

## C6.4 Finite Element Methods for PDEs

Patrick E. Farrell

University of Oxford

In this lecture I will

- ▶ explain what the finite element method is for;
- ▶ give a sketch of how it works;
- ▶ outline the questions we will address in the rest of the course.

# What is the finite element method for?

The finite element method is a framework for computing numerical approximations to boundary and initial-boundary value problems.

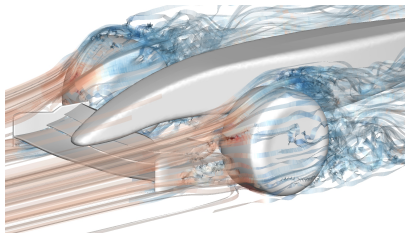
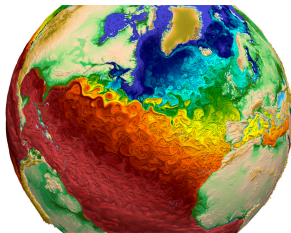
# What is the finite element method for?

The finite element method is a framework for computing numerical approximations to boundary and initial-boundary value problems.

## Example: the Navier–Stokes equations

Find velocity  $u : \Omega \rightarrow \mathbb{R}^3$  and pressure  $p : \Omega \rightarrow \mathbb{R}$  such that

$$-\nabla \cdot \nu \left( \nabla u + (\nabla u)^\top \right) + \nabla \cdot (u \otimes u) + \nabla p = f$$
$$\nabla \cdot u = 0$$





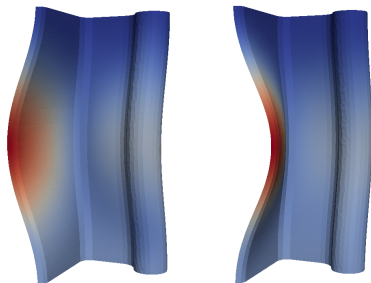
# What is the finite element method for?

The finite element method is a framework for computing numerical approximations to boundary and initial-boundary value problems.

## Example: the equations of elasticity

Find displacement  $u : \Omega \rightarrow \mathbb{R}^3$  such that

$$-\mu \nabla^2 u - (\mu + \gamma) \nabla \nabla \cdot u = f$$



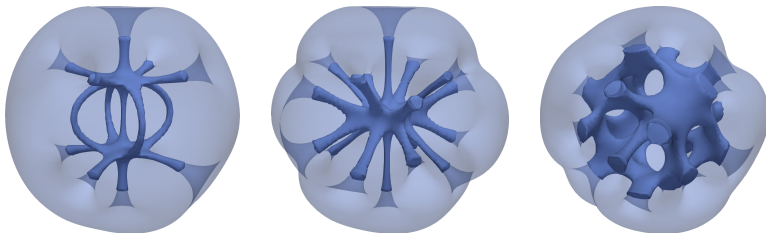
# What is the finite element method for?

The finite element method is a framework for computing numerical approximations to boundary and initial-boundary value problems.

## Example: the nonlinear Schrödinger equation

Find wave function  $\psi : \Omega \times (0, T] \rightarrow \mathbb{C}$  such that

$$i \frac{\partial \psi}{\partial t} = -\frac{1}{2} \nabla^2 \psi + |\psi|^2 \psi + (x^2 + y^2 + z^2) \psi$$

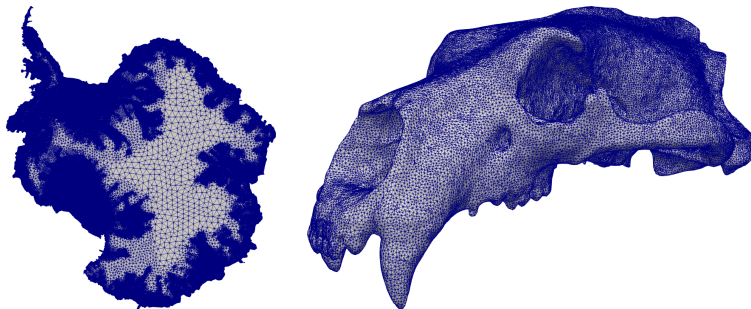


# How does it work?

The finite element method converts PDE problems into algebraic ones.

$$Lu = f \rightsquigarrow Ax = b$$

The discrete system arises by breaking up the domain  $\Omega$  into a mesh:



# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)

# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)

# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)
- ✓ Rooted in the modern variational theory of PDE

# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)
- ✓ Rooted in the modern variational theory of PDE
- ✓ Can design *structure-preserving* discretisations

# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)
- ✓ Rooted in the modern variational theory of PDE
- ✓ Can design *structure-preserving* discretisations
- ✗ Algebraic convergence (vs exponential with spectral methods)



# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)
- ✓ Rooted in the modern variational theory of PDE
- ✓ Can design *structure-preserving* discretisations
- ✗ Algebraic convergence (vs exponential with spectral methods)
- ✗ Can take more work per degree-of-freedom (vs e.g. finite differences)

# Comparing the finite element method

There are many different techniques used to compute numerical approximations of PDE. How does FEM compare?

- ✓ Able to handle complicated geometries (vs e.g. finite differences)
- ✓ Able to handle very general PDE (vs e.g. boundary element methods)
- ✓ Rooted in the modern variational theory of PDE
- ✓ Can design *structure-preserving* discretisations
- ✗ Algebraic convergence (vs exponential with spectral methods)
- ✗ Can take more work per degree-of-freedom (vs e.g. finite differences)
- ✓ The most popular approach used in science and industry

## Section 2

# Basic steps of the finite element method

We illustrate the essential steps of the finite element method with the problem: for given  $f : \Omega \rightarrow \mathbb{R}$ , find  $u : \Omega \rightarrow \mathbb{R}$  such that

$$\begin{aligned} -\nabla^2 u &:= -\nabla \cdot \nabla u = f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega. \end{aligned}$$

For  $\Omega \subset \mathbb{R}^2$ , this is

$$-\left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}\right) = f,$$

while for  $\Omega \subset \mathbb{R}^3$ , this is

$$-\left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2}\right) = f.$$

We illustrate the essential steps of the finite element method with the problem: for given  $f : \Omega \rightarrow \mathbb{R}$ , find  $u : \Omega \rightarrow \mathbb{R}$  such that

$$\begin{aligned} -\nabla^2 u &:= -\nabla \cdot \nabla u = f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega. \end{aligned}$$

The steps are:

- ▶ write as a variational problem over a function space  $V$ ;
- ▶ formulate over a finite-dimensional subspace  $V_h \subset V$ ;
- ▶ construct  $V_h$  + basis with a mesh of  $\Omega$ ;
- ▶ assemble and solve the resulting linear system of equations.

We test the equation with a *test function*  $v$  and integrate:

$$-\int_{\Omega} v \nabla^2 u \, dx = \int_{\Omega} v f \, dx.$$

Just as two surfaces are the same if they look the same from all directions, we demand the LHS and RHS match for “all” test functions (tbd).

We test the equation with a *test function*  $v$  and integrate:

$$-\int_{\Omega} v \nabla^2 u \, dx = \int_{\Omega} v f \, dx.$$

Just as two surfaces are the same if they look the same from all directions, we demand the LHS and RHS match for “all” test functions (tbd).

We can now integrate by parts to shift derivatives onto  $v$ :

$$-\int_{\Omega} v \nabla^2 u \, dx = \int_{\Omega} \nabla v \cdot \nabla u \, dx - \int_{\partial\Omega} v \nabla u \cdot n \, ds = \int_{\Omega} v f \, dx.$$

We test the equation with a *test function*  $v$  and integrate:

$$-\int_{\Omega} v \nabla^2 u \, dx = \int_{\Omega} v f \, dx.$$

Just as two surfaces are the same if they look the same from all directions, we demand the LHS and RHS match for “all” test functions (tbd).

We can now integrate by parts to shift derivatives onto  $v$ :

$$-\int_{\Omega} v \nabla^2 u \, dx = \int_{\Omega} \nabla v \cdot \nabla u \, dx - \int_{\partial\Omega} v \nabla u \cdot n \, ds = \int_{\Omega} v f \, dx.$$

Since we know  $u$  on the boundary, there's no need to test there. Let's fix  $v = 0$  on the boundary:

$$\int_{\Omega} \nabla v \cdot \nabla u \, dx = \int_{\Omega} v f \, dx.$$

This is the *variational* or *weak* formulation of the Poisson equation.



Our problem is to find the *trial function*  $u \in V$  such that

$$\int_{\Omega} \nabla v \cdot \nabla u \, dx = \int_{\Omega} f v \, dx$$

for all *test functions*  $v \in V$ . What function space  $V$  should we look in?

Our problem is to find the *trial function*  $u \in V$  such that

$$\int_{\Omega} \nabla v \cdot \nabla u \, dx = \int_{\Omega} f v \, dx$$

for all *test functions*  $v \in V$ . What function space  $V$  should we look in?

We need

- ▶  $u$  to have all first-order derivatives;
- ▶ the gradient of  $u$  to be square-integrable;
- ▶  $u$  to satisfy  $u = 0$  on the boundary.

The set of such functions is called  $H_0^1(\Omega) =: V$ .

Our problem is to find the *trial function*  $u \in V$  such that

$$\int_{\Omega} \nabla v \cdot \nabla u \, dx = \int_{\Omega} f v \, dx$$

for all *test functions*  $v \in V$ . What function space  $V$  should we look in?

We need

- ▶  $u$  to have all first-order derivatives;
- ▶ the gradient of  $u$  to be square-integrable;
- ▶  $u$  to satisfy  $u = 0$  on the boundary.

The set of such functions is called  $H_0^1(\Omega) =: V$ .

## Key advantage of variational formulation

Classical formulation required  $u \in C^2(\Omega)$ ; we now only require  $u$  to have *first* derivatives.

$V$  is infinite-dimensional; it is too big to search in! So to compute we look instead at *finite-dimensional subspaces*  $V_h \subset V$ . The *Galerkin projection* of our problem onto  $V_h$  is: find  $u_h \in V_h$  such that

$$\int_{\Omega} \nabla v_h \cdot \nabla u_h \, dx = \int_{\Omega} v_h f \, dx$$

for all  $v_h \in V_h$ . Here  $h$  represents the resolution of our discretisation, i.e. the maximal diameter of an element in the mesh.

$V$  is infinite-dimensional; it is too big to search in! So to compute we look instead at *finite-dimensional subspaces*  $V_h \subset V$ . The *Galerkin projection* of our problem onto  $V_h$  is: find  $u_h \in V_h$  such that

$$\int_{\Omega} \nabla v_h \cdot \nabla u_h \, dx = \int_{\Omega} v_h f \, dx$$

for all  $v_h \in V_h$ . Here  $h$  represents the resolution of our discretisation, i.e. the maximal diameter of an element in the mesh.

In this simple case the well-posedness of the discrete problem follows from the continuous one for any  $V_h$ . Some discretisations will be convergent ( $u_h \rightarrow u$  as  $h \rightarrow 0$  in a suitable norm) and some will not. We want discretisations that are well-posed and which converge quickly.

The finite element method is a way to construct a  $V_h$  (+ a basis) with good approximation properties and convenient computational properties.

The finite element method is a way to construct a  $V_h$  (+ a basis) with good approximation properties and convenient computational properties.

### Key idea

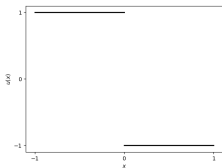
Represent a function with a polynomial on each cell of the mesh.

The finite element method is a way to construct a  $V_h$  (+ a basis) with good approximation properties and convenient computational properties.

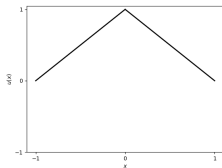
## Key idea

Represent a function with a polynomial on each cell of the mesh.

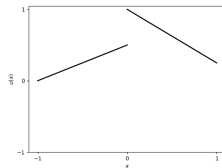
... with a specified degree of continuity.



discont., degree = 0



cont., degree = 1



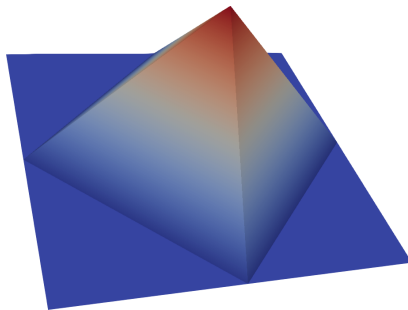
discont., degree = 1



The finite element method is a way to construct a  $V_h$  (+ a basis) with good approximation properties and convenient computational properties.

## Key idea

Represent a function with a polynomial on each cell of the mesh.



Suppose we now have  $V_h = \text{span}\{\phi_1, \dots, \phi_N\}$ . For brevity we write

$$a(u, v) := \int_{\Omega} \nabla v \cdot \nabla u \, dx, \quad F(v) := \int_{\Omega} v f \, dx.$$

Suppose we now have  $V_h = \text{span}\{\phi_1, \dots, \phi_N\}$ . For brevity we write

$$a(u, v) := \int_{\Omega} \nabla v \cdot \nabla u \, dx, \quad F(v) := \int_{\Omega} v f \, dx.$$

Let's expand  $u_h$  and  $v_h$  in terms of our basis. First write

$$v_h = \sum_{i=1}^N V_i \phi_i$$

Suppose we now have  $V_h = \text{span}\{\phi_1, \dots, \phi_N\}$ . For brevity we write

$$a(u, v) := \int_{\Omega} \nabla v \cdot \nabla u \, dx, \quad F(v) := \int_{\Omega} v f \, dx.$$

Let's expand  $u_h$  and  $v_h$  in terms of our basis. First write

$$v_h = \sum_{i=1}^N V_i \phi_i$$

Then our problem becomes:

$$\begin{aligned} a(u_h, v_h) &= F(v_h) \\ \implies a(u_h, \sum_i V_i \phi_i) &= F(\sum_i V_i \phi_i) \\ \implies \sum_i V_i a(u_h, \phi_i) &= \sum_i V_i F(\phi_i). \end{aligned}$$

Suppose we now have  $V_h = \text{span}\{\phi_1, \dots, \phi_N\}$ . For brevity we write

$$a(u, v) := \int_{\Omega} \nabla v \cdot \nabla u \, dx, \quad F(v) := \int_{\Omega} v f \, dx.$$

Let's expand  $u_h$  and  $v_h$  in terms of our basis. First write

$$v_h = \sum_{i=1}^N V_i \phi_i$$

Then our problem becomes:

$$\begin{aligned} a(u_h, v_h) &= F(v_h) \\ \implies a(u_h, \sum_i V_i \phi_i) &= F(\sum_i V_i \phi_i) \\ \implies \sum_i V_i a(u_h, \phi_i) &= \sum_i V_i F(\phi_i). \end{aligned}$$

As this has to hold for all possible values of  $V_i$ , this is equivalent to

$$a(u_h, \phi_i) = F(\phi_i) \text{ for } i = 1, \dots, N.$$

Each test function  $\phi_i$  will yield *one row* of the resulting matrix.

Now expand  $u_h$  as

$$u_h = \sum_{j=1}^N U_j \phi_j$$

Now expand  $u_h$  as

$$u_h = \sum_{j=1}^N U_j \phi_j$$

Substituting, we find

$$\begin{aligned} a\left(\sum_j U_j \phi_j, \phi_i\right) &= F(\phi_i) \\ \implies \sum_j a(\phi_j, \phi_i) U_j &= F(\phi_i) \end{aligned}$$

or in matrix notation

$$AU = b,$$

where

$$A_{ij} = a(\phi_j, \phi_i), \quad b_i = F(\phi_i).$$

Our numerical approximation is computed by solving this linear system.

```
1  from firedrake import *
2
3  mesh = UnitSquareMesh(128, 128, quadrilateral=True)
4  V = FunctionSpace(mesh, "CG", 1)
5  (x, y) = SpatialCoordinate(mesh)
6
7  f = sin(10*pi*x) * sin(5*pi*y)
8  bc = DirichletBC(V, 0, "on_boundary")
9
10 u = Function(V)
11 v = TestFunction(V)
12 G = inner(grad(u), grad(v))*dx - inner(f, v)*dx
13
14 solve(G == 0, u, bc)
15 File("output/poisson.pvd").write(u)
```





## Section 3

### Outlook

This sketch gives the core idea of how the finite element method works. However, each of these steps gets more complicated for real problems:

- ▶ There are different variational formulations for the same PDE, with different advantages and disadvantages.

This sketch gives the core idea of how the finite element method works. However, each of these steps gets more complicated for real problems:

- ▶ There are different variational formulations for the same PDE, with different advantages and disadvantages.
- ▶ For most problems the subspace  $V_h$  must be chosen carefully to achieve a convergent method; any old choice won't work.

This sketch gives the core idea of how the finite element method works. However, each of these steps gets more complicated for real problems:

- ▶ There are different variational formulations for the same PDE, with different advantages and disadvantages.
- ▶ For most problems the subspace  $V_h$  must be chosen carefully to achieve a convergent method; any old choice won't work.
- ▶ Fast solvers for the resulting linear systems must exploit the PDE structure.

The main questions we will address in the remainder of this course are:

- ▶ How do we formulate problems variationally? [Lec 2, 3, 5, 6, 12]

The main questions we will address in the remainder of this course are:

- ▶ How do we formulate problems variationally? [Lec 2, 3, 5, 6, 12]
- ▶ Are the continuous and discrete variational problems well-posed? [Lec 4, 13, 14]

The main questions we will address in the remainder of this course are:

- ▶ How do we formulate problems variationally? [Lec 2, 3, 5, 6, 12]
- ▶ Are the continuous and discrete variational problems well-posed? [Lec 4, 13, 14]
- ▶ What error is incurred in the approximation? [Lec 6, 7, 8, 11, 13]



The main questions we will address in the remainder of this course are:

- ▶ How do we formulate problems variationally? [Lec 2, 3, 5, 6, 12]
- ▶ Are the continuous and discrete variational problems well-posed? [Lec 4, 13, 14]
- ▶ What error is incurred in the approximation? [Lec 6, 7, 8, 11, 13]
- ▶ How do we implement the finite element method? [Lec 8, 9, 10]

## C6.4 Finite Element Methods for PDEs

### Lecture 2: Lebesgue spaces

Patrick E. Farrell

University of Oxford

## Lecture 2: Lebesgue spaces

Our goal for the next three lectures is to prove the *Lax–Milgram* theorem about the well-posedness of the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

where  $a : V \times V \rightarrow \mathbb{R}$  is bilinear and  $F : V \rightarrow \mathbb{R}$  is linear.

## Lecture 2: Lebesgue spaces

Our goal for the next three lectures is to prove the *Lax–Milgram* theorem about the well-posedness of the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

where  $a : V \times V \rightarrow \mathbb{R}$  is bilinear and  $F : V \rightarrow \mathbb{R}$  is linear.

When it applies, the Lax–Milgram theorem will prove the well-posedness of both the continuous problem and its Galerkin approximation.

## Lecture 2: Lebesgue spaces

Our goal for the next three lectures is to prove the *Lax–Milgram* theorem about the well-posedness of the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

where  $a : V \times V \rightarrow \mathbb{R}$  is bilinear and  $F : V \rightarrow \mathbb{R}$  is linear.

When it applies, the Lax–Milgram theorem will prove the well-posedness of both the continuous problem and its Galerkin approximation.

Moreover, we can bound the error in Galerkin approximation in terms of the constants arising in the statement of Lax–Milgram.

## Lecture 2: Lebesgue spaces

Our goal for the next three lectures is to prove the *Lax–Milgram* theorem about the well-posedness of the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

where  $a : V \times V \rightarrow \mathbb{R}$  is bilinear and  $F : V \rightarrow \mathbb{R}$  is linear.

When it applies, the Lax–Milgram theorem will prove the well-posedness of both the continuous problem and its Galerkin approximation.

Moreover, we can bound the error in Galerkin approximation in terms of the constants arising in the statement of Lax–Milgram.

Before getting to Lax–Milgram, we must first understand the *Lebesgue* and *Sobolev* spaces  $V$  in which we look for solutions.

## Definition (normed vector space)

A *normed vector space*  $X$  is a vector space equipped with a norm  $\|\cdot\| : X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $\|x\| \geq 0$ , and  $\|x\| = 0 \iff x = 0$ ;
- ▶  $\|\alpha x\| = |\alpha| \|x\|$  for any scalar  $\alpha \in \mathbb{R}$ ;
- ▶  $\|x + y\| \leq \|x\| + \|y\|$ .

## Definition (normed vector space)

A *normed vector space*  $X$  is a vector space equipped with a norm  $\| \cdot \| : X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $\|x\| \geq 0$ , and  $\|x\| = 0 \iff x = 0$ ;
- ▶  $\|\alpha x\| = |\alpha| \|x\|$  for any scalar  $\alpha \in \mathbb{R}$ ;
- ▶  $\|x + y\| \leq \|x\| + \|y\|$ .

## Definition (Banach space)

A Banach space is a complete normed vector space.



## Definition (normed vector space)

A *normed vector space*  $X$  is a vector space equipped with a norm  $\|\cdot\| : X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $\|x\| \geq 0$ , and  $\|x\| = 0 \iff x = 0$ ;
- ▶  $\|\alpha x\| = |\alpha| \|x\|$  for any scalar  $\alpha \in \mathbb{R}$ ;
- ▶  $\|x + y\| \leq \|x\| + \|y\|$ .

## Definition (Banach space)

A Banach space is a complete normed vector space.

Recall that *completeness* of a normed vector space  $X$  means that all Cauchy sequences in  $X$  converge in  $X$ . A Cauchy sequence  $(x_n)$  is one where  $\forall \varepsilon > 0 \exists N > 0 \forall m, n > N \|x_n - x_m\| < \varepsilon$ .

## Example

Euclidean space  $\mathbb{R}^n$  equipped with the 1-norm, the 2-norm, and the supremum norm

$$\|x\|_1 := \sum_{i=1}^n |x_i|, \quad \|x\|_2 := \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2}, \quad \|x\|_\infty := \max_i |x_i|$$

are all Banach spaces.

## Example

Euclidean space  $\mathbb{R}^n$  equipped with the 1-norm, the 2-norm, and the supremum norm

$$\|x\|_1 := \sum_{i=1}^n |x_i|, \quad \|x\|_2 := \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2}, \quad \|x\|_\infty := \max_i |x_i|$$

are all Banach spaces.

## Example

The space of continuous functions from a domain  $\Omega$  to  $\mathbb{R}$  equipped with the supremum norm

$$\|f\|_\infty = \sup\{|f(x)| : x \in \Omega\}$$

is a Banach space.

## Definition (inner product space)

An *inner product space*  $X$  is a vector space equipped with an inner product  $(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $(u, v) = (v, u)$ ;
- ▶  $(\alpha u + \beta v, w) = \alpha(u, w) + \beta(v, w)$  for  $\alpha, \beta \in \mathbb{R}$ ;
- ▶  $(u, u) \geq 0$  with  $(u, u) = 0 \iff u = 0$ .

## Definition (inner product space)

An *inner product space*  $X$  is a vector space equipped with an inner product  $(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $(u, v) = (v, u)$ ;
- ▶  $(\alpha u + \beta v, w) = \alpha(u, w) + \beta(v, w)$  for  $\alpha, \beta \in \mathbb{R}$ ;
- ▶  $(u, u) \geq 0$  with  $(u, u) = 0 \iff u = 0$ .

An inner product induces a norm  $\|u\| = \sqrt{(u, u)}$ .

## Definition (inner product space)

An *inner product space*  $X$  is a vector space equipped with an inner product  $(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$  that satisfies the following properties:

- ▶  $(u, v) = (v, u)$ ;
- ▶  $(\alpha u + \beta v, w) = \alpha(u, w) + \beta(v, w)$  for  $\alpha, \beta \in \mathbb{R}$ ;
- ▶  $(u, u) \geq 0$  with  $(u, u) = 0 \iff u = 0$ .

An inner product induces a norm  $\|u\| = \sqrt{(u, u)}$ .

## Definition (Hilbert space)

A *Hilbert space* is a complete inner product space.

## Example

The canonical example of a Hilbert space is  $\mathbb{R}^n$  with inner product

$$(u, v)_{\mathbb{R}^n} = u^\top v.$$

## Example

The space of square-integrable functions on a domain  $L^2(\Omega)$  with inner product

$$(u, v)_{L^2(\Omega)} = \int_{\Omega} uv \, dx.$$

## Example

The space  $H_0^1(\Omega)$  of square-integrable functions that are zero on the boundary and that have square-integrable derivatives is a Hilbert space with inner product

$$(u, v)_{H_0^1(\Omega)} = \int_{\Omega} \nabla u \cdot \nabla v \, dx.$$

If  $(u, u)_{H_0^1(\Omega)} = 0$  then  $u$  must be constant; but the only constant function in  $H_0^1(\Omega)$  is the zero function, because of the boundary conditions.



## Example

The space  $H^1(\Omega) = H(\text{grad}, \Omega)$  of square-integrable functions that have square-integrable gradient is a Hilbert space with inner product

$$(u, v)_{H^1(\Omega)} = \int_{\Omega} uv + \nabla u \cdot \nabla v \, dx.$$

## Example

The space  $H^1(\Omega) = H(\text{grad}, \Omega)$  of square-integrable functions that have square-integrable gradient is a Hilbert space with inner product

$$(u, v)_{H^1(\Omega)} = \int_{\Omega} uv + \nabla u \cdot \nabla v \, dx.$$

## Example

The space  $H(\text{div}, \Omega)$  of square-integrable vector-valued functions that have square-integrable divergence is a Hilbert space with inner product

$$(u, v)_{H(\text{div}, \Omega)} = \int_{\Omega} u \cdot v + \nabla \cdot u \nabla \cdot v \, dx.$$

## Example

The space  $H^1(\Omega) = H(\text{grad}, \Omega)$  of square-integrable functions that have square-integrable gradient is a Hilbert space with inner product

$$(u, v)_{H^1(\Omega)} = \int_{\Omega} uv + \nabla u \cdot \nabla v \, dx.$$

## Example

The space  $H(\text{div}, \Omega)$  of square-integrable vector-valued functions that have square-integrable divergence is a Hilbert space with inner product

$$(u, v)_{H(\text{div}, \Omega)} = \int_{\Omega} u \cdot v + \nabla \cdot u \nabla \cdot v \, dx.$$

## Example

For  $\Omega \subset \mathbb{R}^3$ , the space  $H(\text{curl}, \Omega)$  of square-integrable vector-valued functions that have square-integrable curl is a Hilbert space with inner product

$$(u, v)_{H(\text{curl}, \Omega)} = \int_{\Omega} u \cdot v + \nabla \times u \cdot \nabla \times v \, dx.$$

## Theorem (Cauchy–Schwarz inequality)

For a Hilbert space  $X$  and any  $u, v \in X$ ,

$$|(u, v)_X| \leq \|u\|_X \|v\|_X.$$

### Proof.

Let  $\lambda \in \mathbb{R}$ . Then

$$\begin{aligned} 0 \leq \|u + \lambda v\|_X^2 &= (u + \lambda v, u + \lambda v)_X \\ &= (u, u) + (u, \lambda v) + (\lambda v, u) + (\lambda v, \lambda v) \\ &= \|u\|_X^2 + 2\lambda(u, v) + \lambda^2\|v\|_X^2. \end{aligned}$$

The right-hand side is a quadratic polynomial in  $\lambda$  with real coefficients, and it is non-negative for all  $\lambda \in \mathbb{R}$ . Therefore its discriminant is non-positive; it can only be zero or negative. Thus,

$$|2(u, v)_X|^2 - 4\|u\|_X^2\|v\|_X^2 \leq 0.$$

## Section 2

# Dual of a Hilbert space

## Definition (Linear functional on a Hilbert space)

Given a Hilbert space  $X$ , a *linear functional*  $j$  on  $X$  is a function  $j : X \rightarrow \mathbb{R}$  that satisfies

$$j(\alpha u + \beta v) = \alpha j(u) + \beta j(v).$$

## Definition (Linear functional on a Hilbert space)

Given a Hilbert space  $X$ , a *linear functional*  $j$  on  $X$  is a function  $j : X \rightarrow \mathbb{R}$  that satisfies

$$j(\alpha u + \beta v) = \alpha j(u) + \beta j(v).$$

## Example

Integration over a fixed domain  $\Omega$ , evaluation at a fixed point  $x$ , and evaluation of the derivative at a point  $x$  in a fixed direction  $v$  are all examples of linear functionals (when they are defined!).

## Definition (Linear functional on a Hilbert space)

Given a Hilbert space  $X$ , a *linear functional*  $j$  on  $X$  is a function  $j : X \rightarrow \mathbb{R}$  that satisfies

$$j(\alpha u + \beta v) = \alpha j(u) + \beta j(v).$$

## Example

Integration over a fixed domain  $\Omega$ , evaluation at a fixed point  $x$ , and evaluation of the derivative at a point  $x$  in a fixed direction  $v$  are all examples of linear functionals (when they are defined!).

## Example

Drag over a wing, compliance of a structure, average global temperature.



## Definition (Bounded linear functional)

A *bounded* linear functional  $j : X \rightarrow \mathbb{R}$  is one for which there exists  $L \in [0, \infty)$  such that

$$|j(v)| \leq L\|v\|_X \quad \forall v \in X.$$

## Definition (Bounded linear functional)

A *bounded* linear functional  $j : X \rightarrow \mathbb{R}$  is one for which there exists  $L \in [0, \infty)$  such that

$$|j(v)| \leq L\|v\|_X \quad \forall v \in X.$$

## Lemma (Boundedness and continuity)

*Boundedness is equivalent to continuity.*

Proof.

See notes, Lemma 2.3.4. □

## Example

Given any  $g \in X$ , we can construct

$$j(v) = (g, v)_X.$$

This is bounded, since by Cauchy–Schwarz

$$|j(v)| = |(g, v)_X| \leq \|g\|_X \|v\|_X.$$

## Definition (Dual of a Hilbert space)

The *dual*  $X^*$  of a Hilbert space  $X$  is the space of all bounded linear functionals on  $X$ . This has a natural norm induced by the norm on the underlying space:

$$\|j\|_{X^*} := \sup_{\|v\|_X=1} |j(v)|.$$

This gives the “tightest  $L$ ” in the definition of boundedness.

## Definition (Dual of a Hilbert space)

The *dual*  $X^*$  of a Hilbert space  $X$  is the space of all bounded linear functionals on  $X$ . This has a natural norm induced by the norm on the underlying space:

$$\|j\|_{X^*} := \sup_{\|v\|_X=1} |j(v)|.$$

This gives the “tightest  $L$ ” in the definition of boundedness.

Given a  $j \in X^*$ , denote the action of  $j$  on  $v$  by

$$\langle j, v \rangle := j(v).$$

This is called the *duality pairing*.

## Theorem (Riesz Representation Theorem)

*Any bounded linear functional  $j \in X^*$  can be uniquely represented by a  $g \in X$ , via*

$$\langle j, v \rangle = (g, v) \quad \text{for all } v \in X.$$

*Moreover, the norms agree:  $\|j\|_{X^*} = \|g\|_X$ .*

## Theorem (Riesz Representation Theorem)

*Any bounded linear functional  $j \in X^*$  can be uniquely represented by a  $g \in X$ , via*

$$\langle j, v \rangle = (g, v) \quad \text{for all } v \in X.$$

*Moreover, the norms agree:  $\|j\|_{X^*} = \|g\|_X$ .*

This defines a canonical linear map, the *Riesz map*  $\mathcal{R} : X^* \rightarrow X$ , that maps  $j \mapsto g$ . This Riesz map is an isometric isomorphism.

