s_i^T A s_i$ , the formula for exact line search. Thus the minimizer of  $F$  over  $u_0 + W_{i+1}$  is indeed  $u_i + \lambda_i s_i = u_{i+1}$ .  $\square$

Any method which uses  $A$ -orthogonal (also called “conjugate”) search directions has the nice property of the theorem. However it is not so easy to construct such directions. By far the most useful method is the method of conjugate gradients, or the CG method, which defines the search directions by  $A$ -orthogonalizing the residuals  $r_i = f - A u_i$ :

- $s_0 = r_0$
- $s_i = r_i - \sum_{j=0}^{i-1} \frac{s_j^T A r_i}{s_j^T A s_j} s_j.$

This sequence of search directions, together with the exact line search choice of step length (3.5) defines the conjugate gradient. The last formula (which is just the Gram-Schmidt procedure) appears to be quite expensive to implement and to involve a lot of storage, but fortunately we shall see that it may be greatly simplified.

- LEMMA 3.7. (1)  $W_i = \text{span}[s_0, \dots, s_{i-1}] = \text{span}[r_0, \dots, r_{i-1}].$   
 (2) *The residuals are  $l_2$ -orthogonal:  $r_i^T r_j = 0$  for  $i \neq j.$*   
 (3) *There exists  $m \leq n$  such that  $W_1 \subsetneq W_2 \subsetneq \dots \subsetneq W_m = W_{m+1} = \dots$  and  $u_0 \neq u_1 \neq \dots \neq u_m = u_{m+1} = \dots = u.$*   
 (4) *For  $i \leq m$ ,  $\{s_0, \dots, s_{i-1}\}$  is an  $A$ -orthogonal basis for  $W_i$  and  $\{r_0, \dots, r_{i-1}\}$  is an  $l_2$ -orthogonal basis for  $W_i.$*   
 (5)  $s_i^T r_j = s_i^T r_i = r_i^T r_i$  for  $0 \leq j \leq i.$

PROOF. The first statement comes directly from the definitions. To verify the second statement, note that, for  $0 \leq j < i$ ,  $F(u_i + tr_j)$  is minimal when  $t = 0$ , which gives  $r_j^T (A u_i - f) = 0$ , which is the desired orthogonality. For the third statement, certainly there is a least integer  $m \in [1, n]$  so that  $W_m = W_{m+1}$ . Then  $r_m = 0$  since it both belongs to  $W_m$  and is orthogonal to  $W_m$ . This implies that  $u_m = u$  and that  $s_m = 0$ . Since  $s_m = 0$   $u_{m+1} = u_m = u$ . Therefore  $r_{m+1} = 0$ , which implies that  $s_{m+1} = 0$ ,  $u_{m+2} = u$ , etc.

The fourth statement is an immediate consequence of the preceding ones. For the last statement, we use the orthogonality of the residuals to see that  $s_i^T r_i = r_i^T r_i$ . But, if  $0 \leq j \leq i$ , then

$$s_i^T r_j - s_i^T r_0 = s_i^T A(u_0 - u_j) = 0,$$

since  $u_0 - u_j \in W_i$ . □

Since  $s_i \in W_{i+1}$  and the  $r_j$ ,  $j \leq i$  are an orthogonal basis for that space for  $i < m$ , we have

$$s_i = \sum_{j=0}^i \frac{s_i^T r_j}{r_j^T r_j} r_j.$$

In view of part 5 of the lemma, we can simplify

$$s_i = r_i^T r_i \sum_{j=0}^i \frac{r_j}{r_j^T r_j} = r_i + r_i^T r_i \sum_{j=0}^{i-1} \frac{r_j}{r_j^T r_j},$$

whence

$$s_i = r_i + \frac{r_i^T r_i}{r_{i-1}^T r_{i-1}} s_{i-1}.$$

This is the formula which is used to compute the search direction. In implementing this formula it is useful to compute the residual from the formula  $r_{i+1} = r_i - \lambda_i A s_i$  (since  $u_{i+1} = u_i + \lambda_i s_i$ ). Putting things together we obtain the following implementation of CG:

choose initial iterate  $u_0$ , set  $s_0 = r_0 = f - Au_0$

**for**  $i = 0, 1, \dots$

$$\lambda_i = \frac{r_i^T r_i}{s_i^T A s_i}$$

$$u_{i+1} = u_i + \lambda_i s_i$$

$$r_{i+1} = r_i - \lambda_i A s_i$$

$$s_{i+1} = r_{i+1} + \frac{r_{i+1}^T r_{i+1}}{r_i^T r_i} s_i$$

**end**

At each step we have to perform one multiplication of a vector by  $A$ , two dot-products, and three SAXPYs, very similar to steepest descents (one more SAXPY). Here is the algorithm written out in full in pseudocode:

choose initial iterate  $u$

$$r \leftarrow f - Au$$

$$r2 \leftarrow r^T r$$

$$s \leftarrow r$$

**for**  $i = 0, 1, \dots$

$$t \leftarrow As \quad (\text{matrix multiplication})$$

$$s2 \leftarrow s^T t \quad (\text{dot product})$$

$$\lambda \leftarrow r2/s2$$

$$u \leftarrow u + \lambda s \quad (\text{SAXPY})$$

$$r2old \leftarrow r2$$

$$r \leftarrow r - \lambda t \quad (\text{SAXPY})$$

$$r2 \leftarrow r^T r \quad (\text{dot product})$$

$$s \leftarrow r + (r2/r2old)s \quad (\text{SAXPY})$$

**end**

The conjugate gradient method gives the exact solution in  $n$  iterations, but it is most commonly used as an iterative method and terminated with far fewer operations. A typical stopping criterion would be to test if  $r2$  is below a given tolerance. To justify this, we shall show that the method is linearly convergence and we shall establish the rate of convergence. For analytical purposes, it is most convenient to use the vector norm  $\|u\|_A := (u^T A u)^{1/2}$ , and its associated matrix norm.

We start with a third characterization of  $W_i = \text{span}[s_0, \dots, s_{i-1}] = \text{span}[r_0, \dots, r_{i-1}]$ .

LEMMA 3.8.  $W_i = \text{span}[r_0, Ar_0, \dots, A^{i-1}r_0]$  for  $i = 1, 2, \dots, m$ .

PROOF. Since  $\dim W_i = i$ , it is enough to show that  $W_i \subset \text{span}[r_0, Ar_0, \dots, A^{i-1}r_0]$ , which we do by induction. This is certainly true for  $i = 1$ . Assume it holds for some  $i$ . Then, since  $u_i \in u_0 + W_i$ ,  $r_i = f - Au_i \in r_0 + AW_i \in \text{span}[r_0, Ar_0, \dots, A^i r_0]$ , and therefore  $W_{i+1}$ , which is spanned by  $W_i$  and  $r_i$  belongs to  $\text{span}[r_0, Ar_0, \dots, A^i r_0]$ , which completes the induction.  $\square$

The space  $\text{span}[r_0, Ar_0, \dots, A^{i-1}r_0]$  is called the *Krylov space* generated by the matrix  $A$  and the vector  $r_0$ . Note that we have as well

$$W_i = \text{span}[r_0, Ar_0, \dots, A^{i-1}r_0] = \{p(A)r_0 \mid p \in \mathcal{P}_{i-1}\} = \{q(A)(u - u_0) \mid q \in \mathcal{P}_i, q(0) = 0\}.$$

Here  $\mathcal{P}_i$  denotes the space of polynomials of degree at most  $i$ . Since  $r_i$  is  $l_2$ -orthogonal to  $W_i$ ,  $u - u_i$  is  $A$ -orthogonal to  $W_i$ , so

$$\|u - u_i\|_A = \inf_{w \in W_i} \|u - u_i + w\|_A.$$

Since  $u_i - u_0 \in W_i$ ,

$$\inf_{w \in W_i} \|u - u_i + w\|_A = \inf_{w \in W_i} \|u - u_0 + w\|_A.$$

Combining the last three equations, we get

$$\|u - u_i\|_A = \inf_{\substack{q \in \mathcal{P}_i \\ q(0)=0}} \|u - u_0 + q(A)(u - u_0)\|_A = \inf_{\substack{p \in \mathcal{P}_i \\ p(0)=1}} \|p(A)(u - u_0)\|_A.$$

Applying the obvious bound  $\|p(A)(u - u_0)\|_A \leq \|p(A)\|_A \|u - u_0\|_A$  we see that we can obtain an error estimate for the conjugate gradient method by estimating

$$K = \inf_{\substack{p \in \mathcal{P}_i \\ p(0)=1}} \|p(A)\|_A.$$

Now if  $0 < \rho_1 < \dots < \rho_n$  are the eigenvalues of  $A$ , then the eigenvalues of  $p(A)$  are  $p(\rho_j)$ ,  $j = 1, \dots, n$ , and  $\|p(A)\|_A = \max_j |p(\rho_j)|$ . Thus<sup>1</sup>

$$K = \inf_{\substack{p \in \mathcal{P}_i \\ p(0)=1}} \max_j |p(\rho_j)| \leq \inf_{\substack{p \in \mathcal{P}_i \\ p(0)=1}} \max_{\rho_1 \leq \rho \leq \rho_n} |p(\rho)|.$$

The final infimum can be calculated explicitly, as will be explained below. Namely, for any  $0 < a < b$ , and integer  $n > 0$ ,

$$(3.6) \quad \min_{\substack{p \in \mathcal{P}_n \\ p(0)=1}} \max_{x \in [a, b]} |p(x)| = \frac{2}{\left(\frac{\sqrt{b/a+1}}{\sqrt{b/a-1}}\right)^n + \left(\frac{\sqrt{b/a-1}}{\sqrt{b/a+1}}\right)^n}.$$

This gives

$$K \leq \frac{2}{\left(\frac{\sqrt{\kappa+1}}{\sqrt{\kappa-1}}\right)^i + \left(\frac{\sqrt{\kappa-1}}{\sqrt{\kappa+1}}\right)^i} \leq 2 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^i,$$

where  $\kappa = \rho_n/\rho_1$  is the condition number of  $A$ . (To get the right-hand side, we suppressed the second term in the denominator of the left-hand side, which is less than 1 and tends to zero with  $i$ , and kept only the first term, which is greater than 1 and tends to infinity with  $i$ .) We have thus proven that

$$\|u - u_i\|_A \leq 2 \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^i \|u - u_0\|_A,$$

<sup>1</sup>Here we bound  $\max_j |p(\rho_j)|$  by  $\max_{\rho_1 \leq \rho \leq \rho_n} |p(\rho)|$  simply because we can minimize the latter quantity explicitly. However this does not necessarily lead to the best possible estimate, and the conjugate gradient method is often observed to converge faster than the result derived here. Better bounds can sometimes be obtained by taking into account the distribution of the spectrum of  $A$ , rather than just its minimum and maximum.

which is linear convergence with rate

$$r = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}.$$

Note that  $r \sim 1 - 2/\sqrt{\kappa}$  for large  $\kappa$ . So the convergence deteriorates when the condition number is large. However, this is still a notable improvement over the classical iterations. For the discrete Laplacian, where  $\kappa = O(h^{-2})$ , the convergence rate is bounded by  $1 - ch$ , not  $1 - ch^2$ .

The above analysis yields a convergence estimate for the method of steepest descent as well. Indeed, the first step of conjugate gradients coincides with steepest descents, and so, for steepest descents,

$$\|u - u_1\|_A \leq \frac{2}{\frac{\sqrt{\kappa+1}}{\sqrt{\kappa-1}} + \frac{\sqrt{\kappa-1}}{\sqrt{\kappa+1}}} \|u - u_0\|_A = \frac{\kappa - 1}{\kappa + 1} \|u - u_0\|_A.$$

Of course, the same result holds if we replace  $u_0$  by  $u_i$  and  $u_1$  by  $u_{i+1}$ . Thus steepest descents converges linearly, with rate  $(\kappa - 1)/(\kappa + 1)$  (just like Richardson iteration with the optimal parameter). Notice that the estimates indicate that a large value of  $\kappa$  will slow the convergence of both steepest descents and conjugate gradients, but, since the dependence for conjugate gradients is on  $\sqrt{\kappa}$  rather than  $\kappa$ , the convergence of conjugate gradients will usually be much faster.

The figure shows a plot of the norm of the residual versus the number of iterations for the conjugate gradient method and the method of steepest descents applied to a matrix of size 233 arising from a finite element simulation. The matrix is irregular, but sparse (averaging about 6 nonzero elements per row), and has a condition number of about 1,400. A logarithmic scale is used on the  $y$ -axis so the near linearity of the graph reflects linear convergence behavior. For conjugate gradients, the observed rate of linear convergence is between .7 and .8, and it takes 80 iterations to reduce the initial residual by a factor of about  $10^6$ . The convergence of steepest descents is too slow to be useful: in 400 iterations the residual is not even reduced by a factor of 2.

REMARK. There are a variety of conjugate-gradient-like iterative methods that apply to matrix problems  $Au = f$  where  $A$  is either indefinite, non-symmetric, or both. Many share the idea of approximation of the solution in a Krylov space.

Our analysis of conjugate gradients and steepest descents depended on the explicit solution of the minimization problem given in (3.6). Here we outline the proof of this result, leaving the details as an exercise.

The Chebyshev polynomials are defined by the recursion

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x) \text{ for } n = 1, 2, \dots,$$

so  $T_n$  is a polynomial of degree  $n$ . From this follows two explicit formulas for  $T_n$ :

$$T_n(x) = \cos(n \arccos x), \quad T_n(x) = \frac{1}{2}[(x + \sqrt{x^2 - 1})^n + (x - \sqrt{x^2 - 1})^n],$$

with the first equation valid for  $|x| \leq 1$  and the second valid for  $|x| \geq 1$ .

The polynomial  $T_n$  satisfies  $|T_n(x)| \leq 1$  on  $[-1, 1]$  with equality holding for  $n + 1$  distinct numbers in  $[-1, 1]$ . This can be used to establish the following: for any  $\alpha < -1$ , there does

FIGURE 3.3. Convergence of conjugate gradients for solving a finite element system of size 233. On the left 300 iterations are shown, on the right the first 50. Steepest descents is shown for comparison.

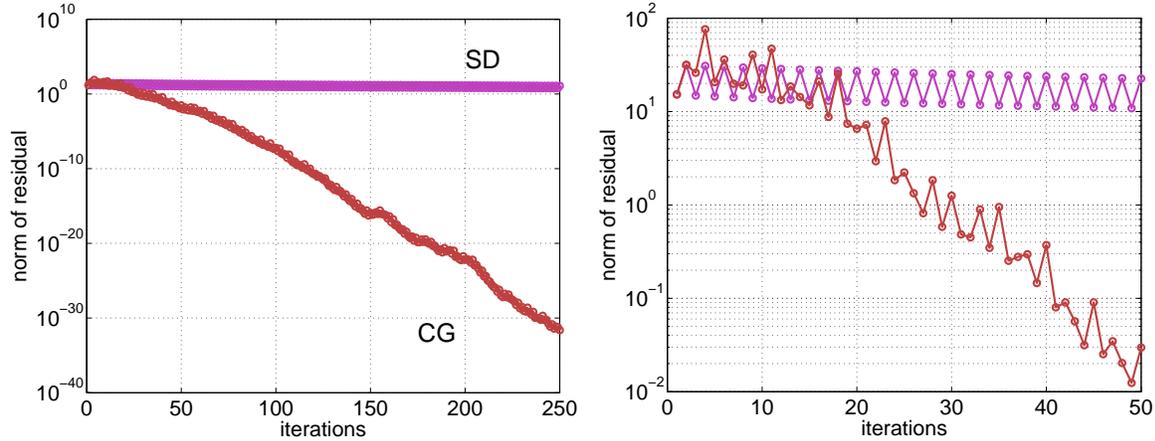
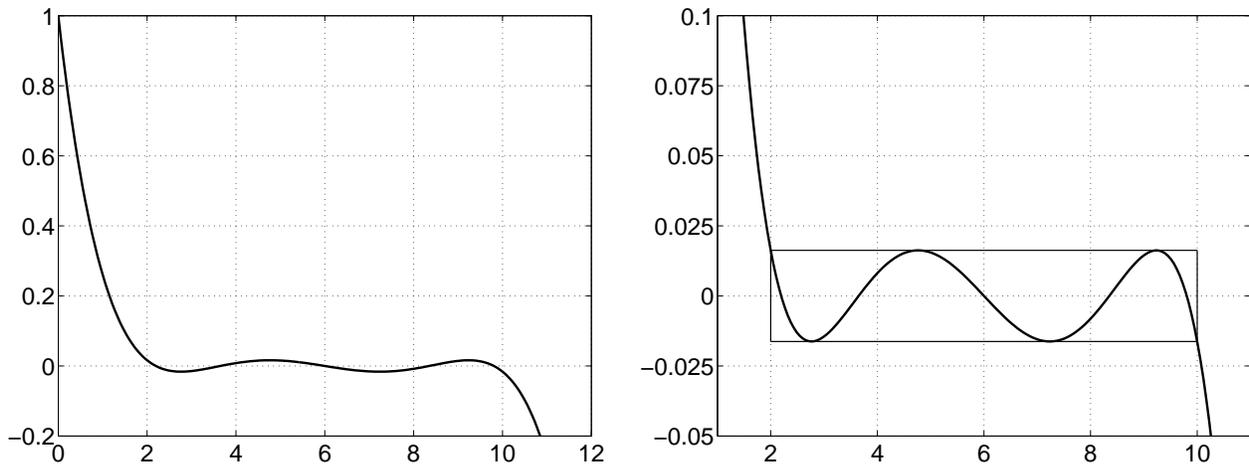


FIGURE 3.4. The quintic polynomial equal to 1 at 0 with the smallest  $L^\infty$  norm on  $[2, 10]$ . This is a scaled Chebyshev polynomial, and so the norm can be computed exactly.



not exist any polynomial  $q \in \mathcal{P}_n$  with  $q(\alpha) = T_n(\alpha)$  and  $|q(x)| < 1$  on  $[-1, 1]$ . In other words,  $T_n$  minimizes  $\max_{x \in [-1, 1]} |p(x)|$  over all polynomials in  $\mathcal{P}_n$  which take the value  $T_n(\alpha)$  at  $\alpha$ .

Scaling this result we find that

$$p(x) = \left[ T_n \left( -\frac{b+a}{b-a} \right) \right]^{-1} T_n \left( \frac{2x-b-a}{b-a} \right)$$

solves the minimization problem (3.6) and gives the minimum value claimed. This polynomial is plotted for  $n = 5$ ,  $a = 2$ ,  $b = 10$  in Figure 3.4.

**2.3. Preconditioning.** The idea is we choose a matrix  $M \approx A$  such that the system  $Mz = c$  is relatively easy to solve. We then consider the *preconditioned system*  $M^{-1}Ax = M^{-1}b$ . The new matrix  $M^{-1}A$  is SPD with respect to the  $M$  inner product, and we solve the preconditioned system using conjugate gradients but using the  $M$ -inner product in place of the  $l_2$ -inner product. Thus to obtain the preconditioned conjugate gradient algorithm, or PCG, we substitute  $M^{-1}A$  for  $A$  everywhere and change expressions of the form  $x^T y$  into  $x^T M y$ . Note that the  $A$ -inner product  $x^T A y$  remains invariant under these two changes. Thus we obtain the algorithm:

```

choose initial iterate  $u_0$ , set  $s_0 = \bar{r}_0 = M^{-1}f - M^{-1}Au_0$ 
for  $i = 0, 1, \dots$ 
   $\lambda_i = \frac{\bar{r}_i^T M \bar{r}_i}{s_i^T A s_i}$ 
   $u_{i+1} = u_i + \lambda_i s_i$ 
   $\bar{r}_{i+1} = \bar{r}_i - \lambda_i M^{-1} A s_i$ 
   $s_{i+1} = \bar{r}_{i+1} + \frac{\bar{r}_{i+1}^T M \bar{r}_{i+1}}{\bar{r}_i^T M \bar{r}_i} s_i$ 
end

```

Note that term  $s_i^T A s_i$  arises as the  $M$ -inner product of  $s_i$  with  $M^{-1} A s_i$ . The quantity  $\bar{r}_i$  is the residual in the preconditioned equation, which is related to the regular residual,  $r_i = f - Au_i$  by  $r_i = M \bar{r}_i$ . Writing PCG in terms of  $r_i$  rather than  $\bar{r}_i$  we get

```

choose initial iterate  $u_0$ , set  $r_0 = f - Au_0$ ,  $s_0 = M^{-1}r_0$ 
for  $i = 0, 1, \dots$ 
   $\lambda_i = \frac{r_i^T M^{-1} r_i}{s_i^T A s_i}$ 
   $u_{i+1} = u_i + \lambda_i s_i$ 
   $r_{i+1} = r_i - \lambda_i A s_i$ 
   $s_{i+1} = M^{-1} r_{i+1} + \frac{r_{i+1}^T M^{-1} r_{i+1}}{r_i^T M^{-1} r_i} s_i$ 
end

```

Thus we need to compute  $M^{-1}r_i$  at each iteration. Otherwise the work is essentially the same as for ordinary conjugate gradients. Since the algorithm is just conjugate gradients for the preconditioned equation we immediately have an error estimate:

$$\|u_i - u\|_A \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^i \|u_0 - u\|_A,$$

where  $\kappa$  now is the ratio of the largest to the least eigenvalue of  $M^{-1}A$ . To the extent that  $M$  approximates  $A$ , this ratio will be close to 1 and so the algorithm will converge quickly.

The matrix  $M$  is called the *preconditioner*. A good preconditioner should have two properties. First, it must be substantially easier to solve systems with the matrix  $M$  than with

the original matrix  $A$ , since we will have to solve such a system at each step of the preconditioned conjugate gradient algorithm. Second, the matrix  $M^{-1}A$  should be substantially better conditioned than  $A$ , so that PCG converges faster than ordinary CG. In short,  $M$  should be near  $A$ , but much easier to invert. Note that these conditions are similar to those we look for in defining a classical iteration via residual correction. If  $u_{i+1} = u_i + B(f - Au_i)$  is an iterative method for which  $B$  is SPD, then we might use  $M = B^{-1}$  as a preconditioner. For example, the Jacobi method suggests taking  $M$  to be the diagonal matrix with the same diagonal entries as  $A$ . When we compute  $M^{-1}r_i$  in the preconditioned conjugate gradient algorithm, we are simply applying one Jacobi iteration. Similarly we could use symmetric Gauss-Seidel to get a preconditioner.

In fact, we can show that conjugate gradients preconditioned by some SPD approximate inverse always converges faster than the corresponding classical iterative method. For if  $\lambda$  is an eigenvalue of  $BA$ , then  $-\rho \leq 1 - \lambda \leq \rho$  where  $\rho$  is the spectral radius of  $I - BA$ , and so

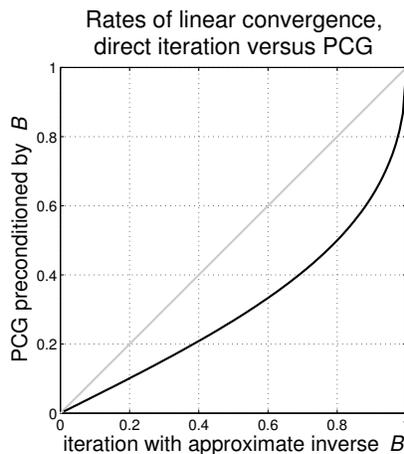
$$\lambda_{\min}(BA) \geq 1 - \rho, \quad \lambda_{\max}(BA) \leq 1 + \rho, \quad \kappa(BA) \leq \frac{1 + \rho}{1 - \rho}.$$

Thus the rate of convergence for the PCG method is at most

$$\frac{\sqrt{\kappa(BA)} - 1}{\sqrt{\kappa(BA)} + 1} \leq \frac{\sqrt{\frac{1+\rho}{1-\rho}} - 1}{\sqrt{\frac{1+\rho}{1-\rho}} + 1} = \frac{1 - \sqrt{1 - \rho^2}}{\rho}.$$

The last quantity is strictly less than  $\rho$  for all  $\rho \in (0, 1)$ ; see Figure 3.5. (For  $\rho$  small it is about  $\rho/2$ , while for the important case of  $\rho \approx 1 - \epsilon$  with  $\epsilon$  small, it is approximately  $1 - \sqrt{2\epsilon}$ .) Thus the rate of convergence of PCG with  $B$  as a preconditioner is better than that of the classical iteration with  $B$  as approximate inverse.