## Theorem (Riesz Representation Theorem)

*Any bounded linear functional  $j \in X^*$  can be uniquely represented by a  $g \in X$ , via*

$$\langle j, v \rangle = (g, v) \quad \text{for all } v \in X.$$

*Moreover, the norms agree:  $\|j\|_{X^*} = \|g\|_X$ .*

This defines a canonical linear map, the *Riesz map*  $\mathcal{R} : X^* \rightarrow X$ , that maps  $j \mapsto g$ . This Riesz map is an isometric isomorphism.

## Example

Let  $X = L^2(\Omega)$  and let

$$j(v) = \langle j, v \rangle = \int_{\Omega} v \, dx.$$

Then its  $L^2(\Omega)$  Riesz representation is the constant function  $g(x) = 1$ .



## Section 3

# Lebesgue spaces

### Definition (Lebesgue $p$ -norm, $p \in [1, \infty)$ )

Let  $p \in [1, \infty)$ . The  $L^p(\Omega)$  norm is defined by

$$\|u\|_{L^p(\Omega)} = \left( \int_{\Omega} |u|^p \, dx \right)^{1/p}.$$

### Definition (Lebesgue $p$ -norm, $p = \infty$ )

The  $L^\infty(\Omega)$  norm is defined by

$$\|u\|_{L^\infty(\Omega)} = \inf\{C \geq 0 : |u(x)| \leq C \text{ almost everywhere}\}.$$

This is the *essential supremum* of  $|u|$ .

## Definition (Lebesgue space)

For  $p \in [1, \infty]$ , consider the definition

$$L^p(\Omega) = \{u : \Omega \rightarrow \mathbb{R} : \|u\|_{L^p(\Omega)} < \infty\}.$$

These are Banach spaces for all  $p$ , a Hilbert space for  $p = 2$ .

## Definition (Lebesgue space)

For  $p \in [1, \infty]$ , consider the definition

$$L^p(\Omega) = \{u : \Omega \rightarrow \mathbb{R} : \|u\|_{L^p(\Omega)} < \infty\}.$$

These are Banach spaces for all  $p$ , a Hilbert space for  $p = 2$ .

## Important remark

The Lebesgue integral ignores differences on a set of measure zero. Two functions  $f$  and  $g$  that differ only on a set of measure zero will have  $\|f - g\|_{L^p} = 0$ . In order to fix this, we actually take elements of  $L^p$  to be *equivalence classes* of functions that differ up to sets of measure zero.

## Definition (Lebesgue space)

For  $p \in [1, \infty]$ , consider the definition

$$L^p(\Omega) = \{u : \Omega \rightarrow \mathbb{R} : \|u\|_{L^p(\Omega)} < \infty\}.$$

These are Banach spaces for all  $p$ , a Hilbert space for  $p = 2$ .

## Important remark

The Lebesgue integral ignores differences on a set of measure zero. Two functions  $f$  and  $g$  that differ only on a set of measure zero will have  $\|f - g\|_{L^p} = 0$ . In order to fix this, we actually take elements of  $L^p$  to be *equivalence classes* of functions that differ up to sets of measure zero.

## Consequence

For  $f \in L^p(\Omega)$ , to evaluate  $f(x)$ , we have to prove that there is a *continuous* function in the equivalence class and evaluate *that*.

Fix  $\Omega \subset \mathbb{R}^n$  to have finite measure.

## Example

The function  $f(x) = 1$  is in  $L^p(\Omega)$  for all  $p$ :

$$\begin{aligned}\|1\|_{L^p(\Omega)} &= \left( \int_{\Omega} 1^p \, dx \right)^{1/p} \\ &= |\Omega|^{1/p} < \infty.\end{aligned}$$

Fix  $\Omega \subset \mathbb{R}^n$  to have finite measure.

## Example

The function  $f(x) = 1$  is in  $L^p(\Omega)$  for all  $p$ :

$$\begin{aligned}\|1\|_{L^p(\Omega)} &= \left( \int_{\Omega} 1^p \, dx \right)^{1/p} \\ &= |\Omega|^{1/p} < \infty.\end{aligned}$$

## Example

Let  $\Omega = (0, 1)$  and let

$$f_q(x) = x^{-q}.$$

Then  $f_q \in L^p(\Omega) \iff q < 1/p$ . That is,  $\frac{1}{x} \notin L^1(\Omega)$ , but  $\frac{1}{x^{0.999}} \in L^1(\Omega)$ , and  $\frac{1}{\sqrt{x}} \notin L^2(\Omega)$ , but  $\frac{1}{x^{0.4999}} \in L^2(\Omega)$ , etc.

In other words, the larger the  $p$ , the slower the allowed rate of blow-up at singularities. For  $L^\infty(\Omega)$  no blow-up whatsoever is allowed.

## Theorem (Hölder's inequality)

Let  $p, q \in [1, \infty]$  such that

$$\frac{1}{p} + \frac{1}{q} = 1.$$

*The elements of such a pair are called Hölder conjugates. By convention here, 1 and  $\infty$  are conjugate.*

*If  $f \in L^p(\Omega)$  and  $g \in L^q(\Omega)$ , then  $fg \in L^1(\Omega)$  and*

$$\|fg\|_{L^1(\Omega)} \leq \|f\|_{L^p(\Omega)} \|g\|_{L^q(\Omega)}.$$



## Theorem (Hölder's inequality)

Let  $p, q \in [1, \infty]$  such that

$$\frac{1}{p} + \frac{1}{q} = 1.$$

The elements of such a pair are called Hölder conjugates. By convention here, 1 and  $\infty$  are conjugate.

If  $f \in L^p(\Omega)$  and  $g \in L^q(\Omega)$ , then  $fg \in L^1(\Omega)$  and

$$\|fg\|_{L^1(\Omega)} \leq \|f\|_{L^p(\Omega)} \|g\|_{L^q(\Omega)}.$$

## Theorem (Inclusion of Lebesgue spaces)

Let  $\Omega$  be bounded. Let  $1 \leq p < q \leq \infty$ . If  $f \in L^q(\Omega)$ , then  $f \in L^p(\Omega)$ .

Proof.

See notes, Theorem 2.5.10. □

## C6.4 Finite Element Methods for PDEs

### Lecture 3: Sobolev spaces

Patrick E. Farrell

University of Oxford

## Lecture 3: Sobolev spaces

In the previous lecture, we studied Lebesgue spaces, which capture how *integrable* a function is.

## Lecture 3: Sobolev spaces

In the previous lecture, we studied Lebesgue spaces, which capture how *integrable* a function is.

In this lecture, we now study *Sobolev* spaces, which also capture how differentiable a function is. But first, we must generalise our notion of taking a derivative.

## Definition (Classical/strong differentiability in one dimension)

$$f'(x) := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

## Definition (Classical/strong differentiability in one dimension)

$$f'(x) := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

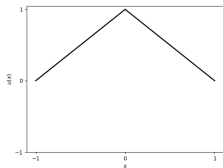
## Remark

This is based on pointwise evaluation. But we've seen that pointwise evaluation isn't a native concept to Lebesgue functions. We can't always do it!

Recall that the variational formulation of the Poisson equation requires evaluating terms like

$$\int_{\Omega} \nabla u_h \cdot \nabla v_h \, dx.$$

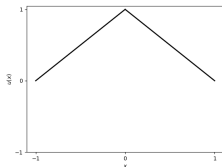
How do we do this if  $u_h$  is a  $C^0$  piecewise polynomial?



Recall that the variational formulation of the Poisson equation requires evaluating terms like

$$\int_{\Omega} \nabla u_h \cdot \nabla v_h \, dx.$$

How do we do this if  $u_h$  is a  $C^0$  piecewise polynomial?



## Answer

We will develop a sense of differentiation, called the weak derivative, built on integration by parts.



To motivate the definition, first suppose  $f \in C^1(a, b)$ . Let  $\phi$  be a differentiable function that is zero on the boundary  $\{a, b\}$ . Then integration by parts tells us that

$$\int_a^b f' \phi \, dx = - \int_a^b f \phi' \, dx,$$

i.e. we can swap the differentiation operator onto the test function  $\phi$ . This is how we will *define* the weak derivative  $f'$  in Lebesgue spaces.

## Definition (Compact support in $\Omega$ )

A function  $\phi \in C(\Omega)$  has *compact support* if

$$\text{supp}(\phi) = \text{closure}\{x \in \Omega : \phi(x) \neq 0\}$$

is compact (i.e. is bounded) and is a subset of the *interior* of  $\Omega$ . In particular, this means that  $\phi$  vanishes on  $\partial\Omega$  (and in a neighbourhood of it).

## Definition (Compact support in $\Omega$ )

A function  $\phi \in C(\Omega)$  has *compact support* if

$$\text{supp}(\phi) = \text{closure}\{x \in \Omega : \phi(x) \neq 0\}$$

is compact (i.e. is bounded) and is a subset of the *interior* of  $\Omega$ . In particular, this means that  $\phi$  vanishes on  $\partial\Omega$  (and in a neighbourhood of it).

## Definition (Bump functions)

The set of *bump functions*  $C_0^\infty(\Omega)$  is the set of  $C^\infty(\Omega)$  functions that have compact support in  $\Omega$ .

## Definition (Compact support in $\Omega$ )

A function  $\phi \in C(\Omega)$  has *compact support* if

$$\text{supp}(\phi) = \text{closure}\{x \in \Omega : \phi(x) \neq 0\}$$

is compact (i.e. is bounded) and is a subset of the *interior* of  $\Omega$ . In particular, this means that  $\phi$  vanishes on  $\partial\Omega$  (and in a neighbourhood of it).

## Definition (Bump functions)

The set of *bump functions*  $C_0^\infty(\Omega)$  is the set of  $C^\infty(\Omega)$  functions that have compact support in  $\Omega$ .

## Example

$$\Psi(x) = \begin{cases} \exp\left(-\frac{1}{1-x^2}\right) & \text{if } |x| < 1 \\ 0 & \text{otherwise} \end{cases}$$

What set of functions might have weak derivatives?

### Definition (Locally integrable functions)

Given a domain  $\Omega$ , the set of *locally integrable functions* is defined by

$$L^1_{\text{loc}}(\Omega) = \{f : \Omega \rightarrow \mathbb{R}, f|_K \in L^1(K) \text{ for all compact } K \subset \text{interior } \Omega\}.$$

This set includes  $L^1(\Omega)$  and  $C^0(\Omega)$  as subsets.

## Definition (Weak first derivative)

A function  $f \in L^1_{\text{loc}}(\Omega)$  has a weak  $i^{\text{th}}$  partial derivative  $\partial f / \partial x_i$  if there exists a function  $g \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} g \phi \, dx = - \int_{\Omega} f \frac{\partial \phi}{\partial x_i} \, dx \quad \text{for all } \phi \in C_0^\infty(\Omega).$$

## Definition (Weak first derivative)

A function  $f \in L^1_{\text{loc}}(\Omega)$  has a weak  $i^{\text{th}}$  partial derivative  $\partial f / \partial x_i$  if there exists a function  $g \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} g \phi \, dx = - \int_{\Omega} f \frac{\partial \phi}{\partial x_i} \, dx \quad \text{for all } \phi \in C_0^\infty(\Omega).$$

With this, we can define the weak gradient, curl, and divergence in the obvious way (collect the relevant weak partial derivatives).

## Definition (Weak first derivative)

A function  $f \in L^1_{\text{loc}}(\Omega)$  has a weak  $i^{\text{th}}$  partial derivative  $\partial f / \partial x_i$  if there exists a function  $g \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} g \phi \, dx = - \int_{\Omega} f \frac{\partial \phi}{\partial x_i} \, dx \quad \text{for all } \phi \in C_0^\infty(\Omega).$$

With this, we can define the weak gradient, curl, and divergence in the obvious way (collect the relevant weak partial derivatives).

## Theorem

*Weak derivatives are unique, up to a set of measure zero.*



## Definition (Weak first derivative)

A function  $f \in L^1_{\text{loc}}(\Omega)$  has a weak  $i^{\text{th}}$  partial derivative  $\partial f / \partial x_i$  if there exists a function  $g \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} g \phi \, dx = - \int_{\Omega} f \frac{\partial \phi}{\partial x_i} \, dx \quad \text{for all } \phi \in C_0^\infty(\Omega).$$

With this, we can define the weak gradient, curl, and divergence in the obvious way (collect the relevant weak partial derivatives).

## Theorem

*Weak derivatives are unique, up to a set of measure zero.*

From now on, whenever I take a derivative, I mean a weak derivative!

## Example

Let  $\Omega = (-1, 1)$  and take  $f(x) = |x|$ . Then it has a weak derivative  $f'$  given by

$$f' = \begin{cases} -1 & x < 0 \\ 1 & x > 0. \end{cases}$$

## Example

Let  $\Omega = (-1, 1)$  and take  $f(x) = |x|$ . Then it has a weak derivative  $f'$  given by

$$f' = \begin{cases} -1 & x < 0 \\ 1 & x > 0. \end{cases}$$

To verify this, let  $\phi \in C_0^\infty(\Omega)$ . Then

$$\begin{aligned} \int_{-1}^1 f(x)\phi'(x) \, dx &= \int_{-1}^0 f(x)\phi'(x) \, dx + \int_0^1 f(x)\phi'(x) \, dx \\ &= - \int_{-1}^0 (-1)\phi(x) \, dx + [f\phi]_{-1}^0 - \int_0^1 (+1)\phi(x) \, dx + [f\phi]_0^1 \\ &= - \int_{-1}^1 f'(x)\phi(x) \, dx + ((f\phi)(0^-) - (f\phi)(0^+)) \\ &= - \int_{-1}^1 f'(x)\phi(x) \, dx. \end{aligned}$$

## Example

Let  $\Omega = (-1, 1)$  and take  $f(x) = |x|$ . Then it has a weak derivative  $f'$  given by

$$f' = \begin{cases} -1 & x < 0 \\ 1 & x > 0. \end{cases}$$

## More generally ...

Any continuous piecewise-differentiable function is weakly differentiable, because the boundary terms arising in integration by parts will cancel.

This is important because this is what we will use to approximate the solutions of PDEs!

## Example

A counterexample: take  $\Omega = (-1, 1)$  and take  $f(x) = \text{sign}(x)$ , i.e.

$$f(x) = \begin{cases} -1 & x < 0 \\ 0 & x = 0 \\ 1 & x > 0. \end{cases}$$

This function has no weak derivative.

An informal proof: the only candidate  $f'$  would be  $f' \equiv 0$ , but the discontinuity at  $x = 0$  means that the extra terms arising from integration by parts do not vanish.

To compactly define higher derivatives, we first need to introduce multi-index notation.

### Definition (multi-index notation)

Let  $\Omega \subset \mathbb{R}^n$ . A multi-index  $\alpha$  is a tuple of  $n$  non-negative integers

$$\alpha = (\alpha_1, \dots, \alpha_n), \quad \alpha_i \in \mathbb{N}^+.$$

Given a multi-index  $\alpha$  and  $\phi \in C^\infty(\Omega)$ , define

$$\partial_x^\alpha \phi = \phi^{(\alpha)} = D^\alpha \phi = \left( \frac{\partial}{\partial x_1} \right)^{\alpha_1} \cdots \left( \frac{\partial}{\partial x_n} \right)^{\alpha_n} \phi.$$

The *length* of  $\alpha$  is the order of the derivative,

$$|\alpha| = \sum_{i=1}^n \alpha_i.$$

## Example

The multi-index  $(1, 0)$  corresponds to  $\partial/\partial x_1$ . The multi-index  $(0, 1)$  corresponds to  $\partial/\partial x_2$ . A sum over  $|\alpha| = 1$  means to sum over all first order derivatives.

## Definition (Weak derivative)

Let  $\Omega \subset \mathbb{R}^n$ . We say that a given function  $f \in L^1_{\text{loc}}(\Omega)$  has a weak derivative  $D^\alpha f$  provided there exists a function  $g \in L^1_{\text{loc}}(\Omega)$  such that

$$\int_{\Omega} g \phi \, dx = (-1)^{|\alpha|} \int_{\Omega} f \phi^{(\alpha)} \, dx \quad \text{for all } \phi \in C_0^\infty(\Omega).$$



## Section 2

# Sobolev spaces

## Definition (Sobolev norm)

Let  $f \in L^1_{\text{loc}}(\Omega)$ . Let  $k$  be a non-negative integer.

## Definition (Sobolev norm)

Let  $f \in L^1_{\text{loc}}(\Omega)$ . Let  $k$  be a non-negative integer. Suppose that the weak derivatives  $D^\alpha f$  exist for all  $|\alpha| \leq k$ .

## Definition (Sobolev norm)

Let  $f \in L^1_{\text{loc}}(\Omega)$ . Let  $k$  be a non-negative integer. Suppose that the weak derivatives  $D^\alpha f$  exist for all  $|\alpha| \leq k$ . For  $p \in [1, \infty)$ , define the *Sobolev norm*

$$\|f\|_{W^k_p(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha f\|_{L^p(\Omega)}^p \right)^{1/p}.$$

## Definition (Sobolev norm)

Let  $f \in L^1_{\text{loc}}(\Omega)$ . Let  $k$  be a non-negative integer. Suppose that the weak derivatives  $D^\alpha f$  exist for all  $|\alpha| \leq k$ . For  $p \in [1, \infty)$ , define the *Sobolev norm*

$$\|f\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha f\|_{L^p(\Omega)}^p \right)^{1/p}.$$

In the case  $p = \infty$

$$\|f\|_{W_p^k(\Omega)} = \max_{|\alpha| \leq k} \|D^\alpha f\|_{L^\infty(\Omega)}.$$

## Definition (Sobolev norm)

Let  $f \in L^1_{\text{loc}}(\Omega)$ . Let  $k$  be a non-negative integer. Suppose that the weak derivatives  $D^\alpha f$  exist for all  $|\alpha| \leq k$ . For  $p \in [1, \infty)$ , define the *Sobolev norm*

$$\|f\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha f\|_{L^p(\Omega)}^p \right)^{1/p}.$$

In the case  $p = \infty$

$$\|f\|_{W_p^k(\Omega)} = \max_{|\alpha| \leq k} \|D^\alpha f\|_{L^\infty(\Omega)}.$$

## Definition (Sobolev space)

Define the *Sobolev space*  $W_p^k(\Omega)$  as

$$W_p^k(\Omega) = \{f \in L^1_{\text{loc}}(\Omega) : \|f\|_{W_p^k(\Omega)} < \infty\}.$$

## Theorem

*The Sobolev space  $W_p^k(\Omega)$  is a Banach space.*

## Proof.

See theorem 1.3.2 of Brenner & Scott. □

## Theorem

*The Sobolev spaces with  $p = 2$  are Hilbert spaces. These are denoted by*

$$H^k(\Omega) = W_2^k(\Omega).$$

## Example

The space  $W_p^0(\Omega) = L^p(\Omega)$ . That is, if we ask for no weak derivatives, we just get the  $L^p(\Omega)$  space back.

## Example

Suppose  $l \geq k$ . Then  $W_p^l(\Omega) \subset W_p^k(\Omega)$ ; we're just asking for fewer derivatives.

## Example

Suppose  $1 \leq p \leq q \leq \infty$  and that  $\Omega$  is bounded. Then  $W_q^k(\Omega) \subset W_p^k(\Omega)$ .



There are other inclusions between Sobolev spaces that are less obvious. These will be encoded in *Sobolev's inequality*. However, in order for the result to be true, we will need an additional regularity requirement on the domain  $\Omega$ .

### Definition (Lipschitz domain, informal)

We say  $\Omega$  is a *Lipschitz domain*, or has Lipschitz boundary, if  $\partial\Omega$  is everywhere locally the graph of a Lipschitz continuous function.

This regularity condition is important: without it, the Sobolev inequality is *not true*. Henceforth, we assume that  $\Omega$  is a Lipschitz domain.

There are three numbers describing a Sobolev space:

- ▶  $n$ , the dimension of the domain;
- ▶  $k$ , the number of weak derivatives possessed;
- ▶  $p$ , the integrability of the function and derivatives.

Sobolev's inequality tells us that if you possess enough weak derivatives that are integrable enough, then your function is continuous and bounded.

### Theorem (Sobolev's inequality)

*Let  $p \in [1, \infty)$ . Suppose*

$$\begin{aligned} k &\geq n \text{ when } p = 1 \\ k &> n/p \text{ when } p > 1. \end{aligned}$$

*Then there is a constant  $C$  such that for all  $u \in W_p^k(\Omega)$ ,*

$$\|u\|_{L^\infty(\Omega)} \leq C \|u\|_{W_p^k(\Omega)},$$

*and moreover there is a continuous function in the equivalence class of  $u$ .*

### Example

For  $n = 1$ , the existence of *a single* weak derivative of any integrability is enough to ensure continuity.

### Example

For  $n = 2$ , we have  $W_1^1(\Omega) \not\subset C(\Omega)$ , but  $W_1^2(\Omega) \subset C(\Omega)$ .

### Example

For  $n = 3$ , we have  $W_1^2(\Omega) \not\subset C(\Omega)$ , but  $W_1^3(\Omega) \subset C(\Omega)$ .

## Example

Let's look at the continuity properties of the Hilbert spaces  $H^k(\Omega)$ , i.e.  $p = 2$ . With  $p = 2$ , Sobolev's inequality tells us that we need

$$k > n/2.$$

In one dimension,

$$H^1(\Omega) \subset C(\Omega).$$

For  $n = 2$ , Sobolev's inequality tells us we need  $k > 1$ , i.e.  $k \geq 2$ , so in two dimensions

$$H^1(\Omega) \not\subset C(\Omega), \quad H^2(\Omega) \subset C(\Omega).$$

For  $n = 3$ , Sobolev's inequality tells us we need  $k > 1.5$ , so  $k \geq 2$  is again sufficient.

We are now in a position to see why the space

$$H_0^1(\Omega) = \{u \in H^1(\Omega) : u|_{\partial\Omega} = 0\}$$

is the “right” one for the variational formulation of the Poisson equation: find  $u \in H_0^1(\Omega)$  such that

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx$$

for all  $v \in H_0^1(\Omega)$ .

- ▶ We want  $v \in L^2(\Omega)$  and  $f \in L^2(\Omega)$  to ensure that the RHS is a bounded linear functional of  $v$ .
- ▶ We want  $u|_{\partial\Omega} = 0$  to satisfy the strongly-imposed boundary conditions.
- ▶ We need the first weak derivatives to exist to talk about  $\nabla u$  and  $\nabla v$ .
- ▶ We want  $u$  and  $v$  to have square-integrable weak derivatives, as this guarantees  $a(u, v) < \infty$  (by Cauchy-Schwarz).

## C6.4 Finite Element Methods for PDEs

### Lecture 4: The Lax–Milgram Theorem

Patrick E. Farrell

University of Oxford

In this course, we will see three main theorems regarding the well-posedness of the linear variational problem: for  $F \in V^*$ ,

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

of increasing generality:

In this course, we will see three main theorems regarding the well-posedness of the linear variational problem: for  $F \in V^*$ ,

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

of increasing generality:

- ▶ Riesz Representation Theorem:  $a$  *bounded, coercive, symmetric*



In this course, we will see three main theorems regarding the well-posedness of the linear variational problem: for  $F \in V^*$ ,

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

of increasing generality:

- ▶ Riesz Representation Theorem:  $a$  *bounded, coercive*, symmetric
- ▶ Lax–Milgram Theorem:  $a$  bounded, coercive

In this course, we will see three main theorems regarding the well-posedness of the linear variational problem: for  $F \in V^*$ ,

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

of increasing generality:

- ▶ Riesz Representation Theorem:  $a$  *bounded, coercive*, symmetric
- ▶ Lax–Milgram Theorem:  $a$  bounded, coercive
- ▶ Babuška's Theorem:  $a$  bounded, satisfies an inf-sup condition

In this course, we will see three main theorems regarding the well-posedness of the linear variational problem: for  $F \in V^*$ ,

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

of increasing generality:

- ▶ Riesz Representation Theorem:  $a$  *bounded, coercive*, symmetric
- ▶ Lax–Milgram Theorem:  $a$  bounded, coercive
- ▶ Babuška's Theorem:  $a$  bounded, satisfies an inf-sup condition

In this lecture we will study the first two.

## Definition (Bounded bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *bounded* if there exists  $C \in [0, \infty)$  such that

$$|a(v, w)| \leq C \|v\|_H \|w\|_H \quad \text{for all } v, w \in H.$$

## Definition (Bounded bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *bounded* if there exists  $C \in [0, \infty)$  such that

$$|a(v, w)| \leq C \|v\|_H \|w\|_H \quad \text{for all } v, w \in H.$$

As with linear functionals, this is equivalent to continuity.

## Definition (Bounded bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *bounded* if there exists  $C \in [0, \infty)$  such that

$$|a(v, w)| \leq C \|v\|_H \|w\|_H \quad \text{for all } v, w \in H.$$

As with linear functionals, this is equivalent to continuity.

The best constant  $C$  satisfying the definition is called the continuity constant of  $a$ :

$$C := \sup_{\substack{v \in H \\ v \neq 0}} \sup_{\substack{w \in H \\ w \neq 0}} \frac{|a(v, w)|}{\|v\|_H \|w\|_H}.$$

## Definition (Coercive bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *coercive* on  $V \subset H$  or  $V$ -coercive if there exists  $\alpha > 0$  such that

$$a(v, v) \geq \alpha \|v\|_H^2 \quad \text{for all } v \in V.$$

## Definition (Coercive bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *coercive* on  $V \subset H$  or  $V$ -coercive if there exists  $\alpha > 0$  such that

$$a(v, v) \geq \alpha \|v\|_H^2 \quad \text{for all } v \in V.$$

This is stronger than  $a$  being positive-definite ( $a(u, u) > 0$  for  $u \neq 0$ ).

## Example

Consider the space  $H = \ell_2(\mathbb{R})$ , the space of square-summable sequences. The form

$$a(x, y) = \sum_{m=1}^{\infty} 2^{-m} x_m y_m$$

is positive-definite but not coercive.



## Definition (Coercive bilinear form)

A bilinear form  $a : H \times H \rightarrow \mathbb{R}$  is said to be *coercive* on  $V \subset H$  or  $V$ -coercive if there exists  $\alpha > 0$  such that

$$a(v, v) \geq \alpha \|v\|_H^2 \quad \text{for all } v \in V.$$

As before, the best constant  $\alpha$  satisfying the definition is called the coercivity constant of  $a$ :

$$\alpha := \inf_{\substack{u \in V \\ u \neq 0}} \frac{a(u, u)}{\|u\|_H^2}.$$

Note that we must have  $\alpha \leq C$ , as

$$\alpha \|u\|_H^2 \leq a(u, u) \leq C \|u\|_H^2.$$

Let's assume for now that  $a$  is also symmetric.

## Theorem

*Let  $H$  be a Hilbert space, and suppose  $a : H \times H \rightarrow \mathbb{R}$  is a symmetric bilinear form that is continuous on  $H$  and coercive on a closed subspace  $V \subset H$ . Then  $(V, a(\cdot, \cdot))$  is a Hilbert space.*

Let's assume for now that  $a$  is also symmetric.

## Theorem

*Let  $H$  be a Hilbert space, and suppose  $a : H \times H \rightarrow \mathbb{R}$  is a symmetric bilinear form that is continuous on  $H$  and coercive on a closed subspace  $V \subset H$ . Then  $(V, a(\cdot, \cdot))$  is a Hilbert space.*

We must prove that  $a$  is an inner product on  $V$ , and that  $V$  is complete with respect to the induced norm.

If  $0 = a(v, v) \geq \alpha \|v\|_H^2 \geq 0$ , then  $v = 0$ . Clearly  $a(v, v) \geq 0$  for all  $v \in V$ . Symmetry and linearity are assumed, so  $a(\cdot, \cdot)$  is an inner product on  $V$ .

Denote

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

It remains to show that  $(V, \|\cdot\|_a)$  is complete.

Suppose that  $\{v_n\}$  is a Cauchy sequence in  $(V, \|\cdot\|_a)$ , i.e.

$$\forall \varepsilon > 0 \exists N > 0 \forall m, n > N \|v_n - v_m\|_a < \varepsilon.$$

Since  $\|v\|_H \leq \frac{1}{\sqrt{\alpha}}\|v\|_a$ ,  $\|v_n - v_m\|_H < \varepsilon/\sqrt{\alpha}$  and  $\{v_n\}$  is also Cauchy in  $(H, \|\cdot\|_H)$ .

Since  $H$  is complete, there exists  $v \in H$  such that  $v_n \rightarrow v$  in the  $\|\cdot\|_H$  norm. Since  $V$  is closed in  $H$ ,  $v \in V$ . Now observe that as  $a$  is bounded

$$\|v - v_n\|_a = \sqrt{a(v - v_n, v - v_n)} \leq \sqrt{C\|v - v_n\|_H^2} = \sqrt{C}\|v - v_n\|_H$$

where  $C$  is the continuity constant for  $a$ . Hence  $v_n \rightarrow v$  in the  $\|\cdot\|_a$  norm too, so  $V$  is complete with respect to this norm.

Faster: note that coercivity and continuity guarantee that

$$\alpha \|v\|_H^2 \leq \|v\|_a^2 \leq C \|v\|_H^2 \quad \text{for all } v \in V.$$

So the norms are equivalent, and hence induce the same notion of convergence and completeness.

The well-posedness of the symmetric coercive bounded linear variational problem follows immediately.

## Theorem

*Let  $V$  be a closed subspace of a Hilbert space  $H$ . Let  $a : H \times H \rightarrow \mathbb{R}$  be a symmetric continuous  $V$ -coercive bilinear form, and let  $F \in V^*$ .*

*Consider the variational problem:*

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V.$$

*This problem has a unique stable solution.*



The well-posedness of the symmetric coercive bounded linear variational problem follows immediately.

## Theorem

*Let  $V$  be a closed subspace of a Hilbert space  $H$ . Let  $a : H \times H \rightarrow \mathbb{R}$  be a symmetric continuous  $V$ -coercive bilinear form, and let  $F \in V^*$ .*

*Consider the variational problem:*

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V.$$

*This problem has a unique stable solution.*

## Proof.

Our previous result implies that  $a(\cdot, \cdot)$  is an inner product on  $V$ , and that  $(V, a)$  is a Hilbert space. Apply the Riesz Representation Theorem, that every bounded linear functional  $F \in V^*$  has a unique representative (in this case  $u$ ).

## Proof.

Stability means that we can find a constant  $c$  such that

$$\|u\|_V \leq c \|F\|_{V^*}.$$

By the Riesz representation theorem, the Riesz map is an isomorphism, so this follows for the norms generated by the inner product with  $c = 1$ .  $\square$

## Example

The variational problem

find  $u \in H_0^1(\Omega)$  such that  $\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx$  for all  $v \in H_0^1(\Omega)$

is well-posed, as  $H_0^1(\Omega)$  is a closed subspace of  $H^1(\Omega)$ , and we will show later that the bilinear form is  $H_0^1(\Omega)$ -coercive, symmetric, and bounded.

## Section 3

### The nonsymmetric case

Now let us drop the assumption that  $a(u, v) = a(v, u)$ .

### Theorem (Lax–Milgram)

*Let  $V$  be a closed subspace of a Hilbert space  $H$ . Let  $a : H \times H \rightarrow \mathbb{R}$  be a (not necessarily symmetric) continuous  $V$ -coercive bilinear form, and let  $F \in V^*$ . Consider the variational problem:*

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V.$$

*This problem has a unique stable solution.*

For the proof, it will be more convenient to treat the LVP as an equation in the dual  $V^*$ .

## Lemma

*Let  $a : V \times V \rightarrow \mathbb{R}$  be linear in its second argument and bounded. For any  $u \in V$ , define a functional via  $A : u \mapsto Au$*

$$(Au)(v) := a(u, v) \quad \text{for all } v \in V.$$

*Then  $Au \in V^*$ , i.e.  $A : V \rightarrow V^*$ . Furthermore,  $A$  is itself linear if  $a$  is linear in its first argument.*

For the proof, it will be more convenient to treat the LVP as an equation in the dual  $V^*$ .

## Lemma

*Let  $a : V \times V \rightarrow \mathbb{R}$  be linear in its second argument and bounded. For any  $u \in V$ , define a functional via  $A : u \mapsto Au$*

$$(Au)(v) := a(u, v) \quad \text{for all } v \in V.$$

*Then  $Au \in V^*$ , i.e.  $A : V \rightarrow V^*$ . Furthermore,  $A$  is itself linear if  $a$  is linear in its first argument.*

## Proof.

Linearity is straightforward. For boundedness (so that  $Au \in V^*$ ),

$$\|Au\|_{V^*} = \sup_{v \neq 0} \frac{|Au(v)|}{\|v\|_H} = \sup_{v \neq 0} \frac{|a(u, v)|}{\|v\|_H} \leq C\|u\|_H < \infty.$$



Thus, the variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

is equivalent to

$$\text{find } u \in V \text{ such that } \langle Au, v \rangle = \langle F, v \rangle \quad \text{for all } v \in V.$$

And since equality of two dual objects means exactly that they have the same output on all possible inputs, this is equivalent to

$$\text{find } u \in V \text{ such that } Au = F,$$

where the equality is between dual objects,  $Au \in V^*$  and  $F \in V^*$ .



## Example

In the case of the homogeneous Dirichlet Laplacian operator, we have  $A : H_0^1(\Omega) \rightarrow (H_0^1(\Omega))^*$ . We could symbolically write  $A = -\nabla^2$  and interpret

$$-\nabla^2 u = f$$

as an equation in the dual of  $H_0^1(\Omega)$ . This dual space is denoted

$$H^{-1}(\Omega) := (H_0^1(\Omega))^*$$

and we can regard the Laplacian as a map  $H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$ .

We know from the Riesz Representation Theorem that there is an isometric isomorphism  $\mathcal{R} : V^* \rightarrow V$  from the dual of a Hilbert space  $V^*$  back to  $V$ . By composing these operators, we have the problem

$$\text{find } u \in V \text{ such that } \mathcal{R}Au = \mathcal{R}F,$$

where the equality is between *primal* objects,  $\mathcal{R}Au \in V$  and  $\mathcal{R}F \in V$ .

Proof strategy: we will define a map  $T : V \rightarrow V$  whose fixed point is the solution of our variational problem, and then show it is a contraction, and invoke the Banach contraction mapping theorem.

### Theorem (Contraction mapping theorem)

*Given a nonempty Banach space  $V$  and a mapping  $T : V \rightarrow V$  satisfying*

$$\|Tv_1 - Tv_2\| \leq M\|v_1 - v_2\|$$

*for all  $v_1, v_2 \in V$  and fixed  $M$ ,  $0 \leq M < 1$ , there exists a unique  $u \in V$  such that*

$$u = Tu.$$

*That is, a contraction  $T$  has a unique fixed point  $u$ .*

We now prove the Lax–Milgram Theorem.

Proof.

Cast the variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in V$$

as the primal equality

$$\text{find } u \in V \text{ such that } \mathcal{R}Au = \mathcal{R}F$$

as discussed. For a fixed  $\rho \in (0, \infty)$ , define the affine map  $T : V \rightarrow V$

$$Tv = v - \rho(\mathcal{R}Av - \mathcal{R}F).$$

If  $T$  is a contraction for some  $\rho$ , then there exists a unique fixed point  $u \in V$  such that

$$Tu = u - \rho(\mathcal{R}Au - \mathcal{R}F) = u,$$

i.e. that  $\mathcal{R}Au = \mathcal{R}F$ . We now show that such a  $\rho$  exists.

Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\|Tv_1 - Tv_2\|_H^2 = \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2$$

Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}\|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\ &= \|v - \rho(\mathcal{R}Av)\|_H^2\end{aligned}\quad (\text{lin. of } \mathcal{R}, A)$$

## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)}
 \end{aligned}$$

## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)} \\
 &= \|v\|_H^2 - 2\rho Av(v) + \rho^2 Av(\mathcal{R}Av) && \text{(definition of } \mathcal{R})
 \end{aligned}$$



## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)} \\
 &= \|v\|_H^2 - 2\rho Av(v) + \rho^2 Av(\mathcal{R}Av) && \text{(definition of } \mathcal{R}) \\
 &= \|v\|_H^2 - 2\rho a(v, v) + \rho^2 a(v, \mathcal{R}Av) && \text{(definition of } A)
 \end{aligned}$$

## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)} \\
 &= \|v\|_H^2 - 2\rho Av(v) + \rho^2 Av(\mathcal{R}Av) && \text{(definition of } \mathcal{R}) \\
 &= \|v\|_H^2 - 2\rho a(v, v) + \rho^2 a(v, \mathcal{R}Av) && \text{(definition of } A) \\
 &\leq \|v\|_H^2 - 2\rho\alpha\|v\|_H^2 + \rho^2 C\|v\|_H\|\mathcal{R}Av\|_H && \text{(coerc. \& cont.)}
 \end{aligned}$$

## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)} \\
 &= \|v\|_H^2 - 2\rho Av(v) + \rho^2 Av(\mathcal{R}Av) && \text{(definition of } \mathcal{R}) \\
 &= \|v\|_H^2 - 2\rho a(v, v) + \rho^2 a(v, \mathcal{R}Av) && \text{(definition of } A) \\
 &\leq \|v\|_H^2 - 2\rho\alpha\|v\|_H^2 + \rho^2 C\|v\|_H\|\mathcal{R}Av\|_H && \text{(coerc. \& cont.)} \\
 &\leq (1 - 2\rho\alpha + \rho^2 C^2)\|v\|_H^2 && (A \text{ cts, } \mathcal{R} \text{ isom.})
 \end{aligned}$$

## Proof.

For any  $v_1, v_2 \in V$ , let  $v = v_1 - v_2$ . Then

$$\begin{aligned}
 \|Tv_1 - Tv_2\|_H^2 &= \|v_1 - v_2 - \rho(\mathcal{R}Av_1 - \mathcal{R}Av_2)\|_H^2 \\
 &= \|v - \rho(\mathcal{R}Av)\|_H^2 && \text{(lin. of } \mathcal{R}, A) \\
 &= \|v\|_H^2 - 2\rho(\mathcal{R}Av, v) + \rho^2\|\mathcal{R}Av\|_H^2 && \text{(lin. of i. prod.)} \\
 &= \|v\|_H^2 - 2\rho Av(v) + \rho^2 Av(\mathcal{R}Av) && \text{(definition of } \mathcal{R}) \\
 &= \|v\|_H^2 - 2\rho a(v, v) + \rho^2 a(v, \mathcal{R}Av) && \text{(definition of } A) \\
 &\leq \|v\|_H^2 - 2\rho\alpha\|v\|_H^2 + \rho^2 C\|v\|_H\|\mathcal{R}Av\|_H && \text{(coerc. \& cont.)} \\
 &\leq (1 - 2\rho\alpha + \rho^2 C^2)\|v\|_H^2 && (A \text{ cts, } \mathcal{R} \text{ isom.)} \\
 &= (1 - 2\rho\alpha + \rho^2 C^2)\|v_1 - v_2\|_H^2.
 \end{aligned}$$

## Proof.

Thus, if we can find a  $\rho$  such that

$$1 - 2\rho\alpha + \rho^2 C^2 < 1,$$

i.e. that

$$\rho(\rho C^2 - 2\alpha) < 0,$$

then we are done. If we choose  $\rho \in (0, 2\alpha/C^2)$  then  $T$  is a contraction and a unique solution exists.

## Proof.

It remains to show stability.

$$\|u\|_H^2 \leq \frac{1}{\alpha} a(u, u) = \frac{1}{\alpha} F(u) \leq \frac{1}{\alpha} \|F\|_{V^*} \|u\|_H,$$

and so

$$\|u\|_H \leq \frac{1}{\alpha} \|F\|_{V^*}$$



## C6.4 Finite Element Methods for PDEs

### Lecture 5: More variational formulations

Patrick E. Farrell

University of Oxford

This lecture has three goals:

- ▶ Study the variational formulation of more problems;
- ▶ Look at how to prove well-posedness with Lax–Milgram;
- ▶ See some problems our current theory cannot handle.



We consider the Poisson problem in one dimension with mixed boundary conditions:

$$-u'' = f, \quad u(0) = 0, \quad u'(1) = g.$$

Let's investigate how we can use Lax–Milgram for this.

We consider the Poisson problem in one dimension with mixed boundary conditions:

$$-u'' = f, \quad u(0) = 0, \quad u'(1) = g.$$

Let's investigate how we can use Lax–Milgram for this.

The solution can be determined from  $f$  via two integrations. First of all, by integrating both sides from  $t$  to 1, we can write

$$u'(t) = \int_t^1 f(s) \, ds + g,$$

We consider the Poisson problem in one dimension with mixed boundary conditions:

$$-u'' = f, \quad u(0) = 0, \quad u'(1) = g.$$

Let's investigate how we can use Lax–Milgram for this.

The solution can be determined from  $f$  via two integrations. First of all, by integrating both sides from  $t$  to 1, we can write

$$u'(t) = \int_t^1 f(s) \, ds + g,$$

and integrating again from 0 to  $x$  yields

$$u(x) = \int_0^x \int_t^1 f(s) \, ds \, dt + gx.$$

We consider the Poisson problem in one dimension with mixed boundary conditions:

$$-u'' = f, \quad u(0) = 0, \quad u'(1) = g.$$

Let's investigate how we can use Lax–Milgram for this.

The solution can be determined from  $f$  via two integrations. First of all, by integrating both sides from  $t$  to 1, we can write

$$u'(t) = \int_t^1 f(s) \, ds + g,$$

and integrating again from 0 to  $x$  yields

$$u(x) = \int_0^x \int_t^1 f(s) \, ds \, dt + gx.$$

This shows that the equation is well-posed.

Now consider a variational formulation. We define the space

$$V = \{v \in H^1(0, 1) : v(0) = 0\}.$$

This makes sense, because in one dimension  $H^1$  functions are continuous.

Now consider a variational formulation. We define the space

$$V = \{v \in H^1(0, 1) : v(0) = 0\}.$$

This makes sense, because in one dimension  $H^1$  functions are continuous.

Two kinds of boundary conditions:

- ▶ The Dirichlet condition  $u(0) = 0$  is carried in the definition of the space (strongly enforced).

Now consider a variational formulation. We define the space

$$V = \{v \in H^1(0, 1) : v(0) = 0\}.$$

This makes sense, because in one dimension  $H^1$  functions are continuous.

Two kinds of boundary conditions:

- ▶ The Dirichlet condition  $u(0) = 0$  is carried in the definition of the space (strongly enforced).
- ▶ The Neumann condition  $u'(1) = g$  will appear in the variational formulation (weakly enforced).

Multiplying the equation by  $v \in V$  and integrating, we find

$$\int_0^1 -u''v \, dx = \int_0^1 fv \, dx.$$



Multiplying the equation by  $v \in V$  and integrating, we find

$$\int_0^1 -u''v \, dx = \int_0^1 f v \, dx.$$

We next integrate by parts:

$$\int_0^1 u'v' \, dx - u'(1)v(1) + u'(0)v(0) = \int_0^1 f v \, dx.$$

Multiplying the equation by  $v \in V$  and integrating, we find

$$\int_0^1 -u''v \, dx = \int_0^1 fv \, dx.$$

We next integrate by parts:

$$\int_0^1 u'v' \, dx - u'(1)v(1) + u'(0)v(0) = \int_0^1 fv \, dx.$$

The surface integral term on the left disappears as  $v(0) = 0$ . On the right, we know that  $u'(1) = g$ , and so we have

$$\int_0^1 u'v' \, dx = \int_0^1 fv \, dx + gv(1).$$

Multiplying the equation by  $v \in V$  and integrating, we find

$$\int_0^1 -u''v \, dx = \int_0^1 fv \, dx.$$

We next integrate by parts:

$$\int_0^1 u'v' \, dx - u'(1)v(1) + u'(0)v(0) = \int_0^1 fv \, dx.$$

The surface integral term on the left disappears as  $v(0) = 0$ . On the right, we know that  $u'(1) = g$ , and so we have

$$\int_0^1 u'v' \, dx = \int_0^1 fv \, dx + gv(1).$$

Thus, we have a linear variational problem with

$$a(u, v) = \int_0^1 u'v' \, dx, \quad F(v) = \int_0^1 fv \, dx + gv(1).$$

## Theorem

*The following bilinear form is coercive on  $V$ :*

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

## Theorem

*The following bilinear form is coercive on  $V$ :*

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

## Proof.

The norm on  $H^1(0, 1)$  is

$$\|v\|_{H^1(0,1)}^2 = \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2.$$

We wish to show that there exists a constant  $\alpha > 0$  such that

$$a(v, v) \geq \alpha \|v\|_{H^1(0,1)}^2 \quad \text{for all } v \in V.$$

Expanding definitions, we want to find an  $\alpha$  such that

$$a(v, v) = \|v'\|_{L^2(0,1)}^2 \geq \alpha \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right).$$

Proof.

If we can prove that there exists an  $\alpha' > 0$  such that

$$\|v'\|_{L^2(0,1)}^2 \geq \alpha' \|v\|_{L^2(0,1)}^2$$

## Proof.

If we can prove that there exists an  $\alpha' > 0$  such that

$$\begin{aligned} \|v'\|_{L^2(0,1)}^2 &\geq \alpha' \|v\|_{L^2(0,1)}^2 \\ \iff \|v'\|_{L^2(0,1)}^2 + \alpha' \|v'\|_{L^2(0,1)}^2 &\geq \alpha' \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right) \end{aligned}$$

## Proof.

If we can prove that there exists an  $\alpha' > 0$  such that

$$\begin{aligned} \|v'\|_{L^2(0,1)}^2 &\geq \alpha' \|v\|_{L^2(0,1)}^2 \\ \iff \|v'\|_{L^2(0,1)}^2 + \alpha' \|v'\|_{L^2(0,1)}^2 &\geq \alpha' \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right) \\ \iff \|v'\|_{L^2(0,1)}^2 &\geq \frac{\alpha'}{\alpha' + 1} \|v\|_{H^1(0,1)}^2 \end{aligned}$$



## Proof.

If we can prove that there exists an  $\alpha' > 0$  such that

$$\begin{aligned}\|v'\|_{L^2(0,1)}^2 &\geq \alpha' \|v\|_{L^2(0,1)}^2 \\ \iff \|v'\|_{L^2(0,1)}^2 + \alpha' \|v'\|_{L^2(0,1)}^2 &\geq \alpha' \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right) \\ \iff \|v'\|_{L^2(0,1)}^2 &\geq \frac{\alpha'}{\alpha' + 1} \|v\|_{H^1(0,1)}^2\end{aligned}$$

then we are done with  $\alpha = \frac{\alpha'}{\alpha' + 1}$ .

## Proof.

We can write

$$v(t) = \int_0^t v'(x) \, dx = \int_0^1 v'(x) w_t'(x) \, dx = a(v, w_t),$$

where the function  $w_t \in V$  is defined by

$$w_t(x) = \begin{cases} x & 0 \leq x \leq t, \\ t & x > t. \end{cases}$$

## Proof.

We can write

$$v(t) = \int_0^t v'(x) \, dx = \int_0^1 v'(x) w_t'(x) \, dx = a(v, w_t),$$

where the function  $w_t \in V$  is defined by

$$w_t(x) = \begin{cases} x & 0 \leq x \leq t, \\ t & x > t. \end{cases}$$

This function is not strongly differentiable, but has weak derivative

$$w_t'(x) = \begin{cases} 1 & 0 \leq x \leq t, \\ 0 & x > t, \end{cases}$$

The function  $w_t(x) \in V$  is the  $a$ -Riesz representation of the functional  $j : v \mapsto v(t)$ .

## Proof.

We can invoke Cauchy–Schwarz on  $L^2(0, 1)$  to get

$$|v(t)| = |a(v, w_t)| \leq \|v'\|_{L^2(0,1)} \|w'_t\|_{L^2(0,1)} = \sqrt{t} \|v'\|_{L^2(0,1)},$$

since

$$\|w'_t\|_{L^2(0,1)} = \left( \int_0^t 1^2 \, dx \right)^{1/2} = \sqrt{t}.$$

Thus,

$$\|v\|_{L^2(0,1)}^2 = \int_0^1 v^2(x) \, dx \leq \int_0^1 x \|v'\|_{L^2(0,1)}^2 \, dx = \frac{1}{2} \|v'\|_{L^2(0,1)}^2$$

so in this case we can take  $\alpha' = 2$  and thus  $\alpha = \frac{2}{3}$ . □

Two remarks:

### Remark

Note that if we consider  $a$  over the whole of  $H^1(0, 1)$ , it is *not* coercive:  $v(x) \equiv 1 \in H^1(0, 1)$  with  $a(v, v) = 0$  but  $\|v\| > 0$ .

The boundary condition  $v(0) = 0$  is essential to the coercivity.

Two remarks:

### Remark

Note that if we consider  $a$  over the whole of  $H^1(0, 1)$ , it is *not* coercive:  $v(x) \equiv 1 \in H^1(0, 1)$  with  $a(v, v) = 0$  but  $\|v\| > 0$ .

The boundary condition  $v(0) = 0$  is essential to the coercivity.

### Remark

Notice that the coercivity constant will depend on the length of the domain: for an interval of length  $L$ ,  $\alpha' = \frac{2}{L^2}$  and  $\alpha = \frac{2}{L^2+2}$ .

## Theorem

*The following bilinear form is continuous:*

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

## Theorem

*The following bilinear form is continuous:*

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

## Proof.

$$\begin{aligned} |a(u, v)| &\leq \|u'\|_{L^2(0,1)} \|v'\|_{L^2(0,1)} \\ &\leq \left( \|u\|_{L^2(0,1)}^2 + \|u'\|_{L^2(0,1)}^2 \right)^{\frac{1}{2}} \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right)^{\frac{1}{2}} \\ &= \|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)}. \end{aligned}$$

That is, the bilinear form is continuous with  $C = 1$ . □



## Theorem

*The following bilinear form is continuous:*

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

## Proof.

$$\begin{aligned} |a(u, v)| &\leq \|u'\|_{L^2(0,1)} \|v'\|_{L^2(0,1)} \\ &\leq \left( \|u\|_{L^2(0,1)}^2 + \|u'\|_{L^2(0,1)}^2 \right)^{\frac{1}{2}} \left( \|v\|_{L^2(0,1)}^2 + \|v'\|_{L^2(0,1)}^2 \right)^{\frac{1}{2}} \\ &= \|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)}. \end{aligned}$$

That is, the bilinear form is continuous with  $C = 1$ . □

We can thus apply Lax–Milgram to problems with this bilinear form.

Now consider a nonsymmetric problem: for  $f \in L^2(0, 1)$ , solve

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

Now consider a nonsymmetric problem: for  $f \in L^2(0, 1)$ , solve

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

Since we have no Dirichlet conditions, we set  $V = H^1(0, 1)$ .

Now consider a nonsymmetric problem: for  $f \in L^2(0, 1)$ , solve

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

Since we have no Dirichlet conditions, we set  $V = H^1(0, 1)$ .

Testing against  $v \in V$  and integrating by parts, we find

$$\int_0^1 u'v' \, dx + \int_0^1 u'v \, dx + \int_0^1 uv \, dx = \int_0^1 fv \, dx.$$

Thus, our standard variational problem has

$$a(u, v) = \int_0^1 u'v' + u'v + uv \, dx, \quad F(v) = \int_0^1 fv \, dx.$$

To prove continuity, observe that

$$\begin{aligned} |a(u, v)| &\leq |(u, v)_{H^1(0,1)}| + \left| \int_0^1 u'v \, dx \right| \\ &\leq \|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)} + \|u'\|_{L^2(0,1)} \|v\|_{L^2(0,1)} \\ &\leq 2\|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)} \end{aligned}$$

so we can take our continuity constant  $C = 2$ .

To prove continuity, observe that

$$\begin{aligned}|a(u, v)| &\leq |(u, v)_{H^1(0,1)}| + \left| \int_0^1 u'v \, dx \right| \\ &\leq \|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)} + \|u'\|_{L^2(0,1)} \|v\|_{L^2(0,1)} \\ &\leq 2\|u\|_{H^1(0,1)} \|v\|_{H^1(0,1)}\end{aligned}$$

so we can take our continuity constant  $C = 2$ .

To prove coercivity, observe that

$$\begin{aligned}a(v, v) &= \int_0^1 v'^2 + v'v + v^2 \, dx \\ &= \frac{1}{2} \int_0^1 (v'^2 + v^2) \, dx + \frac{1}{2} \int_0^1 (v' + v)^2 \, dx \\ &\geq \frac{1}{2} \|v\|_{H^1(0,1)}^2.\end{aligned}$$

## Section 2

# Higher dimensions

Break up  $\partial\Omega$  into disjoint  $\Gamma_D$  and  $\Gamma_N$ . Consider

$$\begin{aligned} -\nabla^2 u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \Gamma_D, \\ \nabla u \cdot n &= g \text{ on } \Gamma_N. \end{aligned}$$



Break up  $\partial\Omega$  into disjoint  $\Gamma_D$  and  $\Gamma_N$ . Consider

$$\begin{aligned}-\nabla^2 u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \Gamma_D, \\ \nabla u \cdot n &= g \text{ on } \Gamma_N.\end{aligned}$$

Define the space

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\}.$$

Multiplying by  $v \in V$ , integrating and integrating by parts, we get

$$\begin{aligned}\int_{\Omega} \nabla u \cdot \nabla v \, dx &= \int_{\Omega} f v \, dx + \int_{\partial\Omega} \nabla u \cdot n v \, ds, \\ &= \int_{\Omega} f v \, dx + \int_{\Gamma_N} g v \, ds.\end{aligned}$$

The continuity proof works as in one dimension. For coercivity, we need a result from functional analysis.

The continuity proof works as in one dimension. For coercivity, we need a result from functional analysis.

### Theorem (Poincaré–Friedrichs inequality)

*Let  $\Omega$  be a bounded Lipschitz domain, and suppose  $\Gamma_D \subset \partial\Omega$  is closed and has nonzero measure. Let*

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\}.$$

*Then there is a constant  $K \in \mathbb{R}$  depending only on  $\Omega$  and  $\Gamma_D$  such that*

$$\|u\|_{L^2(\Omega)}^2 = \int_{\Omega} u^2 \, dx \leq K \int_{\Omega} |\nabla u|^2 \, dx = K \|\nabla u\|_{L^2(\Omega)}^2$$

*for all  $u \in V$ . The constant  $K(\Omega, \Gamma_D)$  is called the Poincaré constant for the domain and boundary.*

Write

$$|u|_{H^1(\Omega)}^2 := \int_{\Omega} |\nabla u|^2 \, dx = \|\nabla u\|_{L^2(\Omega)}^2.$$

In general this is a *seminorm*:  $|u|_{H^1(\Omega)} = 0 \not\Rightarrow u = 0$ .

Write

$$|u|_{H^1(\Omega)}^2 := \int_{\Omega} |\nabla u|^2 \, dx = \|\nabla u\|_{L^2(\Omega)}^2.$$

In general this is a *seminorm*:  $|u|_{H^1(\Omega)} = 0 \not\Rightarrow u = 0$ .

On

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\},$$

manipulating the Poincaré–Friedrichs inequality yields

$$\frac{1}{K+1} \|u\|_{H^1(\Omega)}^2 \leq |u|_{H^1(\Omega)}^2 \leq \|u\|_{H^1(\Omega)}^2.$$

So if you have Dirichlet conditions somewhere then  $|u|_{H^1(\Omega)}$  is an *equivalent norm* to  $\|u\|_{H^1(\Omega)}$ .

Write

$$|u|_{H^1(\Omega)}^2 := \int_{\Omega} |\nabla u|^2 \, dx = \|\nabla u\|_{L^2(\Omega)}^2.$$

In general this is a *seminorm*:  $|u|_{H^1(\Omega)} = 0 \not\Rightarrow u = 0$ .

On

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\},$$

manipulating the Poincaré–Friedrichs inequality yields

$$\frac{1}{K+1} \|u\|_{H^1(\Omega)}^2 \leq |u|_{H^1(\Omega)}^2 \leq \|u\|_{H^1(\Omega)}^2.$$

So if you have Dirichlet conditions somewhere then  $|u|_{H^1(\Omega)}$  is an *equivalent norm* to  $\|u\|_{H^1(\Omega)}$ .

## Remark

If  $\Omega$  is contained within an  $n$ -dimensional cube of side  $L$ , then  $L$  provides a (possibly non-optimal) Poincaré constant. (Braess, 2007)

What happens with inhomogeneous Dirichlet conditions?

$$-\nabla^2 u = f \text{ in } \Omega$$

$$u = h \text{ on } \Gamma_D$$

$$\nabla u \cdot n = g \text{ on } \Gamma_N.$$

What happens with inhomogeneous Dirichlet conditions?

$$\begin{aligned}-\nabla^2 u &= f \text{ in } \Omega \\ u &= h \text{ on } \Gamma_D \\ \nabla u \cdot n &= g \text{ on } \Gamma_N.\end{aligned}$$

Consider

$$\hat{u} = u - h.$$

Then  $\hat{u}$  satisfies

$$\begin{aligned}-\nabla^2 \hat{u} &= f + \nabla^2 h \quad \text{in } \Omega \\ \hat{u} &= 0 \quad \text{on } \Gamma_D \\ \nabla \hat{u} \cdot n &= g - \nabla h \cdot n \text{ on } \Gamma_N,\end{aligned}$$

or variationally

$$a(\hat{u}, v) = a(u - h, v) = a(u, v) - a(h, v) = F(v) - a(h, v)$$

for all  $v \in V$ .



What about pure Neumann conditions?

$$-\nabla^2 u = f \text{ in } \Omega$$

$$\nabla u \cdot n = g \text{ on } \partial\Omega.$$

What about pure Neumann conditions?

$$-\nabla^2 u = f \text{ in } \Omega$$

$$\nabla u \cdot n = g \text{ on } \partial\Omega.$$

If  $u$  satisfies the equation, so does  $u + c$ ,  $c \in \mathbb{R}$ !

What about pure Neumann conditions?

$$\begin{aligned} -\nabla^2 u &= f \text{ in } \Omega \\ \nabla u \cdot n &= g \text{ on } \partial\Omega. \end{aligned}$$

If  $u$  satisfies the equation, so does  $u + c$ ,  $c \in \mathbb{R}$ !

If this is to have any solution,  $f$  and  $g$  must be compatible:

$$\int_{\Omega} f \, dx = \int_{\Omega} -\nabla^2 u \, dx = \int_{\Omega} \nabla u \cdot \nabla 1 \, dx - \int_{\partial\Omega} \nabla u \cdot n \, ds = - \int_{\partial\Omega} g \, ds.$$

If we take the variational formulation with  $V = H^1(\Omega)$ , we get

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx + \int_{\partial\Omega} g v \, ds,$$

We cannot apply Lax–Milgram because the bilinear form is not  $V$ -coercive.

If we take the variational formulation with  $V = H^1(\Omega)$ , we get

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx + \int_{\partial\Omega} g v \, ds,$$

We cannot apply Lax–Milgram because the bilinear form is not  $V$ -coercive.

To eliminate the nullspace of constants, consider the solution space that is  $L^2$ -orthogonal to it:

$$V = \{v \in H^1(\Omega) : \int_{\Omega} v \, dx = 0\}$$

The bilinear form *is* coercive over this space (Poincaré–Neumann).

Let us briefly consider two other kinds of boundary conditions. Robin conditions relate  $\nabla u \cdot n$  and  $u$ :

$$\begin{aligned} -\nabla^2 u &= f \text{ in } \Omega, \\ \nabla u \cdot n + \beta u &= g \text{ on } \partial\Omega. \end{aligned}$$

Taking  $V = H^1(\Omega)$ , this yields the variational problem

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx + \beta \int_{\partial\Omega} uv \, ds = \int_{\Omega} fv \, dx + \int_{\partial\Omega} gv \, ds.$$

Proving the properties required for Lax–Milgram requires knowledge of trace theorems, which we haven't discussed.

Before moving on, let's consider an interesting physical model.

In 1879, Josef Stefan realised that the power radiated from the surface of a black body at temperature  $u$  is proportional to the difference of the fourth powers between  $u$  and the ambient temperature  $u_0$ :

$$\begin{aligned} -\nabla^2 u &= f \text{ in } \Omega, \\ \nabla u \cdot n &= c(u^4 - u_0^4) \text{ on } \partial\Omega. \end{aligned}$$

Before moving on, let's consider an interesting physical model.

In 1879, Josef Stefan realised that the power radiated from the surface of a black body at temperature  $u$  is proportional to the difference of the fourth powers between  $u$  and the ambient temperature  $u_0$ :

$$\begin{aligned}-\nabla^2 u &= f \text{ in } \Omega, \\ \nabla u \cdot n &= c(u^4 - u_0^4) \text{ on } \partial\Omega.\end{aligned}$$

Multiplying by a test function and integrating by parts, we get: find  $u \in H^1(\Omega)$  such that

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx - c \int_{\partial\Omega} (u^4 - u_0^4) v \, ds = \int_{\Omega} f v \, dx$$

for all  $v \in H^1(\Omega)$ .



Before moving on, let's consider an interesting physical model.

In 1879, Josef Stefan realised that the power radiated from the surface of a black body at temperature  $u$  is proportional to the difference of the fourth powers between  $u$  and the ambient temperature  $u_0$ :

$$\begin{aligned}-\nabla^2 u &= f \text{ in } \Omega, \\ \nabla u \cdot n &= c(u^4 - u_0^4) \text{ on } \partial\Omega.\end{aligned}$$

Multiplying by a test function and integrating by parts, we get: find  $u \in H^1(\Omega)$  such that

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx - c \int_{\partial\Omega} (u^4 - u_0^4) v \, ds = \int_{\Omega} f v \, dx$$

for all  $v \in H^1(\Omega)$ .

Problem is no longer a *linear* variational problem! To be discussed.

At this point, I may have given you the impression that there's only one variational formulation for a problem. This is not true, and different formulations have advantages and disadvantages.

At this point, I may have given you the impression that there's only one variational formulation for a problem. This is not true, and different formulations have advantages and disadvantages.