FIGURE 3.5. If an iteration achieves a rate of linear convergence  $\rho < 1$ , then the rate of convergence of conjugate gradients using the iteration as a preconditioner is bounded by  $(1 - \sqrt{1 - \rho^2})/\rho$ , which is always smaller.



Diagonal (Jacobi) preconditioning is often inadequate (in the case of the 5-point Laplacian it accomplishes nothing, since the diagonal is constant). Symmetric Gauss-Seidel is somewhat better, but often insufficient as well. A third possibility which is often applied

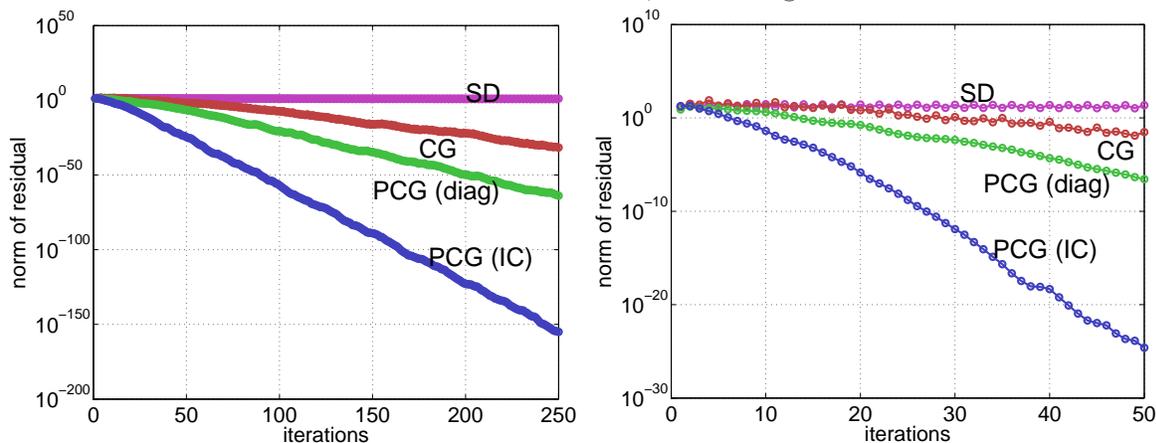
when  $A$  is sparse is to determine  $M$  via the *incomplete Cholesky factorization*. This means that a triangular matrix  $L$  is computed by the Cholesky algorithm applied to  $A$ , except that no fill-in is allowed: only the non-zero elements of  $A$  are altered, and the zero elements left untouched. One then takes  $M = LL^T$ , and, so  $M^{-1}$  is easy to apply. Yet, other preconditioners take into account the source of the matrix problem. For example, if a matrix arises from the discretization of a complex partial differential equation, we might precondition it by the discretization matrix for a simpler related differential equation (if that lead to a linear systems which is easier to solve). In fact the derivation of good preconditioners for important classes of linear systems remain a very active research area.

We close with numerical results for preconditioned conjugate gradients with both the diagonal preconditioner and incomplete Cholesky factorization as preconditioner. In Figure 3.6 we reproduce the results shown in Figure 3.3, together with these preconditioned iterations. By fitting the log of the norm of the residual to a linear polynomial, we can compute the observed rates of linear convergence. They are:

steepest descents	0.997	PCG (diag.)	0.529
conjugate gradients	0.725	PCG (IC)	0.228

The preconditioned methods are much more effective. Diagonal preconditioning reduces the number of iterations needed by conjugate gradients to reduce the initial error by a factor of  $10^{-6}$  from 80 to 44. Incomplete Cholesky preconditioning reduces further to 18 iterations.

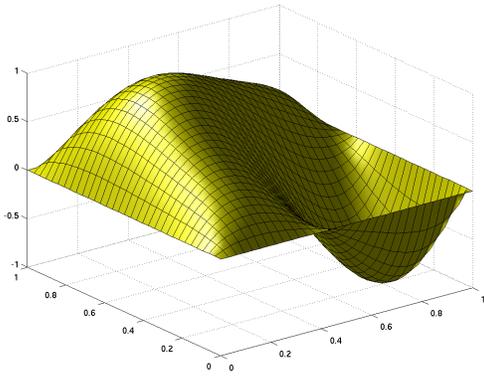
FIGURE 3.6. Convergence of conjugate gradients for solving a finite element system of size 233, unpreconditioned, diagonally preconditioned, and preconditioned by incomplete Cholesky factorization. Steepest descents is shown as well. On the left 300 iterations are shown, on the right the first 50.



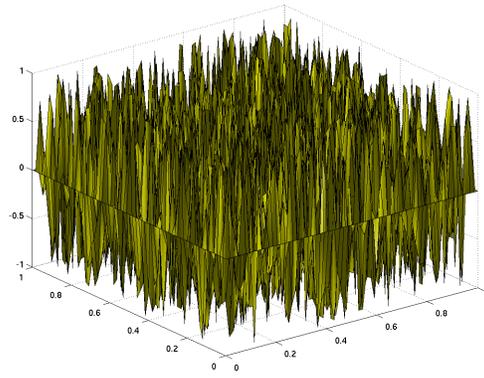
### 3. Multigrid methods

Figure 3.7 shows the result of solving a discrete system of the form  $-\Delta_h u_h = f$  using the Gauss–Seidel iteration. We have take  $h = 64$ , and chosen a smooth right-hand side vector  $f$  which results in the vector  $u_h$  which is shown in the first plot. The initial iterate  $u_0$ , which is shown in the second plot, was chosen at random, and then the iterates  $u_1, u_2, u_{10}, u_{50}$ , and  $u_{500}$  are shown in the subsequent plots. In Figure 3.8, the maximum norm error  $\|u_h - u_i\|/\|u_h\|$  is plotted for  $i = 0, 1, \dots, 50$ .

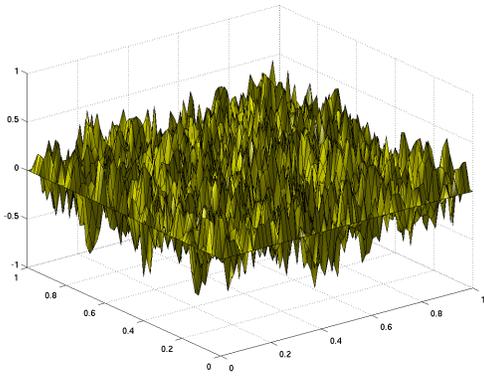
FIGURE 3.7. Iterative solution to  $-\Delta_h u_h = f$ ,  $h = 1/64$ , using Gauss–Seidel. The random initial iterate is rapidly smoothed, but approaches the solution  $u_h$  only very slowly.



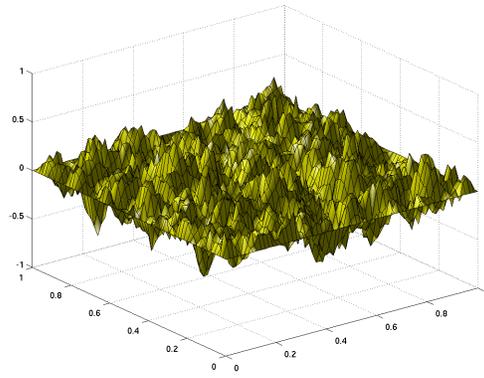
exact solution



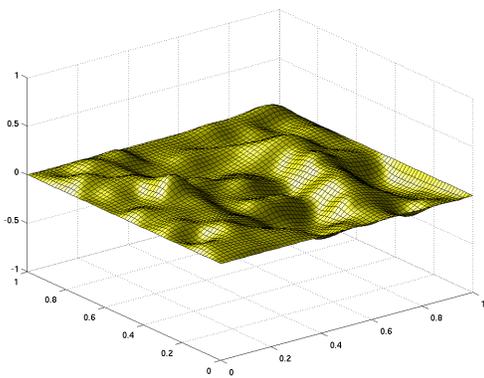
initial iterate



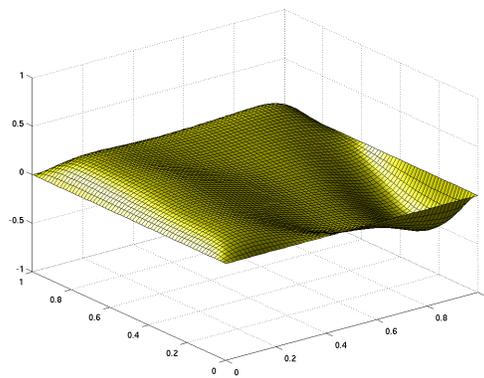
iterate 1



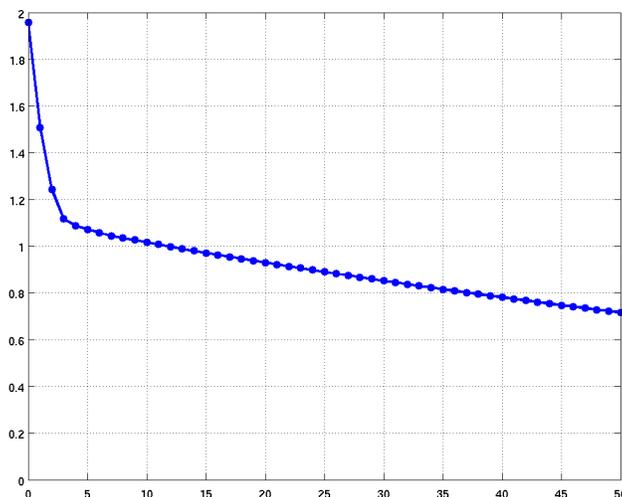
iterate 2



iterate 10



iterate 50

FIGURE 3.8. Error in the Gauss–Seidel iterates 0 through 50 in  $l^\infty$  ( $\bullet$ ).

These numerical experiments illustrate the following qualitative properties, which are typical of the Gauss–Seidel iteration applied to matrices arising from the discretization of elliptic PDEs.

- If we start with a random error, the norm of the error will be reduced fairly quickly for the first few iterations, but the error reduction occurs much more slowly after that.
- After several iterations the error is much smoother, but not much smaller, than initially. Otherwise put, the highly oscillatory modes of the error are suppressed much more quickly by the iteration than the low frequency modes.

The first observation is valid for all the methods we have studied: Richardson, Jacobi, damped Jacobi, and Gauss–Seidel. The second observation—that Gauss–Seidel iteration *smooths* the error—is shared damped Jacobi with  $\alpha < 1$ , but not by Jacobi itself.

If we take the Richardson method with  $\omega = 1/\lambda_{\max}(A)$  for the operator  $A = -D_h^2$ , it is very easy to see how the smoothing property comes about. The initial error can be expanded in terms of the eigenfunctions of  $A$ :  $e_0 = \sum_{m=1}^n c_m \sin m\pi x$ . The  $m$ th component in this expansion is multiplied by  $1 - \lambda_m/\lambda_{\max} = 1 - \lambda_m/\lambda_n$  at each iteration. Thus the high frequency components,  $m \approx n$ , are multiplied by something near to 0 at each iteration, and so are damped very quickly. Even the intermediate eigenvalues,  $\lambda_m \approx \lambda_n/2$  are damped reasonably quickly (by a factor of about 1/2 at each iteration). But the low frequency modes, for which  $\lambda_m \ll \lambda_n$ , decrease very slowly.

This also explains the first observation, that the norm of the error decreases quickly at first, and then more slowly. The norm of the error has contributions from all modes present in the initial error. Those associated to the higher frequency modes disappear in a few iterations, bringing the error down by a significant fraction. But after that the error is dominated by the low frequency modes, and so decays very slowly.

The same analysis applies to damped Jacobi with positive damping, and shows that undamped Jacobi doesn't have the smoothing property: the  $m$ th mode is multiplied by about  $1 - 2\lambda_m/\lambda_n$ , and so convergence is very slow for low frequency modes and also the highest frequency modes  $\lambda_m \approx \lambda_n$ . For the intermediate modes,  $\lambda_m \approx \lambda_n/2$ , convergence is very fast.

Establishing the smoothing property for Gauss–Seidel is more complicated, since the eigenfunctions of the Gauss–Seidel iteration don't coincide with those of  $A$  even for  $A = -D_h^2$ . However both numerical study and careful analysis show that Gauss–Seidel does indeed have the smoothing property for discretized elliptic operators.

The idea behind the multigrid method is to create an iterative method which reduces all components of the residual quickly by putting together two steps. First it applies the approximate inverse from Gauss–Seidel or another classical iterative method with the smoothing property to the residual. This greatly reduces the high frequency components of the residual, but barely reduces the low frequency components. The new residual, being relatively smooth, can then be accurately approximated on a coarser mesh. So, for the second step, the residual is (somehow) transferred to a coarser mesh, and the equation solved there, thus reducing the low frequency components. On the coarser mesh, it is of course less expensive to solve. For simplicity, we assume for now that an exact solver is used on the coarse mesh. Finally this coarse mesh solution to the residual problem is somehow transferred back to the fine mesh where it can be added back to our smoothed approximation.

Thus we have motivated the following rough outline of an algorithm:

- (1) Starting from an initial guess  $u_0$  apply a fine mesh smoothing iteration to get an improved approximation  $\bar{u}$ .
- (2) Transfer the residual in  $\bar{u}$  to a coarser mesh, solve a coarse mesh version of the problem there, transfer the solution back to the fine mesh, and add it back to  $\bar{u}$  to get  $\bar{\bar{u}}$ .

Taking  $\bar{\bar{u}}$  for  $u_1$  and thus have described an iteration to get from  $u_0$  to  $u_1$  (which we can then apply again to get from  $u_1$  to  $u_2$ , and so on). In fact it is much more common to also apply a fine mesh smoothing at the end of the iteration, i.e., to add a third step:

- (3) Starting from  $\bar{\bar{u}}$  apply the smoothing iteration to get an improved approximation  $\bar{\bar{\bar{u}}}$ .

The point of including the third step is that it leads to a multigrid iteration which is symmetric, which is often advantageous (e.g., the iteration can be used as a preconditioner for conjugate gradients). If the approximation inverse  $B$  used for the first smoothing step is not symmetric, we need to apply  $B^T$  (which is also an approximate inverse, since  $A$  is symmetric) to obtain a symmetric iteration.

We have just described a *two-grid* iteration. The true multigrid method will involve not just the original mesh and one coarser mesh, but a whole sequence of meshes. However, once we understand the two-grid iteration, the multigrid iteration will follow easily.

To make the two-grid method more precise we need to explain step 2 more fully, namely (a) how do we transfer the residual from the fine mesh to the coarse mesh?; (b) what problem do we solve on the coarse mesh?; and (c) how do we transfer the solution of that problem from the coarse mesh to the fine mesh? For simplicity, we suppose that  $N = 1/h$  is even and that we are interested in solving  $A_h u = f$  where  $A = -D_h^2$ . Let  $H = 2h = (N/2)^{-1}$ . We will use the mesh of size  $H$  as our coarse mesh. The first step of our multigrid iteration

is then just

$$\bar{u} = u_0 + B_h(f - A_h u_0),$$

where  $B_h$  is just the approximate inverse of  $A_h$  from Gauss–Seidel or some other smoothing iteration. The resulting residual is  $f - A_h \bar{u}$ . This is a function on the fine mesh points  $h, 2h, \dots, (N-1)h$ , and a natural way to transfer it to the coarse mesh is restrict it to the even grid points  $2h, 4h, \dots, (N-2)h = H, 2H, \dots, (N/2-1)H$ , which are exactly the coarse mesh grid points. Denoting this *restriction operator* from fine grid to coarse grid functions (i.e., from  $\mathbb{R}^{N-1} \rightarrow \mathbb{R}^{N/2-1}$ ) by  $P_H$ , we then solve  $A_H e_H = P_H(f - A_h \bar{u}_h)$  where, of course,  $A_H = -D_H^2$  is the 3-point difference operator on the coarse mesh. To transfer the solution  $e_H$ , a coarse grid function, to the fine grid, we need a *prolongation operator*  $Q_H : \mathbb{R}^{N/2-1} \rightarrow \mathbb{R}^{N-1}$ . It is natural to set  $Q_H e_H(jh) = e_H(jh)$  if  $j$  is even. But what about when  $j$  is odd: how should we define  $Q_H e_H$  at the midpoint of two adjacent coarse mesh points? A natural choice, which is simple to implement, is  $Q_H e_H(jh) = [e_H((j-1)h) + e((j+1)h)]/2$ . With these two operators second step is

$$\bar{\bar{u}} = \bar{u} + Q_H A_H^{-1} P_H(f - A_h \bar{u}).$$

And then final post-smoothing step is

$$\bar{\bar{\bar{u}}} = \bar{\bar{u}} + B_h^T(f - A_h \bar{\bar{u}}).$$

Actually this does not give a symmetric iteration. To obtain symmetry we need  $Q_h = cP_H^T$  and that is not the case for the grid transfer operators we defined. We have

$$(3.7) \quad Q_H = \begin{pmatrix} 1/2 & 0 & 0 & 0 & \cdots & 0 \\ 1 & 0 & 0 & 0 & \cdots & 0 \\ 1/2 & 1/2 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 1/2 & 1/2 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1/2 \end{pmatrix},$$

but  $P_H$  as we described it, consists only of 0's and 1's. Therefore one commonly takes a different choice for  $P_H$ , namely  $P_H = (1/2)Q_H^T$ . This means that the transferred coarse grid function doesn't just take the value of the corresponding fine grid function at the coarse grid point, but rather uses a weighted average of the fine grid function's values at the point in question and the fine grid points to the left and right (with weights  $1/4, 1/2, 1/4$ ). With this choice,  $Q_H A_h P_H$  is symmetric; in fact,  $Q_H A_h P_H = A_H$ . This is a useful formula. For operators other than the  $A_h = -D_h^2$ , we can use the same intergrid transfer operators, namely  $Q_H$  given by (3.7) and  $P_H = (1/2)Q_H^T$ , and then define the coarse grid operator by  $A_H = Q_H A_h P_H$ .

REMARK. In a finite element context, the situation is simpler. If the fine mesh is a refinement of the coarse mesh, then a coarse mesh function is already a fine mesh function. Therefore, the operator  $Q_H$  can be taken simply to be the inclusion operator of the coarse mesh space into the fine mesh space. The residual in  $u_0 \in S_h$  is most naturally viewed as a functional on  $S_h$ :  $v \mapsto (f, v) - B(u_0, v)$ . It is then natural to transfer the residual to the coarse mesh simply by restricting the test function  $v$  to  $S_H$ . This operation  $S_h^T \rightarrow S_H^T$  is exactly the adjoint of the inclusion operator  $S_H \rightarrow S_h$ . Thus the second step, solving the

coarse mesh problem for the restricted residual is obvious in the finite element case: we find  $e_H \in S_H$  such that

$$B(e_H, v) = (f, v) - B(\bar{u}, v), \quad v \in S_H,$$

and then we set  $\bar{\bar{u}} = \bar{u} + e_H \in S_h$ .

Returning to the case of finite differences we have arrived at the following two-grid iterative method to solve  $A_h u_h = f_h$ .

---

$u_h = \text{twogrid}(h, A_h, f_h, u_0)$   
*input:*  $h$ , mesh size ( $h = 1/n$  with  $n$  even)  
 $A_h$ , operator on mesh functions  
 $f_h$ , mesh function (right-hand side)  
 $u_0$ , mesh function (initial iterate)  
*output:*  $u_h$ , mesh function (approximate solution)

---

**for**  $i = 0, 1, \dots$  until satisfied  
 1. presmoothing:  $\bar{u} = u_i + B_h(f_h - A_h u_i)$   
 2. coarse grid correction:  
 2.1. residual computation:  $r_h = f_h - A_h \bar{u}$   
 2.2. restriction:  $H = 2h, r_H = P_H r_h, A_H = P_H A_h Q_H$   
 2.3. coarse mesh solve: solve  $A_H e_H = r_H$   
 2.4. prolongation:  $e_h = Q_H e_H$   
 2.5. correction:  $\bar{\bar{u}} = \bar{u} + e_h$   
 3. postsmoothing:  $u_h \leftarrow u_{i+1} = \bar{\bar{u}} + B_h^T(f_h - A_h \bar{\bar{u}})$   
**end**

---

Algorithm 3.1: Two-grid iteration for approximately solving  $A_h u_h = f_h$ .

In the smoothing steps, the matrix  $B_h$  could be, for example,  $(D - L)^{-1}$  where  $D$  is diagonal,  $L$  strictly lower triangular, and  $A_h = D - L - L^T$ . This would be a Gauss–Seidel smoother, but there are other possibilities as well. Besides these steps, the major work is in the coarse mesh solve. To obtain a more efficient algorithm, we may also solve on the coarse mesh using a two-grid iteration, and so involving an even coarser grid. In the following multigrid algorithm, we apply this idea recursively, using multigrid to solve at each mesh level, until we get to a sufficiently coarse mesh,  $h = 1/2$ , at which point we do an exact solve (with a  $1 \times 1$  matrix!).

---

```

 $u_h = \text{multigrid}(h, A_h, f_h, u_0)$ 
  input:  $h$ , mesh size ( $h = 1/n$  with  $n$  a power of 2)
          $A_h$ , operator on mesh functions
          $f_h$ , mesh function (right-hand side)
          $u_0$ , mesh function (initial iterate)
  output:  $u_h$ , mesh function (approximate solution)

```

---

```

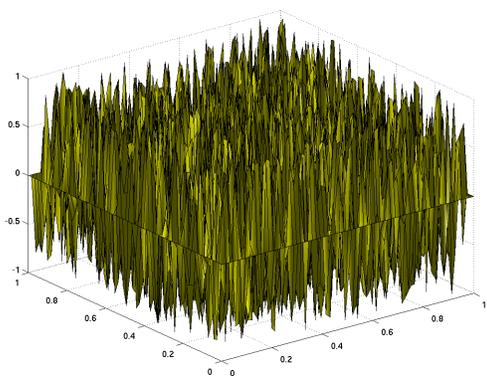
if  $h = 1/2$  then
   $u_h = A_h^{-1} f_h$ 
else
  for  $i = 0, 1, \dots$  until satisfied
    1. presmoothing:  $\bar{u} = u_i + B_h(f - A_h u_i)$ 
    2. coarse grid correction:
      2.1. residual computation:  $r_h = f_h - A_h \bar{u}$ 
      2.2. restriction:  $H = 2h, r_H = P_H r_h, A_H = P_H A_h Q_H$ 
      2.3. coarse mesh solve:  $e_H = \text{multigrid}(H, A_H, r_H, 0)$ 
      2.4. prolongation:  $e_h = Q_H e_H$ 
      2.5. correction:  $\bar{\bar{u}} = \bar{u} + e_h$ 
    3. postsmoothing:  $u_h \leftarrow u_{i+1} = \bar{\bar{u}} + B_h^T(f - A_h \bar{\bar{u}})$ 
  end
end if

```

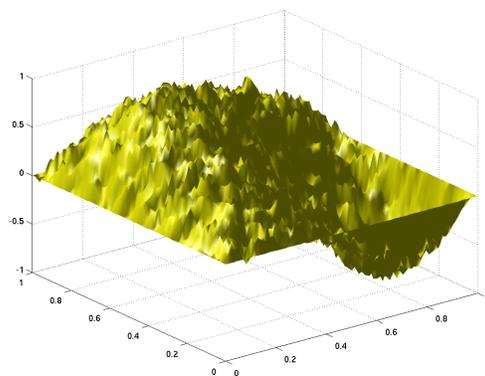
---

Algorithm 3.2: Multigrid iteration for approximately solving  $A_h u_h = f$ .