Suppose we want to know the *flux* in the Poisson equation accurately. We can solve the *mixed* formulation: find  $\sigma : \Omega \rightarrow \mathbb{R}^n$ ,  $u : \Omega \rightarrow \mathbb{R}$  such that

$$\begin{aligned}\sigma &= -\nabla u \text{ in } \Omega, \\ \nabla \cdot \sigma &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega.\end{aligned}$$

Solving this formulation will give an accurate approximation of the flux, and allow for the easy implementation of more complicated constitutive laws.

$$\sigma = -\nabla u \text{ in } \Omega, \quad \nabla \cdot \sigma = f \text{ in } \Omega.$$

Let's multiply the first equation by a vector-valued test function  $v$ , and the second by a scalar-valued function  $w$ :

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx + \int_{\Omega} \nabla u \cdot v &= 0, \\ \int_{\Omega} \nabla \cdot \sigma w \, dx &= \int_{\Omega} f w \, dx. \end{aligned}$$

$$\sigma = -\nabla u \text{ in } \Omega, \quad \nabla \cdot \sigma = f \text{ in } \Omega.$$

Let's multiply the first equation by a vector-valued test function  $v$ , and the second by a scalar-valued function  $w$ :

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx + \int_{\Omega} \nabla u \cdot v &= 0, \\ \int_{\Omega} \nabla \cdot \sigma w \, dx &= \int_{\Omega} f w \, dx. \end{aligned}$$

Since  $\sigma$  needs to have a divergence, and we want  $v$  and  $\sigma$  to come from the same space, let's integrate by parts in the first equation. For symmetry I'll negate the second equation:

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} u \nabla \cdot v + \int_{\partial\Omega} uv \cdot n \, ds &= 0, \\ - \int_{\Omega} \nabla \cdot \sigma w \, dx &= - \int_{\Omega} f w \, dx. \end{aligned}$$

$$\sigma = -\nabla u \text{ in } \Omega, \quad \nabla \cdot \sigma = f \text{ in } \Omega.$$

Let's multiply the first equation by a vector-valued test function  $v$ , and the second by a scalar-valued function  $w$ :

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx + \int_{\Omega} \nabla u \cdot v &= 0, \\ \int_{\Omega} \nabla \cdot \sigma w \, dx &= \int_{\Omega} f w \, dx. \end{aligned}$$

Since  $\sigma$  needs to have a divergence, and we want  $v$  and  $\sigma$  to come from the same space, let's integrate by parts in the first equation. For symmetry I'll negate the second equation:

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} u \nabla \cdot v + \int_{\partial\Omega} uv \cdot n \, ds &= 0, \\ - \int_{\Omega} \nabla \cdot \sigma w \, dx &= - \int_{\Omega} f w \, dx. \end{aligned}$$

Impose the Dirichlet condition weakly!

What function spaces do we need to make sense of

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} u \nabla \cdot v &= 0, \\ - \int_{\Omega} \nabla \cdot \sigma w \, dx &= - \int_{\Omega} f w \, dx. \end{aligned}$$

We don't need any derivatives on  $u$  or  $w$ , so  $u \in L^2(\Omega)$ .

What function spaces do we need to make sense of

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} u \nabla \cdot v &= 0, \\ - \int_{\Omega} \nabla \cdot \sigma w \, dx &= - \int_{\Omega} f w \, dx. \end{aligned}$$

We don't need any derivatives on  $u$  or  $w$ , so  $u \in L^2(\Omega)$ .

For  $\sigma$  and  $v$ , we need  $\sigma \in L^2(\Omega; \mathbb{R}^n)$  and for  $\nabla \cdot \sigma \in L^2(\Omega)$ . This is the space  $H(\text{div}, \Omega)$ :

$$H(\text{div}, \Omega) = \{\sigma \in L^2(\Omega; \mathbb{R}^n) : \nabla \cdot \sigma \in L^2(\Omega)\}.$$



What function spaces do we need to make sense of

$$\begin{aligned} \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} u \nabla \cdot v &= 0, \\ - \int_{\Omega} \nabla \cdot \sigma w \, dx &= - \int_{\Omega} f w \, dx. \end{aligned}$$

We don't need any derivatives on  $u$  or  $w$ , so  $u \in L^2(\Omega)$ .

For  $\sigma$  and  $v$ , we need  $\sigma \in L^2(\Omega; \mathbb{R}^n)$  and for  $\nabla \cdot \sigma \in L^2(\Omega)$ . This is the space  $H(\text{div}, \Omega)$ :

$$H(\text{div}, \Omega) = \{\sigma \in L^2(\Omega; \mathbb{R}^n) : \nabla \cdot \sigma \in L^2(\Omega)\}.$$

Its inner product is

$$(u, v)_{H(\text{div}, \Omega)} = \int_{\Omega} u \cdot v + \nabla \cdot u \nabla \cdot v \, dx.$$

A nice property of variational problems is that we can add the two equations together. The problem is the same as:

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$B(\sigma, u; v, w) := \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot vu - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} fw \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

A nice property of variational problems is that we can add the two equations together. The problem is the same as:

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$B(\sigma, u; v, w) := \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot vu - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} fw \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

Lax–Milgram certainly won't apply:

$$B(0, u; 0, u) = 0 \text{ for all } u \in L^2(\Omega).$$

A nice property of variational problems is that we can add the two equations together. The problem is the same as:

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$B(\sigma, u; v, w) := \int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot vu - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} fw \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

Lax–Milgram certainly won't apply:

$$B(0, u; 0, u) = 0 \text{ for all } u \in L^2(\Omega).$$

We will study the well-posedness of this problem with a more general theory later.

## C6.4 Finite Element Methods for PDEs

### Lecture 6: Differentiation and energy minimisation

Patrick E. Farrell

University of Oxford

In this lecture we will see a fundamental connection between symmetric linear variational problems:

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V,$$

and energy minimisation:

$$u = \operatorname{argmin}_{v \in V} J(v).$$

In this lecture we will see a fundamental connection between symmetric linear variational problems:

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ ,

and energy minimisation:

$$u = \operatorname{argmin}_{v \in V} J(v).$$

This is one reason why the variational formulation is so useful, and will also lead to an insight into Galerkin approximation.

We use derivatives of a function at a point  $u$  to give us a *local model* of how it behaves in the neighbourhood of  $u$ . We want to do the same in Banach spaces.



We use derivatives of a function at a point  $u$  to give us a *local model* of how it behaves in the neighbourhood of  $u$ . We want to do the same in Banach spaces.

Given a function  $J : V \rightarrow W$ ,  $V, W$  Banach spaces, how can we differentiate  $J$ ? As a concrete example, think of  $J : H_0^1(\Omega) \rightarrow \mathbb{R}$  defined by

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx.$$

How will this functional change value if we make a small perturbation  $v$  to the input argument?

We use derivatives of a function at a point  $u$  to give us a *local model* of how it behaves in the neighbourhood of  $u$ . We want to do the same in Banach spaces.

Given a function  $J : V \rightarrow W$ ,  $V, W$  Banach spaces, how can we differentiate  $J$ ? As a concrete example, think of  $J : H_0^1(\Omega) \rightarrow \mathbb{R}$  defined by

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx.$$

How will this functional change value if we make a small perturbation  $v$  to the input argument?

We will introduce progressively stronger notions of differentiation for this.

## Definition (Directional derivative)

Let  $J : V \rightarrow W$ , where  $V$  and  $W$  are Banach spaces. The directional derivative of  $J$  evaluated at  $u \in V$  in the direction  $v \in V$  is

$$J'(u; v) = \lim_{\varepsilon \rightarrow 0^+} \frac{J(u + \varepsilon v) - J(u)}{\varepsilon},$$

if the limit exists.

## Definition (Directional derivative)

Let  $J : V \rightarrow W$ , where  $V$  and  $W$  are Banach spaces. The directional derivative of  $J$  evaluated at  $u \in V$  in the direction  $v \in V$  is

$$J'(u; v) = \lim_{\varepsilon \rightarrow 0^+} \frac{J(u + \varepsilon v) - J(u)}{\varepsilon},$$

if the limit exists.

## Definition (Directionally differentiable)

If the directional derivative of  $J$  at  $u$  in the direction  $v$  exists for all  $v$ , then  $J$  is directionally differentiable at  $u$ .

We will want more than just all directional derivatives exist. We will want

- ▶ that the derivative is *linear and bounded* in the perturbation direction (Gâteaux);
- ▶ that the derivative is *a good approximation* of the function nearby (Fréchet).

- We will want more than just all directional derivatives exist. We will want
- ▶ that the derivative is *linear and bounded* in the perturbation direction (Gâteaux);
  - ▶ that the derivative is a *good approximation* of the function nearby (Fréchet).

### Definition (Gâteaux differentiable)

If  $J$  is directionally differentiable at  $u$ , and there exists a bounded linear map  $J'(u) : V \rightarrow W$  such that

$$J'(u; v) = J'(u)v,$$

then  $J$  is Gâteaux differentiable at  $u$  with derivative  $J'(u)$ .

- We will want more than just all directional derivatives exist. We will want
- ▶ that the derivative is *linear and bounded* in the perturbation direction (Gâteaux);
  - ▶ that the derivative is a *good approximation* of the function nearby (Fréchet).

### Definition (Gâteaux differentiable)

If  $J$  is directionally differentiable at  $u$ , and there exists a bounded linear map  $J'(u) : V \rightarrow W$  such that

$$J'(u; v) = J'(u)v,$$

then  $J$  is Gâteaux differentiable at  $u$  with derivative  $J'(u)$ .

### Example

If  $W = \mathbb{R}$ , then  $J'(u) \in V^*$  for each  $u \in V$ .

## Definition (Fréchet differentiable)

Suppose  $J$  is Gâteaux differentiable at a point  $u \in V$  and that the derivative  $J'$  satisfies

$$\lim_{v \rightarrow 0} \frac{\|J(u+v) - J(u) - J'(u)v\|}{\|v\|} = 0 \quad \text{for all } v \in V.$$

Then  $J$  is Fréchet differentiable at  $u$ .



## Definition (Fréchet differentiable)

Suppose  $J$  is Gâteaux differentiable at a point  $u \in V$  and that the derivative  $J'$  satisfies

$$\lim_{v \rightarrow 0} \frac{\|J(u+v) - J(u) - J'(u)v\|}{\|v\|} = 0 \quad \text{for all } v \in V.$$

Then  $J$  is Fréchet differentiable at  $u$ .

This allows us to approximate what  $J$  does near  $u$ .

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

$$J'(u; v) = \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u + \varepsilon \nabla v|^2 - |\nabla u|^2 \, dx - \frac{1}{\varepsilon} \int_{\Omega} f(u + \varepsilon v - u) \, dx$$

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

$$\begin{aligned} J'(u; v) &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u + \varepsilon \nabla v|^2 - |\nabla u|^2 \, dx - \frac{1}{\varepsilon} \int_{\Omega} f(u + \varepsilon v - u) \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u|^2 + 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 - |\nabla u|^2 \, dx - \int_{\Omega} f v \, dx \end{aligned}$$

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

$$\begin{aligned} J'(u; v) &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u + \varepsilon \nabla v|^2 - |\nabla u|^2 \, dx - \frac{1}{\varepsilon} \int_{\Omega} f(u + \varepsilon v - u) \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u|^2 + 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 - |\nabla u|^2 \, dx - \int_{\Omega} f v \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 \, dx - \int_{\Omega} f v \, dx \end{aligned}$$

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

$$\begin{aligned} J'(u; v) &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u + \varepsilon \nabla v|^2 - |\nabla u|^2 \, dx - \frac{1}{\varepsilon} \int_{\Omega} f(u + \varepsilon v - u) \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u|^2 + 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 - |\nabla u|^2 \, dx - \int_{\Omega} f v \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 \, dx - \int_{\Omega} f v \, dx \\ &= \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\Omega} f v \, dx. \end{aligned}$$

With our example

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 \, dx - \int_{\Omega} f u \, dx,$$

let's calculate  $J'(u; v)$  for a given  $u, v \in H_0^1(\Omega)$ .

$$\begin{aligned} J'(u; v) &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u + \varepsilon \nabla v|^2 - |\nabla u|^2 \, dx - \frac{1}{\varepsilon} \int_{\Omega} f(u + \varepsilon v - u) \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} |\nabla u|^2 + 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 - |\nabla u|^2 \, dx - \int_{\Omega} f v \, dx \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} \int_{\Omega} 2\varepsilon \nabla u \cdot \nabla v + \varepsilon^2 |\nabla v|^2 \, dx - \int_{\Omega} f v \, dx \\ &= \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\Omega} f v \, dx. \end{aligned}$$

So if  $J'(u; v) = 0$  for all  $v \in H_0^1(\Omega)$ ,  $u$  satisfies the weak formulation of the Poisson equation!

For symmetric coercive problems, there is a strong relationship between LVPs and minimisation.

### Theorem (Energy minimisation)

*Suppose  $a$  is symmetric, coercive, and bounded, and  $F \in V^*$ . Let  $u$  be the unique solution to*

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

*Then  $u$  is the unique solution to*

$$u = \underset{v \in V}{\operatorname{argmin}} J(v) := \frac{1}{2}a(v, v) - F(v).$$



## Proof.

Let  $v \in V$ . We want to show  $J(v) \geq J(u)$ . Calculating,

$$J(v) - J(u) = \frac{1}{2}a(v, v) - F(v) - \frac{1}{2}a(u, u) + F(u)$$



## Proof.

Let  $v \in V$ . We want to show  $J(v) \geq J(u)$ . Calculating,

$$\begin{aligned} J(v) - J(u) &= \frac{1}{2}a(v, v) - F(v) - \frac{1}{2}a(u, u) + F(u) \\ &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - F(v - u) \end{aligned}$$



## Proof.

Let  $v \in V$ . We want to show  $J(v) \geq J(u)$ . Calculating,

$$\begin{aligned} J(v) - J(u) &= \frac{1}{2}a(v, v) - F(v) - \frac{1}{2}a(u, u) + F(u) \\ &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - F(v - u) \\ &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - a(u, v - u) \end{aligned}$$



## Proof.

Let  $v \in V$ . We want to show  $J(v) \geq J(u)$ . Calculating,

$$\begin{aligned} J(v) - J(u) &= \frac{1}{2}a(v, v) - F(v) - \frac{1}{2}a(u, u) + F(u) \\ &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - F(v - u) \\ &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - a(u, v - u) \\ &= \frac{1}{2}a(v - u, v - u). \end{aligned}$$



## Proof.

Let  $v \in V$ . We want to show  $J(v) \geq J(u)$ . Calculating,

$$\begin{aligned}
 J(v) - J(u) &= \frac{1}{2}a(v, v) - F(v) - \frac{1}{2}a(u, u) + F(u) \\
 &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - F(v - u) \\
 &= \frac{1}{2}a(v, v) - \frac{1}{2}a(u, u) - a(u, v - u) \\
 &= \frac{1}{2}a(v - u, v - u).
 \end{aligned}$$

Because  $a$  is coercive,

$$J(v) - J(u) \geq \frac{\alpha}{2} \|v - u\|_V^2 \geq 0 \text{ for all } v \in V.$$



The minimiser  $u$  is unique, because if  $\tilde{u}$  also minimises  $J$ , then

$$J(\tilde{u}) - J(u) = 0 \geq \frac{\alpha}{2} \|\tilde{u} - u\|_V^2 \geq 0$$

and hence  $\tilde{u} = u$ .

What about the other way around?

## Theorem

*Let  $u \in V$  be a minimiser of  $J : V \rightarrow \mathbb{R}$ . Then  $u$  is a solution of*

$$J'(u; v) = 0 \text{ for all } v \in V.$$

What about the other way around?

## Theorem

*Let  $u \in V$  be a minimiser of  $J : V \rightarrow \mathbb{R}$ . Then  $u$  is a solution of*

$$J'(u; v) = 0 \text{ for all } v \in V.$$

## Proof.

Since  $u$  is a minimiser of  $J$ ,  $J(u + \varepsilon v) \geq J(u)$  for all  $\varepsilon > 0, v \in V$ . This implies

$$\frac{J(u + \varepsilon v) - J(u)}{\varepsilon} \geq 0.$$

Taking the limit as  $\varepsilon \rightarrow 0^+$  yields  $J'(u; v) \geq 0$ .



What about the other way around?

## Theorem

*Let  $u \in V$  be a minimiser of  $J : V \rightarrow \mathbb{R}$ . Then  $u$  is a solution of*

$$J'(u; v) = 0 \text{ for all } v \in V.$$

## Proof.

Since  $u$  is a minimiser of  $J$ ,  $J(u + \varepsilon v) \geq J(u)$  for all  $\varepsilon > 0, v \in V$ . This implies

$$\frac{J(u + \varepsilon v) - J(u)}{\varepsilon} \geq 0.$$

Taking the limit as  $\varepsilon \rightarrow 0^+$  yields  $J'(u; v) \geq 0$ .

Replacing  $v$  with  $-v$ , we have  $J'(u; -v) = -J'(u; v) \geq 0$ , i.e.  $J'(u; v) \leq 0$ . So  $J'(u; v) = 0$ . □

The variational problem

find  $u \in V$  such that  $J'(u; v) = 0$  for all  $v \in V$

might be linear, nonlinear, coercive, or not! In general it describes *stationary points* of  $J$ , not just minimisers.

The variational problem

find  $u \in V$  such that  $J'(u; v) = 0$  for all  $v \in V$

might be linear, nonlinear, coercive, or not! In general it describes *stationary points* of  $J$ , not just minimisers.

Only symmetric problems can arise from energy minimisation. If a variational problem is nonsymmetric it doesn't enjoy this structure.

## Section 3

# Galerkin approximation

We saw in Lecture 1 that a finite element approximation is a *Galerkin* approximation: given a closed subspace  $V_h \subset V$ , we approximate the solution of

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V,$$

We saw in Lecture 1 that a finite element approximation is a *Galerkin* approximation: given a closed subspace  $V_h \subset V$ , we approximate the solution of

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V,$$

with the Galerkin approximation over  $V_h$

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

Assume that  $a$  is symmetric, coercive, and bounded. Then the variational problems can be restated as (for  $J(v) := \frac{1}{2}a(v, v) - F(v)$ ):

$$u = \operatorname{argmin}_{v \in V} J(v),$$

Assume that  $a$  is symmetric, coercive, and bounded. Then the variational problems can be restated as (for  $J(v) := \frac{1}{2}a(v, v) - F(v)$ ):

$$u = \operatorname{argmin}_{v \in V} J(v),$$

and

$$u_h = \operatorname{argmin}_{v_h \in V_h} J(v_h).$$



Assume that  $a$  is symmetric, coercive, and bounded. Then the variational problems can be restated as (for  $J(v) := \frac{1}{2}a(v, v) - F(v)$ ):

$$u = \operatorname{argmin}_{v \in V} J(v),$$

and

$$u_h = \operatorname{argmin}_{v_h \in V_h} J(v_h).$$

So

$$J(u) \leq J(u_h) \leq J(v_h) \text{ for all } v_h \in V_h!$$

Assume that  $a$  is symmetric, coercive, and bounded. Then the variational problems can be restated as (for  $J(v) := \frac{1}{2}a(v, v) - F(v)$ ):

$$u = \operatorname{argmin}_{v \in V} J(v),$$

and

$$u_h = \operatorname{argmin}_{v_h \in V_h} J(v_h).$$

So

$$J(u) \leq J(u_h) \leq J(v_h) \text{ for all } v_h \in V_h!$$

The finite element method gives you *the best approximation* in this energetic sense, when the problem has a nice quadratic energy functional.

Let's set

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\}$$

and look for minimisers of

$$J(u) = \frac{1}{2} \int_{\Omega} \nabla u \cdot \nabla u \, dx + \frac{1}{2} \int_{\Omega} u^2 \, dx - \int_{\Omega} f u \, dx - \int_{\Gamma_N} g u \, ds.$$

Let's set

$$V = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\}$$

and look for minimisers of

$$J(u) = \frac{1}{2} \int_{\Omega} \nabla u \cdot \nabla u \, dx + \frac{1}{2} \int_{\Omega} u^2 \, dx - \int_{\Omega} f u \, dx - \int_{\Gamma_N} g u \, ds.$$

Calculating, we find

$$J'(u; v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx + \int_{\Omega} u v \, dx - \int_{\Omega} f v \, dx - \int_{\Gamma_N} g v \, ds$$

Setting  $J'(u; v) = 0$ , we have a linear variational problem with

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx + \int_{\Omega} u v \, dx, \quad F(v) = \int_{\Omega} f v \, dx + \int_{\Gamma_N} g v \, ds.$$

This is called the (good) *Helmholtz* problem:

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx + \int_{\Omega} uv \, dx, \quad F(v) = \int_{\Omega} fv \, dx + \int_{\Gamma_N} gv \, ds.$$

This is called the (good) *Helmholtz* problem:

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx + \int_{\Omega} uv \, dx, \quad F(v) = \int_{\Omega} f v \, dx + \int_{\Gamma_N} g v \, ds.$$

In strong form, we have

$$\begin{aligned} -\nabla^2 u + u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \Gamma_D, \\ \nabla u \cdot n &= g \text{ on } \Gamma_N. \end{aligned}$$

Let's consider Lax–Milgram. Here  $a(u, v) = (u, v)_{H^1(\Omega)}$ , so the form is continuous by Cauchy–Schwarz with  $C = 1$ .

Let's consider Lax–Milgram. Here  $a(u, v) = (u, v)_{H^1(\Omega)}$ , so the form is continuous by Cauchy–Schwarz with  $C = 1$ .

Similarly, as

$$a(v, v) = (v, v)_{H^1(\Omega)} = \|v\|_{H^1(\Omega)}^2,$$

the problem is coercive with  $\alpha = 1$ .



Let's look at how to compute minimisers of this Helmholtz problem.

```
1  from firedrake import *
2
3  mesh = UnitSquareMesh(128, 128, quadrilateral=True)
4  (left, right) = (1, 2)
5  V = FunctionSpace(mesh, "CG", 1)
6  (x, y) = SpatialCoordinate(mesh)
7
8  f = Constant(1), g = Constant(0)
9  bc = DirichletBC(V, 0, (left, right))
10
11 u = Function(V)
12
13 J = (0.5 * inner(grad(u), grad(u))*dx + 0.5 * inner(u, u)
14      - inner(f, u)*dx - inner(g, u)*ds)
15 G = derivative(J, u, TestFunction(V))
16
17 solve(G == 0, u, bc)
```

By far the most important problem of this form is the equations of *linear elasticity*. The finite element method was invented by engineers looking to solve linear elasticity; its mathematical foundations and generalisations came much later.

By far the most important problem of this form is the equations of *linear elasticity*. The finite element method was invented by engineers looking to solve linear elasticity; its mathematical foundations and generalisations came much later.

Let  $\Omega \subset \mathbb{R}^n$  be an open bounded Lipschitz domain; its closure  $\bar{\Omega}$  is referred to as the reference configuration. We seek to characterise its shape upon loading via a mapping  $\phi : \bar{\Omega} \rightarrow \mathbb{R}^n$  via

$$\tilde{\Omega} = \phi(\bar{\Omega}).$$

By far the most important problem of this form is the equations of *linear elasticity*. The finite element method was invented by engineers looking to solve linear elasticity; its mathematical foundations and generalisations came much later.

Let  $\Omega \subset \mathbb{R}^n$  be an open bounded Lipschitz domain; its closure  $\bar{\Omega}$  is referred to as the reference configuration. We seek to characterise its shape upon loading via a mapping  $\phi : \bar{\Omega} \rightarrow \mathbb{R}^n$  via

$$\tilde{\Omega} = \phi(\bar{\Omega}).$$

It is useful to write the deformation  $\phi$  as the sum of the identity map plus a displacement:

$$\phi(x) = x + u(x),$$

where  $u(x) : \bar{\Omega} \rightarrow \mathbb{R}^n$ .

For an isotropic homogeneous body, the displacement  $u : \bar{\Omega} \rightarrow \mathbb{R}^n$  minimises the potential energy

$$J(u) = \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) + \lambda (\nabla \cdot u)^2 \, dx - \int_{\Omega} f \cdot u \, dx - \int_{\Gamma_N} g \cdot u \, ds.$$

Here  $f$  is the body loading (e.g. gravity),  $g$  is the surface traction, and

$$\varepsilon(u) = \frac{1}{2} \left( \nabla u + (\nabla u)^{\top} \right),$$

i.e. the symmetric part of the gradient (Jacobian)  $\nabla u : \Omega \rightarrow \mathbb{R}^{n \times n}$  of the displacement. The parameters  $\mu, \lambda > 0$  are material-dependent.

For an isotropic homogeneous body, the displacement  $u : \bar{\Omega} \rightarrow \mathbb{R}^n$  minimises the potential energy

$$J(u) = \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) + \lambda (\nabla \cdot u)^2 \, dx - \int_{\Omega} f \cdot u \, dx - \int_{\Gamma_N} g \cdot u \, ds.$$

Here  $f$  is the body loading (e.g. gravity),  $g$  is the surface traction, and

$$\varepsilon(u) = \frac{1}{2} \left( \nabla u + (\nabla u)^{\top} \right),$$

i.e. the symmetric part of the gradient (Jacobian)  $\nabla u : \Omega \rightarrow \mathbb{R}^{n \times n}$  of the displacement. The parameters  $\mu, \lambda > 0$  are material-dependent.

This yields a linear variational problem with

$$a(u, v) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) + \lambda \nabla \cdot u \nabla \cdot v \, dx, \quad F(v) = \int_{\Omega} f \cdot v \, dx + \int_{\Gamma_N} g \cdot v \, ds.$$

In strong form, we have

$$\begin{aligned} -2\mu \nabla \cdot \varepsilon(u) - \lambda \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \Gamma_D, \\ 2\mu \varepsilon(u) \cdot n + \lambda (\operatorname{tr} \varepsilon(u)) n &= g \text{ on } \Gamma_N. \end{aligned}$$

In strong form, we have

$$\begin{aligned} -2\mu \nabla \cdot \varepsilon(u) - \lambda \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \Gamma_D, \\ 2\mu \varepsilon(u) \cdot n + \lambda (\operatorname{tr} \varepsilon(u)) n &= g \text{ on } \Gamma_N. \end{aligned}$$

The bilinear form is continuous and coercive. In  $n$  dimensions, the continuity constant is

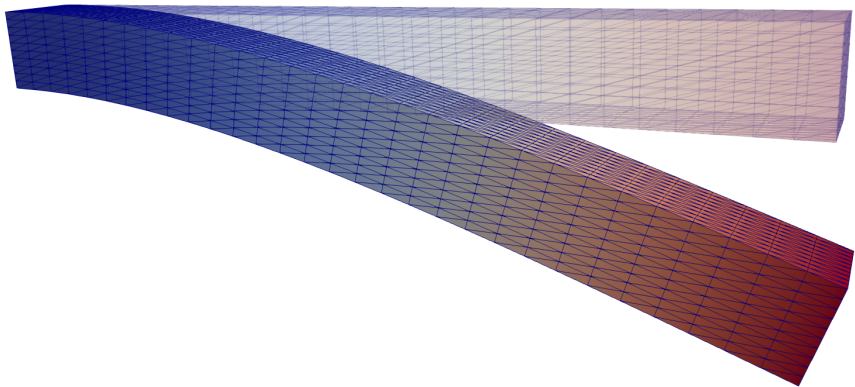
$$C = 2\mu + n\lambda.$$

Coercivity is guaranteed by *Korn's lemma*; if  $\Gamma_D = \partial\Omega$ , the coercivity constant is

$$\alpha = \mu.$$



```
1  from firedrake import *
2
3  mesh = BoxMesh(30, 10, 10, 10, 1, 1)
4  V = VectorFunctionSpace(mesh, "CG", 1)
5
6  g = Constant((0, 0, -5e7))
7  bc = DirichletBC(V, 0, 1)
8
9  u = Function(V)
10 (mu, lam) = (27.4e9, 64.0e9)
11
12 J = (0.5 * 2*mu * inner(sym(grad(u)), sym(grad(u)))*dx
13      + 0.5 * lam * inner(div(u), div(u)) * dx
14      - inner(g, u)*ds(2))
15 G = derivative(J, u, TestFunction(V))
16
17 solve(G == 0, u, bc)
```



## C6.4 Finite Element Methods for PDEs

### Lecture 7: Galerkin approximation

Patrick E. Farrell

University of Oxford

Given a linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V,$$

we form its Galerkin approximation over a closed subspace  $V_h \subset V$

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

Given a linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V,$$

we form its Galerkin approximation over a closed subspace  $V_h \subset V$

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

We first consider its approximation properties over *arbitrary* subspaces  $V_h$ , then in subsequent lectures consider  $V_h$  constructed via finite elements.

## Corollary

*Let  $a$  and  $F$  satisfy the hypothesis of the Lax–Milgram Theorem. Then the Galerkin approximation is well-posed for any closed subspace  $V_h \subset V$ .*

## Corollary

*Let  $a$  and  $F$  satisfy the hypothesis of the Lax–Milgram Theorem. Then the Galerkin approximation is well-posed for any closed subspace  $V_h \subset V$ .*

## Proof.

As  $V_h \subset V$ ,  $a : V_h \times V_h \rightarrow \mathbb{R}$  is bounded and coercive on  $V_h$ , with the same continuity and coercivity constants.

## Corollary

*Let  $a$  and  $F$  satisfy the hypothesis of the Lax–Milgram Theorem. Then the Galerkin approximation is well-posed for any closed subspace  $V_h \subset V$ .*

## Proof.

As  $V_h \subset V$ ,  $a : V_h \times V_h \rightarrow \mathbb{R}$  is bounded and coercive on  $V_h$ , with the same continuity and coercivity constants.  $F : V_h \rightarrow \mathbb{R}$  is linear and bounded.



## Corollary

*Let  $a$  and  $F$  satisfy the hypothesis of the Lax–Milgram Theorem. Then the Galerkin approximation is well-posed for any closed subspace  $V_h \subset V$ .*

## Proof.

As  $V_h \subset V$ ,  $a : V_h \times V_h \rightarrow \mathbb{R}$  is bounded and coercive on  $V_h$ , with the same continuity and coercivity constants.  $F : V_h \rightarrow \mathbb{R}$  is linear and bounded. Thus, by Lax–Milgram, the variational problem defining the Galerkin approximation is well-posed. □

## Corollary

*Let  $a$  and  $F$  satisfy the hypothesis of the Lax–Milgram Theorem. Then the Galerkin approximation is well-posed for any closed subspace  $V_h \subset V$ .*

## Proof.

As  $V_h \subset V$ ,  $a : V_h \times V_h \rightarrow \mathbb{R}$  is bounded and coercive on  $V_h$ , with the same continuity and coercivity constants.  $F : V_h \rightarrow \mathbb{R}$  is linear and bounded. Thus, by Lax–Milgram, the variational problem defining the Galerkin approximation is well-posed. □

For coercive problems, well-posedness is inherited. *This is not true for noncoercive problems.* This makes discretising noncoercive problems much harder.

Once we choose a basis  $\{\phi_i\}$  of  $V_h$ , the linear system is

$$Ax = b,$$

where

$$u_h = \sum_i x_i \phi_i, \quad b_i = F(\phi_i), \quad A_{ji} = a(\phi_i, \phi_j).$$

Once we choose a basis  $\{\phi_i\}$  of  $V_h$ , the linear system is

$$Ax = b,$$

where

$$u_h = \sum_i x_i \phi_i, \quad b_i = F(\phi_i), \quad A_{ji} = a(\phi_i, \phi_j).$$

The linear system we must solve for our Galerkin approximation inherits useful properties of the underlying problem.

Once we choose a basis  $\{\phi_i\}$  of  $V_h$ , the linear system is

$$Ax = b,$$

where

$$u_h = \sum_i x_i \phi_i, \quad b_i = F(\phi_i), \quad A_{ji} = a(\phi_i, \phi_j).$$

The linear system we must solve for our Galerkin approximation inherits useful properties of the underlying problem.

If  $a$  is symmetric, so is  $A$ :

$$A_{ji} = a(\phi_i, \phi_j) = a(\phi_j, \phi_i) = A_{ij}.$$

Once we choose a basis  $\{\phi_i\}$  of  $V_h$ , the linear system is

$$Ax = b,$$

where

$$u_h = \sum_i x_i \phi_i, \quad b_i = F(\phi_i), \quad A_{ji} = a(\phi_i, \phi_j).$$

The linear system we must solve for our Galerkin approximation inherits useful properties of the underlying problem.

If  $a$  is symmetric, so is  $A$ :

$$A_{ji} = a(\phi_i, \phi_j) = a(\phi_j, \phi_i) = A_{ij}.$$

If  $a$  is coercive (hence positive-definite), so is  $A$ :

$$c^\top A c = a \left( \sum_i c_i \phi_i, \sum_i c_i \phi_i \right) \geq 0.$$

We know that the solution  $u$  satisfies

$$a(u, v) = F(v) \quad \text{for all } v \in V,$$

and thus in particular

$$a(u, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

We know that the solution  $u$  satisfies

$$a(u, v) = F(v) \quad \text{for all } v \in V,$$

and thus in particular

$$a(u, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

The Galerkin approximation  $u_h \in V_h$  satisfies

$$a(u_h, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$



We know that the solution  $u$  satisfies

$$a(u, v) = F(v) \quad \text{for all } v \in V,$$

and thus in particular

$$a(u, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

The Galerkin approximation  $u_h \in V_h$  satisfies

$$a(u_h, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

Subtracting, we find

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

We know that the solution  $u$  satisfies

$$a(u, v) = F(v) \quad \text{for all } v \in V,$$

and thus in particular

$$a(u, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

The Galerkin approximation  $u_h \in V_h$  satisfies

$$a(u_h, v_h) = F(v_h) \quad \text{for all } v_h \in V_h \subset V.$$

Subtracting, we find

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

This is called *Galerkin orthogonality*.

Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For any  $v_h \in V_h$ ,

$$\alpha \|u - u_h\|_V^2 \leq a(u - u_h, u - u_h)$$



Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For any  $v_h \in V_h$ ,

$$\begin{aligned} \alpha \|u - u_h\|_V^2 &\leq a(u - u_h, u - u_h) \\ &= a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h) \end{aligned}$$



Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For any  $v_h \in V_h$ ,

$$\begin{aligned} \alpha \|u - u_h\|_V^2 &\leq a(u - u_h, u - u_h) \\ &= a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h) \\ &= a(u - u_h, u - v_h) \end{aligned}$$



Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For any  $v_h \in V_h$ ,

$$\begin{aligned} \alpha \|u - u_h\|_V^2 &\leq a(u - u_h, u - u_h) \\ &= a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h) \\ &= a(u - u_h, u - v_h) \\ &\leq C \|u - u_h\|_V \|u - v_h\|_V. \end{aligned}$$



Let's assume that  $a$  is coercive and bounded, but not symmetric.

## Lemma (Céa's Lemma)

*The Galerkin approximation  $u_h \in V_h$  to  $u \in V$  is quasi-optimal, in that it satisfies*

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For any  $v_h \in V_h$ ,

$$\begin{aligned} \alpha \|u - u_h\|_V^2 &\leq a(u - u_h, u - u_h) \\ &= a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h) \\ &= a(u - u_h, u - v_h) \\ &\leq C \|u - u_h\|_V \|u - v_h\|_V. \end{aligned}$$

Dividing by  $\alpha$  and minimising over  $v_h \in V$ , we obtain the result. □



## Remark

This quasi-optimality result relates (the error in the PDE approximation) with (the approximating power of the space  $V_h$ ). This decouples the error analysis from the specific PDE and turns the focus to constructing  $V_h$  with good approximation properties.

## Remark

This quasi-optimality result relates (the error in the PDE approximation) with (the approximating power of the space  $V_h$ ). This decouples the error analysis from the specific PDE and turns the focus to constructing  $V_h$  with good approximation properties.

This leads to the question: given  $u \in V$ , what is

$$\min_{v_h \in V_h} \|u - v_h\|_V?$$

In the finite element context, the answer will depend on the smoothness of  $u$ , the mesh size  $h$ , and the polynomial degree  $p$ .

## Remark

This quasi-optimality result relates (the error in the PDE approximation) with (the approximating power of the space  $V_h$ ). This decouples the error analysis from the specific PDE and turns the focus to constructing  $V_h$  with good approximation properties.

This leads to the question: given  $u \in V$ , what is

$$\min_{v_h \in V_h} \|u - v_h\|_V?$$

In the finite element context, the answer will depend on the smoothness of  $u$ , the mesh size  $h$ , and the polynomial degree  $p$ .

## Remark

The ratio  $C/\alpha$  is crucial. If  $C/\alpha = 5$ , things are fine. But if  $C/\alpha = 1000$ , our discretisation won't be very useful.

Now let's also assume that  $a$  is symmetric.

Now let's also assume that  $a$  is symmetric.

Recall that  $a$  defines a norm  $\|v\|_a := \sqrt{a(v, v)}$  on  $V$ , with

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

where the continuity and coercivity constants are measured in the  $V$  norm.

Now let's also assume that  $a$  is symmetric.

Recall that  $a$  defines a norm  $\|v\|_a := \sqrt{a(v, v)}$  on  $V$ , with

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

where the continuity and coercivity constants are measured in the  $V$  norm.

When we measure the continuity and coercivity constants *in the energy norm*, we get that  $C = 1$  (by Cauchy–Schwarz) and  $\alpha = 1$  (by definition).

Apply Céa's Lemma in the energy norm:

$$\begin{aligned}\|u - u_h\|_a &\leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_a \\ &= \min_{v_h \in V_h} \|u - v_h\|_a.\end{aligned}$$

Apply Céa's Lemma in the energy norm:

$$\begin{aligned}\|u - u_h\|_a &\leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_a \\ &= \min_{v_h \in V_h} \|u - v_h\|_a.\end{aligned}$$

Since  $u_h \in V_h$ , we must have equality, and thus *the error is optimal in the norm induced by the problem*:

$$\|u - u_h\|_a = \min_{v_h \in V_h} \|u - v_h\|_a.$$



Apply Céa's Lemma in the energy norm:

$$\begin{aligned}\|u - u_h\|_a &\leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_a \\ &= \min_{v_h \in V_h} \|u - v_h\|_a.\end{aligned}$$

Since  $u_h \in V_h$ , we must have equality, and thus *the error is optimal in the norm induced by the problem*:

$$\|u - u_h\|_a = \min_{v_h \in V_h} \|u - v_h\|_a.$$

The Galerkin approximation  $u_h$  is the *projection* of  $u$  onto  $V_h$  in the  $a$ -inner product!

What if we want to measure our error for a symmetric problem in the  $V$ -norm?

What if we want to measure our error for a symmetric problem in the  $V$ -norm?

Using the equivalences

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

we have

$$\|u - u_h\|_V \leq \frac{1}{\sqrt{\alpha}} \|u - u_h\|_a$$

What if we want to measure our error for a symmetric problem in the  $V$ -norm?

Using the equivalences

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

we have

$$\begin{aligned} \|u - u_h\|_V &\leq \frac{1}{\sqrt{\alpha}} \|u - u_h\|_a \\ &= \frac{1}{\sqrt{\alpha}} \min_{v_h \in V_h} \|u - v_h\|_a \end{aligned}$$

What if we want to measure our error for a symmetric problem in the  $V$ -norm?

Using the equivalences

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

we have

$$\begin{aligned} \|u - u_h\|_V &\leq \frac{1}{\sqrt{\alpha}} \|u - u_h\|_a \\ &= \frac{1}{\sqrt{\alpha}} \min_{v_h \in V_h} \|u - v_h\|_a \\ &\leq \sqrt{\frac{C}{\alpha}} \min_{v_h \in V_h} \|u - v_h\|_V \end{aligned}$$

What if we want to measure our error for a symmetric problem in the  $V$ -norm?

Using the equivalences

$$\alpha \|v\|_V^2 \leq \|v\|_a^2 \leq C \|v\|_V^2,$$

we have

$$\begin{aligned} \|u - u_h\|_V &\leq \frac{1}{\sqrt{\alpha}} \|u - u_h\|_a \\ &= \frac{1}{\sqrt{\alpha}} \min_{v_h \in V_h} \|u - v_h\|_a \\ &\leq \sqrt{\frac{C}{\alpha}} \min_{v_h \in V_h} \|u - v_h\|_V \end{aligned}$$

so we improve the constant of quasi-optimality by a square root!

## Section 5

# Linear elasticity: the nearly incompressible case

Let's consider linear elasticity again for a problem on  $\Omega \subset \mathbb{R}^n$ . We seek

$$u = \operatorname{argmin}_{v \in [H_0^1(\Omega)]^n} \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) + \lambda (\nabla \cdot u)^2 \, dx - \int_{\Omega} f \cdot u \, dx.$$



Let's consider linear elasticity again for a problem on  $\Omega \subset \mathbb{R}^n$ . We seek

$$u = \operatorname{argmin}_{v \in [H_0^1(\Omega)]^n} \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) + \lambda (\nabla \cdot u)^2 \, dx - \int_{\Omega} f \cdot u \, dx.$$

When engineers implemented the finite element method for this problem, they observed something puzzling: it worked well for steel and concrete, but did not work for rubber. Why not?

Let's consider linear elasticity again for a problem on  $\Omega \subset \mathbb{R}^n$ . We seek

$$u = \operatorname{argmin}_{v \in [H_0^1(\Omega)]^n} \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) + \lambda (\nabla \cdot u)^2 \, dx - \int_{\Omega} f \cdot u \, dx.$$

When engineers implemented the finite element method for this problem, they observed something puzzling: it worked well for steel and concrete, but did not work for rubber. Why not?

We claimed in Lecture 6 that its coercivity and continuity constants are

$$C = 2\mu + n\lambda, \quad \alpha = \mu.$$

Let's investigate the practical consequences of this.

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55
rubber A	0.018	0.9	2.75	0.018	12.2

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55
rubber A	0.018	0.9	2.75	0.018	12.2
rubber B	0.018	9.0	27.5	0.018	38.7

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55
rubber A	0.018	0.9	2.75	0.018	12.2
rubber B	0.018	9.0	27.5	0.018	38.7
rubber C	0.018	90	275	0.018	122.4

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55
rubber A	0.018	0.9	2.75	0.018	12.2
rubber B	0.018	9.0	27.5	0.018	38.7
rubber C	0.018	90	275	0.018	122.4

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

As  $\lambda \rightarrow \infty$ , the Galerkin approximation breaks down. This is because  $C = C(\mu, \lambda)$ , while  $\alpha = \alpha(\mu)$ . This is called *locking*.



Consider the following different materials for  $n = 3$ :

Material	$\mu$	$\lambda$	$C$	$\alpha$	$\sqrt{C/\alpha}$
steel	75	112	486	75	2.55
concrete	18	27	117	18	2.55
rubber A	0.018	0.9	2.75	0.018	12.2
rubber B	0.018	9.0	27.5	0.018	38.7
rubber C	0.018	90	275	0.018	122.4

**Table:** Lamé parameters for different materials. All units for  $\mu, \lambda, C, \alpha$  are multiplied by  $10^9$ .

As  $\lambda \rightarrow \infty$ , the Galerkin approximation breaks down. This is because  $C = C(\mu, \lambda)$ , while  $\alpha = \alpha(\mu)$ . This is called *locking*.

The parameter  $\lambda$  penalises  $\|\nabla \cdot u\|_{L^2(\Omega)}^2$ ; as  $\lambda \rightarrow \infty$ , the displacement is not allowed to change the volume. The different rubber samples are becoming *nearly incompressible*.

Note that this problem of locking is not specific to any particular discretisation; it is that *the formulation of the problem is becoming ill-conditioned*.

Note that this problem of locking is not specific to any particular discretisation; it is that *the formulation of the problem is becoming ill-conditioned*.

What can we do?

Note that this problem of locking is not specific to any particular discretisation; it is that *the formulation of the problem is becoming ill-conditioned*.

What can we do? Use a different formulation!

We introduce an auxiliary variable

$$p = \lambda \nabla \cdot u$$

which in weak form becomes

$$\int_{\Omega} q \nabla \cdot u \, dx = \frac{1}{\lambda} \int_{\Omega} p q \, dx.$$

We introduce an auxiliary variable

$$p = \lambda \nabla \cdot u$$

which in weak form becomes

$$\int_{\Omega} q \nabla \cdot u \, dx = \frac{1}{\lambda} \int_{\Omega} p q \, dx.$$

We then consider: find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\begin{aligned} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ -\frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx &= 0, \end{aligned}$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

We introduce an auxiliary variable

$$p = \lambda \nabla \cdot u$$

which in weak form becomes

$$\int_{\Omega} q \nabla \cdot u \, dx = \frac{1}{\lambda} \int_{\Omega} p q \, dx.$$

We then consider: find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\begin{aligned} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ -\frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx &= 0, \end{aligned}$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

This remains uniformly well-posed as  $\lambda \rightarrow \infty$ , even for  $\lambda = \infty$ !

Does the mixed formulation have an energetic structure?

$$\begin{aligned} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx \\ -\frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx &= 0 \end{aligned}$$



Does the mixed formulation have an energetic structure?

$$\begin{aligned} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx \\ -\frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx &= 0 \end{aligned}$$

are the Euler–Lagrange stationarity conditions for the Lagrangian

$$L(u, p) = \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) \, dx + \int_{\Omega} p \nabla \cdot u \, dx - \frac{1}{2\lambda} \int_{\Omega} p^2 \, dx - \int_{\Omega} f \cdot u \, dx.$$

This becomes the familiar linear elastic energy under  $p \mapsto \lambda \nabla \cdot u$ .

Does the mixed formulation have an energetic structure?

$$\begin{aligned} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx \\ -\frac{1}{\lambda} \int_{\Omega} pq \, dx + \int_{\Omega} q \nabla \cdot u \, dx &= 0 \end{aligned}$$

are the Euler–Lagrange stationarity conditions for the Lagrangian

$$L(u, p) = \frac{1}{2} \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(u) \, dx + \int_{\Omega} p \nabla \cdot u \, dx - \frac{1}{2\lambda} \int_{\Omega} p^2 \, dx - \int_{\Omega} f \cdot u \, dx.$$

This becomes the familiar linear elastic energy under  $p \mapsto \lambda \nabla \cdot u$ .

We will see later that the solution  $(u, p)$  is a saddle point of  $L$ .

Note that the bilinear form

$$B(u, p; v, q) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx + \int_{\Omega} q \nabla \cdot u \, dx - \frac{1}{\lambda} \int_{\Omega} pq \, dx$$

Note that the bilinear form

$$B(u, p; v, q) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx + \int_{\Omega} q \nabla \cdot u \, dx - \frac{1}{\lambda} \int_{\Omega} pq \, dx$$

is not coercive:

$$B(0, p; 0, p) = -\frac{1}{\lambda} \|p\|_{L^2(\Omega)}^2 < 0.$$

Note that the bilinear form

$$B(u, p; v, q) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx + \int_{\Omega} q \nabla \cdot u \, dx - \frac{1}{\lambda} \int_{\Omega} pq \, dx$$

is not coercive:

$$B(0, p; 0, p) = -\frac{1}{\lambda} \|p\|_{L^2(\Omega)}^2 < 0.$$

We will see in subsequent lectures how to analyse such systems. We have exchanged a simple-to-discretise but unstable formulation for a stable but harder-to-discretise one.

```

1  from firedrake import *
2  mesh = BoxMesh(30, 10, 10, 10, 1, 1)
3  V = VectorFunctionSpace(mesh, "CG", 2)
4  Q = FunctionSpace(mesh, "CG", 1)
5  Z = V*Q
6
7  g = Constant((0, 0, -5e7))
8  bc = DirichletBC(Z.sub(0), 0, 1)
9
10 z = Function(Z)
11 (u, p) = split(z)
12 (mu, lam) = (27.4e9, 64.0e9)
13
14 L = (mu * inner(sym(grad(u)), sym(grad(u)))*dx
15      + p * div(u) * dx - 1/(2*lam) * p**2 * dx
16      - inner(g, u)*ds(2))
17 G = derivative(L, z, TestFunction(Z))
18 solve(G == 0, z, bc)

```

## C6.4 Finite Element Methods for PDEs

### Lecture 8: Constructing function spaces with finite elements

Patrick E. Farrell

University of Oxford

In the last lecture, we saw Céa's Lemma: for a coercive bounded  $a$ , the error in Galerkin approximation is bounded by

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$



In the last lecture, we saw Céa's Lemma: for a coercive bounded  $a$ , the error in Galerkin approximation is bounded by

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

How do we construct discrete spaces  $V_h \subset V$  with good approximating properties on general domains?

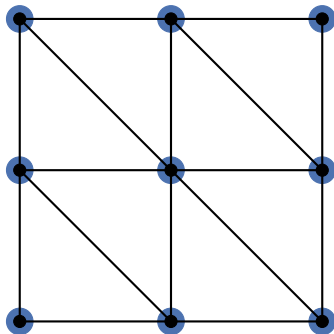
In the last lecture, we saw Céa's Lemma: for a coercive bounded  $a$ , the error in Galerkin approximation is bounded by

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V.$$

How do we construct discrete spaces  $V_h \subset V$  with good approximating properties on general domains?

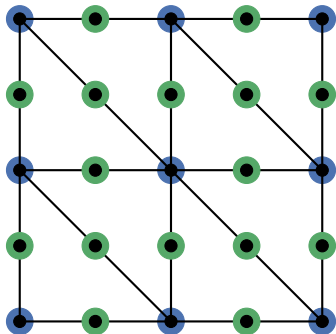
The finite element method!

Key idea: use piecewise polynomials on a mesh of  $\Omega$ .



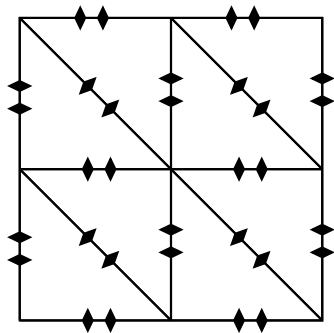
Data stored to represent a piecewise linear function,  $V_h \subset H^1(\Omega)$ .

Key idea: use piecewise polynomials on a mesh of  $\Omega$ .



Data stored to represent a piecewise quadratic function,  $V_h \subset H^1(\Omega)$ .

Key idea: use piecewise polynomials on a mesh of  $\Omega$ .



Data stored to represent a piecewise linear vector function,  $V_h \subset H(\text{div}, \Omega)$ .

Let's consider what happens on a single cell first, then stitch them together to enforce the continuity properties we need to conform to the Sobolev space we want.

Let's consider what happens on a single cell first, then stitch them together to enforce the continuity properties we need to conform to the Sobolev space we want.

### Definition (Finite element)

A finite element is a triple  $(K, \mathcal{V}, \mathcal{L})$  where

- ▶ The cell  $K$  is a bounded, closed subset of  $\mathbb{R}^n$  with nonempty connected interior and piecewise smooth boundary;

Let's consider what happens on a single cell first, then stitch them together to enforce the continuity properties we need to conform to the Sobolev space we want.

### Definition (Finite element)

A finite element is a triple  $(K, \mathcal{V}, \mathcal{L})$  where

- ▶ The cell  $K$  is a bounded, closed subset of  $\mathbb{R}^n$  with nonempty connected interior and piecewise smooth boundary;
- ▶ The space  $\mathcal{V} = \mathcal{V}(K)$  is a finite dimensional function space on  $K$  of dimension  $d$ ;

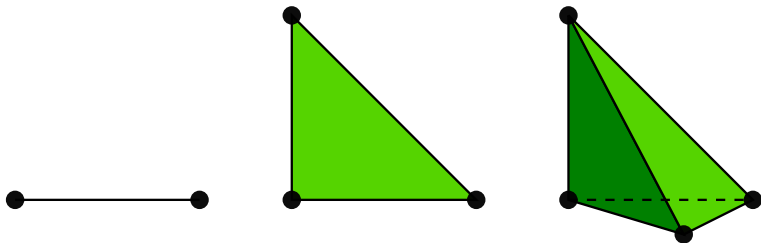


Let's consider what happens on a single cell first, then stitch them together to enforce the continuity properties we need to conform to the Sobolev space we want.

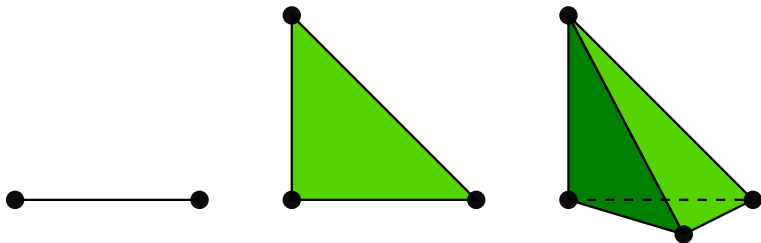
### Definition (Finite element)

A finite element is a triple  $(K, \mathcal{V}, \mathcal{L})$  where

- ▶ The cell  $K$  is a bounded, closed subset of  $\mathbb{R}^n$  with nonempty connected interior and piecewise smooth boundary;
- ▶ The space  $\mathcal{V} = \mathcal{V}(K)$  is a finite dimensional function space on  $K$  of dimension  $d$ ;
- ▶ The set of degrees of freedom  $\mathcal{L} = \{\ell_1, \dots, \ell_d\}$  is a basis for  $\mathcal{V}^*$ , the dual space of  $\mathcal{V}$ .



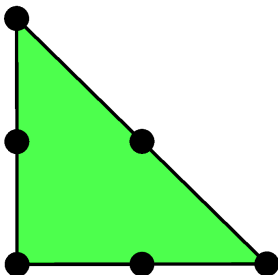
The linear Lagrange finite element  $CG_1$  in one, two and three dimensions. The black circles denote pointwise evaluation.



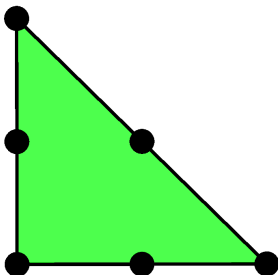
The linear Lagrange finite element  $CG_1$  in one, two and three dimensions. The black circles denote pointwise evaluation.

### Example (2D)

$$K = \triangle, \mathcal{V} = \text{span}(1, x, y), \mathcal{L} = \{\ell_1, \ell_2, \ell_3\}, \ell_i : v \mapsto v(x_i).$$



The quadratic Lagrange finite element  $CG_2$  in two dimensions.



The quadratic Lagrange finite element  $\text{CG}_2$  in two dimensions.

### Example

$K = \triangle$ ,  $\mathcal{V} = \text{span}(1, x, y, x^2, y^2, xy)$ ,  $\mathcal{L} = \{\ell_1, \dots, \ell_6\}$ , each  $\ell_i$  evaluates the function at a vertex or edge midpoint.

## Definition (Polynomial spaces)

Denote the space of polynomials of *total* degree  $q$  on  $K \subset \mathbb{R}^n$ :

$$\mathcal{P}_q(K) = \text{span} \left\{ x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n} |_K : \sum_{i=1}^n \alpha_i \leq q, \alpha_i \geq 0 \text{ for all } i = 1, \dots, n \right\},$$

of *maximal* degree  $q$  on  $K$ :

$$\mathcal{Q}_q(K) = \text{span} \left\{ x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n} |_K : \alpha_i \leq q, \alpha_i \geq 0 \text{ for all } i = 1, \dots, n \right\},$$

## Definition (Polynomial spaces)

Denote the space of polynomials of *total* degree  $q$  on  $K \subset \mathbb{R}^n$ :

$$\mathcal{P}_q(K) = \text{span} \left\{ x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n} |_K : \sum_{i=1}^n \alpha_i \leq q, \alpha_i \geq 0 \text{ for all } i = 1, \dots, n \right\},$$

of *maximal* degree  $q$  on  $K$ :

$$\mathcal{Q}_q(K) = \text{span} \left\{ x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n} |_K : \alpha_i \leq q, \alpha_i \geq 0 \text{ for all } i = 1, \dots, n \right\},$$

## Example

For  $\triangle, \square \subset \mathbb{R}^2$ ,

$$\mathcal{P}_2(\triangle) = \text{span}\{1, x, x^2, y, y^2, xy\},$$

$$\mathcal{Q}_2(\square) = \text{span}\{1, x, x^2, y, y^2, xy, x^2y, xy^2, x^2y^2\}.$$

## Intuition

The degrees of freedom are what we need to store to specify a particular  $v \in \mathcal{V}$ .



## Intuition

The degrees of freedom are what we need to store to specify a particular  $v \in \mathcal{V}$ .

For fixed  $x$ , consider the functional  $\ell_x(v) = v(x)$ . Since  $\ell_x \in \mathcal{V}^*$ , we can write

$$\ell_x = \alpha_1 \ell_1 + \cdots + \alpha_d \ell_d$$

for some coefficients  $\alpha$ . So if we know  $\ell_i(v)$ ,  $i = 1, \dots, d$ , then we know  $v(x)$  at every  $x \in K$ .

## Intuition

The degrees of freedom are what we need to store to specify a particular  $v \in \mathcal{V}$ .

For fixed  $x$ , consider the functional  $\ell_x(v) = v(x)$ . Since  $\ell_x \in \mathcal{V}^*$ , we can write

$$\ell_x = \alpha_1 \ell_1 + \cdots + \alpha_d \ell_d$$

for some coefficients  $\alpha$ . So if we know  $\ell_i(v)$ ,  $i = 1, \dots, d$ , then we know  $v(x)$  at every  $x \in K$ .

The main work in verifying that something is a finite element is in checking that  $\mathcal{L}$  is indeed a basis for  $\mathcal{V}^*$ .

## Lemma (Verifying finite elements)

*Let  $\mathcal{V}$  be a  $d$ -dimensional vector space and let  $\mathcal{L} = \{\ell_1, \dots, \ell_d\}$  be a subset of the dual space  $\mathcal{V}^*$ . Then the following two statements are equivalent:*

- (a)  $\mathcal{L}$  is a basis for  $\mathcal{V}^*$ ;*
- (b) Given  $v \in \mathcal{V}$  with  $\ell_i(v) = 0$  for  $i = 1, \dots, d$ , then  $v \equiv 0$ .*

## Lemma (Verifying finite elements)

*Let  $\mathcal{V}$  be a  $d$ -dimensional vector space and let  $\mathcal{L} = \{\ell_1, \dots, \ell_d\}$  be a subset of the dual space  $\mathcal{V}^*$ . Then the following two statements are equivalent:*

- (a)  $\mathcal{L}$  is a basis for  $\mathcal{V}^*$ ;*
- (b) Given  $v \in \mathcal{V}$  with  $\ell_i(v) = 0$  for  $i = 1, \dots, d$ , then  $v \equiv 0$ .*

This means that we just need to verify condition (b), which is much easier; we set the degrees of freedom to be zero and show that the only element of  $\mathcal{V}$  that satisfies this is the zero element.

## Lemma (Verifying finite elements)

Let  $\mathcal{V}$  be a  $d$ -dimensional vector space and let  $\mathcal{L} = \{\ell_1, \dots, \ell_d\}$  be a subset of the dual space  $\mathcal{V}^*$ . Then the following two statements are equivalent:

- (a)  $\mathcal{L}$  is a basis for  $\mathcal{V}^*$ ;
- (b) Given  $v \in \mathcal{V}$  with  $\ell_i(v) = 0$  for  $i = 1, \dots, d$ , then  $v \equiv 0$ .

This means that we just need to verify condition (b), which is much easier; we set the degrees of freedom to be zero and show that the only element of  $\mathcal{V}$  that satisfies this is the zero element.

## Definition

We say that  $\mathcal{L}$  *determines*  $\mathcal{V}$  if given  $v \in \mathcal{V}$ ,  $\ell_i(v) = 0 \ \forall i \implies v = 0$ . We also say that  $\mathcal{L}$  is *unisolvent*.

## Example

For  $\text{CG}_1(\triangle)$ , if  $v$  is zero at each vertex, then  $v$  must be zero everywhere as a plane is uniquely determined by its values at three non-collinear points. Thus, the linear Lagrange element on a triangle is indeed a finite element.

Having fixed  $\mathcal{L}$ , the usual choice for a basis of  $\mathcal{V}$  is the *nodal basis*.

### Definition (nodal basis)

The basis  $\{\phi_1, \dots, \phi_d\}$  of  $\mathcal{V}$  dual to  $\mathcal{L}$ , i.e. with the property that

$$\ell_i(\phi_j) = \delta_{ij}$$

is called the *nodal basis* for  $\mathcal{V}$ .

### Example (CG<sub>1</sub> in one dimension)

Let  $K = [0, 1]$ ,  $\mathcal{V} = \mathcal{P}_1(K)$ , and  $\mathcal{L}$  be pointwise evaluation at the endpoints. Then the nodal basis is given by

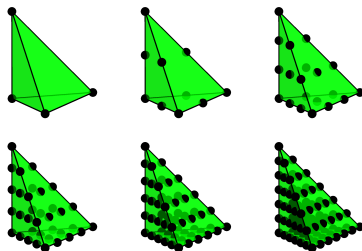
$$\phi_1(x) = 1 - x, \quad \phi_2(x) = x.$$

### Example ( $\text{CG}_1$ in two dimensions)

Let  $K$  be the triangle with vertices at  $(0, 0)$ ,  $(0, 1)$ ,  $(1, 0)$ . Let  $\mathcal{V} = \mathcal{P}_1(K)$ , and  $\mathcal{L}$  be pointwise evaluation at the vertices. Then the nodal basis is given by

$$\phi_1(x) = 1 - x_1 - x_2, \quad \phi_2(x) = x_1, \quad \phi_3(x) = x_2.$$





The Lagrange  $\text{CG}_q$  elements on tetrahedra for  $q = 1, \dots, 6$ .

## Definition (Lagrange element on a simplex)

The Lagrange element  $\text{CG}_q$  of dimension  $n$  and degree  $q \geq 1$  is

- ▶  $K$  is an  $n$ -dimensional simplex (interval, triangle, tetrahedron),
- ▶  $\mathcal{V} = \mathcal{P}_q(K)$ ,
- ▶  $\ell_i : v \mapsto v(x_i)$ ,  $i = 1, \dots, f(q)$ ,

where  $x_i$ ,  $i = 1, \dots, f(q)$  is an enumeration of points in the element.

One of the main things we will do with finite elements is interpolate functions onto them.

One of the main things we will do with finite elements is interpolate functions onto them.

### Definition (Interpolant on an element)

Let  $(K, \mathcal{V}, \mathcal{L})$  be a finite element. For a suitable function space  $H$ , define the interpolant  $\mathcal{I}_K : H \rightarrow \mathcal{V}$  via

$$\begin{aligned}\mathcal{I}_K : u &\mapsto \mathcal{I}_K u \\ \ell_i(\mathcal{I}_K u) &= \ell_i(u) \quad \text{for all } \ell_i \in \mathcal{L}.\end{aligned}$$

That is, the interpolant matches the function being interpolated at the degrees of freedom.

One of the main things we will do with finite elements is interpolate functions onto them.

### Definition (Interpolant on an element)

Let  $(K, \mathcal{V}, \mathcal{L})$  be a finite element. For a suitable function space  $H$ , define the interpolant  $\mathcal{I}_K : H \rightarrow \mathcal{V}$  via

$$\begin{aligned}\mathcal{I}_K : u &\mapsto \mathcal{I}_K u \\ \ell_i(\mathcal{I}_K u) &= \ell_i(u) \quad \text{for all } \ell_i \in \mathcal{L}.\end{aligned}$$

That is, the interpolant matches the function being interpolated at the degrees of freedom.