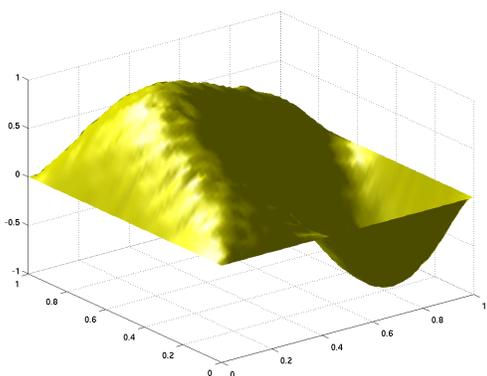
Figure 3.9 shows 5 iterations of this multigrid algorithm for solving the system  $-\Delta_h u_h = f$ ,  $h = 1/64$ , considered at the beginning of this section, starting from a random initial guess (we would get even better results starting from a zero initial guess). Compare with Figure 3.7. The fast convergence of the multigrid algorithm is remarkable. Indeed, for the multigrid method discussed here it is possible to show that the iteration is linearly convergent with a rate independent of the mesh size (in this example, it is roughly 0.2). This means that the number of iterations needed to obtain a desired accuracy remains bounded independent of  $h$ . It is also easy to count the number of operations per iteration. Each iteration involves two applications of the smoothing iteration, plus computation of the residual, restriction, prolongation, and correction on the finest mesh level. All those procedures cost  $O(n)$  operations. But then, during the coarse grid solve, the same procedures are applied on the grid of size  $2h$ , incurring an additional cost of  $O(n/2)$ . Via the recursion the work will be incurred for each mesh size  $h, 2h, 4h, \dots$ . Thus the total work per iteration will be  $O(n + n/2 + n/4 + \dots + 1) = O(n)$  (since the geometric series sums to  $2n$ ). Thus the total work to obtain the solution of the discrete system to any desired accuracy is itself  $O(n)$ , i.e., optimal.

FIGURE 3.9. Iterative solution to  $-\Delta_h u_h = f$ ,  $h = 1/64$ , using multigrid.

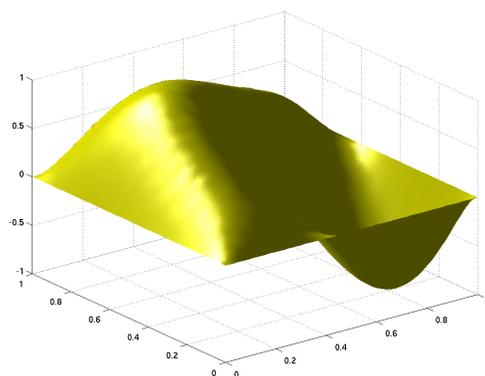
initial iterate



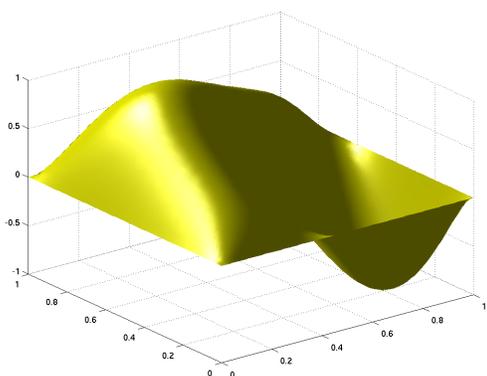
iterate 1



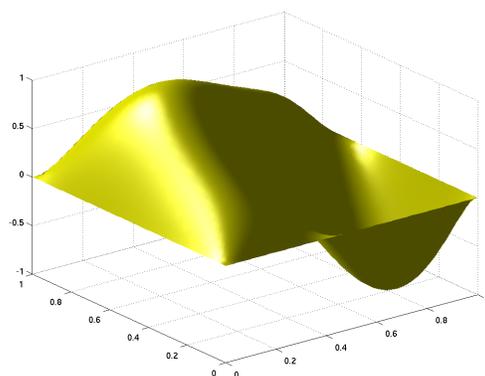
iterate 2



iterate 3



iterate 4



iterate 5



## CHAPTER 4

### Finite element methods for elliptic equations

#### 1. Weak and variational formulations

Model PDE:  $-\operatorname{div} a \operatorname{grad} u + cu = f$  in  $\Omega$

Here  $\Omega$  is a bounded domain in  $\mathbb{R}^n$ ;  $0 < \underline{a} \leq a(x) \leq \bar{a}$ ,  $0 \leq c(x) \leq \bar{c}$

First consider the homogeneous Dirichlet BC:  $u = 0$  on  $\partial\Omega$ .

Assuming that  $a \in C^1(\bar{\Omega})$ ,  $c \in C(\bar{\Omega})$ ,  $u \in C^2(\bar{\Omega})$  satisfies the PDE and BC (a *strong solution*), then it also satisfies the *weak formulation*:

Find  $u \in \mathring{H}^1(\Omega)$  such that

$$\int (a \operatorname{grad} u \cdot \operatorname{grad} v + cuv) = \int f v, \quad v \in \mathring{H}^1(\Omega),$$

A solution of the weak formulation need not belong to  $C^2(\bar{\Omega})$ , but if it does, then it is a strong solution.

The *variational formulation* is completely equivalent to the weak formulation

$$u = \operatorname{argmin}_{v \in \mathring{H}^1(\Omega)} \int_{\Omega} \left[ \frac{1}{2} (a \operatorname{grad} v \cdot \operatorname{grad} v + cv^2) - f v \right]$$

Extensions: Neumann BC, Robin BC, mixed BC, inhomogeneous Dirichlet BC. First order term to the PDE (then the problem is not symmetric and there is no variational formulation, but weak formulation is fine).

All these problems can be put in the weak form: Find  $u \in V$  such that

$$(4.1) \quad b(u, v) = F(v), \quad v \in V,$$

where  $V$  is a Hilbert space ( $H^1$  or  $\mathring{H}^1$ ),  $b : V \times V \rightarrow \mathbb{R}$  is a bilinear form,  $F : V \rightarrow \mathbb{R}$  is a linear form. (The inhomogeneous Dirichlet problem takes this form if we solve for  $u - u_g$  where  $u_g$  is a function satisfying the inhomogeneous Dirichlet BC  $u_g = g$  on  $\partial\Omega$ .) For symmetric problems (no first order term), the bilinear form  $b$  is symmetric, and the weak form is equivalent to the variational form:

$$u = \operatorname{argmin}_{v \in V} \left[ \frac{1}{2} b(v, v) - F(v) \right].$$

## 2. Galerkin method and finite elements

Let  $V_h$  be a finite dimensional subspace of  $V$ . If we replace the  $V$  in the weak formulation with  $V_h$  we get a discrete problem: Find  $u_h \in V_h$  such that

$$(4.2) \quad b(u_h, v) = F(v), \quad v \in V_h.$$

This is called the *Galerkin method*. For symmetric problems it is equivalent to the *Rayleigh–Ritz* method, which replaces  $V$  by  $V_h$  in the variational formulation:

$$u_h = \operatorname{argmin}_{v \in V_h} \left[ \frac{1}{2} b(v, v) - F(v) \right].$$

The Galerkin solution can be reduced to a set of  $n$  linear equations in  $n$  unknowns where  $n = \dim V_h$  by choosing a basis. Adopting terminology from elasticity, the matrix is called the *stiffness matrix* and the right hand side is the *load vector*.

Comparing (4.1) and (4.2), we find that the error in the Galerkin method  $u - u_h$  satisfies

$$(4.3) \quad b(u - u_h, v) = 0, \quad v \in V_h.$$

This relation, known as *Galerkin orthogonality*, is key to the analysis of Galerkin methods.

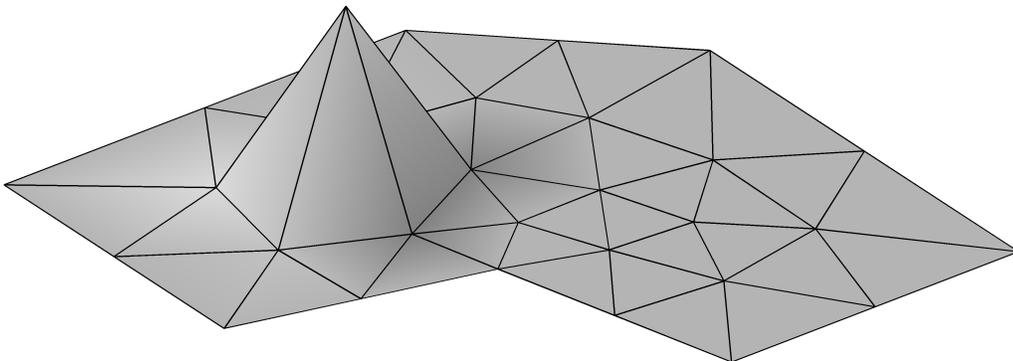
To define a simple finite element method, we suppose that  $\Omega$  is a polygon in  $\mathbb{R}^2$  and let  $\mathcal{T}_h$  be a simplicial decomposition of  $\Omega$  (covering of  $\bar{\Omega}$  by closed triangles so that the intersection of any two distinct elements of  $\mathcal{T}_h$  is either empty or a common edge or vertex). Let

$$M_0^1(\mathcal{T}_h) = \{ v \in C(\Omega) \mid V|_T \in \mathcal{P}_1(T) \forall T \in \mathcal{T}_h \} = \{ v \in H^1(\Omega) \mid V|_T \in \mathcal{P}_1(T) \forall T \in \mathcal{T}_h \},$$

and  $\mathring{M}_0^1(\mathcal{T}_h) = \mathring{H}^1(\Omega) \cap M_0^1(\mathcal{T}_h)$ . The  $\mathcal{P}_1$  finite element method for the Dirichlet problem is the Galerkin method with  $V_h = \mathring{M}_0^1(\mathcal{T}_h)$ .

We can use the Lagrange (hat function) basis for  $V_h$  to ensure that (1) the matrix is sparse, and (2) the integrals entering into the stiffness matrix and load vector are easy to compute.

FIGURE 4.1. A hat function basis element for  $M_0^1(\mathcal{T}_h)$ .



In the special case where  $\Omega$  is the unit square, and  $\mathcal{T}_h$  is obtained from a uniform  $m \times m$  partition into subsquares, each bisected by its SW-NE diagonal (so  $n = (m - 1)^2$ ), the resulting stiffness matrix is exactly the same as the matrix of the 5-point Laplacian.

### 3. Lagrange finite elements

This section is written mostly for 2D, although extending to  $n$  dimensions is straightforward.

A finite element space is a space of piecewise polynomials with respect to a given triangulation (simplicial decomposition)  $\mathcal{T}_h$ , but not just any space of piecewise polynomials. It is constructed by specifying the following things for each  $T \in \mathcal{T}_h$ :

- *Shape functions*: a finite dimensional space  $V(T)$  consisting of polynomial functions on  $T$ .
- *Degrees of freedom*: a finite set of linear functionals  $V(T) \rightarrow \mathbb{R}$  which are *unisolvent* on  $V(T)$ . This means that real values can be assigned arbitrarily to each DOF, and these determine one and only one element of  $V(T)$ . In other words, the DOF form a basis for the dual space of  $V(T)$ .

We further assume that each degree of freedom on  $T$  is associated to a subsimplex of  $T$ , i.e., to a vertex, an edge, or  $T$  itself (in 2D). Moreover, if a subsimplex is shared by two different triangles in  $T_1$  and  $T_2$  in  $\mathcal{T}_h$ , the DOFs for  $T_1$  and  $T_2$  associated to the subsimplex are in 1-to-1 correspondence.

When all this is specified, the *assembled finite element space* is defined as all functions  $v \in L^2(\Omega)$  such that

- $v|_T \in V(T)$  for all  $T \in \mathcal{T}_h$
- The DOFs are single-valued in the sense that whenever  $q$  is a subsimplex shared by  $T_1$  and  $T_2$ , then the corresponding DOFs on applied to  $v|_{T_1}$  and  $v|_{T_2}$  take on the same value.

Note that we do not specify the interelement continuity explicitly. It is determined by the fact that the shared DOFs are single-valued.

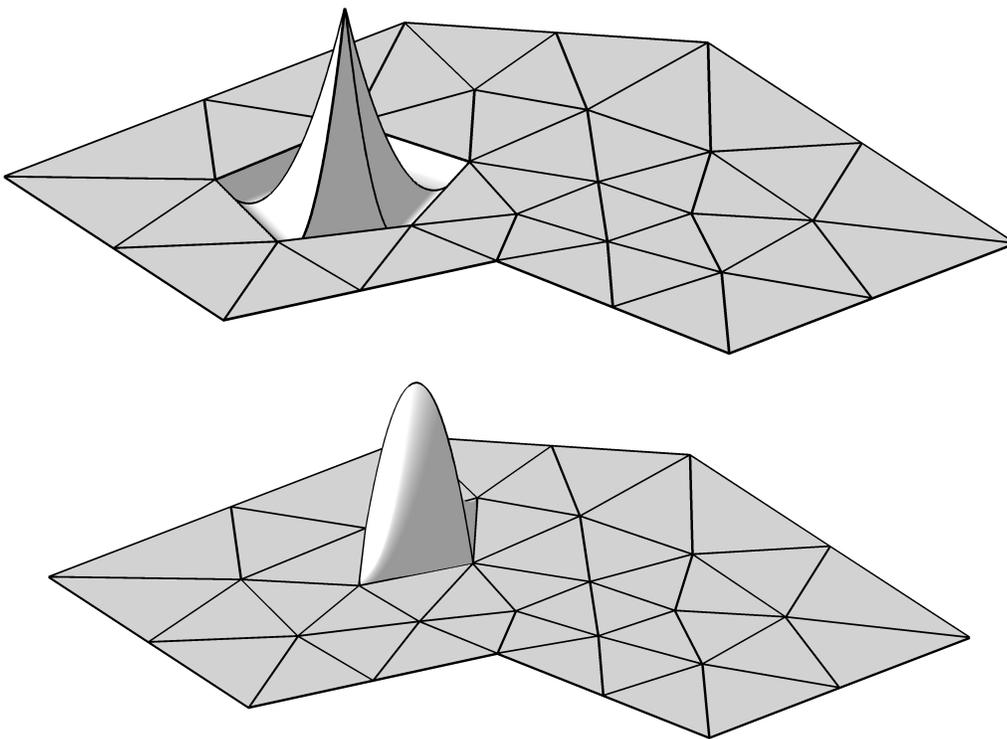
The reason for this definition is that it is easy to construct and compute with piecewise polynomial spaces defined in this way. First of all, we immediately obtain a set of *global degrees of freedom*, by considering all the degrees of freedom associated with all the subsimplices of the triangulation. An element of the FE space is uniquely determined by an arbitrary assignment of values to the global degrees of freedom. Thus the dimension of the FE space is the sum over the subsimplices of the number of degrees of freedom associated to the subsimplex. A basis for the FE space is obtained by setting one of the global DOFs to 1 and all the rest to zero. The resulting basis function is supported in the union of the triangles which contain the subsimplex. Thus we have a *local basis* (small supports), and will obtain a sparse stiffness matrix.

The simplest example is the  $\mathcal{P}_1$  element, or Lagrange element of degree 1, discussed above. Then the shape functions are simply the linear polynomials:  $V(T) = \mathcal{P}_1(T)$  (dimension equals 3 in 2D). The degrees of freedom on  $T$  are the evaluation functionals associated to the 3 vertices. These DOFs are certainly unisolvent: a linear function in 2D is determined by its value at any 3 non-colinear points. Clearly any continuous piecewise linear function belongs to the FE space, since it can be specified by assigning its vertex values. Conversely, if  $v$  is an element of the assembled FE space and two triangles  $T_1$  and  $T_2$  share a common edge  $e$ , then  $v|_{T_1}$  and  $v|_{T_2}$  must agree on all of  $e$ , since on  $e$  they are both linear functions, and they agree at the two end points of  $e$  (a linear function in 1D is determined by its value at any 2 distinct points). This shows that the assembled FE space consists precisely of the

continuous piecewise linears. The global degrees of freedom are the vertex values, and the corresponding local basis consists of the hat functions.

For the Lagrange element of degree 2, the shape functions are  $V(T) = \mathcal{P}_2(T)$ . There is one DOF associated to each vertex (evaluation at the vertex), and one associated to each edge (evaluation at the midpoint of the edge). Let us check unisolvence. Since  $\dim V(T) = 6$  and there are 6 DOFs on  $T$ , we must show that if all 6 DOFs vanish for some  $v \in V(T)$ , then  $v \equiv 0$ . For this choose a numbering of the vertices of  $T$ , and let  $\lambda_i \in \mathcal{P}_1(\mathbb{R}^2)$ ,  $i = 1, 2, 3$ , denote the linear function that is 1 on the  $i$ th vertex and which vanishes on the opposite edge,  $e$  of  $T$ . Thus the equation  $\lambda_i = 0$  is the equation of the line through the edge  $e$ . Since  $v$  vanishes at the two endpoints and the midpoint of  $e$ ,  $v|_e$  vanishes (a quadratic in 1D which vanishes at 3 points is zero). Therefore  $v$  is divisible by the linear polynomial  $\lambda_i$ , for each  $i$ . Thus  $v$  is a multiple of  $\lambda_1\lambda_2\lambda_3$ . But  $v$  is quadratic and this product is cubic, so the only possibility is that  $v \equiv 0$ . It is easy to check that the assembled space is exactly the space  $M_0^2(\mathcal{T}_h)$  of continuous piecewise quadratic functions. There is one basis function associated to each vertex and one to each edge.

FIGURE 4.2. Basis functions for  $M_0^2(\mathcal{T}_h)$ .

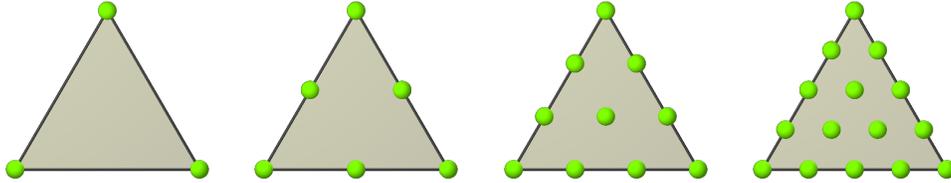


Note: the linear functions  $\lambda_i$  are the barycentric coordinate functions on  $T$ . They satisfy  $\lambda_1 + \lambda_2 + \lambda_3 \equiv 1$ .

Higher degree Lagrange elements are defined similarly.  $V(T) = \mathcal{P}_r(T)$ .  $\dim V(T) = (r+1)(r+2)/2$ . There is 1 DOF at each vertex,  $r-1$  on each edge, and  $(r-2)(r-1)/2$  in each triangle. Note that  $3 \times 1 + 3 \times (r-1) + (r-2)(r-1)/2 = (r+1)(r+2)/2$ . The DOFs

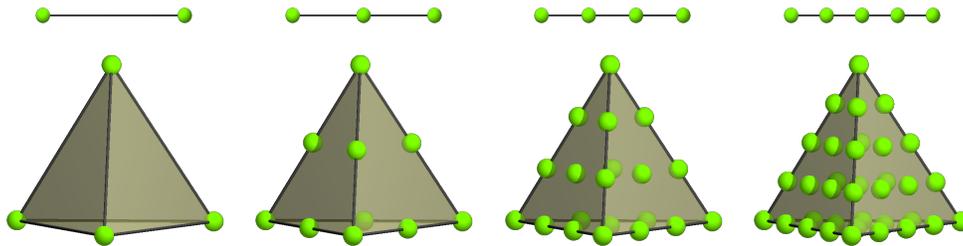
are the evaluation functionals at the points with barycentric coordinates all of the form  $i/r$  with  $0 \leq i \leq r$  and integer. See Figure 4.3.

FIGURE 4.3. Lagrange finite elements of degree 1, 2, and 3.



Lagrange elements can be defined in a similar way in  $n$ -dimensions.

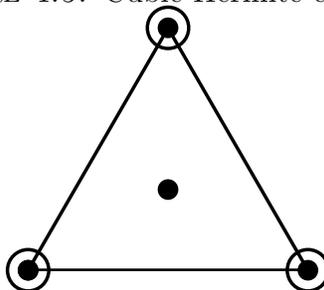
FIGURE 4.4. Lagrange finite elements of degree 1, 2, and 3 in 1-D and 3-D.



The restriction of a Lagrange element of degree  $r$  to a face or an edge is a Lagrange element of degree  $r$  on the face or edge. This leads to an inductive proof of unisolvence valid for all dimensions.

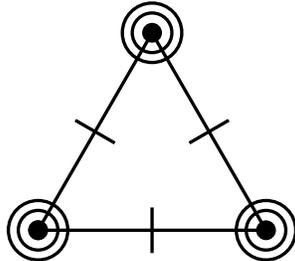
Other finite element spaces. Cubic Hermite finite elements. Shape functions are  $\mathcal{P}_3(T)$  for a triangle  $T$ . Three DOFs for each vertex  $v$ :  $u \mapsto u(v)$ ,  $u \mapsto (\partial u / \partial x)(v)$ ,  $u \mapsto (\partial u / \partial y)(v)$ ; and one DOF associated to the interior  $u \mapsto \int_T v$  (alternatively, evaluation at the barycenter). Proof of unisolvence. Note that the assembled finite element space belongs to  $C^0$  and  $H^1$ , but not to  $C^1$  and  $H^2$ .

FIGURE 4.5. Cubic Hermite element.



Quintic Hermite finite elements (Argyris elements). Shape functions are  $\mathcal{P}_5(T)$  for a triangle  $T$ . Six DOFs for each vertex  $v$ : evaluation of  $u$ , both 1st partials, and all three 2nd partials of  $u$  at  $v$ . One DOF for each edge: evaluation of  $\partial u / \partial n$  at midpoint. Unisolvent,  $H^2$ .

FIGURE 4.6. Quintic Hermite element.



#### 4. Finite element assembly

Now we consider how to efficiently compute the stiffness matrix in a finite element computation. We consider the Neumann problem

$$-\operatorname{div} a \operatorname{grad} u + b \cdot \operatorname{grad} u v + c u v = f \text{ in } \Omega, \quad a \frac{\partial u}{\partial n} = 0 \text{ on } \partial\Omega.$$

so the finite element method seeks  $u_h \in V_h$  such that

$$b(u_h, v) = F(v), \quad v \in V_h,$$

where

$$b(u, v) := \int_{\Omega} (a \operatorname{grad} u \cdot \operatorname{grad} v + b \cdot \operatorname{grad} u v + c u v) dx, \quad F(v) = \int_{\Omega} f v dx.$$

For simplicity suppose we use  $V_h = M_0^1(\mathcal{T}_h)$ , the Lagrange finite element space of degree 1 for some triangulation in 2-D. Let  $p_1, \dots, p_{n_{\text{vert}}}$  denote the vertices of the triangulation  $\mathcal{T}_h$  (where  $n_{\text{vert}}$  is the number of vertices, which is equal to  $\operatorname{div} V_h$ ), and let  $\phi_1, \dots, \phi_{n_{\text{vert}}}$  denote the corresponding hat function basis functions (such as the one pictured in Figure 4.1). The stiffness matrix, which we need to compute, is given by

$$A_{ij} = b(\phi_j, \phi_i), \quad i, j = 1, \dots, n_{\text{vert}}.$$

This might suggest as an algorithm

```

for  $i = 1, \dots, n_{\text{vert}}$ 
  for  $j = 1, \dots, n_{\text{vert}}$ 
    compute  $A_{ij} = b(\phi_j, \phi_i)$ 
  end
end

```

but, in fact, such an algorithm is very inefficient and *should never be used*.

To define an efficient algorithm, we introduce the local vertices, basis functions, and stiffness matrix. For each triangle  $T$ , let  $p_k^T$ ,  $k = 1, 2, 3$  denote the three vertices of  $T$ , and let  $\phi_k^T$  denote the three local basis functions on  $T$ . The local stiffness matrix on  $T$  is the  $3 \times 3$  matrix

$$A_{ij}^T = b^T(\phi_j^T, \phi_i^T), \quad i, j = 1, 2, 3,$$

where

$$b^T(v, w) = \int_T (a \operatorname{grad} v \cdot \operatorname{grad} w + b \cdot \operatorname{grad} vw + cvw) dx.$$

Note that the integral is only over the triangle  $T$ . We need to compute the quantity  $A_{ij}^T = b^T(\phi_j^T, \phi_i^T)$  in order to compute the (global) stiffness matrix. It will be part of exactly one element of the stiffness matrix. Fortunately, this quantity is easily computable. The three functions  $\phi_j^T$  and their gradients can be easily expressed analytically in terms of the coordinates of the vertices  $p_1^T, p_2^T, p_3^T$  of  $T$ . Therefore, if the coefficients  $a$ ,  $b$ , and  $c$  are constant on the element  $T$ , it is straightforward to give an arithmetic expression for  $A_{ij}^T$ . If the coefficients are variable, one commonly evaluates them at one or a few quadrature points in  $T$  and computes  $b^T(\phi_j^T, \phi_i^T)$  through a quadrature rule.