In the nodal basis, the interpolation operator is particularly simple:

$$\mathcal{I}_K u = \sum_{i=1}^d \ell_i(u) \phi_i.$$

## Section 2

# Meshes and the local-to-global mapping

To define a global function space and basis

$$V_h = \text{span}\{\phi_1, \dots, \phi_N\} \subset V,$$

we need to decompose  $\Omega$  into cells, define a finite element on each, and then specify how the local function spaces are to be stitched together.

To define a global function space and basis

$$V_h = \text{span}\{\phi_1, \dots, \phi_N\} \subset V,$$

we need to decompose  $\Omega$  into cells, define a finite element on each, and then specify how the local function spaces are to be stitched together.

Assume that  $\Omega$  is polytopic, so that it can be decomposed into cells exactly. (Otherwise we have to worry about geometric approximation errors also.)

## Definition (mesh)

A *mesh*  $\mathcal{M}$  is a geometric decomposition of a domain  $\Omega$  into a finite set of *cells*  $\mathcal{M} = \{K_i\}$  such that

1.  $\cup_i K_i = \overline{\Omega}$ .



## Definition (mesh)

A *mesh*  $\mathcal{M}$  is a geometric decomposition of a domain  $\Omega$  into a finite set of *cells*  $\mathcal{M} = \{K_i\}$  such that

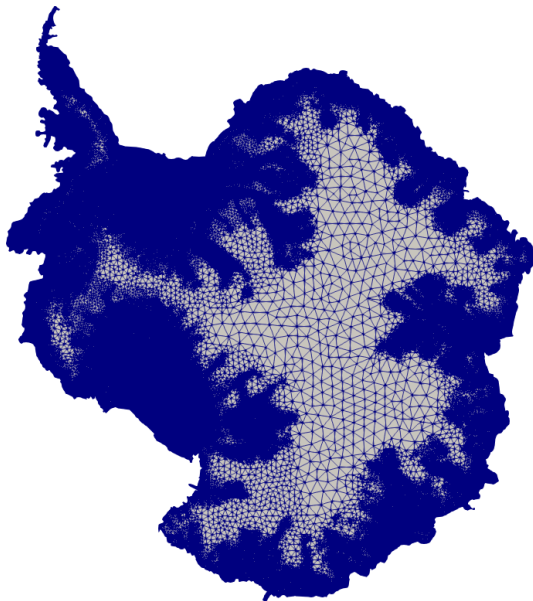
1.  $\cup_i K_i = \overline{\Omega}$ .
2. If  $K_i \cap K_j$  for  $i \neq j$  is exactly one point, it is a common vertex of  $K_i$  and  $K_j$ .

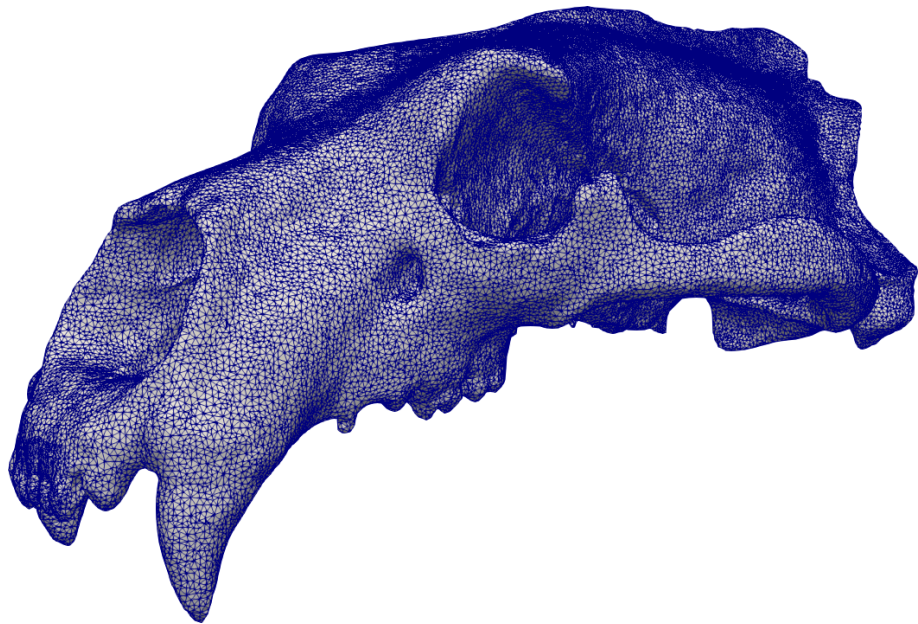
## Definition (mesh)

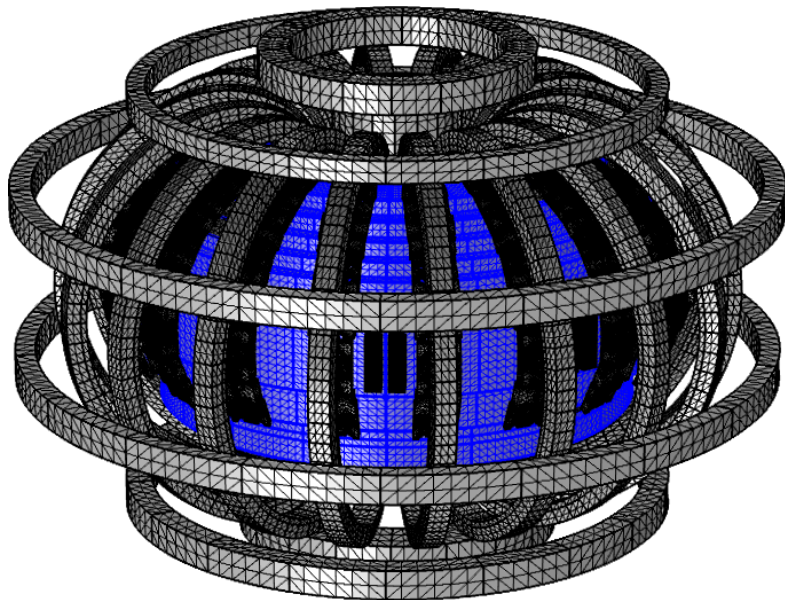
A *mesh*  $\mathcal{M}$  is a geometric decomposition of a domain  $\Omega$  into a finite set of *cells*  $\mathcal{M} = \{K_i\}$  such that

1.  $\cup_i K_i = \overline{\Omega}$ .
2. If  $K_i \cap K_j$  for  $i \neq j$  is exactly one point, it is a common vertex of  $K_i$  and  $K_j$ .
3. If  $K_i \cap K_j$  for  $i \neq j$  is not exactly one point, it is a common facet of  $K_i$  and  $K_j$  (edge in two dimensions, face in three dimensions).

Meshing is a huge subject of computational geometry in its own right.







We equip each cell  $K \in \mathcal{M}$  with a finite element, so we have a set of finite elements  $\{(K, \mathcal{V}_K, \mathcal{L}_K) : K \in \mathcal{M}\}$ . We'll assume we always equip each cell with the same type of element to gloss over technicalities.

We equip each cell  $K \in \mathcal{M}$  with a finite element, so we have a set of finite elements  $\{(K, \mathcal{V}_K, \mathcal{L}_K) : K \in \mathcal{M}\}$ . We'll assume we always equip each cell with the same type of element to gloss over technicalities.

We can thus give our first specification of a finite element space. Suppose we solve a variational problem over  $V$ . Then we take

$$V_h = \{v \in V : v|_K \in \mathcal{V}_K \ \forall \ K \in \mathcal{M}\}.$$

We equip each cell  $K \in \mathcal{M}$  with a finite element, so we have a set of finite elements  $\{(K, \mathcal{V}_K, \mathcal{L}_K) : K \in \mathcal{M}\}$ . We'll assume we always equip each cell with the same type of element to gloss over technicalities.

We can thus give our first specification of a finite element space. Suppose we solve a variational problem over  $V$ . Then we take

$$V_h = \{v \in V : v|_K \in \mathcal{V}_K \ \forall \ K \in \mathcal{M}\}.$$

Let's consider  $V = H^1(\Omega)$ . To enforce that  $V_h \subset H^1(\Omega)$ , we need to make sure that functions in the space are continuous. How do we do that?



We specify how the elements fit together with the *local-to-global mapping*. For each cell  $K \in \mathcal{M}$ , we must specify a local-to-global map

$$\iota_K : \{1, \dots, d(K)\} \rightarrow \{1, \dots, N\}$$

which specifies how the *local* degrees of freedom  $\ell_i^K(v)$  relate to the *global* degrees of freedom. Each local degree of freedom corresponds to a global degree of freedom, under the action of the local-to-global map:

$$\ell_{\iota_K(i)}(v) = \ell_i^K(v|_K), \quad i = 1, \dots, d(K).$$

We specify how the elements fit together with the *local-to-global mapping*. For each cell  $K \in \mathcal{M}$ , we must specify a local-to-global map

$$\iota_K : \{1, \dots, d(K)\} \rightarrow \{1, \dots, N\}$$

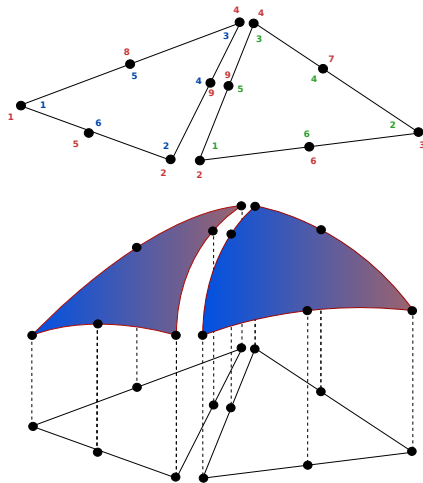
which specifies how the *local* degrees of freedom  $\ell_i^K(v)$  relate to the *global* degrees of freedom. Each local degree of freedom corresponds to a global degree of freedom, under the action of the local-to-global map:

$$\ell_{\iota_K(i)}(v) = \ell_i^K(v|_K), \quad i = 1, \dots, d(K).$$

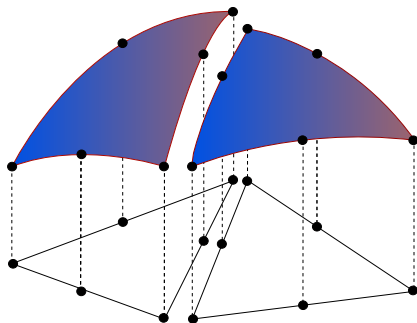
If two different degrees of freedom  $\ell^K, \ell^{K'}$  on two different cells  $K, K'$  both map to the same global degree of freedom, we demand

$$\ell^K(v|_K) = \ell^{K'}(v|_{K'}),$$

i.e. matching degrees of freedom agree.



The local-to-global mapping for a mesh of two triangles, both equipped with  $\text{CG}_2$ . By mapping matching local degrees of freedom at the common edge to the same global degree of freedom, the local-to-global map ensures the  $C^0$ -continuity of the approximation:  $V_h \subset H^1(\Omega)$ .



By *not* mapping matching local degrees of freedom at the common edge to the same global degree of freedom, a discontinuous approximation results:

$$V_h \subset L^2(\Omega).$$

The matching properties of the local-to-global map determine the global continuity of the function space, and hence which Sobolev space it *conforms* to.

### Definition (conforming approximation)

Suppose the continuous variational problem is posed over a Hilbert space  $V$ . If  $V_h \subset V$ , the approximation is *conforming*; if  $V_h \not\subset V$ , then the approximation is said to be nonconforming.

The matching properties of the local-to-global map determine the global continuity of the function space, and hence which Sobolev space it *conforms* to.

### Definition (conforming approximation)

Suppose the continuous variational problem is posed over a Hilbert space  $V$ . If  $V_h \subset V$ , the approximation is *conforming*; if  $V_h \not\subset V$ , then the approximation is said to be nonconforming.

In this course we will only consider conforming discretisations, although nonconforming ones are important, common, and useful.

Once we have the local-to-global map, we gather all of the global degrees of freedom

$$\mathcal{L} = \{\ell_1, \dots, \ell_N\} = \bigcup_{K \in \mathcal{M}} \{\ell_{\iota_K(i)}, i = 1, \dots, d(K)\}.$$

Once we have the local-to-global map, we gather all of the global degrees of freedom

$$\mathcal{L} = \{\ell_1, \dots, \ell_N\} = \bigcup_{K \in \mathcal{M}} \{\ell_{\iota_K(i)}, i = 1, \dots, d(K)\}.$$

and use its associated nodal basis our basis for  $V_h$ :

$$V_h = \text{span}\{\phi_1, \dots, \phi_N\}, \quad \ell_i(\phi_j) = \delta_{ij}, \quad \phi_i|_K \in \mathcal{V}_K \quad \forall K \in \mathcal{M}.$$

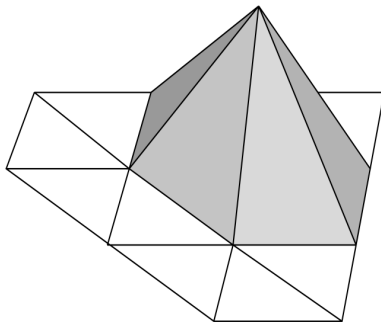


Once we have the local-to-global map, we gather all of the global degrees of freedom

$$\mathcal{L} = \{\ell_1, \dots, \ell_N\} = \bigcup_{K \in \mathcal{M}} \{\ell_{\iota_K(i)}, i = 1, \dots, d(K)\}.$$

and use its associated nodal basis as our basis for  $V_h$ :

$$V_h = \text{span}\{\phi_1, \dots, \phi_N\}, \quad \ell_i(\phi_j) = \delta_{ij}, \quad \phi_i|_K \in \mathcal{V}_K \quad \forall K \in \mathcal{M}.$$



Reproduced from Braess (2007).

Now that we have a global function space, we can construct a global interpolation operator.

### Definition (global interpolation operator)

Let  $V_h$  be a finite element function space constructed by equipping a mesh  $\mathcal{M}$  with finite elements. Then the interpolation operator  $\mathcal{I}_h : H \rightarrow V_h$  is defined by

$$(\mathcal{I}_h u)|_K = \mathcal{I}_K u,$$

and that  $\mathcal{I}_h u$  satisfies any necessary continuity requirements.

Now that we have a global function space, we can construct a global interpolation operator.

### Definition (global interpolation operator)

Let  $V_h$  be a finite element function space constructed by equipping a mesh  $\mathcal{M}$  with finite elements. Then the interpolation operator  $\mathcal{I}_h : H \rightarrow V_h$  is defined by

$$(\mathcal{I}_h u)|_K = \mathcal{I}_K u,$$

and that  $\mathcal{I}_h u$  satisfies any necessary continuity requirements.

### Example

If we take  $\text{CG}_q$  in 2D or 3D, we have  $V = H^1(\Omega)$  and  $H = H^2(\Omega)$ .

## C6.4 Finite Element Methods for PDEs

### Lecture 9: Local and global assembly

Patrick E. Farrell

University of Oxford

We are solving the Galerkin approximation

find  $u_h \in V_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in V_h$

over  $V_h := \text{span}\{\phi_1, \dots, \phi_N\}$ .

In this lecture we study the central algorithm executed by a finite element code, the *assembly* algorithm for computing  $A$  and  $b$  of

$$Ax = b.$$

We are solving the Galerkin approximation

find  $u_h \in V_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in V_h$

over  $V_h := \text{span}\{\phi_1, \dots, \phi_N\}$ .

In this lecture we study the central algorithm executed by a finite element code, the *assembly* algorithm for computing  $A$  and  $b$  of

$$Ax = b.$$

Recall that

$$A_{ji} = a(\phi_i, \phi_j), \quad b_j = F(\phi_j).$$

The naïve algorithm for assembly:

```
1: for  $i = 1, \dots, N$  do  
2:   for  $j = 1, \dots, N$  do  
3:     Compute  $A_{ji} = a(\phi_i, \phi_j)$ .  
4:   end for  
5: end for
```

has two major disadvantages:

The naïve algorithm for assembly:

```
1: for  $i = 1, \dots, N$  do  
2:   for  $j = 1, \dots, N$  do  
3:     Compute  $A_{ji} = a(\phi_i, \phi_j)$ .  
4:   end for  
5: end for
```

has two major disadvantages:

- ▶ Each  $\phi_i$  has *local* support. For most pairs  $i, j$ ,  $a(\phi_i, \phi_j) = 0$ .



The naïve algorithm for assembly:

```
1: for  $i = 1, \dots, N$  do  
2:   for  $j = 1, \dots, N$  do  
3:     Compute  $A_{ji} = a(\phi_i, \phi_j)$ .  
4:   end for  
5: end for
```

has two major disadvantages:

- ▶ Each  $\phi_i$  has *local* support. For most pairs  $i, j$ ,  $a(\phi_i, \phi_j) = 0$ .
- ▶ Each evaluation of  $a$  requires integrating over  $\Omega$ , i.e. a loop over cells. The calculations required to integrate over each cell are repeated many times.

A better idea:

- ▶ Loop over each cell of the mesh once.
- ▶ Calculate all contributions of each cell to all entries that it influences.

A better idea:

- ▶ Loop over each cell of the mesh once.
- ▶ Calculate all contributions of each cell to all entries that it influences.

Notation:  $(K, \mathcal{V}_K, \mathcal{L}_K)$ ,  $d = \dim(\mathcal{V}_K)$ ,  $\iota_K$ ,  $\phi_i^K = \phi_{\iota_K(i)}|_K$ .

A better idea:

- ▶ Loop over each cell of the mesh once.
- ▶ Calculate all contributions of each cell to all entries that it influences.

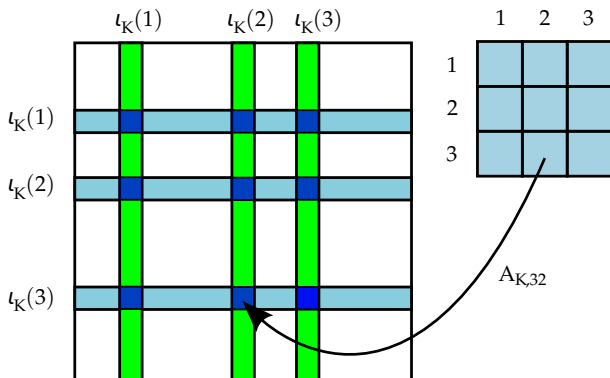
Notation:  $(K, \mathcal{V}_K, \mathcal{L}_K)$ ,  $d = \dim(\mathcal{V}_K)$ ,  $\iota_K$ ,  $\phi_i^K = \phi_{\iota_K(i)}|_K$ .

```

1: for  $K \in \mathcal{M}$  do
2:   Fetch the local-to-global map  $\iota_K$ .

3:   Compute the local tensor  $A_K$ :
4:   for  $i = 1, \dots, d$  do
5:     for  $j = 1, \dots, d$  do
6:       Compute  $(A_K)_{ji} = a(\phi_i^K, \phi_j^K)$  (only on the cell  $K$ ).
7:     end for
8:   end for

9:   Add the local tensor to the global tensor:
10:   $A_{\iota_K, \iota_K} \stackrel{+}{=} A_K$ 
11: end for
  
```



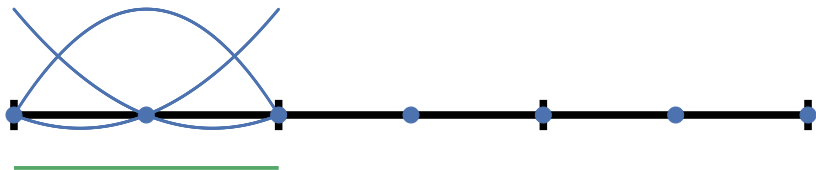
Finite element assembly. We loop over each cell  $K$  of the mesh and assemble the local stiffness matrix  $A_K$  (top right). We add this matrix to the submatrix of the global stiffness matrix  $A$  formed by taking the rows and columns associated with the local-to-global map  $l_K$ .

Assembly of a  $\text{CG}_2$  discretisation in 1D.



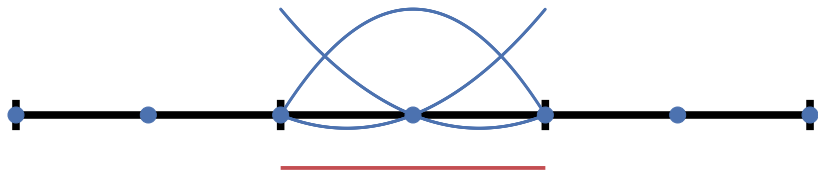
$$\begin{bmatrix} \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times \end{bmatrix}$$

# Assembly of a CG<sub>2</sub> discretisation in 1D.



$$\begin{bmatrix}
 \times & \times & \times & 0 & 0 & 0 & 0 \\
 \times & \times & \times & 0 & 0 & 0 & 0 \\
 \times & \times & \times & \times & \times & 0 & 0 \\
 0 & 0 & \times & \times & \times & 0 & 0 \\
 0 & 0 & \times & \times & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times
 \end{bmatrix}$$

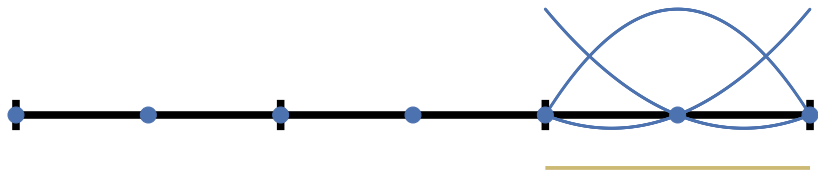
Assembly of a  $CG_2$  discretisation in 1D.



$$\begin{bmatrix}
 \times & \times & \times & 0 & 0 & 0 & 0 \\
 \times & \times & \times & 0 & 0 & 0 & 0 \\
 \times & \times & \times & \times & \times & 0 & 0 \\
 0 & 0 & \times & \times & \times & 0 & 0 \\
 0 & 0 & \times & \times & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times \\
 0 & 0 & 0 & 0 & \times & \times & \times
 \end{bmatrix}$$



Assembly of a  $CG_2$  discretisation in 1D.



$$\begin{bmatrix} \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times \end{bmatrix}$$

## Section 2

### Assembling the local tensor

How do we assemble the local tensor on a cell?

$$\int_K \phi_i(x) \phi_j(x) \, dx \text{ or } \int_K \nabla \phi_i(x) \cdot \nabla \phi_j(x) \, dx \text{ or } \dots$$

How do we assemble the local tensor on a cell?

$$\int_K \phi_i(x) \phi_j(x) \, dx \text{ or } \int_K \nabla \phi_i(x) \cdot \nabla \phi_j(x) \, dx \text{ or } \dots$$

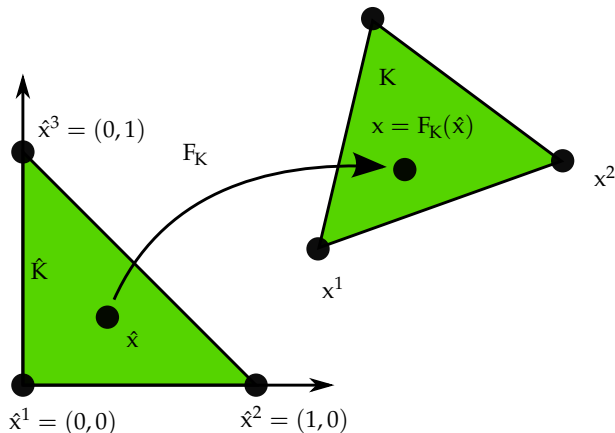
We calculate the integral with a *quadrature rule* on a *reference element*.

Introduce a reference element

$$(\hat{K}, \hat{\mathcal{V}}, \hat{\mathcal{L}})$$

and a set of diffeomorphisms

$$\{F_K : K \in \mathcal{M}\} \text{ such that } K = F_K(\hat{K}) \text{ for all } K \in \mathcal{M}.$$



For each  $K \in \mathcal{M}$ , the map  $F_K$  generates a function space on  $K$  via

$$\mathcal{V}(K) = \{v = \hat{v} \circ F_K^{-1} : \hat{v} \in \hat{\mathcal{V}}\},$$

and a set of degrees of freedom on  $K$  via

$$\mathcal{L}(K) = \{\ell(v) = \hat{\ell}(v \circ F_K) : \hat{\ell} \in \hat{\mathcal{L}}\}.$$

For each  $K \in \mathcal{M}$ , the map  $F_K$  generates a function space on  $K$  via

$$\mathcal{V}(K) = \{v = \hat{v} \circ F_K^{-1} : \hat{v} \in \hat{\mathcal{V}}\},$$

and a set of degrees of freedom on  $K$  via

$$\mathcal{L}(K) = \{\ell(v) = \hat{\ell}(v \circ F_K) : \hat{\ell} \in \hat{\mathcal{L}}\}.$$

By construction, we obtain a nodal basis for  $\mathcal{V}(K)$  from one on  $\hat{\mathcal{V}}$ . Suppose  $\{\hat{\phi}_i\}_{i=1}^d$  satisfies

$$\hat{\ell}_i(\hat{\phi}_j) = \delta_{ij}.$$

Define  $\phi_i^K = \hat{\phi}_i \circ F_K^{-1}$ . Computing, we find

$$\ell_j^K(\phi_i^K) = \hat{\ell}_j(\phi_i^K \circ F_K) = \hat{\ell}_j(\hat{\phi}_i \circ F_K^{-1} \circ F_K) = \hat{\ell}_j(\hat{\phi}_i) = \delta_{ij}.$$

For this simple approach of mapping to a reference element to work, we need

$$\mathcal{V}(K) = \mathcal{V}_K, \quad \mathcal{L}(K) = \mathcal{L}_K,$$

i.e. the finite element constructed via transformation is the same as that constructed directly on the cell.



For this simple approach of mapping to a reference element to work, we need

$$\mathcal{V}(K) = \mathcal{V}_K, \quad \mathcal{L}(K) = \mathcal{L}_K,$$

i.e. the finite element constructed via transformation is the same as that constructed directly on the cell.

This is true for Lagrange finite elements, and more complicated maps make it true for other finite elements we will meet.

Let's consider an example. Suppose we need to calculate

$$\int_K \phi_i(x) \phi_j(x) \, dx.$$

Let's consider an example. Suppose we need to calculate

$$\int_K \phi_i(x) \phi_j(x) \, dx.$$

We transform coordinates in the integral, at the cost of the determinant of the Jacobian of the mapping:

$$\int_K \phi_i(x) \phi_j(x) \, dx = \int_{\hat{K}} \hat{\phi}_i(\hat{x}) \hat{\phi}_j(\hat{x}) |J_K(\hat{x})| \, d\hat{x},$$

where  $J_K(\hat{x})$  is the Jacobian of  $F_K(\hat{x})$ .

Let's consider an example. Suppose we need to calculate

$$\int_K \phi_i(x) \phi_j(x) \, dx.$$

We transform coordinates in the integral, at the cost of the determinant of the Jacobian of the mapping:

$$\int_K \phi_i(x) \phi_j(x) \, dx = \int_{\hat{K}} \hat{\phi}_i(\hat{x}) \hat{\phi}_j(\hat{x}) |J_K(\hat{x})| \, d\hat{x},$$

where  $J_K(\hat{x})$  is the Jacobian of  $F_K(\hat{x})$ .

We then approximate the integral with a *quadrature rule*:

$$\int_{\hat{K}} f(\hat{x}) \, d\hat{x} \approx \sum_{i=1}^q w_i f(\hat{x}_i).$$

## Definition (quadrature rule of degree $m$ )

A quadrature rule over a cell  $\hat{K}$  is a choice of  $q$  quadrature points  $\hat{x}_i \in \hat{K}$  and weights  $w_i$  such that

$$\int_{\hat{K}} f(\hat{x}) \, d\hat{x} \approx \sum_{i=1}^q w_i f(\hat{x}_i).$$

It has *degree of precision* (or degree)  $m$  if the approximation is exact for polynomials of degree  $m$  or less.

## Definition (quadrature rule of degree $m$ )

A quadrature rule over a cell  $\hat{K}$  is a choice of  $q$  quadrature points  $\hat{x}_i \in \hat{K}$  and weights  $w_i$  such that

$$\int_{\hat{K}} f(\hat{x}) \, d\hat{x} \approx \sum_{i=1}^q w_i f(\hat{x}_i).$$

It has *degree of precision* (or degree)  $m$  if the approximation is exact for polynomials of degree  $m$  or less.

In 1D, Gaussian quadrature gives us the optimal choice of weights and quadrature points to maximise the degree of the rule. For  $q$  points in an interval we get  $m = 2q - 1$ . In higher dimensions things are not as simple, and the best quadrature rules are collated in an encyclopaedia.

To summarise. Offline calculations:

- ▶ Quadrature rule on the reference cell
- ▶ Basis functions at the quadrature points of the reference cell

Online calculations:

- ▶ Coordinate transformation & Jacobian

To summarise. Offline calculations:

- ▶ Quadrature rule on the reference cell
- ▶ Basis functions at the quadrature points of the reference cell

Online calculations:

- ▶ Coordinate transformation & Jacobian

In particular we do *not* need to calculate the basis functions for each cell.



Let's look at another example. Consider

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx,$$

which is what we need to calculate for solving Poisson.

Let's look at another example. Consider

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx,$$

which is what we need to calculate for solving Poisson.

Transforming, we get

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx = \int_{\hat{K}} \nabla_x \hat{\phi}_i(\hat{x}) \cdot \nabla_x \hat{\phi}_j(\hat{x}) |J_K(\hat{x})| \, d\hat{x},$$

but this is still not computable because it requires derivatives with respect to the physical coordinate.

Let's look at another example. Consider

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx,$$

which is what we need to calculate for solving Poisson.

Transforming, we get

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx = \int_{\hat{K}} \nabla_{\hat{x}} \hat{\phi}_i(\hat{x}) \cdot \nabla_{\hat{x}} \hat{\phi}_j(\hat{x}) |J_K(\hat{x})| \, d\hat{x},$$

but this is still not computable because it requires derivatives with respect to the physical coordinate.

To replace these, we apply the chain rule:

$$\frac{\partial \phi}{\partial x_k} = \sum_l \frac{\partial \hat{x}_l}{\partial x_k} \frac{\partial \phi}{\partial \hat{x}_l}.$$

Some calculation finds that

$$\nabla_x \hat{\phi}(\hat{x}) = J_K^{-\top}(\hat{x}) \nabla_{\hat{x}} \hat{\phi}(\hat{x}).$$

Thus, finally, we write

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx = \int_{\hat{K}} \left( J_K^{-\top} \nabla_{\hat{x}} \hat{\phi}_i(\hat{x}) \right) \cdot \left( J_K^{-\top} \nabla_{\hat{x}} \hat{\phi}_j(\hat{x}) \right) |J_K(\hat{x})| \, d\hat{x}.$$

Thus, finally, we write

$$\int_K \nabla_x \phi_i(x) \cdot \nabla_x \phi_j(x) \, dx = \int_{\hat{K}} \left( J_K^{-\top} \nabla_{\hat{x}} \hat{\phi}_i(\hat{x}) \right) \cdot \left( J_K^{-\top} \nabla_{\hat{x}} \hat{\phi}_j(\hat{x}) \right) |J_K(\hat{x})| \, d\hat{x}.$$

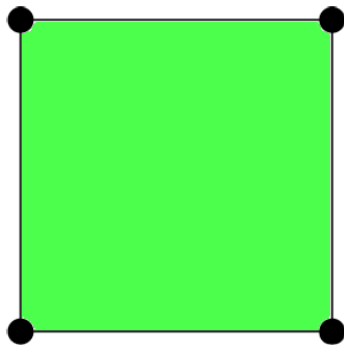
Remarks:

- ▶ Modern finite element software does this for you.
- ▶ We need to calculate the gradients of the basis functions at quadrature points (offline).

## Section 3

# Representing the element map

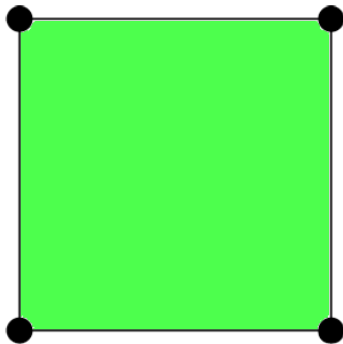
Prelude: some finite elements.



Definition (Lagrange element  $CG_1$  on a quadrilateral)

$K = \square$ ,  $\mathcal{V} = \mathcal{Q}_1(\square)$ ,  $\mathcal{L}$  evaluation at vertices.

Prelude: some finite elements.

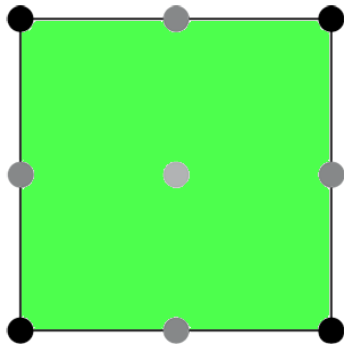


Definition (Vector Lagrange element  $[CG_1]^2$  on a quadrilateral)

$K = \square$ ,  $\mathcal{V} = \mathcal{Q}_1(\square)^2$ ,  $\mathcal{L}$  evaluation of each component at vertices.



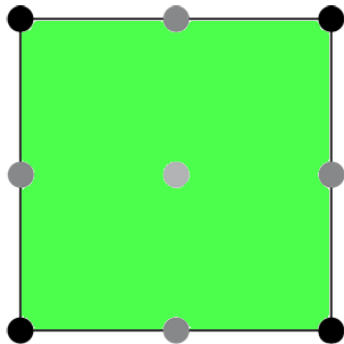
Prelude: some finite elements.



Definition (Lagrange element  $\text{CG}_2$  on a quadrilateral)

$K = \square$ ,  $\mathcal{V} = \mathcal{Q}_2(\square)$ ,  $\mathcal{L}$  evaluation at points shown.

Prelude: some finite elements.

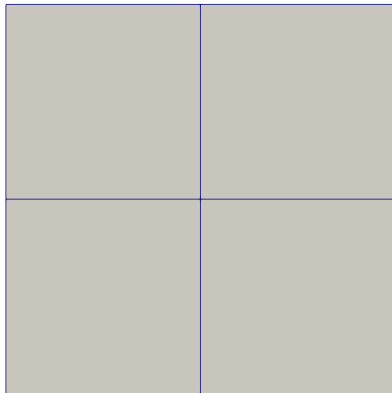


Definition (Vector Lagrange element  $[CG_2]^2$  on a quadrilateral)

$K = \square$ ,  $\mathcal{V} = \mathcal{Q}_2(\square)^2$ ,  $\mathcal{L}$  evaluation of each component at points shown.

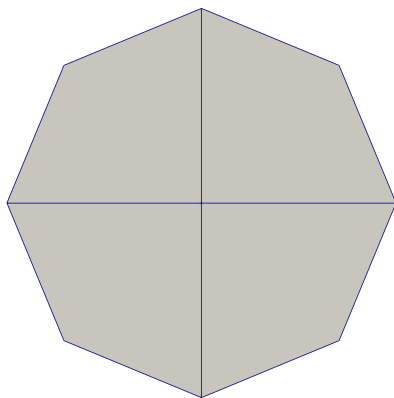
One can consider a purely *topological* mesh: sets of cells and vertices, with connectivity maps between them.

One can consider a purely *topological* mesh: sets of cells and vertices, with connectivity maps between them.



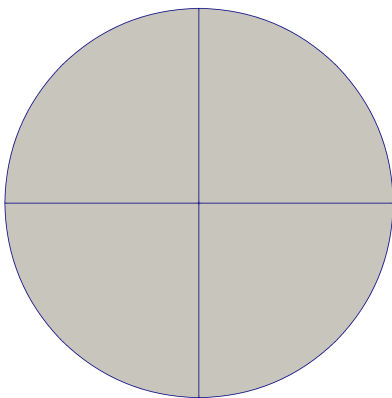
The *geometric* mesh of  $\Omega$  assigns coordinates to each cell  $([CG_1]^2)$ .

One can consider a purely *topological* mesh: sets of cells and vertices, with connectivity maps between them.

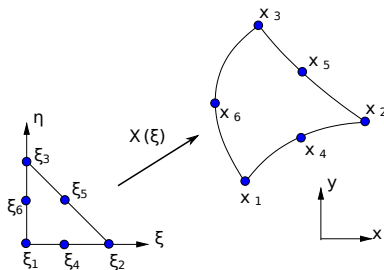


The *geometric* mesh of  $\Omega$  assigns coordinates to each cell  $([CG_1]^2)$ .

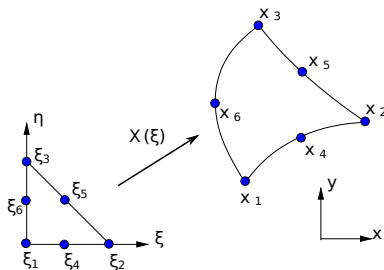
One can consider a purely *topological* mesh: sets of cells and vertices, with connectivity maps between them.



The *geometric* mesh of  $\Omega$  assigns coordinates to each cell  $([CG_2]^2)$ .



We represent the coordinate field with Lagrange elements of arbitrary order, allowing us to bend the mesh. This is useful if  $\Omega$  is not a polygon or polyhedron.



We represent the coordinate field with Lagrange elements of arbitrary order, allowing us to bend the mesh. This is useful if  $\Omega$  is not a polygon or polyhedron.

This means that for each element we can write

$$x = \sum_{i=1}^d x_i \hat{\psi}_i(\hat{x})$$

for (scalar-valued) coefficients  $x_i$  and (vector-valued) basis functions  $\hat{\psi}_i$ . This is an explicit construction for the map  $x = F_K(\hat{x})$ .



## C6.4 Finite Element Methods for PDEs

### Lecture 10: Finite elements beyond Lagrange

Patrick E. Farrell

University of Oxford

We are solving the linear variational problem

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ .

We are solving the linear variational problem

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ .

One of the great advantages of the finite element method is that you can *tailor* the approximation  $V_h$  to the function space  $V$ . For  $V \neq H^1(\Omega)$  we usually want to use different finite elements than Lagrange.

We are solving the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

One of the great advantages of the finite element method is that you can *tailor* the approximation  $V_h$  to the function space  $V$ . For  $V \neq H^1(\Omega)$  we usually want to use different finite elements than Lagrange.

The fundamental Hilbert spaces we have met are related by the *de Rham* complex:

$$H^1 \xrightarrow{\text{grad}} H(\text{curl}) \xrightarrow{\text{curl}} H(\text{div}) \xrightarrow{\text{div}} L^2$$

We are solving the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

One of the great advantages of the finite element method is that you can *tailor* the approximation  $V_h$  to the function space  $V$ . For  $V \neq H^1(\Omega)$  we usually want to use different finite elements than Lagrange.

The fundamental Hilbert spaces we have met are related by the *de Rham* complex:

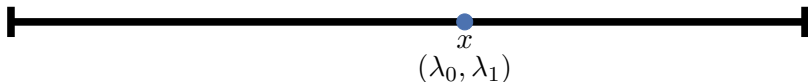
$$\begin{array}{ccccccc} H^1 & \xrightarrow{\text{grad}} & H(\text{curl}) & \xrightarrow{\text{curl}} & H(\text{div}) & \xrightarrow{\text{div}} & L^2 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ \Sigma_h & \xrightarrow{\text{grad}} & V_h & \xrightarrow{\text{curl}} & W_h & \xrightarrow{\text{div}} & Q_h \end{array}$$

We can build finite element spaces for *all* spaces in the  $L^2$  de Rham complex in a structure-preserving way.

On the interval  $\Omega = [0, 1]$ , the obvious way to label a point is with a single coordinate  $x$ .

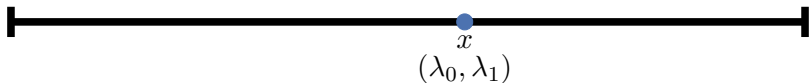


On the interval  $\Omega = [0, 1]$ , the obvious way to label a point is with a single coordinate  $x$ .



It will be convenient for us to label points with *barycentric coordinates*:  $n + 1$  coordinates for a simplex in  $\mathbb{R}^n$  with the constraint that  $\sum_i \lambda_i = 1$ .

On the interval  $\Omega = [0, 1]$ , the obvious way to label a point is with a single coordinate  $x$ .

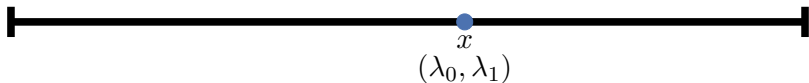


It will be convenient for us to label points with *barycentric coordinates*:  $n + 1$  coordinates for a simplex in  $\mathbb{R}^n$  with the constraint that  $\sum_i \lambda_i = 1$ .

In 1D,  $\lambda_0 = x$ ,  $\lambda_1 = 1 - x$ .



On the interval  $\Omega = [0, 1]$ , the obvious way to label a point is with a single coordinate  $x$ .

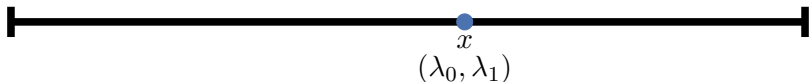


It will be convenient for us to label points with *barycentric coordinates*:  $n + 1$  coordinates for a simplex in  $\mathbb{R}^n$  with the constraint that  $\sum_i \lambda_i = 1$ .

In 1D,  $\lambda_0 = x$ ,  $\lambda_1 = 1 - x$ .

The reason why this is convenient is because it gives us a nice way to describe different geometric parts of our simplex. In 1D, we have the left vertex is given by  $\lambda_0 = 0$ , the right vertex by  $\lambda_1 = 0$ .

On the interval  $\Omega = [0, 1]$ , the obvious way to label a point is with a single coordinate  $x$ .

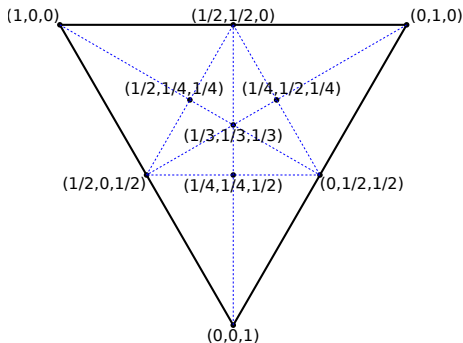


It will be convenient for us to label points with *barycentric coordinates*:  $n + 1$  coordinates for a simplex in  $\mathbb{R}^n$  with the constraint that  $\sum_i \lambda_i = 1$ .

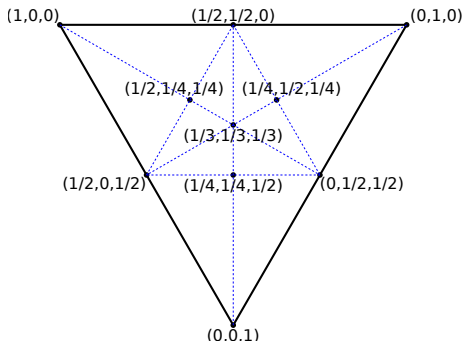
In 1D,  $\lambda_0 = x$ ,  $\lambda_1 = 1 - x$ .

The reason why this is convenient is because it gives us a nice way to describe different geometric parts of our simplex. In 1D, we have the left vertex is given by  $\lambda_0 = 0$ , the right vertex by  $\lambda_1 = 0$ .

Another way to look at it: the barycentric coordinates express a point  $p$  as a convex combination of the vertices.



For a triangle, we represent points  $p = (\lambda_1, \lambda_2, \lambda_3)$ ; three numbers, but only two degrees of freedom because of the summation constraint.



For a triangle, we represent points  $p = (\lambda_1, \lambda_2, \lambda_3)$ ; three numbers, but only two degrees of freedom because of the summation constraint.

The edge opposite vertex  $i$  is described by  $\lambda_i = 0$ .

## Lemma (First factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 1$  that vanishes on a hyperplane  $\{x : L(x) = 0\}$ , where  $L(x)$  is a non-degenerate linear function. Then we can write  $P = LQ$ , where  $Q$  is a polynomial of degree  $d - 1$ .*

## Lemma (First factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 1$  that vanishes on a hyperplane  $\{x : L(x) = 0\}$ , where  $L(x)$  is a non-degenerate linear function. Then we can write  $P = LQ$ , where  $Q$  is a polynomial of degree  $d - 1$ .*

## Example (One dimension)

Suppose  $P$  vanishes on the hyperplane  $\{x : x - r = 0\}$ . Then we can write  $P = (x - r)Q$  for some polynomial  $Q$ .

## Lemma (First factorisation lemma)

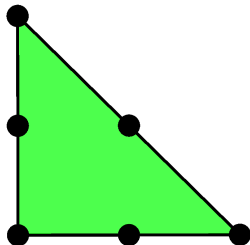
*Let  $P$  be a polynomial of degree  $d \geq 1$  that vanishes on a hyperplane  $\{x : L(x) = 0\}$ , where  $L(x)$  is a non-degenerate linear function. Then we can write  $P = LQ$ , where  $Q$  is a polynomial of degree  $d - 1$ .*

### Example (One dimension)

Suppose  $P$  vanishes on the hyperplane  $\{x : x - r = 0\}$ . Then we can write  $P = (x - r)Q$  for some polynomial  $Q$ .

### Example (Higher dimensions)

Suppose  $P$  vanishes on each edge of a triangle. Then  $P = \lambda_1 \lambda_2 \lambda_3 Q$  for some  $Q$ .



The quadratic Lagrange finite element  $\text{CG}_2$  in two dimensions.

## Unisolvence of $\text{CG}_2$

Suppose  $v \in \mathcal{P}_2(\triangle)$  with all degrees of freedom zero. Restricted to an edge,  $v$  is a quadratic polynomial with three roots, hence  $v = 0$  on each edge. By the factorisation lemma,  $v = \lambda_1 \lambda_2 c$  for a constant  $c \in \mathbb{R}$ . Evaluating both sides on the edge  $\lambda_3 = 0$  shows that  $c = 0$ .



## Lemma (Second factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .*

## Lemma (Second factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .*

## Remark

If  $P$  vanishes on  $\{x : L(x) = 0\}$ , so does  $\nabla P \cdot t$  for a tangent vector  $t$ .

## Lemma (Second factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .*

## Remark

If  $P$  vanishes on  $\{x : L(x) = 0\}$ , so does  $\nabla P \cdot t$  for a tangent vector  $t$ .

## Proof.

Since  $P$  vanishes on  $\{x : L(x) = 0\}$ , we have  $P = L\tilde{Q}$ .

## Lemma (Second factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .*

## Remark

If  $P$  vanishes on  $\{x : L(x) = 0\}$ , so does  $\nabla P \cdot t$  for a tangent vector  $t$ .

## Proof.

Since  $P$  vanishes on  $\{x : L(x) = 0\}$ , we have  $P = L\tilde{Q}$ . Calculating,

$$\nabla P \cdot n = \tilde{Q} \nabla L \cdot n + L \nabla \tilde{Q} \cdot n.$$

## Lemma (Second factorisation lemma)

*Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .*

## Remark

If  $P$  vanishes on  $\{x : L(x) = 0\}$ , so does  $\nabla P \cdot t$  for a tangent vector  $t$ .

## Proof.

Since  $P$  vanishes on  $\{x : L(x) = 0\}$ , we have  $P = L\tilde{Q}$ . Calculating,

$$\nabla P \cdot n = \tilde{Q} \nabla L \cdot n + L \nabla \tilde{Q} \cdot n.$$

Since  $L$  vanishes on the plane, and  $\nabla L$  is normal to the plane (hence colinear with  $n$ ), this forces  $\tilde{Q} = 0$  on  $\{x : L(x) = 0\}$ .

## Lemma (Second factorisation lemma)

Let  $P$  be a polynomial of degree  $d \geq 2$  such that  $P$  and  $\nabla P \cdot n$  vanish on a hyperplane  $\{x : L(x) = 0\}$ , where  $n$  is the normal to  $L$ . Then we can write  $P = L^2Q$ , where  $Q$  is a polynomial of degree  $d - 2$ .

### Remark

If  $P$  vanishes on  $\{x : L(x) = 0\}$ , so does  $\nabla P \cdot t$  for a tangent vector  $t$ .

### Proof.

Since  $P$  vanishes on  $\{x : L(x) = 0\}$ , we have  $P = L\tilde{Q}$ . Calculating,

$$\nabla P \cdot n = \tilde{Q} \nabla L \cdot n + L \nabla \tilde{Q} \cdot n.$$

Since  $L$  vanishes on the plane, and  $\nabla L$  is normal to the plane (hence colinear with  $n$ ), this forces  $\tilde{Q} = 0$  on  $\{x : L(x) = 0\}$ . Hence  $\tilde{Q} = LQ$  for some  $Q$ . □

## Section 3

# The biharmonic problem

The biharmonic equation arises in many areas of physics. It describes equilibrium configurations of clamped plates under transverse loading, the stresses in an elastic body, the stream function in certain flow regimes, and other things. The equation is

$$\begin{aligned}\nabla^4 u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega, \\ \nabla u \cdot n &= 0 \text{ on } \partial\Omega.\end{aligned}$$



The biharmonic equation arises in many areas of physics. It describes equilibrium configurations of clamped plates under transverse loading, the stresses in an elastic body, the stream function in certain flow regimes, and other things. The equation is

$$\begin{aligned}\nabla^4 u &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega, \\ \nabla u \cdot n &= 0 \text{ on } \partial\Omega.\end{aligned}$$

Here  $\nabla^4 = \Delta^2 = \nabla \cdot \nabla \nabla \cdot \nabla$ . More simply, in two dimensions

$$\frac{\partial^4 u}{\partial x^4} + \frac{\partial^4 u}{\partial y^4} + 2 \frac{\partial^4 u}{\partial x^2 \partial y^2} = f.$$

Let's cast this into variational form formally, i.e. not yet specifying  $V$ .  
Multiplying by  $v \in V$  for some  $V$ , we find

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla \cdot \nabla \nabla \cdot \nabla uv \, dx = \int_{\Omega} f v \, dx.$$

Let's cast this into variational form formally, i.e. not yet specifying  $V$ . Multiplying by  $v \in V$  for some  $V$ , we find

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla \cdot \nabla \nabla \cdot \nabla uv \, dx = \int_{\Omega} f v \, dx.$$

We want to invoke Lax–Milgram, so we want  $u \in V$  and  $v \in V$ . Let's integrate by parts once:

$$\int_{\Omega} \nabla^4 uv \, dx = - \int_{\Omega} (\nabla \nabla \cdot \nabla) u \cdot \nabla v \, dx + \int_{\partial\Omega} \nabla (\nabla \cdot \nabla u) \cdot n v \, ds$$

Let's cast this into variational form formally, i.e. not yet specifying  $V$ . Multiplying by  $v \in V$  for some  $V$ , we find

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla \cdot \nabla \nabla \cdot \nabla uv \, dx = \int_{\Omega} f v \, dx.$$

We want to invoke Lax–Milgram, so we want  $u \in V$  and  $v \in V$ . Let's integrate by parts once:

$$\int_{\Omega} \nabla^4 uv \, dx = - \int_{\Omega} (\nabla \nabla \cdot \nabla) u \cdot \nabla v \, dx + \int_{\partial\Omega} \nabla (\nabla \cdot \nabla u) \cdot n v \, ds$$

and again:

$$\begin{aligned} \int_{\Omega} \nabla^4 uv \, dx &= \int_{\Omega} (\nabla \cdot \nabla u) (\nabla \cdot \nabla v) \, dx \\ &\quad - \int_{\partial\Omega} (\nabla \cdot \nabla u) \nabla v \cdot n \, ds \\ &\quad + \int_{\partial\Omega} \nabla (\nabla \cdot \nabla u) \cdot n v \, ds. \end{aligned}$$

Noting that the Laplacian  $\nabla^2 = \nabla \cdot \nabla$ , we rewrite

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla^2 u \nabla^2 v \, dx - \int_{\partial\Omega} \nabla^2 u \nabla v \cdot n \, ds + \int_{\partial\Omega} \nabla (\nabla^2 u) \cdot n v \, ds.$$

Noting that the Laplacian  $\nabla^2 = \nabla \cdot \nabla$ , we rewrite

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla^2 u \nabla^2 v \, dx - \int_{\partial\Omega} \nabla^2 u \nabla v \cdot n \, ds + \int_{\partial\Omega} \nabla (\nabla^2 u) \cdot n v \, ds.$$

We need square-integrable second derivatives, so  $V \subset H^2(\Omega)$ .

Noting that the Laplacian  $\nabla^2 = \nabla \cdot \nabla$ , we rewrite

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla^2 u \nabla^2 v \, dx - \int_{\partial\Omega} \nabla^2 u \nabla v \cdot n \, ds + \int_{\partial\Omega} \nabla (\nabla^2 u) \cdot n v \, ds.$$

We need square-integrable second derivatives, so  $V \subset H^2(\Omega)$ .

We have nowhere convenient to enforce the boundary conditions  $u = \nabla u \cdot n = 0$ . So we should take

$$V = H_0^2(\Omega) := \{v \in H^2(\Omega) : v = 0, \nabla v \cdot n = 0 \text{ on } \partial\Omega\}.$$

Noting that the Laplacian  $\nabla^2 = \nabla \cdot \nabla$ , we rewrite

$$\int_{\Omega} \nabla^4 uv \, dx = \int_{\Omega} \nabla^2 u \nabla^2 v \, dx - \int_{\partial\Omega} \nabla^2 u \nabla v \cdot n \, ds + \int_{\partial\Omega} \nabla (\nabla^2 u) \cdot n v \, ds.$$

We need square-integrable second derivatives, so  $V \subset H^2(\Omega)$ .

We have nowhere convenient to enforce the boundary conditions  $u = \nabla u \cdot n = 0$ . So we should take

$$V = H_0^2(\Omega) := \{v \in H^2(\Omega) : v = 0, \nabla v \cdot n = 0 \text{ on } \partial\Omega\}.$$

With  $v \in V$ , the surface integrals vanish, leaving us with the problem: find  $u \in H_0^2(\Omega)$  such that

$$\int_{\Omega} \nabla^2 u \nabla^2 v \, dx = \int_{\Omega} f v \, dx \text{ for all } v \in H_0^2(\Omega).$$



How do we discretise this problem? If we take  $V_h \sim \text{CG}_p$ ,  $V_h \notin H^2$ !

How do we discretise this problem? If we take  $V_h \sim \text{CG}_p$ ,  $V_h \notin H^2$ !

For a piecewise smooth function  $u$ ,  $u \in H^1(\Omega) \iff u \in C^0(\Omega)$ .

Since  $u \in H^2(\Omega)$  iff  $u$  and all its first derivatives are in  $H^1(\Omega)$ , that means for  $u \in H^2(\Omega)$  we need  $u \in C^1(\Omega)$ .

How do we discretise this problem? If we take  $V_h \sim \text{CG}_p$ ,  $V_h \notin H^2$ !

For a piecewise smooth function  $u$ ,  $u \in H^1(\Omega) \iff u \in C^0(\Omega)$ .

Since  $u \in H^2(\Omega)$  iff  $u$  and all its first derivatives are in  $H^1(\Omega)$ , that means for  $u \in H^2(\Omega)$  we need  $u \in C^1(\Omega)$ .

Two approaches:

- ▶  $C^1$ -continuous finite elements;
- ▶ nonconforming discretisations.

## Section 4

### The Hermite element



## Definition (Hermite finite element)

$K = [0, 1]$ ,  $\mathcal{V} = \mathcal{P}_3(K)$ , and

$$\mathcal{L} = \{v \mapsto v(0), \\ v \mapsto v'(0), \\ v \mapsto v(1), \\ v \mapsto v'(1)\}.$$



## Definition (Hermite finite element)

$K = [0, 1]$ ,  $\mathcal{V} = \mathcal{P}_3(K)$ , and

$$\mathcal{L} = \{v \mapsto v(0), \\ v \mapsto v'(0), \\ v \mapsto v(1), \\ v \mapsto v'(1)\}.$$

This gives a  $C^1$  approximation, because the function value *and derivative* agree across cells (by construction).



## Definition (Hermite finite element)

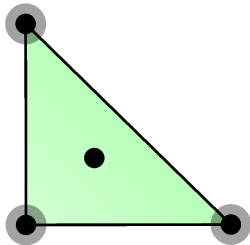
$K = [0, 1]$ ,  $\mathcal{V} = \mathcal{P}_3(K)$ , and

$$\mathcal{L} = \{v \mapsto v(0), \\ v \mapsto v'(0), \\ v \mapsto v(1), \\ v \mapsto v'(1)\}.$$

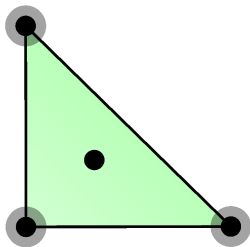
This gives a  $C^1$  approximation, because the function value *and derivative* agree across cells (by construction).

## Unisolvence

Suppose  $v \in \mathcal{P}_3(K)$  with all dofs zero. Then  $v$  is a cubic polynomial with four roots (two double roots), hence zero.

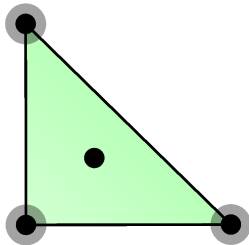






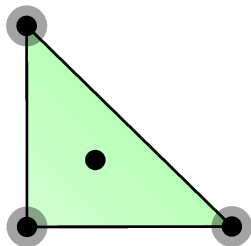
### Definition (Hermite element in 2D)

$K = \triangle, \mathcal{V} = \mathcal{P}_3(\triangle), \mathcal{L}$  shown.



## Lemma (Unisolvence of the triangular Hermite element)

*The Hermite element in two dimensions is unisolvent.*



## Lemma (Unisolvence of the triangular Hermite element)

*The Hermite element in two dimensions is unisolvent.*

### Proof.

Suppose  $u \in \mathcal{P}_3(\triangle)$  with all dofs zero. Along an edge of the triangle,  $u$  is a cubic polynomial with 2 double roots, so  $u = 0$  along each edge. Thus,  $u(x) = c\lambda_1\lambda_2\lambda_3$  for some  $c \in \mathbb{R}$ . Since the value at the barycentre is also zero, and  $\lambda_1 = \lambda_2 = \lambda_3 = 1/3$  at the barycentre, we must have  $c = 0$ .  $\square$

Job done?

Job done? No!

Job done? No!

In the unisolvence proof we saw the dofs on an edge (i.e. those shared with a neighbour) determine  $u$  and thus  $\nabla u \cdot t$ ,  $t$  the tangent vector of the edge. But they do *not* determine  $\nabla u \cdot n$ .

Job done? No!

In the unisolvence proof we saw the dofs on an edge (i.e. those shared with a neighbour) determine  $u$  and thus  $\nabla u \cdot t$ ,  $t$  the tangent vector of the edge. But they do *not* determine  $\nabla u \cdot n$ .

For two elements, take

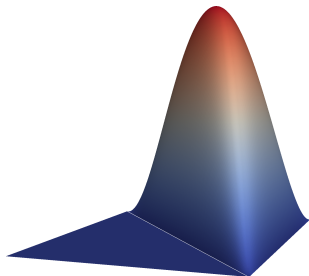
$$p(x) = \begin{cases} \lambda_1 \lambda_2 \lambda_3 & x \in K_1, \\ 0 & x \in K_2, \end{cases}$$

Job done? No!

In the unisolvence proof we saw the dofs on an edge (i.e. those shared with a neighbour) determine  $u$  and thus  $\nabla u \cdot t$ ,  $t$  the tangent vector of the edge. But they do *not* determine  $\nabla u \cdot n$ .

For two elements, take

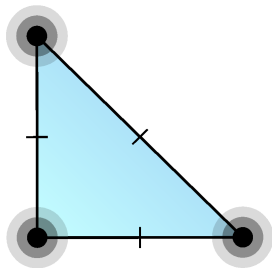
$$p(x) = \begin{cases} \lambda_1 \lambda_2 \lambda_3 & x \in K_1, \\ 0 & x \in K_2, \end{cases}$$

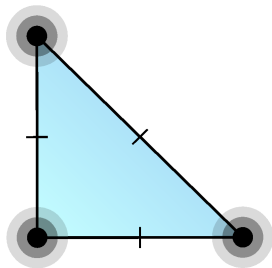




## Section 5

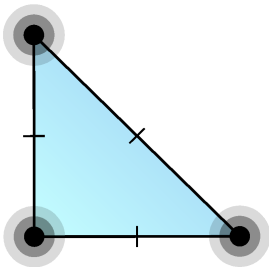
# The Argyris element





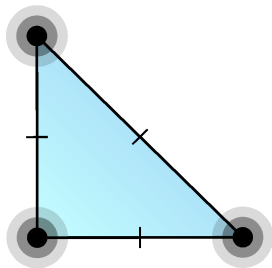
Definition (Hermite element in 2D)

$K = \triangle, \mathcal{V} = \mathcal{P}_5(\triangle), \mathcal{L}$  shown.