To relate the local quantities and the global quantities, we define  $I_k^T$  as the global vertex number of the  $k$ th vertex of  $T$ :

$$p_{I_k^T} = p_k^T, \quad k = 1, 2, 3.$$

The values of  $I_k^T$  can be stored in an integer table with one row for each triangle and 3 columns. This the  $(j, k)$  entry of the table is the global vertex number of the  $k$ th vertex of the  $j$ th triangle of the mesh. This is called the *connectivity table* of the mesh.

The finite element assembly algorithm to compute the stiffness matrix organizes the computation as a loop over the elements, in each element we: (1) compute the local stiffness matrix, (2) add the resulting elements into the appropriate elements on the global stiffness matrix.

```

Initialize A to 0
for T ∈ Th
  compute AijT = bT(ϕjT, ϕiT), i, j = 1, 2, 3
  AIiT IjT = AIiT IjT + AijT, i, j = 1, 2, 3
end

```

This is how a finite element stiffness matrix is computed in practice. Note that the computation is organized as a single loop over elements, rather than as a double loop over vertices.

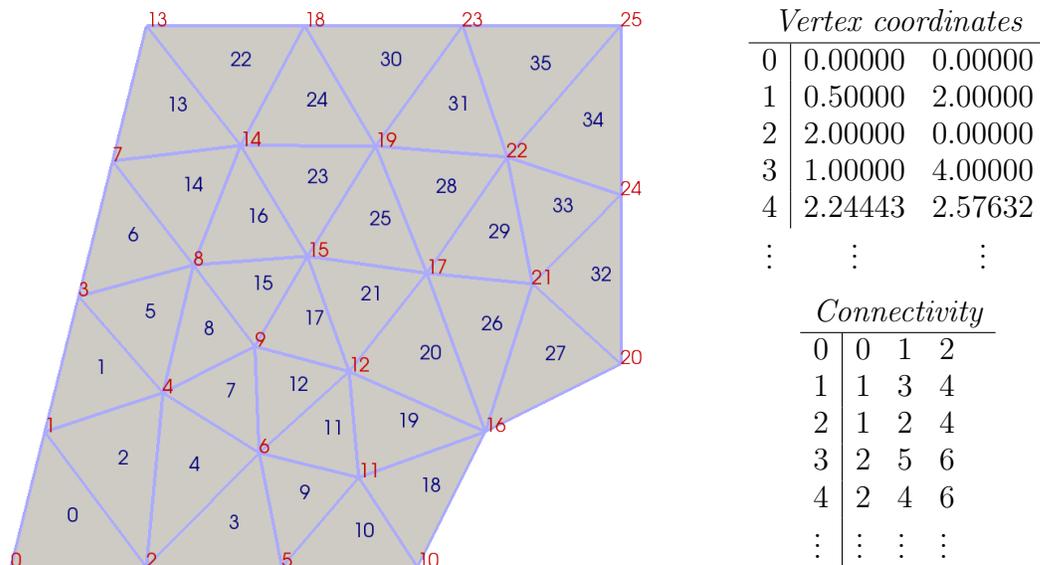
Thus to compute the stiffness matrix, we need two tables which describe the mesh: a real table of size  $n_{\text{vert}} \times 2$  which gives the coordinates of the vertices, and the integer connectivity table of size  $n_{\text{elt}} \times 3$ . These tables are created when the mesh is generated. Figure 4.7 shows a mesh of  $n_{\text{vert}} = 26$  triangles (numbered, in red, from 0 to 25 rather than 1 to 26), and  $n_{\text{elt}} = 36$  elements (numbered in blue from 0 to 35), and the corresponding mesh data tables.

## 5. Coercivity, inf-sup condition, and well-posedness

First we consider this in an abstract framework, in which we are given a Hilbert space  $V$ , a bounded bilinear form  $b : V \times V \rightarrow \mathbb{R}$  and a bounded linear form  $F : V \rightarrow \mathbb{R}$ , and the problem we wish to solve is

$$(4.4) \quad b(u, v) = F(v), \quad v \in V.$$

FIGURE 4.7. Finite element mesh data structure.



We will analyze the Galerkin method using the ideas of stability and consistency we introduced in Chapter 2, § 1.2. Recall that there we studied a problem of the form find  $u \in X$  such that  $Lu = f$  where  $L$  mapped the vector space  $X$  where the solution was sought to the vector space  $Y$  where the data  $f$  belonged. For the problem (4.4), the solution space  $X = V$ , and the data space  $Y = V^*$ , the dual space of  $V$ . The linear operator  $L$  is simply given by

$$Lw(v) = b(w, v) \quad w, v \in V.$$

So the problem (4.4) is simply: given  $F \in V^*$  find  $u \in V$  such that  $Lu = F$ .

First of all, before turning to discretization, we consider the well-posedness of this problem. We have already assumed that the bilinear form  $b$  is bounded, i.e., there exists a constant  $M$  such that

$$|b(w, v)| \leq M\|w\|\|v\|, \quad w, v \in V,$$

where, of course, the norm is the  $V$  norm. This implies that the operator  $L : V \rightarrow V^*$  is a bounded linear operator (with the same bound).

We will now consider hypotheses on the bilinear form  $b$  that ensure that the problem (4.4) is well-posed. We will consider three cases, in increasing generality.

**5.1. The symmetric coercive case.** First we assume that the bilinear form  $b$  is symmetric ( $b(w, v) = b(v, w)$ ) and *coercive*, which means that there exists a constant  $\gamma > 0$ , such that

$$b(v, v) \geq \gamma\|v\|^2, \quad v \in V.$$

In this case  $b(w, v)$  is an inner product on  $V$ , and the associated norm is equivalent to the  $V$  norm:

$$\gamma\|v\|_V^2 \leq b(v, v) \leq M\|v\|_V^2.$$

Thus  $V$  is a Hilbert space when endowed with the  $b$  inner product, and the Riesz Representation Theorem gives that for all  $F \in V^*$  there exists a unique  $u \in V$  such that

$$b(u, v) = F(v), \quad v \in V,$$

i.e., (4.4) has a unique solution for any data, and  $L^{-1} : V^* \rightarrow V$  is well-defined. Taking  $v = u$  and using the coercivity condition we immediately have  $\gamma \|u\|_V^2 \leq F(u) \leq \|F\|_{V^*} \|u\|_V$ , so  $\|u\|_V \leq \gamma^{-1} \|F\|_{V^*}$ , i.e.,  $\|L^{-1}\|_{\mathcal{L}(V^*, V)} \leq \gamma^{-1}$ .

In short, well-posedness is an immediate consequence of the Riesz Representation Theorem in the symmetric coercive case.

As a simple example of the utility of this result, let us consider the Neumann problem

$$-\operatorname{div} a \operatorname{grad} u + cu = f \text{ in } \Omega, \quad a \frac{\partial u}{\partial n} = 0 \text{ on } \partial\Omega$$

where  $0 < \underline{a} \leq a(x) \leq \bar{a}$ ,  $0 < \underline{c} \leq c(x) \leq \bar{c}$ . The weak formulation is: Find  $u \in V$  such that

$$b(u, v) = F(v), \quad v \in V,$$

where  $V = H^1$ ,

$$b(w, v) = \int_{\Omega} (a \operatorname{grad} w \cdot \operatorname{grad} v + c w v), \quad F(v) = \int_{\Omega} f v.$$

Clearly the bilinear form  $b$  is bounded with  $M = \max(\bar{a}, \bar{c})$ , and is coercive with  $\gamma = \min(\underline{a}, \underline{c})$ . It follows that the weak formulation is well-posed. It admits a unique solution and  $\|u\|_{H^1} \leq \gamma^{-1} \|F\|_{(H^1)'} \leq \gamma^{-1} \|f\|_{L^2}$ .

**5.2. The coercive case.** Even if we dispense with the assumption of symmetry, coercivity implies well-posedness. From coercivity we have  $\gamma \|w\|^2 \leq b(w, w) = Lw(w) \leq \|Lw\|_{V^*} \|w\|$ , so

$$(4.5) \quad \|w\| \leq \gamma^{-1} \|Lw\|_{V^*}, \quad w \in V.$$

This immediately leads to three conclusions:

- $L$  is one-to-one.
- If  $L$  is also onto, so  $L^{-1}$  is well-defined, then  $\|L^{-1}\|_{\mathcal{L}(V^*, V)} \leq \gamma^{-1}$ .
- The range  $W = L(V)$  is closed in  $V^*$ .

The first two points are immediate. For the third, suppose that for some  $u_1, u_2, \dots \in V$ ,  $Lu_n$  converges to some  $G \in V^*$ . We must show that  $G = Lu$  for some  $u \in V$ . Since  $Lu_n$  converges in  $V^*$  it forms a Cauchy sequence. Using (4.5) we conclude that  $u_n$  forms a Cauchy sequence in  $V$ , and so converge to some  $u$  in  $V$ . Since  $L$  is bounded  $Lu_n \rightarrow Lu$  in  $V^*$ , so  $Lu = G$ , showing that indeed  $W$  is closed in  $V^*$ .

It remains to show that  $L$  is onto, i.e., the closed subspace  $W = L(V)$  is the whole of  $V^*$ . If  $W$  were a strict closed subspace of  $V^*$  then there would exist a nonzero element  $v \in V$  such that  $G(v) = 0$  for all  $G \in W$ , i.e.,  $b(w, v) = Lw(v) = 0$  for all  $w \in V$  and this particular  $v$ . But, taking  $w = v$  and using coercivity we get a contradiction.

Thus we have shown that for a bounded coercive bilinear form, symmetric or not, the abstract weak formulation (4.4) is well-posed. This result is known as the *Lax-Milgram theorem*.

**5.3. The inf-sup condition.** It turns out to be very useful to consider a much more general case. Suppose that, instead of coercivity, we assume that

- (1) (inf-sup condition) There exists  $\gamma > 0$  such that for all  $0 \neq w \in V$  there exists  $0 \neq v \in V$  such that

$$b(w, v) \geq \gamma \|w\| \|v\|.$$

- (2) (dense range condition) For all  $0 \neq v \in V$  there exists a  $w \in V$  such that  $b(w, v) \neq 0$ .

We shall see that it is easy to adapt the proof of the Lax-Milgram theorem to this case.

Note that the inf-sup condition can be written

$$\inf_{0 \neq w \in V} \sup_{0 \neq v \in V} \frac{b(w, v)}{\|w\| \|v\|} > 0,$$

which explains its name. The dense range condition is equivalent to the condition that the range  $W = L(V)$  is dense in  $V^*$ . Clearly coercivity implies both these conditions (take  $v = w$  for the first and  $w = v$  for the second). In the symmetric case the second condition follows from the first. In any case, using these two conditions it is easy to carry out the above argument, as we do now.

The bound (4.5) follows directly from the inf-sup condition. Again this implies  $L$  is 1-to-1 and that  $W = L(V)$  is *closed* in  $V^*$ , and furnishes a bound on  $\|L^{-1}\|$  if  $L$  is onto. Since  $L$  has dense range, by the second condition, and closed range, it is indeed onto.

This version is in some sense the most general possible. If (4.4) is well-posed, so  $L^{-1} : V^* \rightarrow V$  exists, then it is a simple matter of untangling the definitions to see that

$$\inf_{0 \neq w \in V} \sup_{0 \neq v \in V} \frac{b(w, v)}{\|w\| \|v\|} = \|L^{-1}\|_{\mathcal{L}(V^*, V)}^{-1},$$

and so the inf-sup condition holds with  $\gamma = 1/\|L^{-1}\|_{\mathcal{L}(V^*, V)}^{-1}$ . Thus the inf-sup condition and dense range condition are equivalent to well-posedness.

## 6. Stability, consistency, and convergence

Now we turn to discretization, again using the framework of Chapter 2, § 1.2. First we consider the coercive (but not necessarily symmetric) case. Thus we suppose again that  $b : V \times V \rightarrow \mathbb{R}$  is a bounded, coercive bilinear form, with constants  $M$  and  $\gamma$ . Consider  $V_h$  a finite dimensional subspace of  $V$ . Restricting the bilinear form  $b$  to  $V_h \times V_h$  defines an operator  $L_h : V_h \rightarrow V_h^*$ , and restricting  $F$  to  $V_h$  gives a linear form  $F_h : V_h \rightarrow \mathbb{R}$ . Galerkin's method is just  $L_h u_h = F_h$ . We will show that this method is consistent and stable, and so convergent.

Stability just means that  $L_h$  is invertible with the stability constant given by  $\|L_h^{-1}\|_{\mathcal{L}(V_h^*, V_h)}$ . Since  $b$  is coercive over all of  $V$  it is certainly coercive over  $V_h$ , and so the last section implies stability with stability constant  $\gamma^{-1}$ . In short, *if the bilinear form is coercive, then for any choice of subspace the Galerkin method is stable with the stability constant bounded by the inverse of the coercivity constant.*

To talk about consistency, as in Chapter 2, we need to define a representative  $r_h u$  of the solution  $u$  in  $V_h$ . A natural choice, which we shall make, is that  $r_h : V \rightarrow V_h$  is the orthogonal projection so that

$$\|u - r_h u\| = \inf_{v \in V_h} \|u - v\|.$$

The consistency error is

$$\|L_h r_h u - F_h\|_{V_h^*} = \sup_{0 \neq v \in V_h} \frac{|(L_h r_h u - F_h)(v)|}{\|v\|}.$$

But  $(L_h r_h u - F_h)(v) = b(r_h u, v) - F(v) = b(r_h u - u, v)$ , so  $|(L_h r_h u - F_h)(v)| \leq M \|u - r_h u\| \|v\|$ . Therefore the consistency error is bounded by

$$M \inf_{v \in V_h} \|u - v\|.$$

We therefore obtain the convergence estimate

$$\|r_h u - u_h\| \leq M \gamma^{-1} \inf_{v \in V_h} \|u - v\|.$$

We can then apply the triangle inequality to deduce

$$\|u - u_h\| \leq (1 + M \gamma^{-1}) \inf_{v \in V_h} \|u - v\|.$$

This is the fundamental estimate for finite elements. It shows that finite elements are *quasi-optimal*, i.e., that the error in the finite element solution is no more than a constant multiple of the error in the best possible approximation from the subspace. The constant can be taken to be 1 plus the bound of the bilinear form times the stability constant.

REMARK. We obtained stability using coercivity. From the last section we know that we could obtain stability as well if we had instead of coercivity, a discrete inf-sup condition: There exists  $\gamma > 0$  such that for all  $0 \neq w \in V_h$  there exists  $0 \neq v \in V_h$  such that

$$b(w, v) \geq \gamma \|w\| \|v\|.$$

(In the finite dimensional case the dense range condition follows from the inf-sup condition since an operator from  $V_h$  to  $V_h^*$  which is 1-to-1 is automatically onto) The big difficulty however is that *the fact that  $b$  satisfies the inf-sup condition over  $V$  does not by any means imply that it satisfies the inf-sup condition over a finite dimensional subspace  $V_h$* . In short, for coercive problems stability is automatic, but for more general well-posed problems Galerkin methods may or may not be stable (depending on the choice of subspace), and proving stability can be difficult.

## 7. Finite element approximation theory

In this section we turn to the question of finite element approximation theory, that is of estimating

$$\inf_{v \in V_h} \|u - v\|_1$$

where  $V_h$  is a finite element space. For simplicity, we first consider the case where  $V_h = M_0^1(\mathcal{T}_h)$ , the Lagrange space of continuous piecewise linear functions on a given mesh  $\mathcal{T}_h$  where  $\mathcal{T}_h$  is a simplicial decomposition of  $\Omega$  with mesh size  $h = \max_{T \in \mathcal{T}_h} \text{diam } T$ . (Note: we are overloading the symbol  $h$ . If we were being more careful we would consider a sequence of meshes  $T_i$  with mesh size  $h_i$  tending to zero. But the common practice of using  $h$  as both the index and the mesh size saves writing subscripts and does not lead to confusion.)

First we need some preliminary results on Sobolev spaces: density of smooth functions, Poincaré inequality, Sobolev embedding  $H^s(\Omega) \subset C(\bar{\Omega})$  if  $s > n/2$ .

**THEOREM 4.1** (Poincaré inequality). *Let  $\Omega$  be a bounded domain in  $\mathbb{R}^n$  with Lipschitz boundary (e.g., smooth boundary or a polygon). Then there exists a constant  $c$ , depending only on  $\Omega$ , such that*

$$\|u\|_{L^2(\Omega)} \leq c \|\text{grad } u\|_{L^2(\Omega)}, \quad u \in H^1(\Omega) \text{ such that } \int_{\Omega} u = 0.$$

*The same inequality holds for all  $u \in \mathring{H}^1(\Omega)$ , or even for all  $u \in H^1(\Omega)$  which vanish on a non-empty open subset of the boundary.*

The result for  $u$  of mean value zero is sometimes called the Poincaré–Neumann inequality. In one dimension it is called Wirtinger’s inequality. The result for  $u \in \mathring{H}^1(\Omega)$  is sometimes called the Poincaré–Friedrichs inequality or just the Friedrichs inequality.

One proof of this result is based on Rellich’s theorem that  $H^1(\Omega)$  is compactly embedded in  $L^2(\Omega)$ . Other proofs are more explicit. Here we give a very simple proof of the Poincaré–Neumann inequality in one dimension.

If  $u \in C^1(\bar{I})$  where  $I$  is an interval of length  $L$ , and  $\int u = 0$ , then there exists a point  $x_0 \in I$  such that  $u(x_0) = 0$ . Therefore

$$|u(x)|^2 = \left| \int_{x_0}^x u'(s) ds \right|^2 \leq \left| \int_I |u'(s)| ds \right|^2 \leq L \int_I |u'(s)|^2 ds.$$

Integrating, we get  $\|u\| \leq L\|u'\|$ . This can be extended to  $u \in H^1$  using density of  $C^\infty$  in  $H^1$ .

An alternative proof uses Fourier cosine series:  $u(x) = \sum_{n=1}^{\infty} a_n \cos n\pi x/L$  (where the sum starts at  $n = 1$ , since  $\int u = 0$ ). This gives the result  $\|u\| \leq L/\pi \|u'\|$ , in which the constant is sharp (achieved by  $u(x) = \cos \pi x/L$ ). In fact the result can be proved with the constant  $d/\pi$ ,  $d = \text{diameter of } \Omega$  for any *convex* domain in  $n$ -dimensions (Payne and Weinberger, Arch. Rat. Mech. Anal. 5 1960, pp. 286–292). The dependence of the constant on the domain is more complicated for non-convex domains.

Multi-index notation: In  $n$ -dimensions a multi-index  $\alpha$  is an  $n$ -tuple  $(\alpha_1, \dots, \alpha_n)$  with the  $\alpha_i$  non-negative integers. We write  $|\alpha| = \alpha_1 + \dots + \alpha_n$  for the degree of the multi-index,  $\alpha! = \alpha_1! \dots \alpha_n!$ ,

$$|\alpha| = \alpha_1 + \dots + \alpha_n, \quad \alpha! = \alpha_1! \dots \alpha_n!, \quad x^\alpha = x_1^{\alpha_1} \dots x_n^{\alpha_n}, \quad D^\alpha u = \frac{\partial^{|\alpha|} u}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}}.$$

Thus a general element of  $\mathcal{P}_r(\mathbb{R}^n)$  is  $p(x) = \sum_{|\alpha| \leq r} a_\alpha x^\alpha$ , and a general constant-coefficient linear partial differential operator of degree  $r$  is  $Lu = \sum_{|\alpha| \leq r} a_\alpha D^\alpha u$ . Taylor’s theorem for a smooth function defined in a neighborhood of a point  $x_0 \in \mathbb{R}^n$  is

$$u(x) = \sum_{|\alpha| \leq m} \frac{1}{\alpha!} D^\alpha u(x_0) (x - x_0)^\alpha + O(|x - x_0|^{m+1})$$

We write  $\alpha \leq \beta \iff \alpha_i \leq \beta_i, i = 1, \dots, n$ . We have

$$D^\alpha x^\beta = \begin{cases} \frac{\beta!}{(\beta-\alpha)!} x^{\beta-\alpha}, & \alpha \leq \beta, \\ 0, & \text{otherwise.} \end{cases}$$

In particular  $D^\alpha x^\alpha = \alpha!$ .

Let  $\Omega$  be a bounded domain in  $\mathbb{R}^n$  with Lipschitz boundary (for our applications, it will be a triangle). It is easy to see that the DOFs

$$u \mapsto \int_{\Omega} D^{\alpha} u, \quad |\alpha| \leq r,$$

are unisolvent on  $\mathcal{P}_r(\Omega)$ . Therefore we can define  $P_r : H^r(\Omega) \rightarrow \mathcal{P}_r(\Omega)$  by

$$\int_{\Omega} D^{\alpha} P_r u(x) dx = \int_{\Omega} D^{\alpha} u(x) dx, \quad |\alpha| \leq r.$$

It follows immediately from this definition that  $D^{\beta} P_r u = P_{r-|\beta|} D^{\beta} u$  for  $|\beta| \leq r$ .

REMARK. The  $r$ th Taylor polynomial of  $u$  at  $x_0$  is  $T_r u$  given by

$$D^{\alpha} T_r u(x_0) = D^{\alpha} u(x_0), \quad |\alpha| \leq r.$$

So  $P_r u$  is a sort of averaged Taylor polynomial of  $u$ .

Let  $u \in H^{r+1}(\Omega)$ . Then  $u - P_r u$  has integral zero on  $\Omega$ , so the Poincaré inequality gives

$$\|u - P_r u\| \leq c_1 \sum_{|\alpha|=1} \|D^{\alpha}(u - P_r u)\| = c_1 \sum_{|\alpha|=1} \|D^{\alpha} u - P_{r-1}(D^{\alpha} u)\|,$$

for some constant  $c_1$  depending only on  $\Omega$  (where we use the  $L^2(\Omega)$  norm). Applying the same reasoning to  $D^{\alpha} u - P_{r-1}(D^{\alpha} u)$ , we have  $\|D^{\alpha} u - P_{r-1}(D^{\alpha} u)\|$  is bounded by the sum of the norms of second partial derivatives, so

$$\|u - P_r u\| \leq c_2 \sum_{|\alpha|=2} \|D^{\alpha} u - P_{r-2}(D^{\alpha} u)\|.$$

Continuing in this way we get

$$\|u - P_r u\| \leq c_r \sum_{|\alpha|=r} \|D^{\alpha} u - P_0(D^{\alpha} u)\| \leq C \sum_{|\alpha|=r+1} \|D^{\alpha} u\|.$$

For any  $|\beta| \leq r$  may also apply this result to  $D^{\beta} u \in H^{r+1-|\beta|}$  to get

$$\|D^{\beta} u - P_{r-|\beta|} D^{\beta} u\| \leq C \sum_{|\gamma| \leq r-|\beta|+1} \|D^{\gamma} D^{\beta} u\|$$

so

$$\|D^{\beta}(u - P_r u)\| \leq C \sum_{|\alpha|=r+1} \|D^{\alpha} u\|.$$

Since this holds for all  $|\beta| \leq r$ , we

$$\|u - P_r u\|_{H^r} \leq c |u|_{H^{r+1}}, \quad u \in H^{r+1}(\Omega).$$

We have thus given a constructive proof of the follow important result.

**THEOREM 4.2** (Bramble–Hilbert lemma). *Let  $\Omega$  be a Lipschitz domain and  $r \geq 0$ . Then there exists a constant  $c$  only depending on the domain  $\Omega$  and on  $r$  such that*

$$\inf_{p \in \mathcal{P}_r} \|u - p\|_{H^r} \leq c |u|_{H^{r+1}}, \quad u \in H^{r+1}(\Omega).$$

REMARK. This proof of the Bramble–Hilbert lemma, based on the Poincaré inequality, is due to Verfürth (*A note on polynomial approximation in Sobolev spaces*, M2AN 33, 1999). The method is constructive in that it exhibits a specific polynomial  $p$  satisfying the estimate (namely  $P_r u$ ). Based on classical work of Payne and Weinberger on the dependence of the Poincaré constant on the domain mentioned above, it leads to good explicit bounds on the constant in the Bramble–Hilbert lemma. A much older constructive proof is due to Dupont and Scott and taught in the textbook of Brenner and Scott. However that method is both more complicated and it leads a worse bound on the constant. Many texts (e.g., Braess) give a non-constructive proof of the Bramble–Hilbert lemma based on Rellich’s compactness theorem.

Now we derive an corollary of the Bramble–Hilbert lemma (which is of such importance that sometimes the corollary is itself referred to as the Bramble–Hilbert lemma).

COROLLARY 4.3. *Let  $\Omega$  be a Lipschitz domain and  $r \geq 0$ , and  $\pi : H^{r+1}(\Omega) \rightarrow \mathcal{P}_r(\Omega)$  be a bounded linear projection onto  $\mathcal{P}_r(\Omega)$ . Then there exists a constant  $c$  which only depends on the domain  $\Omega$ ,  $r$ , and the norm of  $\pi$  such that*

$$\|u - \pi u\|_{H^r} \leq c|u|_{H^{r+1}}.$$

Note: the hypothesis means that  $\pi$  is a bounded linear operator mapping  $H^{r+1}(\Omega)$  into  $\mathcal{P}_r(\Omega)$  such that  $\pi u = u$  if  $u \in \mathcal{P}_r(\Omega)$ . Bounded means that  $\|\pi\|_{\mathcal{L}(H^{r+1}, H^r)} < \infty$ . It doesn’t matter what norm we choose on  $\mathcal{P}_r$ , since it is a finite dimensional space.

PROOF.

$$\begin{aligned} \|u - \pi u\|_{H^r} &= \inf_{p \in \mathcal{P}_r} \|(u - p) - \pi(u - p)\|_{H^r} \\ &\leq (1 + \|\pi\|_{\mathcal{L}(H^{r+1}, H^r)}) \inf_{p \in \mathcal{P}_r} \|u - p\|_{H^{r+1}} = c(1 + \|\pi\|)|u|_{H^{r+1}}. \end{aligned}$$

□

We will be applying this Bramble–Hilbert corollary on the individual triangles  $T$  of the finite element mesh. However, if we apply the corollary with  $\Omega = T$ , the unknown constant  $c$  which arises in the corollary will depend on the individual triangle  $T$ , and we will not be able to control it. So instead we will apply the corollary on one fixed reference triangle, and then scale the result from the reference triangle to an arbitrary triangle  $T$ , and determine how the constant is effected. Thus, we let  $\Omega = \hat{T}$  be the unit triangle with vertices  $\hat{v}_0 = (0, 0)$ ,  $\hat{v}_1 = (1, 0)$ , and  $\hat{v}_2 = (0, 1)$ ,  $r = 1$ , and let  $\pi = I_{\hat{T}}$  the linear interpolant:  $I_{\hat{T}}u \in \mathcal{P}_1(\hat{T})$  and  $I_{\hat{T}}u(\hat{v}_i) = u(\hat{v}_i)$ ,  $i = 0, 1, 2$ . The  $I_{\hat{T}}u$  is defined for all  $u \in C(\hat{T})$  and  $\|I_{\hat{T}}u\|_{L^\infty} \leq \|u\|_{L^\infty} \leq c\|u\|_{H^2}$ , where we use the Sobolev embedding theorem in the last step (and  $c$  is some absolute constant). From the corollary we get

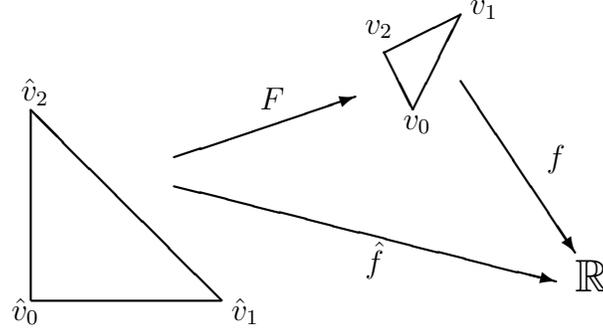
$$(4.6) \quad \|u - I_{\hat{T}}u\|_{H^1(\hat{T})} \leq c|u|_{H^2(\hat{T})}.$$

This result will turn out to be a key step in analyzing piecewise linear interpolation.

The next step is to scale this result from the unit triangle to an arbitrary triangle  $T$ . Suppose the vertices of  $T$  are  $v_0$ ,  $v_1$ , and  $v_2$ . There exists a unique affine map  $F$  taking  $\hat{v}_i$  to  $v_i$ ,  $i = 0, 1, 2$ . Indeed,

$$x = F\hat{x} = v_0 + B\hat{x}, \quad B = (v_1 - v_0 | v_2 - v_0),$$

FIGURE 4.8. Mapping between the reference triangle and an arbitrary triangle.



where the last notation means that  $B$  is the  $2 \times 2$  matrix whose columns are the vectors  $v_1 - v_0$  and  $v_2 - v_0$ . The map  $F$  takes  $\hat{T}$  1-to-1 onto  $T$ . Since the columns of  $B$  are both vectors of length at most  $h_T$ , certainly the four components  $b_{ij}$  of  $B$  are bounded by  $h_T$ , and so, in any convenient norm,  $\|B\| \leq ch_T$  (with  $c$  depending on the norm chosen). Moreover,  $\det B = 2|T|$ , the ratio of the area of  $T$  to the area of  $\hat{T}$ ,  $|\hat{T}| = 1/2$ .

Now to any function  $f$  on  $T$  we may associate the pulled-back function  $\hat{f}$  on  $\hat{T}$  where

$$\hat{f}(\hat{x}) = f(x) \quad \text{with } x = F\hat{x}.$$

I.e.,  $\hat{f} = f \circ F$ . See Figure 4.8.

Next we relate derivatives and norms of a function  $f$  with its pull-back  $\hat{f}$ . For the derivative we simply use the chain rule:

$$\frac{\partial \hat{f}}{\partial \hat{x}_j}(\hat{x}) = \sum_{i=1}^2 \frac{\partial f}{\partial x_i}(x) \frac{\partial x_i}{\partial \hat{x}_j} = \sum_{i=1}^2 b_{ij} \frac{\partial f}{\partial x_i}(x).$$

Similarly,

$$\frac{\partial^2 \hat{f}}{\partial \hat{x}_j \partial \hat{x}_l}(\hat{x}) = \sum_{i=1}^2 \sum_{k=1}^2 b_{ij} b_{kl} \frac{\partial^2 f}{\partial x_i \partial x_k}(x),$$

etc. Thus we have

$$(4.7) \quad \sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})| \leq c \|B\|^r \sum_{|\beta|=r} |D^\beta f(x)|.$$

In the same way we have

$$(4.8) \quad \sum_{|\beta|=r} |D^\beta f(x)| \leq c \|B^{-1}\|^r \sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})|.$$

In (4.7) we may bound  $\|B\|$  by  $ch_T$ . To bound  $\|B^{-1}\|$  in (4.8), we introduce another geometric quantity, namely the diameter  $\rho_T$  of the inscribed disk in  $T$ . Then any vector of length  $\rho_T$  is the difference of two points in  $T$  (two opposite points on the inscribed circle), and these are mapped by  $B^{-1}$  to two points in  $\hat{T}$ , which are at most  $\sqrt{2}$  apart. Thus, using the Euclidean

norm  $\|B^{-1}\| \leq \sqrt{2}/\rho_T$ , i.e.,  $\|B^{-1}\| = O(\rho_T^{-1})$ . We have thus shown

$$\sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})| \leq ch_T^r \sum_{|\beta|=r} |D^\beta f(x)|, \quad \sum_{|\beta|=r} |D^\beta f(x)| \leq c\rho_T^{-r} \sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})|.$$

Now let us consider how norms map under pull-back. First we consider the  $L^2$  norm. Let  $|T|$  denote the area of  $T$ . Changing variables from  $\hat{x}$  to  $x = F\hat{x}$ , we have

$$\|f\|_{L^2(T)}^2 = \int_T |f(x)|^2 dx = 2|T| \int_{\hat{T}} |\hat{f}(\hat{x})|^2 d\hat{x} = 2|T| \|\hat{f}\|_{L^2(\hat{T})}^2.$$

That is,  $\|f\|_{L^2(T)} = \sqrt{2|T|} \|\hat{f}\|_{L^2(\hat{T})}$ . Next consider the  $H^r$  seminorm:

$$|f|_{H^r(T)}^2 = \int_T \sum_{|\beta|=r} |D^\beta f(x)|^2 dx \leq c\rho_T^{-2r} \int_T \sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})|^2 dx = 2c|T|\rho_T^{-2r} \int_{\hat{T}} \sum_{|\alpha|=r} |D^\alpha \hat{f}(\hat{x})|^2 d\hat{x},$$

so

$$|f|_{H^r(T)} \leq c\sqrt{|T|}\rho_T^{-r} |\hat{f}|_{H^r(\hat{T})}.$$

Similarly,

$$|\hat{f}|_{H^r(\hat{T})} \leq c \frac{1}{\sqrt{|T|}} h_T^r |f|_{H^r(T)}.$$

Now let  $u \in H^2(T)$ , and let  $\hat{u} \in H^2(\hat{T})$  be the corresponding function. We saw in (4.6) that

$$\|\hat{u} - I_{\hat{T}}\hat{u}\|_{H^1(\hat{T})} \leq c|\hat{u}|_{H^2(\hat{T})}.$$

Now it is easy to see that the pull-back of  $I_T u$  is  $I_{\hat{T}}\hat{u}$  (both are linear functions which equal  $u(v_i)$  at the vertex  $\hat{v}_i$ ). Therefore  $\hat{u} - I_{\hat{T}}\hat{u} = \widehat{u - I_T u}$ . We then have

$$\|u - I_T u\|_{L^2(T)} \leq c\sqrt{|T|} \|\hat{u} - I_{\hat{T}}\hat{u}\|_{L^2(\hat{T})} \leq c\sqrt{|T|} |\hat{u}|_{H^2(\hat{T})} \leq ch_T^2 |u|_{H^2(T)},$$

and

$$|u - I_T u|_{H^1(T)} \leq c\sqrt{|T|}\rho_T^{-1} |\hat{u} - I_{\hat{T}}\hat{u}|_{H^1(\hat{T})} \leq c\sqrt{|T|}\rho_T^{-1} |\hat{u}|_{H^2(\hat{T})} \leq ch_T^2/\rho_T |u|_{H^2(T)}.$$

If the triangle  $T$  is not too distorted, then  $\rho_T$  is not much smaller than  $h_T$ . Let us define  $\sigma_T = h_T/\rho_T$ , the *shape constant* of  $T$ . We have proved:

**THEOREM 4.4.** *Let  $T$  be a triangle with diameter  $h_T$ , and let  $I_T$  be the linear interpolant at the vertices of  $T$ . Then there exists an absolute constant  $c$  such that*

$$\|u - I_T u\|_{L^2(T)} \leq ch_T^2 |u|_{H^2(T)}, \quad u \in H^2(T).$$

Moreover there exists a constant  $c'$  depending only on the shape constant for  $T$  such that

$$|u - I_T u|_{H^1(T)} \leq c'h_T |u|_{H^2(T)}, \quad u \in H^2(T).$$

Now we have analyzed linear interpolation on a single but arbitrary triangle, we can just add over the triangles to analyze piecewise linear interpolation.

**THEOREM 4.5.** *Suppose we have a sequence of triangulations  $\mathcal{T}_h$  with mesh size  $h$  tending to 0. For  $u$  a continuous function on  $\bar{\Omega}$ , let  $I_h u$  denote the continuous piecewise linear interpolant of  $u$  on the mesh  $\mathcal{T}_h$ . Then there exists an absolute constant  $c$  such that*

$$\|u - I_h u\|_{L^2(\Omega)} \leq ch^2 |u|_{H^2(\Omega)}, \quad u \in H^2(\Omega).$$

*If the mesh sequence is shape regular (i.e., the shape constant is uniformly bounded), then there exists a constant  $c'$  depending only on a bound for the shape constant such that*

$$|u - I_h u|_{H^1(\Omega)} \leq c'h |u|_{H^2(\Omega)}, \quad u \in H^2(\Omega).$$

In a similar fashion, for the space of Lagrange finite elements of degree  $r$  we can analyze the interpolant defined via the degrees of freedom.

**THEOREM 4.6.** *Suppose we have a sequence of triangulations  $\mathcal{T}_h$  with mesh size  $h$  tending to 0. Let  $V_h$  be the space of Lagrange finite elements of degree  $r$  with respect to the mesh, and for  $u$  a continuous function on  $\bar{\Omega}$ , let  $I_h u$  denote the interpolant of  $u$  into  $V_h$  defined through the degrees of freedom. Then there exists an absolute constant  $c$  such that*

$$\|u - I_h u\|_{L^2(\Omega)} \leq ch^s |u|_{H^s(\Omega)}, \quad u \in H^s(\Omega), \quad 2 \leq s \leq r + 1.$$

*If the mesh sequence is shape regular (i.e., the shape constant is uniformly bounded), then there exists a constant  $c'$  depending only on a bound for the shape constant such that*

$$|u - I_h u|_{H^1(\Omega)} \leq c'h^{s-1} |u|_{H^s(\Omega)}, \quad u \in H^s(\Omega), \quad 2 \leq s \leq r + 1.$$

Thus for smooth  $u$  (more precisely,  $u \in H^{r+1}$ ), we obtain the rate of convergence  $O(h^{r+1})$  in  $L^2$  and  $O(h^r)$  in  $H^1$  when we approximation with Lagrange elements of degree  $r$ .

The proof of this result is just the Bramble–Hilbert lemma and scaling. Note that we must assume  $s \geq 2$  so that  $u \in H^s$  is continuous and the interpolant is defined. On the other hand we are limited to a rate of  $O(h^{r+1})$  in  $L^2$  and  $O(h^r)$  in  $H^1$ , since the interpolant is exact on polynomials of degree  $r$ , but not higher degree polynomials.

## 8. Error estimates for finite elements

**8.1. Estimate in  $H^1$ .** To be concrete, consider the Dirichlet problem

$$-\operatorname{div} a \operatorname{grad} u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega,$$

with a coefficient  $a$  bounded above and below by positive constants of  $\Omega$  and  $f \in L^2$ . The weak formulation is: find  $u \in V = \dot{H}^1(\Omega)$  such that

$$(4.9) \quad b(u, v) = F(v), \quad v \in V,$$

where  $b(u, v) = \int_{\Omega} \operatorname{grad} u \cdot \operatorname{grad} v \, dx$ ,  $F(v) = \int_{\Omega} f v \, dx$ . Clearly  $b$  is bounded:  $|b(u, v)| \leq M \|u\|_1 \|v\|_1$ , (with  $M = \sup a$ ). By Poincaré's inequality, Theorem 4.1,  $b$  is coercive:  $b(v, v) \geq \gamma \|v\|_1^2$ .

Now suppose that  $\Omega$  is a polygon and let  $V_h$  be the space of Lagrange finite elements of degree  $r$  vanishing on the boundary with respect to a mesh of  $\Omega$  of mesh size  $h$ , and define  $u_h$  to be the finite element solution:  $u_h \in V_h$ ,

$$b(u_h, v) = F(v), \quad v \in V_h.$$

By the fundamental estimate for finite element methods, proven in Section 6,

$$\|u - u_h\|_1 \leq c \inf_{v \in V_h} \|u - v\|_1,$$

(where  $c = 1 + M\gamma^{-1}$ ). Then we may apply the finite element approximation theory summarized in Theorem 4.6, and conclude that

$$(4.10) \quad \|u - u_h\|_1 \leq ch^r \|u\|_{r+1}$$

as long as the solution  $u$  belongs to  $H^{r+1}$ . If  $u$  is less smooth, the rate of convergence will be decreased accordingly.

In short, one proves the error estimate in  $H^1$  by using quasi-optimality in  $H^1$  (which comes from coercivity), and then finite element approximation theory.

**8.2. Estimate in  $L^2$ .** Now let  $g \in L^2(\Omega)$  be a given function, and consider the computation of  $G(u) := \int_{\Omega} ug \, dx$ , which is a functional of the solution  $u$  of our Dirichlet problem. We ask how accurately  $G(u_h)$  approximates  $G(u)$ . To answer this, we define an auxiliary function  $\phi \in V$  by

$$b(w, \phi) = G(w), \quad w \in V.$$

This is simply a weak formulation of the Dirichlet problem

$$-\operatorname{div} a \operatorname{grad} \phi = g \text{ in } \Omega, \quad \phi = 0 \text{ on } \partial\Omega.$$

We will assume that this problem satisfies  $H^2$  regularity, i.e., the solution  $\phi \in H^2(\Omega)$  and satisfies

$$\|\phi\|_2 \leq c\|g\|_0.$$

This is true, for example, if  $\Omega$  is either a convex Lipschitz domain or a smooth domain and  $a$  is a smooth coefficient.

REMARK. Note that we write  $b(w, \phi)$  with the trial function  $\phi$  second and the test function  $w$  first, the opposite as for the original problem (4.9). Since the bilinear form we are considering is symmetric, this is not a true distinction. But if we started with a nonsymmetric bilinear form, we would still define the auxiliary function  $\phi$  in this way. In short  $\phi$  satisfies a boundary value problem for the *adjoint* differential equation.

Now consider the error in  $G(u)$ :

$$G(u) - G(u_h) = \int_{\Omega} (u - u_h)g \, dx = b(u - u_h, \phi) = b(u - u_h, \phi - v)$$

for any  $v \in V_h$ , where the second equality comes from the definition of the auxiliary function  $\phi$  and the third from Galerkin orthogonality (4.3). Therefore

$$|G(u) - G(u_h)| \leq M \|u - u_h\|_1 \inf_{v \in V_h} \|\phi - v\|_1.$$

Now finite element approximation theory and 2-regularity tell us

$$\inf_{v \in V_h} \|\phi - v\|_1 \leq ch \|\phi\|_2 \leq ch \|g\|_0.$$

Thus

$$|G(u) - G(u_h)| \leq ch \|u - u_h\|_1 \|g\|_0 \leq ch^{r+1} \|u\|_{r+1} \|g\|_0.$$

In short, if  $g \in L^2(\Omega)$ , the error in  $G(u) = \int_{\Omega} ug \, dx$  is  $O(h^{r+1})$ , one power of  $h$  higher than the  $H^1$  error.

A very important special case is when  $g = u - u_h$ . Then  $G(u) - G(u_h) = \|u - u_h\|_0^2$ , so we have

$$\|u - u_h\|_0^2 \leq ch \|u - u_h\|_1 \|u - u_h\|_0,$$

or

$$\|u - u_h\|_0 \leq ch \|u - u_h\|_1 \leq ch^{r+1} \|u\|_{r+1}.$$

That is, the  $L^2$  error in the finite element method is one power of  $h$  higher than the  $H^1$  error.

REMARK. The idea of introducing an auxilliary function  $\phi$ , so we can express  $G(u - u_h)$  or  $\|u - u_h\|_0^2$  as  $b(u - u_h, \phi)$  and estimate it using Galerkin orthogonality is the Aubin–Nitsche duality method. If we use it to estimate  $G(u - u_h)$  where  $g$  is smoother than  $L^2$  and we have higher order elliptic regularity, we can get even higher order estimates, so called negative-norm estimates.

## 9. A posteriori error estimates and adaptivity

The error estimate (4.10) is a typical *a priori* error estimate for the finite element method. It indicates that, as long as we know a priori that the unknown solution of our problem belongs to  $H^{r+1}$ , then the error  $\|u - u_h\|_1$  will converge to zero as  $O(h^r)$ . By contrast an *a posteriori* error estimate attempts to bound the error in terms of  $u_h$ , allowing the error in the finite element solution to be approximated once the finite element solution itself has been calculated. One important use of a posteriori error estimates is in estimating how accurate the computed solution is. Another relates to the fact that the some a posteriori error estimates give a way of attributing the error to the different elements of the mesh. Therefore they suggest how the mesh might be refined to most effectively decrease the error (basically by subdividing the elements which are contributing a lot to the error). This is the basic idea of *adaptivity*, which we shall discuss below.

**9.1. The Clément interpolant.** First we need a new tool from finite element approximation theory. Suppose we are given a polygonal domain  $\Omega$  and a mesh of mesh size  $h$ . Let  $V_h$  be the usual Lagrange finite element space of degree  $r$ . Given a continuous function  $u$  on  $\bar{\Omega}$ , we may define the interpolant  $I_h u$  of  $u$  into  $V_h$  through the usual degrees of freedom. Then we have the error estimate

$$\|u - I_h u\|_t \leq ch^{s-t} \|u\|_s, \quad u \in H^s(\Omega),$$

valid for integers  $0 \leq t \leq 1$ ,  $2 \leq s \leq r + 1$ . See Theorem 4.6 (the constant  $c$  here depends only on  $r$  and the shape regularity of the mesh). We proved this result element-by-element, using the Bramble–Hilbert lemma and scaling. Of course this result implies that

$$\inf_{v \in V_h} \|u - v\|_t \leq ch^{s-t} \|u\|_s, \quad u \in H^s(\Omega),$$

for the same ranges of  $t$  and  $s$ . The restriction  $t \leq 1$  is needed, since otherwise the functions in  $V_h$ , being continuous but not generally  $C^1$ , do not belong to  $H^t(\Omega)$ . Here we are concerned with weakening the restriction  $s \geq 2$ , so we can say something about the approximation by piecewise polynomials of a function  $u$  that does not belong to  $H^2(\Omega)$ . We might hope for example that

$$\inf_{v \in V_h} \|u - v\|_0 \leq ch \|u\|_1, \quad u \in H^1(\Omega).$$

In fact, this estimate is true, and is important to the development of a posteriori error estimates and adaptivity. However it can not be proven using the usual interpolant  $I_h u$ , because  $I_h u$  is not defined unless the function  $u$  has well-defined point values at the node points of the mesh, and this is not true for a general function  $u \in H^1$ . (In 2- and 3-dimensions the Sobolev embedding theorem implies the existence of point values for function in  $H^2$ , but not in  $H^1$ .)

The way around this is through a different operator than  $I_h$ , called the Clément interpolant, or quasi-interpolant. For each polynomial degree  $r \geq 1$  and each mesh, the Clément interpolant  $\Pi_h : L^2(\Omega) \rightarrow V_h$  is a bounded linear operator. Its approximation properties are summarized in the following theorem.

**THEOREM 4.7** (Clément interpolant). *Let  $\Omega$  be a domain in  $\mathbb{R}^n$  furnished with a simplicial triangulation with shape constant  $\gamma$  and maximum element diameter  $h$ , let  $r$  be a positive integer, and let  $V_h$  denote the Lagrange finite element space of continuous piecewise polynomials of degree  $r$ . Then there exists a bounded linear operator  $\Pi_h : L^2(\Omega) \rightarrow V_h$  and a constant  $c$  depending only on  $\gamma$  and  $r$  such that*

$$\|u - \Pi_h u\|_t \leq ch^{s-t} \|u\|_s, \quad u \in H^s(\Omega),$$

for all  $0 \leq t \leq s \leq r + 1$ ,  $t \leq 1$ .

Now we define the Clément interpolant. Let  $\mu_i : C(\bar{\Omega}) \rightarrow \mathbb{R}$ ,  $i = 1, \dots, N$ , be the usual DOFs for  $V_h$  and  $\phi_i$  the corresponding basis functions. Thus the usual interpolant is

$$I_h u = \sum_i \mu_i(u) \phi_i, \quad u \in C(\bar{\Omega}).$$

To define the Clément interpolant we let  $S_i$  denote the support of  $\phi_i$ , i.e., the union of triangles where  $\phi_i$  is not identically zero (if  $\mu_i$  is a vertex degree of freedom this is the union of the elements with that vertex, if an edge degree of freedom, the union of the triangles with that edge, etc.). Denote by  $P_i : L^2(S_i) \rightarrow \mathcal{P}_r(S_i)$  the  $L^2$ -projection. Then we set

$$\Pi_h u = \sum_i \mu_i(P_i u) \phi_i, \quad u \in L^2(\Omega).$$

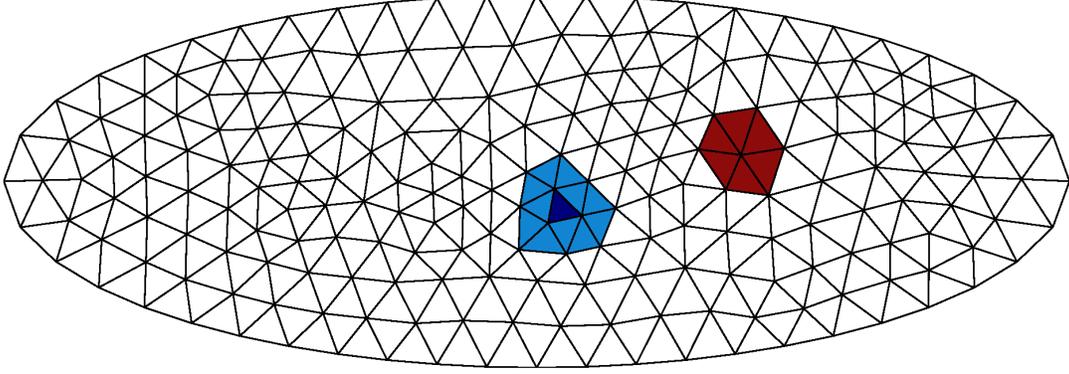
The usual interpolant  $I_h u$  is completely local in the sense that if  $u$  vanishes on a particular triangle  $T$ , then  $I_h u$  also vanishes on  $T$ . The Clément interpolation operator is not quite so local, but is nearly local in the following sense. If  $u$  vanishes on the set  $\tilde{T}$ , defined to be the union of the triangles that share at least one vertex with  $T$  (see Figure 4.9), then  $\Pi_h u$  vanishes on  $T$ . In fact, for any  $0 \leq t \leq s \leq r + 1$ ,

$$(4.11) \quad \|u - \Pi_h u\|_{H^t(T)} \leq ch_T^{s-t} \|u\|_{H^s(\tilde{T})}, \quad u \in H^s(\tilde{T}),$$

where the constant depends only on the shape regularity of the mesh and  $r$ . From (4.11), using the shape regularity, the estimate of Theorem 3.7 easily follows.

To avoid too much technicality, we shall prove (4.11) in the case of linear elements,  $r = 1$ . Thus we are interested in the case  $t = 0$  or  $1$  and  $t \leq s \leq 2$ . Let  $T$  be a particular triangle and let  $z_i$ ,  $\mu_i$ ,  $\phi_i$ ,  $S_i$  denote its vertices and corresponding DOFs, basis functions, and their supports, for  $i = 1, 2, 3$ . Note that it is easy to see that

$$(4.12) \quad \|\phi_i\|_{L^2(T)} \leq |T|^{1/2} \leq ch_T, \quad \|\text{grad } \phi_i\|_{L^2(T)} \leq ch_T^{-1} |T|^{1/2} \leq c,$$

FIGURE 4.9. Shown in brown  $S_z$  for a vertex  $z$  and in blue  $\tilde{T}$  for a triangle  $T$ .

where the constants may depend on the shape constant for  $T$ . Next, using the Bramble–Hilbert lemma and scaling on  $S_i$ , we have for  $0 \leq t \leq s \leq 2$ ,

$$(4.13) \quad \|u - P_i u\|_{H^t(S_i)} \leq ch_T^{s-t} \|u\|_{H^s(S_i)}.$$

In performing the scaling, we map the whole set  $S_i$  to the corresponding set  $\hat{S}_i$  using the affine map  $F$  which takes one of the triangles  $T$  in  $S_i$  to the unit triangle  $\hat{T}$ . The Bramble–Hilbert lemma is applied on the scaled domain  $\hat{S}_i$ . Although there is not just a single domain  $S_i$ —it depends on the number of triangles meeting at the vertex  $z_i$  and their shapes—the constant which arises when applying the Bramble–Hilbert lemma on the scaled domain can be bounded on the scaled domain in terms only of the shape constant for the triangulation (this can be established using compactness).

We also need one other bound. Let  $i$  and  $j$  denote the indices of two different vertices of  $T$ . For  $u \in L^2(\tilde{T})$ , both  $P_i u$  and  $P_j u$  are defined on  $T$ . If  $\hat{u}$  denotes the corresponding function on the scaled domain  $\hat{T}$ , then we have

$$\begin{aligned} \|P_i u - P_j u\|_{L^\infty(T)} &= \|\hat{P}_i \hat{u} - \hat{P}_j \hat{u}\|_{L^\infty(\hat{T})} \leq c \|\hat{P}_i \hat{u} - \hat{P}_j \hat{u}\|_{L^2(\hat{T})} \\ &\leq c (\|\hat{P}_i \hat{u} - \hat{u}\|_{L^2(\hat{T})} + \|\hat{P}_j \hat{u} - \hat{u}\|_{L^2(\hat{T})}) \\ &\leq c (\|\hat{P}_i \hat{u} - \hat{u}\|_{L^2(\hat{S}_i)} + \|\hat{P}_j \hat{u} - \hat{u}\|_{L^2(\hat{S}_j)}) \end{aligned}$$

where the first inequality comes from equivalence of norms on the finite dimensional space  $\mathcal{P}_1(\hat{T})$  and the second from the triangle inequality. Both of the terms on the right-hand side can be bounded using the Bramble–Hilbert lemma and scaled back to  $S_i$ , just as for (4.13). In this way we obtain the estimate

$$(4.14) \quad \|P_i u - P_j u\|_{L^\infty(T)} \leq ch_T^s |T|^{-1/2} \|u\|_{H^s(\tilde{T})}.$$

Therefore also

$$(4.15) \quad |\mu_i(P_i u - P_j u)| \leq ch_T^s |T|^{-1/2} \|u\|_{H^s(\tilde{T})}.$$

Now on the triangle  $T$  with vertices numbered  $z_1, z_2, z_3$  for simplicity,

$$\Pi_h u = \sum_{i=1}^3 \mu_i(P_i u) \phi_i.$$

Since  $P_1u$  is a linear polynomial on  $T$ ,

$$u - \Pi_h u = (u - P_1u) - \sum_{i=2}^3 \mu_i (P_i u - P_1 u) \phi_i.$$

The first term on the right hand side is bounded using (4.13):

$$\|u - P_1u\|_{H^t(T)} \leq ch_T^{s-t} \|u\|_{H^s(\tilde{T})}.$$

For the second term we have

$$\|\mu_i (P_i u - P_1 u) \phi_i\|_{H^t(T)} \leq \|P_i u - P_1 u\|_{L^\infty(T)} \|\phi_i\|_{H^t(T)},$$

which satisfies the desired bound by (4.12) and (4.15).

For the next section an important special case of (4.11) is

$$(4.16) \quad \|u - \Pi_h u\|_{L^2(T)} \leq ch_T \|u\|_{H^1(\tilde{T})}.$$

Another case is  $H^1$  boundedness:

$$(4.17) \quad \|u - \Pi_h u\|_{H^1(T)} \leq c \|u\|_{H^1(\tilde{T})}.$$

We draw one more conclusion, which we will need below. Let  $\hat{T}$  denote the unit triangle, and  $\hat{e}$  one edge of it. The trace theorem then tells us that

$$\|\hat{u}\|_{L^2(\hat{e})}^2 \leq c(\|\hat{u}\|_{L^2(\hat{T})}^2 + \|\text{grad } \hat{u}\|_{L^2(\hat{T})}^2), \quad \hat{u} \in H^1(\hat{T}).$$

If we use linear scaling to an arbitrary triangle, we get

$$\|u\|_{L^2(e)}^2 \leq c(h_T^{-1} \|u\|_{L^2(T)}^2 + h_T \|\text{grad } u\|_{L^2(T)}^2), \quad u \in H^1(T),$$

where the constant depends only on the shape constant of  $T$ . If we now apply this with  $u$  replaced by  $u - \Pi_h u$  and use (4.16) and (4.17), we get this bound for the Clément interpolant:

$$(4.18) \quad \|u - \Pi_h u\|_{L^2(e)} \leq ch_e^{1/2} \|u\|_{H^1(\tilde{T})},$$

where  $h_e$  is the length of the edge  $e$ ,  $T$  is a triangle containing  $e$ , and  $c$  depends only on the shape constant for the mesh.

**9.2. The residual and the error.** Consider our usual model problem

$$-\text{div } a \text{ grad } u = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega,$$

with a continuous positive coefficient  $a$  on  $\bar{\Omega}$  and  $f \in L^2$ . The weak formulation is to find  $u \in V$  satisfying

$$b(u, v) = F(v), \quad v \in V,$$

where  $V = \mathring{H}^1(\Omega)$  and

$$b(w, v) = \int a \text{ grad } w \cdot \text{grad } v \, dx, \quad F(v) = \int f v \, dx, \quad w, v \in V.$$

The bilinear form is bounded and coercive on  $\mathring{H}^1$ :

$$|b(w, v)| \leq M \|w\|_1 \|v\|_1, \quad b(v, v) \geq \gamma \|v\|_1^2, \quad w, v \in V.$$

Now we suppose that we have computed an approximation  $U$  of the solution  $u$  and we wish to assess the norm of the error  $u - U$  (we are most interested in the case  $U = u_h$ , the finite element solution). Just as for linear algebra problems, in which we find an

approximate solution to a linear system, we shall approach the error through the residual, which is computable. But what do we mean by the *residual* in the solution to such a weakly formulated equation? As discussed at the start of Section 5, the weak formulation may be viewed as an operator equation  $Lu = F$ , where  $L$  is a linear operator for  $V$  to  $V^*$ . In our case,  $V = \dot{H}^1$  and its dual,  $V^*$  is generally denoted  $H^{-1}$ , with the dual norm denoted by  $\|\cdot\|_{-1}$ . Thus the residual  $R(U) = F - LU \in H^{-1}$ , i.e., it is a linear functional on  $\dot{H}^1$ . Specifically,

$$R(U)w = F(w) - b(U, w), \quad w \in V.$$

Clearly  $R(U)w = b(u - U, w)$ . It follows immediately that  $|R(U)w| \leq M\|u - U\|_1\|w\|_1$  for all  $w \in \dot{H}^1$ , or, equivalently, that  $\|R(U)\|_{-1} \leq M\|u - U\|_1$ . On the other hand, taking  $w = u - U$  and using the coercivity, we get  $\gamma\|u - U\|_1^2 \leq R(U)(u - U) \leq \|R(U)\|_{-1}\|u - U\|_1$ . Thus

$$M^{-1}\|R(U)\|_{-1} \leq \|u - U\|_1 \leq \gamma^{-1}\|R(U)\|_{-1}.$$

In short, *the  $H^{-1}$  norm of the residual  $R(U)$  is equivalent to the  $H^1$  norm of the error  $u - U$ .*

**9.3. Estimating the residual.** Now let  $V_h$  be the Lagrange finite element subspace of  $V = \dot{H}^1$  corresponding to some mesh  $\mathcal{T}_h$  and some polynomial degree  $r$ , and let  $u_h$  be the corresponding finite element solution. We have just seen that we may estimate the error  $\|u - u_h\|_1$  by estimating the  $H^{-1}$  error in the residual

$$R(u_h)w = F(w) - b(u_h, w), \quad w \in V.$$

This quantity does not involve the unknown solution  $u$ , so we may hope to compute it *a posteriori*, i.e., after we have computed  $u_h$ .

We start by integrating by parts on each element  $T$  of the mesh to rewrite  $R(u_h)w$ :

$$\begin{aligned} R(u_h)w &= \sum_{T \in \mathcal{T}_h} \int_T (fw - a \operatorname{grad} u_h \cdot \operatorname{grad} w) dx \\ &= \sum_T \int_T (f + \operatorname{div} a \operatorname{grad} u_h)w dx - \sum_T \int_{\partial T} a \frac{\partial u_h}{\partial n_T} w ds. \end{aligned}$$

Consider the final sum. We can split each integral over  $\partial T$  into the sum of the integrals of the three edges of  $T$ . Each edge  $e$  which is not contained in the boundary comes in twice. For such an edge, let  $T_-$  and  $T_+$  be the two triangles which contain  $e$  and set

$$R_e(u_h) = -a \left( \frac{\partial u_h|_{T_-}}{\partial n_{T_-}} + \frac{\partial u_h|_{T_+}}{\partial n_{T_+}} \right) \in L^2(e)$$

on  $e$ . Since  $n_{T_-} = -n_{T_+}$  the term in parenthesis is the *jump* in the normal derivative of  $u_h$  across the edge  $e$ . Also, for  $T \in \mathcal{T}_h$ , we set  $R_T(u_h) = f + \operatorname{div} a \operatorname{grad} u_h \in L^2(T)$ . Then

$$(4.19) \quad R(u_h)w = \sum_T \int_T R_T(u_h)w dx + \sum_{e \in \mathcal{E}} \int_e R_e(u_h)w ds.$$

Next we use *Galerkin orthogonality*: since  $u_h$  is the finite element solution, we have

$$b(u - u_h, v) = 0, \quad v \in V_h.$$

In terms of the residual this says that

$$R(u_h)w = R(u_h)(w - v), \quad v \in V_h.$$

In particular, we may choose  $v = \Pi_h w$ , the Clément interpolant in this equation. Combining with (4.19) (with  $w$  replaced by  $w - \Pi_h w$ ) we get

$$\begin{aligned}
(4.20) \quad R(u_h)w &= R(u_h)(w - \Pi_h w) \\
&= \sum_T \int_T R_T(u_h)(w - \Pi_h w) dx + \sum_{e \in \mathring{\mathcal{E}}} \int_e R_e(u_h)(w - \Pi_h w) ds \\
&= \sum_T \left[ \int_T R_T(u_h)(w - \Pi_h w) dx + \frac{1}{2} \sum_{\substack{e \in \mathring{\mathcal{E}} \\ e \subset T}} \int_e R_e(u_h)(w - \Pi_h w) ds \right],
\end{aligned}$$

where in the last step we used the fact that each  $e \in \mathring{\mathcal{E}}$  belongs to 2 triangles. Next, we bound the terms in the brackets on the right hand side of (4.20) for  $w \in \mathring{H}^1$ . First we use (4.16) to get

$$\int_T R_T(u_h)(w - \Pi_h w) dx \leq \|R_T(u_h)\|_{L^2(T)} \|w - \Pi_h w\|_{L^2(T)} \leq ch_T \|R_T(u_h)\|_{L^2(T)} \|w\|_{H^1(\tilde{T})}.$$

In a similar way, but using (4.18), we obtain for  $e \subset T$ ,

$$\left| \int_e R_e(u_h)(w - \Pi_h w) ds \right| \leq ch_e^{1/2} \|R_e(u_h)\|_{L^2(e)} \|w\|_{H^1(\tilde{T})}.$$

Combining the last three estimates, we get

$$\begin{aligned}
|R(u_h)w| &\leq c \sum_T [h_T \|R_T(u_h)\|_{L^2(T)} + \sum_{e \subset T} h_e^{1/2} \|R_e(u_h)\|_{L^2(e)}] \|w\|_{H^1(\tilde{T})} \\
&\leq c \left\{ \sum_T [h_T^2 \|R_T(u_h)\|_{L^2(T)}^2 + \sum_{e \subset T} h_e \|R_e(u_h)\|_{L^2(e)}^2] \right\}^{1/2} \left( \sum_T \|w\|_{H^1(\tilde{T})}^2 \right)^{1/2} \\
&\leq c \left\{ \sum_T [h_T^2 \|R_T(u_h)\|_{L^2(T)}^2 + \sum_{e \subset T} h_e \|R_e(u_h)\|_{L^2(e)}^2] \right\}^{1/2} \|w\|_1,
\end{aligned}$$

where we invoked the shape regularity in the last step. Since this estimate holds for all  $w \in \mathring{H}^1$ , we have shown that

$$\|R(u_h)\|_{-1} \leq c \left\{ \sum_T [h_T^2 \|R_T(u_h)\|_{L^2(T)}^2 + \sum_{e \subset T} h_e \|R_e(u_h)\|_{L^2(e)}^2] \right\}^{1/2}.$$

In view of the equivalence of the  $H^{-1}$  norm of the residual and the  $H^1$  norm of the error established in the preceding subsection, this gives us the a posteriori error estimate

$$(4.21) \quad \|u - u_h\|_1 \leq c \left\{ \sum_T [h_T^2 \|R_T(u_h)\|_{L^2(T)}^2 + \sum_{e \subset T} h_e \|R_e(u_h)\|_{L^2(e)}^2] \right\}^{1/2}.$$

This is a key result. First of all, it gives us a bound on the norm of the error of the finite element solution in terms of quantities that can be explicitly computed (except for the unknown constant  $c$ ). Second, the error bound is the square root of a sum of terms associated to the individual triangles. Thus, we have a way of assigning portions of the error to the various elements. This will enable us to base our adaptive strategy on refining those triangles for which the corresponding portion of the error for either the triangle itself or for one of its edges is relatively large.

**9.4. A posteriori error indicators and adaptivity.** Specifically, we associate to each triangle an *error indicator*:

$$\eta_T^2 := h_T^2 \|R_T(u_h)\|_{L^2(T)}^2 + \frac{1}{2} \sum_{e \subset \partial T} h_e \|R_e(u_h)\|_{L^2(e)}^2$$

The factor of 1/2 is usually used, to account for the fact that each edge belongs to two triangles. In terms of the error indicators, we can rewrite the a posteriori estimate as

$$\|u - u_h\|_1 \leq c \left( \sum_T \eta_T^2 \right)^{1/2}.$$

Our basic adaptive strategy then proceeds via the following SOLVE-ESTIMATE-MARK-REFINE loop:

- SOLVE: Given a mesh, compute  $u_h$
- ESTIMATE: For each triangle  $T$  compute  $\eta_T$ . If  $(\sum_T \eta_T^2)^{1/2} \leq tol$ , quit.
- MARK: Mark those elements  $T$  for which  $\eta_T$  is too large for refinement.
- REFINE: Create a new mesh with the marked elements refined.

We have already seen how to carry out the SOLVE and ESTIMATE steps. There are a number of possible strategies for choosing which elements to mark. One of the simplest is *maximal marking*. We pick a number  $\rho$  between 0 and 1, compute  $\eta_{\max} = \max_T \eta_T$ , and refine those elements  $T$  for  $\eta_T \geq \rho \eta_{\max}$ . Another approach, which is usually preferred, imposes the *Dörfler marking* criterion, which requires that some collection of elements  $\mathcal{S}$  is marked so that  $\sum_{T \in \mathcal{S}} \eta_T^2 \geq \rho^2 \sum_{T \in \mathcal{T}_h} \eta_T^2$ , i.e., we mark enough elements that they account for a given portion  $\rho$  of the total error. The program on the next page shows one way to implement this (there are others).

Once we have marked the elements, there is the question of how to carry out the refinement to be sure that all the marked elements are refined and there is not too much additional refinement. In 2-dimensions this is quite easy. Most schemes are based either on dividing each triangle in two, or dividing the marked triangles into 4 congruent triangles. Generally, after refining the marked elements, additional elements have to be refined to avoid hanging nodes in which a vertex of an element fall in the interior of the edge of a neighboring element. In 3-dimensions things are more complicated, but good refinement schemes (which retain shape regularity and avoid hanging nodes) are known.

**9.5. Examples of adaptive finite element computations.** On the next page we present a bare-bones adaptive Poisson solver written in FEniCS, displayed on the next page. This code uses Lagrange piecewise linear finite elements to solve the Dirichlet problem

$$-\Delta u = 1 \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega,$$

with  $\Omega$  an L-shaped domain and  $f \equiv 1$ , with the error indicators and marking strategy described above. The solution behaves like  $r^{2/3} \sin(2\theta/3)$  in a neighborhood of the reentrant corner, and so is not in  $H^2$ . The results can be seen in Figure 4.10. The final mesh has 6,410 elements, all right isocles triangles, with hypotenuse length  $h$  ranging from 0.044 to 0.002. If we used a uniform mesh with the smallest element size, this would require over 3 million elements. Figure 4.11 displays an adaptive mesh in 3D.

```

"""
Adaptive Poisson solver using a residual-based energy-norm error
estimator

eta_h**2 = sum_T eta_T**2

with

eta_T**2 = h_T**2 ||R_T||_T**2 + c h_T ||R_dT||_dT**2

where

R_T = f + div grad u_h
R_dT = 2 avg(grad u_h * n) (2*avg is jump, since n switches sign across edges)

and a Dorfler marking strategy

Adapted by Douglas Arnold from code of Marie Rognes
"""

from dolfin import *
from sys import stdin
from numpy import zeros

# Stop when sum of eta_T**2 < tolerance or max_iterations is reached
tolerance = 0.04
max_iterations = 20

# Create initial mesh
mesh = Mesh("l-shape-mesh.xml")
mesh.order()
figure(0) # reuse plotting window

# Define boundary and boundary value for Dirichlet conditions
u0 = Constant(0.0)
def boundary(x, on_boundary):
    return on_boundary

# SOLVE - ESTIMATE - MARK - REFINER loop
for i in range(max_iterations):

    # *** SOLVE step
    # Define variational problem and boundary condition
    # Solve variational problem on current mesh
    V = FunctionSpace(mesh, "CG", 1)
    u = TrialFunction(V)
    v = TestFunction(V)
    f = Constant(1.0)
    a = inner(grad(u), grad(v))*dx
    L = f*v*dx
    u_h = Function(V)
    solve(a==L, u_h, DirichletBC(V, u0, boundary))

```

— continued on next page —

```

# *** ESTIMATE step
# Define cell and edge residuals
R_T = f + div(grad(u_h))
# get the normal to the cells
n = V.cell().n
R_dT = 2*avg(dot(grad(u_h), n))
# Will use space of constants to localize indicator form
Constants = FunctionSpace(mesh, "DG", 0)
w = TestFunction(Constants)
h = CellSize(mesh)
# Assemble squared error indicators, eta_T^2, and store into a numpy array
eta2 = assemble(h**2*R_T**2*w*dx + 4.*avg(h)*R_dT**2*avg(w)*dS) # dS is integral over interior edges only
eta2 = eta2.array()
# compute maximum and sum (which is the estimate for squared H1 norm of error)
eta2_max = max(eta2)
sum_eta2 = sum(eta2)
# stop error estimate is less than tolerance
if sum_eta2 < tolerance:
    print "Mesh %g: %d triangles, %d vertices, hmax = %g, hmin = %g, errest = %g" \
          % (i, mesh.num_cells(), mesh.num_vertices(), mesh.hmax(), mesh.hmin(), sqrt(sum_eta2))
    print "\nTolerance achieved. Exiting."
    break

# *** MARK step
# Mark cells for refinement for which eta > frac eta_max for frac = .95, .90, ...;
# choose frac so that marked elements account for a given part of total error
frac = .95
delfrac = .05
part = .5
marked = zeros(eta2.size, dtype='bool') # marked starts as False for all elements
sum_marked_eta2 = 0. # sum over marked elements of squared error indicators
while sum_marked_eta2 < part*sum_eta2:
    new_marked = (~marked) & (eta2 > frac*eta2_max)
    sum_marked_eta2 += sum(eta2[new_marked])
    marked += new_marked
    frac -= delfrac
# convert marked array to a MeshFunction
cells_marked = MeshFunction("bool", mesh, mesh.topology().dim())
cells_marked.array[:] = marked

# *** REFINE step
mesh = refine(mesh, cells_marked)
plot(mesh, title="Mesh q" + str(i))
print "Mesh %g: %d triangles, %d vertices, hmax = %g, hmin = %g, errest = %g" \
      % (i, mesh.num_cells(), mesh.num_vertices(), mesh.hmax(), mesh.hmin(), sqrt(sum_eta2))
stdin.readline()

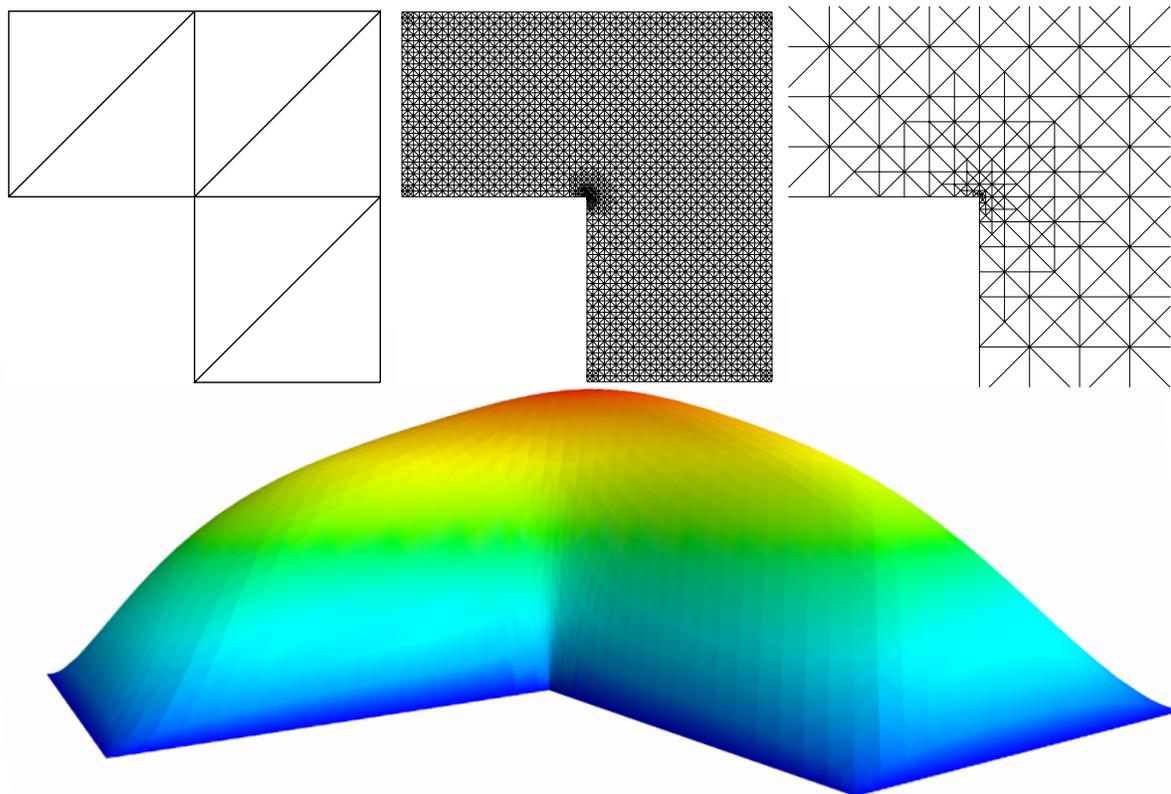
plot(mesh)
interactive()

```

**9.6. Nonlinear problems.** So far we have only discussed linear PDE. In many situations in which PDE models are applied, linear PDE are a simplification, which often is not sufficiently accurate. For example, in our model problem  $-\operatorname{div} a \operatorname{grad} u = f$  modeling a steady-state temperature distribution, the thermal conductivity  $a$  might depend on the temperature giving a nonlinear equation  $-\operatorname{div} a(u(x)) \operatorname{grad} u(x) = 0$ . Dependence on the temperature gradient is possible as well, and we might have convection and source terms, which might also depend on the temperature or gradient, leading to an equation of the form

$$(4.22) \quad -\operatorname{div}[a(u(x), \operatorname{grad} u(x)) \operatorname{grad} u(x)] - f(u(x), \operatorname{grad} u(x)) = 0,$$

FIGURE 4.10. Adaptive solution of Poisson's equation by the FEniCS program on the preceding page. Shown are the input mesh, the computed adaptive mesh, and a blow-up of that mesh near the re-entrant corner, as well as the final solution.



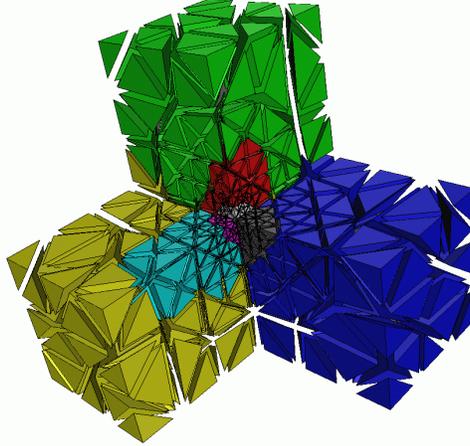
where  $a$  and  $f$  are functions of  $n+1$  variables in  $n$  dimensions (or  $2n+1$ : they might depend explicitly on  $x$  as well). A PDE of the form (4.22) is called *quasilinear*. If  $a$  is independent of  $u$  and  $\text{grad } u$ , so the only nonlinearity arises in the lower order terms in  $f$ , then it is called *semilinear*. As long as the coefficient  $a$  is everywhere positive (or a symmetric positive definite matrix), (4.22) is elliptic and we can hope to treat it by the sorts of finite element methods discussed heretofore. Some simple examples of such PDE are

$$-\Delta u + \lambda e^u = 0, \quad -\text{div}(1 + e^{-u^2}) \text{grad } u = 0, \quad -\text{div} \frac{1}{\sqrt{1 + |\text{grad } u|^2}} \text{grad } u = 0.$$

The first of these, called Bratu's problem, is semilinear and arises in combustion modeling. The other two are quasilinear. The third is the minimal surface equation, satisfied by functions whose graphs are minimal area surfaces subject to their boundary conditions (like soap bubbles on a frame).

The theory needed to analyze nonlinear PDE, e.g., to prove existence, uniqueness, continuous dependence, and various qualitative behaviors, is extensive, diverse, often complex, and an area of active research. Many different approaches have been developed in order to

FIGURE 4.11. An adaptive mesh in 3-dimensions produced by Michael Holst using his MC code. (The colors relate to partitioning among processors for parallel computation.)



address different equations. We will not consider these at all, but briefly consider how one might devise numerical methods to compute the solution to a nonlinear elliptic PDE, assuming that a locally unique solution exists (“locally unique” means that in some neighborhood of the solution no other solution exists). In this case, the general approach to computation is to approximate the solution of the nonlinear problem by solving a sequence of linear problems.

Consider the quasilinear PDE (4.22) subject (for example) to the Dirichlet boundary condition  $u = g$  on  $\partial\Omega$ . We obtain a weak formulation just as in the linear case, by multiplying the equation by a test function  $v$  satisfying homogeneous Dirichlet boundary conditions and integrating over  $\Omega$  by parts. Thus we obtain the nonlinear weak formulation: find  $u$  which is equal to  $g$  on  $\partial\Omega$  and such that

$$(4.23) \quad F(u, v) := \int a(u, \text{grad } u) \text{grad } u \cdot \text{grad } v \, dx - \int f(u, \text{grad } u) v \, dx = 0,$$

for all test functions  $v$  vanishing on  $\partial\Omega$ . We might use the space  $H^1(\Omega)$  as the space for the trial and test functions, as in the linear case, although, depending on the nonlinearity, more complicated function spaces may be needed to insure that the integrals all exist. This is an issue for the analysis of the numerical method, but need not concern us here where we will only discuss the formulation of algorithms. Note that the bivariate form  $F(u, v)$  is not bilinear. It remains linear in  $v$ , but is nonlinear in  $u$ .

To solve the nonlinear problem, we use Galerkin’s method as in the linear case. Thus we choose a finite dimensional space  $V_h$  for the trial and test functions (satisfying the boundary conditions), such as a finite element space based on a mesh and Lagrange finite elements. Then the Galerkin method seeks  $u_h \in V_h$  such that

$$(4.24) \quad F(u_h, v) = 0 \text{ for all } v \in V_h.$$

If we choose a basis  $\phi_i$ ,  $i = 1, \dots, n$ , for  $V_h$  and expand  $u_h = \sum_{j=1}^n U_j \phi_j$ , then the coefficients  $U_j$  may be determined from the system of equations

$$(4.25) \quad F\left(\sum_j U_j \phi_j, \phi_i\right) = 0, \quad i = 1, \dots, n,$$

which is a (nonlinear) system of  $n$  equations in  $n$  unknowns.

9.6.1. *Picard iteration.* Let

$$b(w; u, v) = \int a(w, \text{grad } w) \text{grad } u \cdot \text{grad } v \, dx - \int f(w, \text{grad } w) v \, dx$$

so the form in (4.23) is  $F(u, v) = b(u; u, v)$ . If we fix some  $w \in V_h$  the problem of finding  $u_h(w) \in V_h$  such that

$$b(w; u_h(w), v) = 0 \text{ for all } v \in V_h,$$

is a standard linear finite element problem. This defines a mapping  $w \mapsto u_h$  from  $V_h$  to itself. If  $u_h$  is a fixed point of this map, then it satisfies  $b(u_h; u_h, v) = 0$  for all  $v \in V_h$  which is the desired Galerkin equation (4.24). Thus we may try to solve (4.24) by fixed point iteration, with each iteration requiring the solution of a linear finite element system. This approach is called *Picard iteration*. Thus the basic iteration takes the form:

```
choose initial iterate  $u_h^0 \in V_h$ 
for  $i = 0, 1, \dots$ 
  find  $u_h^{i+1} \in V_h$  such that
```

$$\int a(u_h^i, \text{grad } u_h^i) \text{grad } u_h^{i+1} \cdot \text{grad } v \, dx = \int f(u_h^i, \text{grad } u_h^i) v \, dx, \quad v \in V_h.$$

```
end
```

Thus, in each iteration we solve a linear finite element problem, where the coefficients depend on the previously computed iterate. For example, for the minimal surface equation, the nonlinear weak formulation is

$$\int \frac{1}{\sqrt{1 + |\text{grad } u|^2}} \text{grad } u \cdot \text{grad } v \, dx = 0, \quad v \in \mathring{H}^1,$$

so the Picard iteration seeks  $u_h^{i+1} \in V_h$  satisfying the Dirichlet boundary conditions and

$$\int \frac{1}{\sqrt{1 + |\text{grad } u_h^i|^2}} \text{grad } u_h^{i+1} \cdot \text{grad } v \, dx = 0, \quad v \in \mathring{V}_h.$$

which is a standard linear finite element system.

The Picard iteration does not always converge, but it often does, especially if the initial guess is reasonably close to the exact solution. When it does converge, it typically does so with a linear rate of convergence to the solution of the nonlinear Galerkin equations. This can be quite slow, requiring many linear solves.

9.6.2. *Newton iteration.* We start by recalling what it means to linearize a system of  $n$  algebraic equations in  $n$  unknowns. We write the system as  $G(u) = 0$  where  $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is some function (supposed smooth), and the solution  $u$  is sought in  $\mathbb{R}^n$ . To linearize around some  $u^0 \in \mathbb{R}^n$ , not too far from the solution  $u$ , we replace  $u$  in the equations with a perturbation  $u^0 + \delta u$  of  $u^0$ , where  $\delta u \in \mathbb{R}^n$  is expected to be small. Thus we want  $G(u^0 + \delta u)$  to vanish, or nearly so. Expanding this quantity via Taylor's theorem gives

$$G(u^0 + \delta u) = G(u^0) + DG(u^0)\delta u + \cdots,$$

where  $DG(u^0)$  is a linear operator from  $\mathbb{R}^n$  to  $\mathbb{R}^n$  (its matrix is  $\partial G_i / \partial u_j$  evaluated at  $u^0$ ), and the dots indicate terms that are quadratic in  $\delta u$ . If we drop the quadratic terms and set the result equal to zero, we get a linear system of  $n$  equations in  $n$  unknowns to solve for  $\delta u$ :

$$DG(u^0)\delta u = -G(u^0).$$

Newton's method defines  $u^1$  to be  $u^0 + \delta u$ , which is hopefully an improved approximation of the solution  $u$ . It continues by linearizing  $G(u) = 0$  about  $u^1$  to find  $u^2$ , etc. The typical behavior of Newton's method is that it converges if the initial iterate  $u^0$  is chosen close enough to the solution  $u$ , and in that case the convergence is very fast—quadratic. It may, however, be difficult to find a suitably close initial iterate.

Just as for a nonlinear algebraic system, Newton's method may be applied to a nonlinear PDE, or to the finite element discretization of a nonlinear PDE. We can apply it directly to the nonlinear PDE (4.22), and then solve the resulting linear PDE by converting it into weak form, and then discretizing it with Galerkin's method. Alternatively, we can start with the nonlinear weak formulation (4.23), linearize that, obtaining a linear weak formulation, which we can discretize by Galerkin's method. A third possibility is to start with the nonlinear algebraic system (4.24) obtained by discretizing the nonlinear weak formulation, and implement Newton's method for this nonlinear algebraic system. It turns out that all three approaches are equivalent: the final discrete iterates computed are the same for all three. (After reading this section, try to prove this—it is a good test of your understanding.) We prefer to describe the middle alternative: linearization of the nonlinear weak formulation.

For simplicity, consider the case of (4.23) in which  $f$  vanishes, so the nonlinear weak formulation seeks  $u \in H^1$  with given Dirichlet boundary values such that

$$(4.26) \quad F(u, v) := \int a(u, \text{grad } u) \text{grad } u \cdot \text{grad } v \, dx = 0, \quad v \in \mathring{H}^1.$$

Here we are assuming that the coefficient  $a = a(y, z)$  is a smooth real-valued function of a scalar variable  $y$  and vector variable  $z$ . Now suppose we have an approximation  $u^0 \in H^1$  to  $u$ , which we suppose satisfies the boundary conditions. Now

$$a(u^0 + \delta u, \text{grad}(u^0 + \delta u)) = a + \frac{\partial a}{\partial y} \delta u + \sum_{j=1}^n \frac{\partial a}{\partial z_j} \frac{\partial \delta u}{\partial x_j} + \cdots,$$

where on the right-hand side  $a$  and its partial derivatives are evaluated at  $(u^0, \text{grad } u^0)$ , and the dots represent terms which are quadratic or higher in  $\delta u$  and its derivatives. Thus the

integrand of (4.26) becomes

$$\begin{aligned} a(u, \text{grad } u) \text{grad } u \cdot \text{grad } v &= a(u^0 + \delta u, \text{grad}(u^0 + \delta u)) \text{grad}(u^0 + \delta u) \cdot \text{grad } v \\ &= a \text{grad } u^0 \cdot \text{grad } v + a \text{grad } \delta u \cdot \text{grad } v + \frac{\partial a}{\partial y} \delta u \text{grad } u^0 \cdot \text{grad } v + \sum_{j=1}^n \frac{\partial a}{\partial z_j} \frac{\partial \delta u}{\partial x_j} \text{grad } u^0 \cdot \text{grad } v + \dots \\ &\quad A \text{grad } \delta u \cdot \text{grad } v + \delta u B \cdot \text{grad } v + a \text{grad } u^0 \cdot \text{grad } v + \dots, \end{aligned}$$

where  $A$  is the matrix-valued functions

$$A_{ij} = a + \frac{\partial a}{\partial z_j} \frac{\partial u^0}{\partial x_i},$$

and  $B = (\partial a / \partial y) \text{grad } u^0$ , both evaluated at  $(u^0, \text{grad } u^0)$ . Thus  $\delta u \in \mathring{H}^1$  is determined by the problem

$$(4.27) \quad \int (A \text{grad } \delta u \cdot \text{grad } v + \delta u B \cdot \text{grad } v) dx = - \int a \text{grad } u^0 \cdot \text{grad } v dx, \quad v \in \mathring{H}^1,$$

which is the linear weak formulation of a PDE (with homogeneous Dirichlet boundary conditions, even though the nonlinear problem had inhomogeneous boundary conditions).

To implement Newton's method for the finite element method, we start with an approximation  $u_h^0$  in the finite element space, satisfying the Dirichlet boundary conditions, we then define  $\delta u_h \in \mathring{V}^h$  by the linear weak formulation (4.27) with the test function  $v$  restricted to  $\mathring{V}_h$ , and set  $u_h^1 = u_h^0 + \delta u_h \in V_h$ , as the next iterate.

## Time-dependent problems

So far we have considered the numerical solution of elliptic PDEs. In this chapter we will consider some parabolic and hyperbolic PDEs.

### 1. Finite difference methods for the heat equation

In this section we consider the Dirichlet problem for the heat equation: find  $u : \bar{\Omega} \times [0, T] \rightarrow \mathbb{R}$  such that

$$(5.1) \quad \frac{\partial u}{\partial t}(x, t) - \Delta u(x, t) = f(x, t), \quad x \in \Omega, \quad 0 \leq t \leq T,$$

$$(5.2) \quad u(x, t) = 0, \quad x \in \partial\Omega, \quad 0 \leq t \leq T.$$

Since this is a time-dependent problem, we also need an initial condition:

$$u(x, 0) = u_0(x), \quad x \in \Omega.$$

For simplicity, we will assume  $\Omega = (0, 1) \times (0, 1)$  is the unit square in  $\mathbb{R}^2$  (or the unit interval in  $\mathbb{R}$ ). Let us consider first discretization in space only, which we have already studied. Following the notations we used in Chapter 2, we use a mesh with spacing  $h = 1/N$ , and let  $\Omega_h$  be the set of interior mesh points,  $\Gamma_h$  the set of boundary mesh points, and  $\bar{\Omega}_h = \Omega_h \cup \Gamma_h$ . The semidiscrete finite difference method is: find  $u_h : \bar{\Omega}_h \times [0, T] \rightarrow \mathbb{R}$  such that

$$\frac{\partial u_h}{\partial t}(x, t) - \Delta_h u_h(x, t) = f(x, t), \quad x \in \Omega_h, \quad 0 \leq t \leq T,$$

$$u_h(x, t) = 0, \quad x \in \partial\Omega_h, \quad 0 \leq t \leq T.$$

$$u_h(x, 0) = u_0(x), \quad x \in \Omega_h.$$

If we let  $U_{mn}(t) = u_h((mh, nh), t)$ , then we may write the first equation as

$$U'_{mn}(t) - \frac{U_{m+1,n}(t) + U_{m-1,n}(t) + U_{m,n+1}(t) + U_{m,n-1}(t) - 4U_{mn}(t)}{h^2} = f_{mn}(t),$$

$$0 < m, n < N, \quad 0 \leq t \leq T.$$

Thus we have a system of  $(n-1)^2$  ordinary differential equations, with given initial conditions.

One could feed this system of ODEs to an ODE solver. But we shall consider simple ODE solution schemes, which are, after all, themselves finite difference schemes, and analyze them directly. We shall focus on three simple schemes, although there are much more sophisticated possibilities.

For a system of ODEs, find  $u : [0, T] \rightarrow \mathbb{R}^m$  such that

$$u'(t) = f(t, u(t)), \quad 0 \leq t \leq T, \quad u(0) = u_0,$$

(where  $f : [0, T] \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ ,  $u_0 \in \mathbb{R}^m$ ) the simplest discretization is Euler's method. For a given timestep  $k > 0$ , let  $t_j = jk$ ,  $j = 0, 1, \dots$ , and define  $u_j = u_h(t_j) \in \mathbb{R}^m$  for  $j = 0, 1, \dots$  by  $u_h(0) = u(0)$ , and

$$\frac{u_{j+1} - u_j}{k} = f(t_j, u_j), \quad j = 0, 1, \dots$$

Explicitly,

$$u_{j+1} = u_j + kf(t_j, u_j), \quad j = 0, 1, \dots$$

An alternative method is the backward Euler method or implicit Euler method

$$\frac{u_{j+1} - u_j}{k} = f(t_{j+1}, u_{j+1}), \quad j = 0, 1, \dots$$

This method involves solving the algebraic system

$$u_{j+1} - kf(t_{j+1}, u_{j+1}) = u_j, \quad j = 0, 1, \dots$$

This is a linear or nonlinear algebraic system according to whether  $f$  is linear or nonlinear in  $u$ , i.e., according to whether the original ODE system is linear or nonlinear.

**1.1. Forward differences in time.** Now consider the application of Euler's method to the semidiscretized heat equation. We take the timestep  $k = T/M$  for some integer  $M$ , so that the discrete time values are  $0, k, 2k, \dots, T = Mk$ . Writing  $U_{mn}^j$  for  $u_h((mh, nh), jk)$  we get the explicit method

$$(5.3) \quad \frac{U_{mn}^{j+1} - U_{mn}^j}{k} - (\Delta_h U)_{mn}^j = f_{mn}^j,$$

i.e.,

$$U_{mn}^{j+1} = U_{mn}^j + k[(\Delta_h U)_{mn}^j + f_{mn}^j], \quad 0 < m, n < N, \quad j = 0, 1, \dots$$

This is called the forward-centered difference method for the heat equation, because it uses forward differences in time and centered differences in space.

We shall analyze the forward-centered scheme (5.3) as usual, by establishing consistency and stability. Let  $u_{mn}^j = u((mh, nh), t_j)$  denote the restriction of the exact solution, the consistency error is just

$$(5.4) \quad E_{mn}^j := \frac{u_{mn}^{j+1} - u_{mn}^j}{k} - (\Delta_h u)_{mn}^j - f_{mn}^j$$

$$(5.5) \quad = \left[ \frac{u_{mn}^{j+1} - u_{mn}^j}{k} - (\Delta_h u)_{mn}^j \right] - \left[ \left( \frac{\partial u}{\partial t} - \Delta u \right) ((mh, nh), jk) \right].$$

In Chapter 2 we used Taylor's expansion to get

$$|(\Delta_h u)_{mn}^j - \Delta u((mh, nh), jk)| \leq c_1 h^2,$$

where

$$c_1 = (\|\partial^4 u / \partial x^4\|_{L^\infty(\bar{\Omega} \times [0, T])} + \|\partial^4 u / \partial y^4\|_{L^\infty(\bar{\Omega} \times [0, T])}) / 12.$$

Even easier is

$$\left| \frac{u_{mn}^{j+1} - u_{mn}^j}{k} - \frac{\partial u}{\partial t}((mh, nh), jk) \right| \leq c_2 k,$$

where  $c_2 = \|\partial^2 u / \partial t^2 u\|_{L^\infty} / 2$ . Thus

$$|E_{mn}^j| \leq c(k + h^2),$$

with  $c = \max(c_1, c_2)$ .

Next we establish a stability result. Suppose that a mesh function  $U_{mn}^j$  satisfies (5.3). We want to bound an appropriate norm of the mesh function in terms of an appropriate norm of the function  $f_{mn}^j$  on the right hand side. For the norm, we use the max norm:

$$\|U\|_{L^\infty} = \max_{0 \leq j \leq M} \max_{0 \leq m, n \leq N} |U_{mn}^j|.$$

Write  $K^j = \max_{0 \leq m, n \leq N} |U_{mn}^j|$ , and  $F^j = \max_{0 \leq m, n \leq N} |f_{mn}^j|$ . From (5.3), we have

$$U_{mn}^{j+1} = \left(1 - \frac{4k}{h^2}\right)U_{mn}^j + \frac{k}{h^2}(U_{m-1,n}^j + U_{m+1,n}^j + U_{m,n-1}^j + U_{m,n+1}^j) + kf_{mn}^j.$$

Now we make the assumption that  $4k/h^2 \leq 1$ . Then the 5 coefficients of  $U$  on the right hand side are all nonnegative numbers which add to 1, so it easily follows that

$$K^{j+1} \leq K^j + kF^j.$$

Therefore  $K^1 \leq K^0 + kF^0$ ,  $K^2 \leq K^0 + k(F^0 + F^1)$ , etc. Thus  $\max_{0 \leq j \leq M} K^j \leq K^0 + T \max_{0 \leq j < M} F^j$ , where we have used that  $kM = T$ . Thus, if  $U$  satisfies (5.3), then

$$(5.6) \quad \|U\|_{L^\infty} \leq \|U^0\|_{L^\infty(\Omega_h)} + T\|f\|_{L^\infty},$$

which is a stability result. We have thus shown stability under the condition that  $k \leq h^2/4$ . We say that the forward-centered difference method for the heat equation is *conditionally stable*.

Now we apply this stability result to the error  $e_{mn}^j = u_{mn}^j - U_{mn}^j$ , which satisfies

$$\frac{e_{mn}^{j+1} - e_{mn}^j}{k} - (\Delta_h e)_{mn}^j = E_{mn}^j,$$

with  $E$  the consistency error (so  $\|E\|_{L^\infty} \leq c(k + h^2)$ ). Note that  $E_{mn}^0 = 0$ , so the stability result gives

$$\|e\|_{L^\infty} \leq T\|E\|_{L^\infty}.$$

Using our estimate for the consistency error, we have proven the following theorem.

**THEOREM 5.1.** *Let  $u$  solve the heat equation (5.1), and let  $U_{mn}^j$  be determined by the forward-centered finite difference method with mesh size  $h = 1/N$  and timestep  $k = T/M$ . Suppose that  $k \leq h^2/4$ . Then*

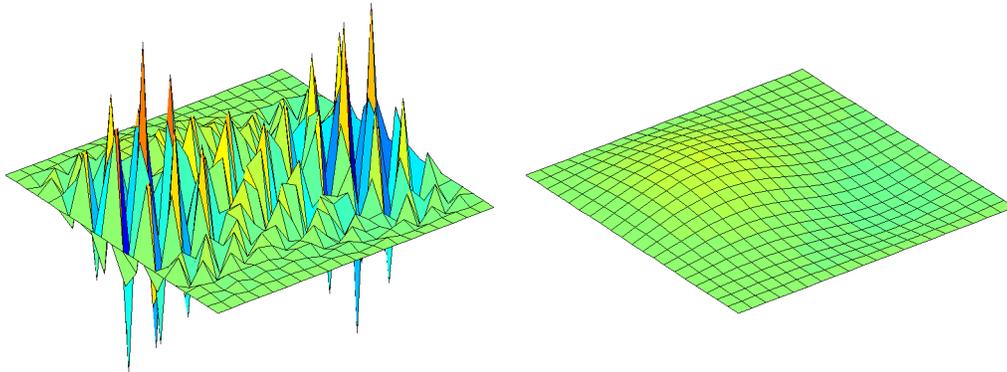
$$\max_{0 \leq j \leq M} \max_{0 \leq m, n \leq N} |u((mh, nh), jk) - U_{mn}^j| \leq C(k + h^2),$$

where  $C = cT$  with  $c$  as above.

In short,  $\|u - u_h\|_{L^\infty} = O(k + h^2)$ .

The requirement that  $4k/h^2 \leq 1$  is not needed for consistency. But it is required to prove stability, and a simple example shows that it is necessary for convergence. See Figure 5.1. An even simpler example would be the 1D case, for which we obtain stability under the condition  $k \leq h^2/2$ .

FIGURE 5.1. Centered differences in space and forward differences in time for the Dirichlet problem for the heat equation on the unit square. The mesh size is  $h = 1/20$ , and the timestep is  $k = 1/800$  on the left and  $k = 1/1600$  on the right. On the left we show the result after 18 time steps, i.e.,  $t = 18/800 = .0225$ . On the right we show the result after 36 time steps (the same time). The computation on the right remains stable for long times.



**1.2. Backward differences in time.** Next we consider of backward differences in time, i.e., the backward Euler method to solve the semidiscrete system. Then (5.3) becomes

$$(5.7) \quad \frac{U_{mn}^{j+1} - U_{mn}^j}{k} - (\Delta_h U)_{mn}^{j+1} = f_{mn}^{j+1}.$$

This is now an *implicit* method. Instead of writing down  $u_h^{j+1}$  explicitly in terms of  $u_h^j$ , we need to solve for the vector of its values,  $U_{mn}^{j+1}$ , from a system of linear equations:

$$U_{mn}^{j+1} - k\Delta_h U_{mn}^{j+1} = U_{mn}^j + kf_{mn}^{j+1}, \quad 0 < m, n < N,$$

or, written out,

$$(5.8) \quad (1 + 4\mu)U_{mn}^{j+1} - \mu(U_{m+1,n}^{j+1} + U_{m-1,n}^{j+1} + U_{m,n+1}^{j+1} + U_{m,n-1}^{j+1}) = U_{mn}^j + kf_{mn}^{j+1},$$

with  $\mu = k/h^2$ . This is a sparse system of  $(N-1)^2$  equations in  $(N-1)^2$  unknowns with the same sparsity pattern as  $\Delta_h$ . Since  $-\Delta_h$  is symmetric and positive definite, this system, whose matrix is  $I - k\Delta_h$ , is as well.

REMARK. The computational difference between explicit and implicit methods was very significant before the advent of fast solvers, like multigrid. Since such solvers reduce the computational work to the order of the number of unknowns, they are not so much slower than explicit methods.

Now suppose that (5.8) holds. Let  $K^j$  again denote the maximum of  $|U_{mn}^j|$ . Then there exists  $m, n$  such that  $K^{j+1} = sU_{mn}^{j+1}$  where  $s = \pm 1$ . Therefore  $sU_{mn}^{j+1} \geq sU_{m+1,n}^{j+1}$  and similarly for each of the other three neighbors. For this particular  $m, n$ , we multiply (5.8) by  $s$  and obtain

$$K^{j+1} = sU_{mn}^{j+1} \leq sU_{mn}^j + skf_{mn}^{j+1} \leq K^j + k\|f^{j+1}\|_{L^\infty}.$$

From this we get the stability result (5.6) as before, but now the stability is unconditional: it holds for any  $h, k > 0$ . We immediately obtain the analogue of the convergence theorem Theorem 5.1 for the backward-centered method.

**THEOREM 5.2.** *Let  $u$  solve the heat equation (5.1), and let  $U_{mn}^j$  be determined by the backward-centered finite difference method with mesh size  $h = 1/N$  and timestep  $k = T/M$ . Then*

$$\max_{0 \leq j \leq M} \max_{0 \leq m, n \leq N} |u((mh, nh), jk) - U_{mn}^j| \leq C(k + h^2),$$

where  $C = cT$  with  $c$  as above.

**1.3. Fourier analysis.** As we did for the Poisson problem, we can use Fourier analysis to analyze difference schemes for the heat equation. Recall that for a mesh function  $v$  on  $\Omega_h$  we defined the norm

$$\|u\|_h^2 = h^2 \sum_{m=1}^{N-1} \sum_{n=1}^{N-1} |u(mh, nh)|^2,$$

and the corresponding inner product. We then showed that  $-\Delta_h$  had an orthogonal basis of eigenfunctions  $\phi_{mn}$  with corresponding eigenvalues satisfying

$$2\pi^2 \approx \lambda_{1,1} \leq \lambda_{mn} \leq \lambda_{N-1, N-1} < 8/h^2.$$

Consider now the forward-centered difference equations for the heat equation  $\partial u / \partial t = \Delta u$ , with homogeneous Dirichlet boundary data, and given initial data  $u_0$  (for simplicity we assume  $f \equiv 0$ ). We have

$$U_{mn}^{j+1} = GU_{mn}^j,$$

where the  $G = I + k\Delta_h$ . By iteration, we have

$$\|U^j\| \leq \|G\|^j \|U^0\|.$$

Thus we have  $L^2$  stability if and only if the spectral radius of  $G$  is bounded by 1. Now the eigenvalues  $1 - k\lambda_{mn}$ , so they satisfy  $1 - 8k/h^2 < \mu < 1$ , so the spectral radius condition is satisfied if  $1 - 8k/h^2 \geq -1$ , i.e.,  $k \leq h^2/4$ . In this way, Fourier analysis leads us to the same conditional stability condition we obtained above.

If we consider instead the backward-centered difference scheme, then the corresponding operator  $G$  is  $G = (I + k\Delta_h)^{-1}$  with eigenvalues  $(1 + k\lambda_{mn})^{-1}$ , which has spectral radius bounded by 1 for any  $k > 0$ . Thus we obtain unconditional stability.

**1.4. Crank–Nicolson.** Although we are free to choose the timestep and space type as we like for the backward-centered method, accuracy considerations still indicate that we should take  $k = O(h^2)$ , so very small timesteps. It is natural to seek a method which is second order in time as well as space. If we use the trapezoidal method for time discretization, the resulting fully discrete scheme is called the Crank–Nicolson method. It is an implicit method given by

$$\frac{U_{mn}^{j+1} - U_{mn}^j}{k} - \frac{1}{2}[(\Delta_h U)_{mn}^{j+1} + (\Delta_h U)_{mn}^j] = \frac{1}{2}[f_{mn}^j + f_{mn}^{j+1}].$$

We leave it as an exercise to show that the Crank–Nicolson scheme is unconditionally stable scheme with respect to the  $L^2$  norm, with error  $O(h^2 + k^2)$ .

## 2. Finite element methods for the heat equation

In this section we consider the initial boundary value problem for the heat equation

$$\begin{aligned} \frac{\partial u}{\partial t}(x, t) - \operatorname{div} a(x) \operatorname{grad} u(x, t) &= f(x, t), \quad x \in \Omega, \quad 0 \leq t \leq T, \\ u(x, t) &= 0, \quad x \in \partial\Omega, \quad 0 \leq t \leq T, \\ u(x, 0) &= u_0(x), \quad x \in \Omega. \end{aligned}$$

We have allowed the thermal conductivity  $a$  to be variable, assuming only that it is bounded above and below by positive constants, since this will cause no additional complications. We could easily generalize further, allowing a variable specific heat (coefficient of  $\partial u/\partial t$ ), lower order terms, and different boundary conditions.

Now we consider finite elements for spatial discretization. To derive a weak formulation, we multiply the heat equation (5.1) by a test function  $v(x) \in \dot{H}^1(\Omega)$  and integrate over  $\Omega$ . This gives

$$\int_{\Omega} \frac{\partial u}{\partial t}(x, t) v(x) dx + \int_{\Omega} a(x) \operatorname{grad} u(x, t) \cdot \operatorname{grad} v(x) dx = \int_{\Omega} f(x, t) v(x) dx, \quad 0 \leq t \leq T.$$

Writing  $\langle \cdot, \cdot \rangle$  for the  $L^2(\Omega)$  inner product and  $b(w, v) = \int_{\Omega} a \operatorname{grad} w \cdot \operatorname{grad} v dx$ , we may write the weak formulation as

$$\left\langle \frac{\partial u}{\partial t}, v \right\rangle + b(u, v) = \langle f, v \rangle, \quad v \in \dot{H}^1(\Omega), \quad 0 \leq t \leq T.$$

For the finite element method it is often useful to think of  $u(x, t)$  as a function of  $t$  taking values in functions of  $x$ . Specifically, we may think of  $u$  mapping  $t \in [0, T]$  to  $u(\cdot, t) \in \dot{H}^1(\Omega)$ . Specifically, we may seek the solution  $u \in C^1([0, T], \dot{H}^1(\Omega))$ , which means that  $u(\cdot, t) \in \dot{H}^1(\Omega)$  for each  $t$ , that  $u$  is differentiable with respect to  $t$ , and that  $\partial u(\cdot, t)/\partial t \in \dot{H}^1(\Omega)$  for all  $t$ . Thus when we write  $u(t)$ , we mean the function  $u(\cdot, t) \in \dot{H}^1(\Omega)$ .

Now let  $V_h \subset \dot{H}^1(\Omega)$  denote the usual space of Lagrange finite elements of degree  $r$  with respect to a triangulation of  $\Omega$ . For a *semidiscrete finite element approximation* we seek  $u_h$  mapping  $[0, T]$  into  $V_h$ , i.e.,  $u_h \in C^1([0, T], V_h)$ , satisfying

$$(5.9) \quad \left\langle \frac{\partial u_h}{\partial t}, v \right\rangle + b(u_h, v) = \langle f, v \rangle, \quad v \in V_h, \quad 0 \leq t \leq T.$$

We also need to specify an initial condition for  $u_h$ . For this we choose some  $u_h^0 \in V_h$  which approximates  $u_0$  (common choices are the  $L^2$  projection or the interpolant).

We now show that this problem may be viewed as a system of ODEs. For this let  $\phi_i$ ,  $1 \leq i \leq D$  be a basis for  $V_h$  (for efficiency we will choose a local basis). We may then write

$$u_h(x, t) = \sum_{j=1}^D \alpha_j(t) \phi_j(x).$$

Plugging this into (5.9) and taking the test function  $v = \phi_i$ , we get

$$\sum_j \langle \phi_j, \phi_i \rangle \alpha_j'(t) + \sum_j b(\phi_j, \phi_i) \alpha_j(t) = \langle f, \phi_i \rangle.$$

In terms of the *mass matrix*, stiffness matrix, and load vector:

$$M_{ij} = \langle \phi_j, \phi_i \rangle, \quad A_{ij} = \langle \phi_j, \phi_i \rangle, \quad F_i(t) = \langle f, \phi_i \rangle,$$

this can be written

$$M\alpha'(t) + A\alpha(t) = F(t), \quad 0 \leq t \leq T.$$

This is a system of linear ODEs for the unknown coefficients  $\alpha = (\alpha_j(t))$ . The initial condition  $u_h(0) = u_h^0$  can be written  $\alpha(0) = \alpha^0$  where  $u_h^0 = \sum_j \alpha_j^0 \phi_j$ .

We now turn to *fully discrete approximation* using the finite element method for discretization in space, and finite differences for discretization in time. Consider first using Euler's method for time discretization. This leads to the system

$$M \frac{\alpha^{j+1} - \alpha^j}{k} + A\alpha^j = F^j,$$

or

$$M\alpha^{j+1} = M\alpha^j + k(-A\alpha^j + F^j).$$

Notice that for finite elements this method is not truly explicit, since we have to solve an equation involving the mass matrix at each time step.

The backward Euler's method

$$M \frac{\alpha^{j+1} - \alpha^j}{k} + A\alpha^{j+1} = F^{j+1},$$

leads to a different linear system at each time step:

$$(M + kA)\alpha^{j+1} = M\alpha^j + kF^{j+1},$$

while Crank–Nicolson would give

$$\left(M + \frac{k}{2}A\right)\alpha^{j+1} = \left(M - \frac{k}{2}A\right)\alpha^j + \frac{k}{2}(F^j + F^{j+1}).$$

**2.1. Analysis of the semidiscrete finite element method.** Before analyzing a fully discrete finite element scheme, we analyze the convergence of the semidiscrete scheme, since it is less involved. The key to the analysis of the semidiscrete finite element method is to compare  $u_h$  not directly to  $u$ , but rather to an appropriate representative  $w_h \in C^1([0, T], V_h)$ . For  $w_h$  we choose the *elliptic projection* of  $u$ , defined by

$$(5.10) \quad b(w_h, v) = b(u, v), \quad v \in V_h, \quad 0 \leq t \leq T.$$

From our study of the finite element method for elliptic problems, we have the  $L^2$  estimate

$$(5.11) \quad \|u(t) - w_h(t)\| \leq ch^{r+1} \|u(t)\|_{r+1}, \quad 0 \leq t \leq T.$$

If we differentiate (5.10), we see that  $\partial w_h / \partial t$  is the elliptic projection of  $\partial u / \partial t$ , so

$$\left\| \frac{\partial u}{\partial t}(t) - \frac{\partial w_h}{\partial t}(t) \right\| \leq ch^{r+1} \left\| \frac{\partial u}{\partial t}(t) \right\|_{r+1}, \quad 0 \leq t \leq T.$$

Now

$$(5.12) \quad \begin{aligned} \left\langle \frac{\partial w_h}{\partial t}, v \right\rangle + b(w_h, v) &= \left\langle \frac{\partial w_h}{\partial t}, v \right\rangle + b(u, v) \\ &= \left\langle \frac{\partial(w_h - u)}{\partial t}, v \right\rangle + \langle f, v \rangle, \quad v \in V_h, \quad 0 \leq t \leq T. \end{aligned}$$

Let  $y_h = w_h - u_h$ . Subtracting (5.9) from (5.12), we get

$$\left\langle \frac{\partial y_h}{\partial t}, v \right\rangle + b(y_h, v) = \left\langle \frac{\partial(w_h - u)}{\partial t}, v \right\rangle, \quad v \in V_h, \quad 0 \leq t \leq T.$$

Now, for each  $t$  we choose  $v = y_h(t) \in V_h$ . Note that for any function  $y \in C^1([0, T]; L^2(\Omega))$ ,

$$\|y\| \frac{d}{dt} \|y\| = \frac{1}{2} \frac{d}{dt} \|y\|^2 = \left\langle \frac{\partial y}{\partial t}, y \right\rangle.$$

Thus we get

$$(5.13) \quad \|y_h\| \frac{d}{dt} \|y_h\| + b(y_h, y_h) = \left\langle \frac{\partial(w_h - u)}{\partial t}, y_h \right\rangle \leq \left\| \frac{\partial(w_h - u)}{\partial t} \right\| \|y_h\|,$$

so

$$\frac{d}{dt} \|y_h\| \leq \left\| \frac{\partial(w_h - u)}{\partial t} \right\| \leq ch^{r+1} \left\| \frac{\partial u}{\partial t}(t) \right\|_{r+1}.$$

This holds for each  $t$ . Integrating over  $[0, t]$ , we get

$$\|y_h(t)\| \leq \|y_h(0)\| + ch^{r+1} \left\| \frac{\partial u}{\partial t} \right\|_{L^1([0, T]; H^{r+1}(\Omega))}.$$

For  $y_h(0)$  we have

$$\|y_h(0)\| = \|w_h(0) - u_h(0)\| \leq \|w_h(0) - u(0)\| + \|u_0 - u_h(0)\| \leq ch^{r+1} \|u_0\|_{r+1} + \|u_0 - u_h(0)\|.$$

Thus, assuming that the exact solution is sufficiently smooth and the initial data  $u_h(0)$  is chosen so that  $\|u_0 - u_h(0)\| = O(h^{r+1})$ , we have

$$\|y_h\|_{L^\infty([0, T]; L^2(\Omega))} = O(h^{r+1}).$$

Combining this estimate with the elliptic estimate (5.11) we get an estimate on the error

$$\|u - u_h\|_{L^\infty([0, T]; L^2(\Omega))} = O(h^{r+1}).$$

**REMARK.** We can put this analysis into the framework of consistency and stability introduced in Section 3. We take our discrete solution space  $X_h$  as  $C^1([0, T]; V_h)$ , and the discrete operator  $L_h : X_h \rightarrow Y_h := C([0, T]; V_h^*)$  is

$$(L_h u_h)(v) = \left\langle \frac{\partial u_h}{\partial t}, v \right\rangle + b(u_h, v), \quad u_h \in X_h, \quad v \in V_h, \quad 0 \leq t \leq T.$$

Thus our numerical method is to find  $u_h \in X_h$  such that  $L_h u_h = F_h$ , where

$$F_h(v) = \int f v \, dx, \quad v \in V_h, \quad 0 \leq t \leq T.$$

As a representative  $r_h u \in X_h$  of the exact solution  $u$  we use the elliptic projection  $w_h$ . Then the consistency error is given by

$$E(v) := \left\langle \frac{\partial w_h}{\partial t}, v \right\rangle + b(w_h, v) - \langle f, v \rangle, \quad v \in V_h, \quad 0 \leq t \leq T.$$

In the first part of our analysis we showed that

$$E(v) = \left\langle \frac{\partial(w_h - u)}{\partial t}, v \right\rangle,$$

so  $\|E\| = O(h^{r+1})$ , where the norm we use on  $Y_h$  is

$$\|E\| = \int_0^T \sup_{0 \neq v \in V_h} \frac{|E(v)|}{\|v\|} dt.$$

The second part of the analysis was a stability result. Essentially we showed that if  $u_h \in X_h$  and  $F_h \in Y_h$  satisfy  $L_h u_h = F_h$ , then

$$\max_{0 \leq t \leq T} \|u_h\| \leq \|u_h(0)\| + \|F_h\|.$$

REMARK. In the case of finite elements for elliptic problems, we first got an estimate in  $H^1$ , then an estimate in  $L^2$ , and I mentioned that there are others possible. In the case of the parabolic problem, there are many other estimates we could derive in different norms in space or time or both. For example, by integrating (5.13) in time we get that  $\|y_h\|_{L^2([0,T];H^1(\Omega))} = O(h^{r+1})$ . For the elliptic projection we have  $\|u - w_h\|_{H^1(\Omega)} = O(h^r)$  for each  $t$ , so the triangle inequality gives  $\|u - u_h\|_{L^2([0,T];H^1(\Omega))} = O(h^r)$ .

**2.2. Analysis of a fully discrete finite element method.** Now we turn to the analysis of a fully discrete scheme: finite elements in space and backward Euler in time. Writing  $u_h^j$  for  $u_h(\cdot, jk)$  (with  $k$  the time step), the scheme is

$$(5.14) \quad \left\langle \frac{u_h^{j+1} - u_h^j}{k}, v \right\rangle + b(u^{j+1}, v) = \langle f^{j+1}, v \rangle, \quad v \in V_h, \quad j = 0, 1, \dots$$

We initialize the iteration by choosing  $u_h^0 \in V_h$  to be, e.g., the interpolant,  $L^2$  projection, or elliptic projection. Notice that, at each time step, we have to solve the linear system

$$(M + kA)\alpha^{j+1} = M\alpha^j + kF^{j+1},$$

where  $\alpha^j$  is the vector of coefficients of  $u_h^j$  with respect to a basis, and  $M$ ,  $A$ , and  $F$ , are the mass matrix, stiffness matrix, and load vector respectively.

To analyze this scheme, we proceed as we did for the semidiscrete scheme, with some extra complications coming from the time discretization. In particular, we continue to use the elliptic projection  $w_h$  as a representative of  $u$ . Thus the consistency error is given by

$$\begin{aligned} & \left\langle \frac{w_h^{j+1} - w_h^j}{k}, v \right\rangle + b(w_h^{j+1}, v) - \langle f^{j+1}, v \rangle \\ &= \left\langle \frac{u^{j+1} - u^j}{k}, v \right\rangle + b(u^{j+1}, v) - \langle f^{j+1}, v \rangle + \left\langle \frac{(w_h^{j+1} - u^{j+1}) - (w_h^j - u^j)}{k}, v \right\rangle \\ &= \left\langle \frac{u^{j+1} - u^j}{k} - \frac{\partial u^{j+1}}{\partial t}, v \right\rangle + \left\langle \frac{(w_h^{j+1} - u^{j+1}) - (w_h^j - u^j)}{k}, v \right\rangle = \langle z^j, v \rangle, \end{aligned}$$

where the last line gives the definition of  $z^j$ . Next we estimate the two terms that comprise  $z^j$ , in  $L^2$ . First we have

$$\left\| \frac{u^{j+1} - u^j}{k} - \frac{\partial u^{j+1}}{\partial t} \right\| \leq \frac{k}{2} \left\| \frac{\partial^2 u}{\partial t^2} \right\|_{L^\infty(L^2)},$$

by Taylor's theorem. Next,

$$\frac{(w_h^{j+1} - u^{j+1}) - (w_h^j - u^j)}{k} = \frac{1}{k} \int_{jk}^{(j+1)k} \frac{\partial}{\partial t} [w_h(s) - u(s)] ds,$$

so

$$\left\| \frac{(w_h^{j+1} - u^{j+1}) - (w_h^j - u^j)}{k} \right\| \leq ch^{r+1} \left\| \frac{\partial u}{\partial t} \right\|_{L^\infty([jk, (j+1)k]; H^{r+1}(\Omega))}.$$

Thus we have obtained a bound on the consistency error:

$$\left\langle \frac{w_h^{j+1} - w_h^j}{k}, v \right\rangle + b(w_h^{j+1}, v) - \langle f^{j+1}, v \rangle = \langle z^j, v \rangle, \quad v \in V_h, \quad j = 0, 1, \dots$$

with

$$\|z^j\| \leq c(k \left\| \frac{\partial^2 u}{\partial t^2} \right\|_{L^\infty([0, T]; L^2(\Omega))} + h^{r+1} \left\| \frac{\partial u}{\partial t} \right\|_{L^\infty([0, T]; H^{r+1}(\Omega))}) =: E, \quad j = 0, 1, \dots$$

Combining with the scheme (5.14), we get (for  $y_h = w_h - u_h$ )

$$\left\langle \frac{y_h^{j+1} - y_h^j}{k}, v \right\rangle + b(y_h^{j+1}, v) = \langle z^j, v \rangle, \quad v \in V_h.$$

We conclude the argument with a stability argument. Choose  $v = y_h^{j+1} \in V_h$ . This becomes:

$$\|y_h^{j+1}\|^2 + kb(y_h^{j+1}, y_h^{j+1}) = \langle y_h^j + kz^j, y_h^{j+1} \rangle,$$

so

$$\|y_h^{j+1}\| \leq \|y_h^j\| + kE,$$

and, by iteration,

$$\max_{0 \leq j \leq M} \|y_h^j\| \leq \|y_h^0\| + TE.$$

In this way we prove that

$$\max_{0 \leq j \leq M} \|u^j - w_h^j\| = O(k + h^{r+1}).$$

Exercise for the reader: analyze the convergence of the fully discrete finite element method using Crank–Nicolson for time discretization. In the stability argument, you will want to use the test function  $v = (y_h^{j+1} + y_h^j)/2$ .