## Lemma (Unisolvence of the triangular Argyris element)

*The Argyris element in two dimensions is unisolvent.*



## Proof.

Suppose  $u \in \mathcal{P}_5(\triangle)$  with all dofs zero. Along an edge,  $u$  is a quintic polynomial with 2 treble roots, so  $u = 0$  along each edge. Moreover,  $\nabla u \cdot n$  is a quartic polynomial with 2 double roots and a single root, hence zero. Thus,  $u$  is divisible by  $\lambda_1^2 \lambda_2^2 \lambda_3^2$ , which is of degree 6. Thus  $u = 0$ .  $\square$

```
1  from firedrake import *
2
3  mesh = UnitSquareMesh(20, 20)
4  V = FunctionSpace(mesh, "Argyris", 5)
5
6  u = Function(V)
7  v = TestFunction(V)
8
9  (x, y) = SpatialCoordinate(mesh)
10 f = sin(2*pi*x) * sin(2*pi*y)
11
12 F = (inner(div(grad(u)), div(grad(v)))*dx
13      + inner(u, v)*dx
14      - inner(f, v)*dx)
15 solve(F == 0, u)
```

dim	lowest $p$ for $V_h \subset H^2(\Omega)$
1	3
2	5
3	9

dim	lowest $p$ for $V_h \subset H^2(\Omega)$
1	3
2	5
3	9

An alternative: *nonconforming* discretisations. Use Lagrange elements and *weakly* enforce  $C^1$  continuity by penalising the jump in  $\nabla u \cdot n$  across edges.



dim	lowest $p$ for $V_h \subset H^2(\Omega)$
1	3
2	5
3	9

An alternative: *nonconforming* discretisations. Use Lagrange elements and *weakly* enforce  $C^1$  continuity by penalising the jump in  $\nabla u \cdot n$  across edges.

$$a_h(u, v) = \sum_{K \in \mathcal{M}} \int_K \nabla^2 u \nabla^2 v \, dx + \gamma \sum_{E \in \mathcal{E}_h} \frac{1}{|E|} \int_E \llbracket \nabla u \cdot n \rrbracket \llbracket \nabla v \cdot n \rrbracket \, ds + \dots$$

dim	lowest $p$ for $V_h \subset H^2(\Omega)$
1	3
2	5
3	9

An alternative: *nonconforming* discretisations. Use Lagrange elements and *weakly* enforce  $C^1$  continuity by penalising the jump in  $\nabla u \cdot n$  across edges.

$$a_h(u, v) = \sum_{K \in \mathcal{M}} \int_K \nabla^2 u \nabla^2 v \, dx + \gamma \sum_{E \in \mathcal{E}_h} \frac{1}{|E|} \int_E \llbracket \nabla u \cdot n \rrbracket \llbracket \nabla v \cdot n \rrbracket \, ds + \dots$$

where the jump over cells  $K_+$ ,  $K_-$  is

$$\llbracket \nabla u \cdot n \rrbracket := (\nabla u)_+ \cdot n_+ + (\nabla u)_- \cdot n_-.$$

Research aside: in recent work coauthors & I have proposed the energy functional

$$\begin{aligned} F(\delta\rho, Q) = \int_{\Omega} & \left[ \frac{a}{2} (\delta\rho)^2 + \frac{b}{3} (\delta\rho)^3 + \frac{c}{4} (\delta\rho)^4 \right. \\ & + B \left| \mathcal{D}^2 \delta\rho + q^2 \left( Q + \frac{I_d}{d} \right) \delta\rho \right|^2 \\ & \left. + \frac{K}{2} |\nabla Q|^2 + f_n(Q) \right] \end{aligned}$$

to model smectic liquid crystals.

Research aside: in recent work coauthors & I have proposed the energy functional

$$\begin{aligned}
 F(\delta\rho, Q) = \int_{\Omega} & \left[ \frac{a}{2} (\delta\rho)^2 + \frac{b}{3} (\delta\rho)^3 + \frac{c}{4} (\delta\rho)^4 \right. \\
 & + B \left| \mathcal{D}^2 \delta\rho + q^2 \left( Q + \frac{I_d}{d} \right) \delta\rho \right|^2 \\
 & \left. + \frac{K}{2} |\nabla Q|^2 + f_n(Q) \right]
 \end{aligned}$$

to model smectic liquid crystals.

This requires  $\delta\rho \in H^2(\Omega)$  because of the Hessian  $\mathcal{D}^2 \delta\rho$  in the energy.

Research aside: in recent work coauthors & I have proposed the energy functional

$$\begin{aligned}
 F(\delta\rho, Q) = \int_{\Omega} & \left[ \frac{a}{2} (\delta\rho)^2 + \frac{b}{3} (\delta\rho)^3 + \frac{c}{4} (\delta\rho)^4 \right. \\
 & + B \left| \mathcal{D}^2 \delta\rho + q^2 \left( Q + \frac{I_d}{d} \right) \delta\rho \right|^2 \\
 & \left. + \frac{K}{2} |\nabla Q|^2 + f_n(Q) \right]
 \end{aligned}$$

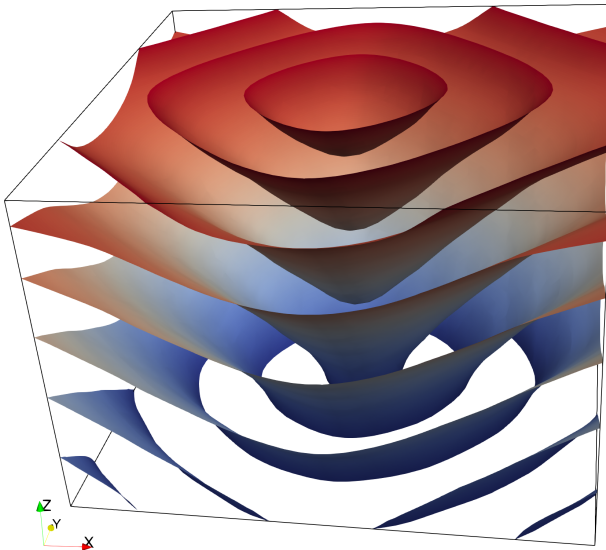
to model smectic liquid crystals.

This requires  $\delta\rho \in H^2(\Omega)$  because of the Hessian  $\mathcal{D}^2 \delta\rho$  in the energy.

We discretise this with Lagrange elements by adding

$$\sum_{E \in \mathcal{E}} \int_E \frac{\gamma}{|E|^3} (\llbracket \nabla \delta\rho \cdot n \rrbracket)^2$$

to the energy.



A *toroidal focal conic domain*, computed for the first time as a minimiser of an energy functional.

## Section 7

Elements for  $H(\text{div})$  and  $H(\text{curl})$

Let's consider another model problem, the  $H(\text{div})$  Riesz map:

$$\begin{aligned}u - \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u \cdot n &= 0 \text{ on } \partial\Omega.\end{aligned}$$



Let's consider another model problem, the  $H(\text{div})$  Riesz map:

$$\begin{aligned}u - \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u \cdot n &= 0 \text{ on } \partial\Omega.\end{aligned}$$

Casting to variational form, we have

$$\int_{\Omega} u \cdot v \, dx - \int_{\Omega} (\nabla \nabla \cdot u) \cdot v \, dx = \int_{\Omega} f \cdot v \, dx,$$

Let's consider another model problem, the  $H(\text{div})$  Riesz map:

$$\begin{aligned} u - \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u \cdot n &= 0 \text{ on } \partial\Omega. \end{aligned}$$

Casting to variational form, we have

$$\int_{\Omega} u \cdot v \, dx - \int_{\Omega} (\nabla \nabla \cdot u) \cdot v \, dx = \int_{\Omega} f \cdot v \, dx,$$

and integrating by parts yields

$$\int_{\Omega} u \cdot v \, dx + \int_{\Omega} \nabla \cdot u \nabla \cdot v \, dx - \int_{\Omega} \nabla \cdot uv \cdot n \, ds = \int_{\Omega} f \cdot v \, dx.$$

Let's consider another model problem, the  $H(\text{div})$  Riesz map:

$$\begin{aligned}u - \nabla \nabla \cdot u &= f \text{ in } \Omega, \\ u \cdot n &= 0 \text{ on } \partial\Omega.\end{aligned}$$

Casting to variational form, we have

$$\int_{\Omega} u \cdot v \, dx - \int_{\Omega} (\nabla \nabla \cdot u) \cdot v \, dx = \int_{\Omega} f \cdot v \, dx,$$

and integrating by parts yields

$$\int_{\Omega} u \cdot v \, dx + \int_{\Omega} \nabla \cdot u \nabla \cdot v \, dx - \int_{\Omega} \nabla \cdot u v \cdot n \, ds = \int_{\Omega} f \cdot v \, dx.$$

The base space is  $H(\text{div})$ , and we need to enforce BCs:

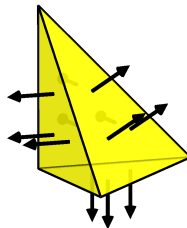
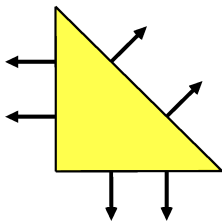
$$V = H_0(\text{div}) := \{v \in H(\text{div}) : v \cdot n = 0 \text{ on } \partial\Omega\},$$

with final formulation: find  $u \in V$  such that

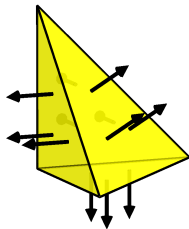
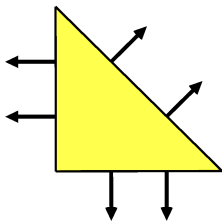
$$\int_{\Omega} u \cdot v \, dx + \int_{\Omega} \nabla \cdot u \nabla \cdot v \, dx = \int_{\Omega} f \cdot v \, dx \text{ for all } v \in V.$$

On problem sheet 1 we learn that a pw smooth function is in  $H(\text{div})$  iff its normal component is continuous.

On problem sheet 1 we learn that a pw smooth function is in  $H(\text{div})$  iff its normal component is continuous.



On problem sheet 1 we learn that a pw smooth function is in  $H(\text{div})$  iff its normal component is continuous.



Definition (Lowest order Brezzi–Douglas–Marini element)

$K = \triangle$  or  $K = \triangleleft$ ,  $\mathcal{V} = \mathcal{P}_1(K)^n$ ,  $\mathcal{L}$  eval normal component on facets.

Similarly, for the  $H(\text{curl})$  Riesz map

$$\begin{aligned}u + \nabla \times \nabla \times u &= f \text{ in } \Omega, \\u \times n &= 0 \text{ on } \partial\Omega,\end{aligned}$$

Similarly, for the  $H(\text{curl})$  Riesz map

$$\begin{aligned} u + \nabla \times \nabla \times u &= f \text{ in } \Omega, \\ u \times n &= 0 \text{ on } \partial\Omega, \end{aligned}$$

we end up with the variational formulation over

$$V := \{v \in H(\text{curl}) : v \times n = 0\}$$

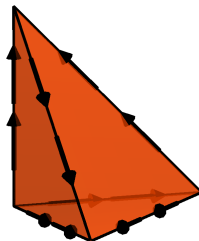
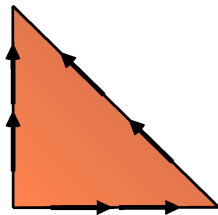
to find  $u \in V$  such that

$$\int_{\Omega} u \cdot v \, dx + \int_{\Omega} \nabla \times u \cdot \nabla \times v \, dx = \int_{\Omega} f \cdot v \, dx \text{ for all } v \in V.$$

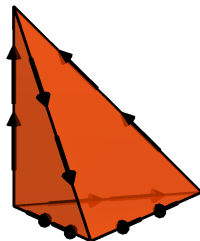
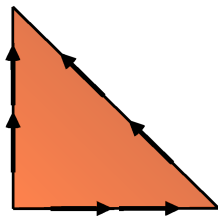


On problem sheet 1 we learn that a pw smooth function is in  $H(\text{curl})$  iff its tangential components are continuous.

On problem sheet 1 we learn that a pw smooth function is in  $H(\text{curl})$  iff its tangential components are continuous.



On problem sheet 1 we learn that a pw smooth function is in  $H(\text{curl})$  iff its tangential components are continuous.



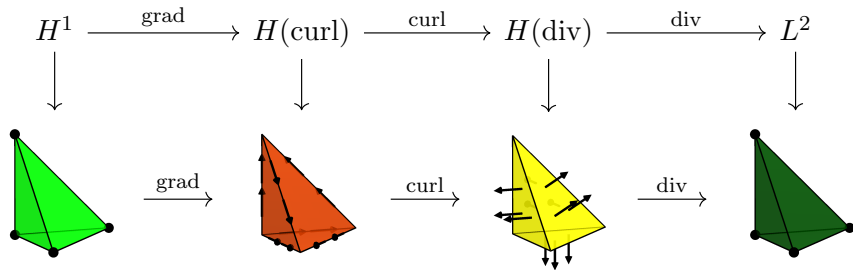
Definition (Lowest order Nédélec element of the second kind)

$K = \triangle$  or  $K = \diamond$ ,  $\mathcal{V} = \mathcal{P}_1(K)^n$ ,  $\mathcal{L}$  eval tangential component on edges.

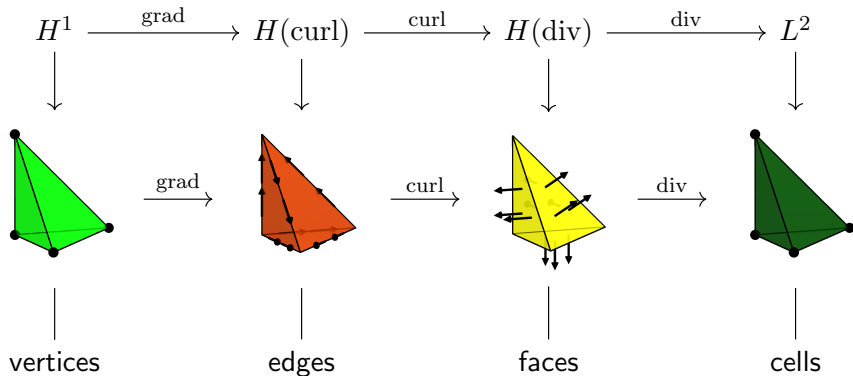
Consider again the de Rham complex:

$$H^1 \xrightarrow{\text{grad}} H(\text{curl}) \xrightarrow{\text{curl}} H(\text{div}) \xrightarrow{\text{div}} L^2$$

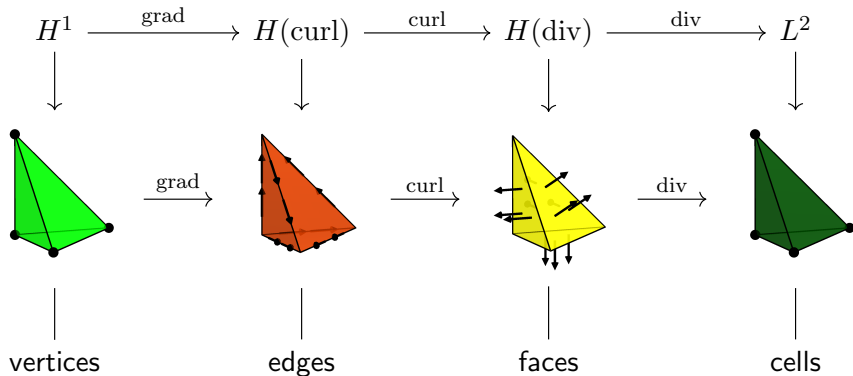
Consider again the de Rham complex:



Consider again the de Rham complex:



Consider again the de Rham complex:



The finite element method has deep connections to algebraic topology and differential geometry: the *finite element exterior calculus*.

## C6.4 Finite Element Methods for PDEs

### Lecture 11: Interpolation error estimates

Patrick E. Farrell

University of Oxford



We are approximating the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

with the solution of

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

We are approximating the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

with the solution of

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

In Lecture 7 we saw Céa's Lemma for coercive, bounded  $a$ :

$$\|u - u_h\|_V \leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V$$

We are approximating the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

with the solution of

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

In Lecture 7 we saw Céa's Lemma for coercive, bounded  $a$ :

$$\begin{aligned} \|u - u_h\|_V &\leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V \\ &\leq \frac{C}{\alpha} \|u - \mathcal{I}_h u\|_V \end{aligned}$$

where  $\mathcal{I}_h : V \rightarrow V_h$  is the finite element interpolation operator.

We are approximating the linear variational problem

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

with the solution of

$$\text{find } u_h \in V_h \text{ such that } a(u_h, v_h) = F(v_h) \text{ for all } v_h \in V_h.$$

In Lecture 7 we saw Céa's Lemma for coercive, bounded  $a$ :

$$\begin{aligned} \|u - u_h\|_V &\leq \frac{C}{\alpha} \min_{v_h \in V_h} \|u - v_h\|_V \\ &\leq \frac{C}{\alpha} \|u - \mathcal{I}_h u\|_V \\ &\leq \frac{C}{\alpha} f(h, p) \end{aligned}$$

where  $\mathcal{I}_h : V \rightarrow V_h$  is the finite element interpolation operator.

We seek a bound in terms of parameters we control:  $h, p$ .

In lecture 3, we saw that the Sobolev space  $W_p^k(\Omega)$  for  $p < \infty$  is equipped with

$$\|u\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

In lecture 3, we saw that the Sobolev space  $W_p^k(\Omega)$  for  $p < \infty$  is equipped with

$$\|u\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

At this point it is convenient to introduce the Sobolev *seminorm*

$$|u|_{W_p^k(\Omega)} = \left( \sum_{|\alpha|=k} \|D^\alpha u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

In lecture 3, we saw that the Sobolev space  $W_p^k(\Omega)$  for  $p < \infty$  is equipped with

$$\|u\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \leq k} \|D^\alpha u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

At this point it is convenient to introduce the Sobolev *seminorm*

$$|u|_{W_p^k(\Omega)} = \left( \sum_{|\alpha|=k} \|D^\alpha u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

## Example

$$\|u\|_{H^1(\Omega)}^2 = \|u\|_{L^2(\Omega)}^2 + |u|_{H^1(\Omega)}^2.$$

We want to bound the interpolation error in terms of the mesh size  $h \rightarrow 0$  and polynomial degree  $p \rightarrow \infty$ . How should we characterise the size of the mesh?



We want to bound the interpolation error in terms of the mesh size  $h \rightarrow 0$  and polynomial degree  $p \rightarrow \infty$ . How should we characterise the size of the mesh?

The quantity we will use on each cell is its *diameter*.

### Definition (diameter of a cell)

$$h_K = \text{diam}(K) = \sup\{\|x - y\| : x, y \in K\}.$$

For a triangle or tetrahedron, this resolves to the length of its longest edge.

We want to bound the interpolation error in terms of the mesh size  $h \rightarrow 0$  and polynomial degree  $p \rightarrow \infty$ . How should we characterise the size of the mesh?

The quantity we will use on each cell is its *diameter*.

### Definition (diameter of a cell)

$$h_K = \text{diam}(K) = \sup\{\|x - y\| : x, y \in K\}.$$

For a triangle or tetrahedron, this resolves to the length of its longest edge.

Over the whole mesh we take a pessimistic view.

### Definition (mesh size)

Given a mesh  $\mathcal{M}$ , its mesh size  $h$  is given by

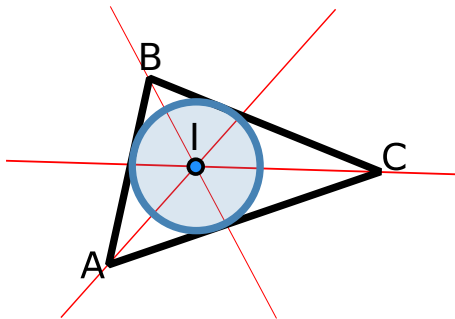
$$h = \max_{K \in \mathcal{M}} \text{diam}(K).$$

We consider a sequence of meshes  $(\mathcal{M}_h)_h$  indexed by mesh size  $h \rightarrow 0$ .  
We need a technical condition on the sequence of meshes.

We consider a sequence of meshes  $(\mathcal{M}_h)_h$  indexed by mesh size  $h \rightarrow 0$ . We need a technical condition on the sequence of meshes.

### Definition (incircle diameter of a cell)

The incircle diameter  $\rho_K$  of a cell  $K$  is the diameter of the largest hyperdisc (i.e. disc in two dimensions, ball in three dimensions) that is completely contained within  $K$ .



We consider a sequence of meshes  $(\mathcal{M}_h)_h$  indexed by mesh size  $h \rightarrow 0$ . We need a technical condition on the sequence of meshes.

### Definition (incircle diameter of a cell)

The incircle diameter  $\rho_K$  of a cell  $K$  is the diameter of the largest hyperdisc (i.e. disc in two dimensions, ball in three dimensions) that is completely contained within  $K$ .

### Definition (shape regularity of mesh sequence $(\mathcal{M}_h)_h$ )

A sequence of meshes  $(\mathcal{M}_h)_h$  is *shape regular* if there exists a constant  $\sigma$  such that

$$\sup_h \max_{K \in \mathcal{M}_h} \frac{h_K}{\rho_K} \leq \sigma.$$

## Section 3

# Interpolation error for Lagrange elements

First, consider  $\text{CG}_p$  elements applied to problems in  $H^1(\Omega)$ .

## Theorem

*Let  $(V_h)_h$  be the function spaces constructed with continuous Lagrange elements of order  $p$  on a shape-regular sequence of meshes  $(\mathcal{M}_h)_h$  indexed by mesh size  $h$ . Let  $u \in H^{p+1}(\Omega)$ , and let  $\mathcal{I}_h : H^{p+1}(\Omega) \rightarrow V_h$  be the interpolation operator associated with each  $V_h$ . Then there exists a constant  $D < \infty$  independent of  $u$  such that*

$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq Dh^p |u|_{H^{p+1}(\Omega)}.$$

$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq Dh^p |u|_{H^{p+1}(\Omega)}.$$



$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq Dh^p |u|_{H^{p+1}(\Omega)}.$$

## Remark

For  $p = 1$ , the interpolation error depends on the second derivatives (curvature). If  $|u|_{H^2(\Omega)}$  is zero, i.e. for a linear function, the interpolant is exact.

$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq D h^p |u|_{H^{p+1}(\Omega)}.$$

### Remark

For  $p = 1$ , the interpolation error depends on the second derivatives (curvature). If  $|u|_{H^2(\Omega)}$  is zero, i.e. for a linear function, the interpolant is exact.

### Remark

Error scales like  $h^p$ . If solutions are smooth, increasing  $p$  is better. If not, decreasing  $h$  is better. These can be combined in  $hp$ -adaptive schemes.

$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq Dh^p |u|_{H^{p+1}(\Omega)}.$$

### Remark

For  $p = 1$ , the interpolation error depends on the second derivatives (curvature). If  $|u|_{H^2(\Omega)}$  is zero, i.e. for a linear function, the interpolant is exact.

### Remark

Error scales like  $h^p$ . If solutions are smooth, increasing  $p$  is better. If not, decreasing  $h$  is better. These can be combined in  $hp$ -adaptive schemes.

### Remark

Notice that we require  $u \in H^2(\Omega)$  for  $\text{CG}_1$ , as before.

How do we know if  $u \in H^2(\Omega)$ ? An elliptic regularity result.

How do we know if  $u \in H^2(\Omega)$ ? An elliptic regularity result.

### Theorem (Example elliptic regularity result)

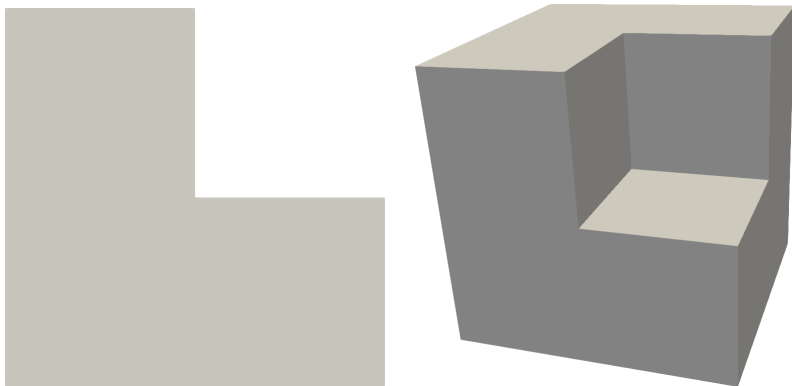
*Let  $\Omega$  be  $C^\infty$ -smooth, i.e.  $\partial\Omega$  possesses a local parametrisation by  $C^\infty$  functions. Then the solution  $u \in H_0^1(\Omega)$  to the Poisson equation is an element of  $H^2(\Omega)$  and satisfies*

$$\|u\|_{H^2(\Omega)} \leq c \|f\|_{L^2(\Omega)}$$

*for some constant  $c$ .*

The requirement that  $\Omega$  has some smoothness is indispensable.

The requirement that  $\Omega$  has some smoothness is indispensable.



With a re-entrant corner,  $u \in H^1(\Omega) \setminus H^2(\Omega)$ .

## Section 4

# Changing norms: Aubin–Nitsche duality



The interpolation error bound also depends on the norm used.

The interpolation error bound also depends on the norm used.

The  $L^2(\Omega)$  norm is *weaker* than the  $H^1(\Omega)$  norm: it only measures how good your approximation of the function is, while the  $H^1(\Omega)$  norm also measures the function and its derivative.

The interpolation error bound also depends on the norm used.

The  $L^2(\Omega)$  norm is *weaker* than the  $H^1(\Omega)$  norm: it only measures how good your approximation of the function is, while the  $H^1(\Omega)$  norm also measures the function and its derivative.

When you measure the interpolation error in a weaker norm, the convergence rate improves:

## Theorem

*Let  $(V_h)_h$  be the function spaces constructed with continuous Lagrange elements of order  $p$  on a shape-regular sequence of meshes  $(\mathcal{M}_h)_h$  indexed by mesh size  $h$ . Let  $u \in H^{p+1}(\Omega)$ , and let  $\mathcal{I}_h : H^{p+1}(\Omega) \rightarrow V_h$  be the interpolation operator associated with each  $V_h$ . Then there exists a constant  $D < \infty$  independent of  $u$  such that*

$$\|u - \mathcal{I}_h u\|_{L^2(\Omega)} \leq D h^{p+1} |u|_{H^{p+1}(\Omega)}.$$

Let's consider solving the Poisson equation. We're not primarily interested in  $\|u - \mathcal{I}_h u\|_{L^2(\Omega)}$ ! We're interested in  $\|u - u_h\|_{L^2(\Omega)}$ .

Let's consider solving the Poisson equation. We're not primarily interested in  $\|u - \mathcal{I}_h u\|_{L^2(\Omega)}$ ! We're interested in  $\|u - u_h\|_{L^2(\Omega)}$ .

Céa's Lemma only tells us about  $\|u - u_h\|_{H^1(\Omega)}$  in terms of  $\|u - \mathcal{I}_h u\|_{H^1(\Omega)}$ . How can we get a hold of estimates of  $\|u - u_h\|_{L^2(\Omega)}$ ?

Let's consider solving the Poisson equation. We're not primarily interested in  $\|u - \mathcal{I}_h u\|_{L^2(\Omega)}$ ! We're interested in  $\|u - u_h\|_{L^2(\Omega)}$ .

Céa's Lemma only tells us about  $\|u - u_h\|_{H^1(\Omega)}$  in terms of  $\|u - \mathcal{I}_h u\|_{H^1(\Omega)}$ . How can we get a hold of estimates of  $\|u - u_h\|_{L^2(\Omega)}$ ?

The Aubin–Nitsche duality argument.

Consider the variational problem

find  $u \in H_0^1(\Omega)$  such that  $a(u, v) = (f, v)_{L^2(\Omega)}$  for all  $v \in H_0^1(\Omega)$ ,

where  $\Omega$  is  $C^\infty$ -smooth and

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

Consider the variational problem

find  $u \in H_0^1(\Omega)$  such that  $a(u, v) = (f, v)_{L^2(\Omega)}$  for all  $v \in H_0^1(\Omega)$ ,

where  $\Omega$  is  $C^\infty$ -smooth and

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

We know that this has a unique solution  $u \in H_0^1(\Omega)$  by the Riesz representation theorem, and that  $u \in H^2(\Omega)$  by elliptic regularity.



Consider the variational problem

find  $u \in H_0^1(\Omega)$  such that  $a(u, v) = (f, v)_{L^2(\Omega)}$  for all  $v \in H_0^1(\Omega)$ ,

where  $\Omega$  is  $C^\infty$ -smooth and

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

We know that this has a unique solution  $u \in H_0^1(\Omega)$  by the Riesz representation theorem, and that  $u \in H^2(\Omega)$  by elliptic regularity.

We thus know that its Galerkin approximation with  $\text{CG}_1$  elements satisfies

$$\begin{aligned} \|u - u_h\|_{H^1(\Omega)} &\leq CD\alpha^{-1}h|u|_{H^2(\Omega)} \\ &\leq cCD\alpha^{-1}h\|f\|_{L^2(\Omega)}. \end{aligned}$$

Consider the error  $e = u - u_h \in H_0^1(\Omega)$ . Given any element of a Hilbert space, we can construct its associated dual element as

$$e^*(v) = (u - u_h, v)_{L^2(\Omega)}.$$

and this makes sense as data for an auxiliary problem (‘adjoint’ or ‘dual’):

find  $w \in H_0^1(\Omega)$  such that  $a(w, v) = e^*(v)$  for all  $v \in H_0^1(\Omega)$ .

Consider the error  $e = u - u_h \in H_0^1(\Omega)$ . Given any element of a Hilbert space, we can construct its associated dual element as

$$e^*(v) = (u - u_h, v)_{L^2(\Omega)}.$$

and this makes sense as data for an auxiliary problem (‘adjoint’ or ‘dual’):

$$\text{find } w \in H_0^1(\Omega) \text{ such that } a(w, v) = e^*(v) \text{ for all } v \in H_0^1(\Omega).$$

By the same arguments we know that this has a unique solution and that there exists  $c$  such that

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

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\|u - u_h\|_{L^2(\Omega)}^2 = (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h)$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}\|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\ &= a(w, u - u_h)\end{aligned}$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}\|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\ &= a(w, u - u_h) \\ &= a(u - u_h, w - \mathcal{I}_h w) \quad (\text{sym, GO})\end{aligned}$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}
 \|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\
 &= a(w, u - u_h) \\
 &= a(u - u_h, w - \mathcal{I}_h w) && \text{(sym, GO)} \\
 &\leq C \|u - u_h\|_{H^1(\Omega)} \|w - \mathcal{I}_h w\|_{H^1(\Omega)} && (a \text{ cts})
 \end{aligned}$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}
 \|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\
 &= a(w, u - u_h) \\
 &= a(u - u_h, w - \mathcal{I}_h w) && \text{(sym, GO)} \\
 &\leq C \|u - u_h\|_{H^1(\Omega)} \|w - \mathcal{I}_h w\|_{H^1(\Omega)} && (a \text{ cts}) \\
 &\leq CDh \|u - u_h\|_{H^1(\Omega)} |w|_{H^2(\Omega)} && (\text{interp})
 \end{aligned}$$



Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}
 \|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\
 &= a(w, u - u_h) \\
 &= a(u - u_h, w - \mathcal{I}_h w) && \text{(sym, GO)} \\
 &\leq C \|u - u_h\|_{H^1(\Omega)} \|w - \mathcal{I}_h w\|_{H^1(\Omega)} && (a \text{ cts}) \\
 &\leq CDh \|u - u_h\|_{H^1(\Omega)} |w|_{H^2(\Omega)} && (\text{interp}) \\
 &\leq C^2 D^2 \alpha^{-1} h^2 |u|_{H^2(\Omega)} |w|_{H^2(\Omega)} && (\text{err est})
 \end{aligned}$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}
 \|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\
 &= a(w, u - u_h) \\
 &= a(u - u_h, w - \mathcal{I}_h w) && \text{(sym, GO)} \\
 &\leq C \|u - u_h\|_{H^1(\Omega)} \|w - \mathcal{I}_h w\|_{H^1(\Omega)} && (a \text{ cts}) \\
 &\leq CDh \|u - u_h\|_{H^1(\Omega)} |w|_{H^2(\Omega)} && (\text{interp}) \\
 &\leq C^2 D^2 \alpha^{-1} h^2 |u|_{H^2(\Omega)} |w|_{H^2(\Omega)} && (\text{err est}) \\
 &\leq C^2 D^2 \alpha^{-1} ch^2 |u|_{H^2(\Omega)} \|u - u_h\|_{L^2(\Omega)} && (\text{ell reg})
 \end{aligned}$$

Now consider  $\|u - u_h\|_{L^2(\Omega)}^2$ , the quantity we wish to bound. We have

$$\begin{aligned}
 \|u - u_h\|_{L^2(\Omega)}^2 &= (u - u_h, u - u_h)_{L^2(\Omega)} = e^*(u - u_h) \\
 &= a(w, u - u_h) \\
 &= a(u - u_h, w - \mathcal{I}_h w) && \text{(sym, GO)} \\
 &\leq C \|u - u_h\|_{H^1(\Omega)} \|w - \mathcal{I}_h w\|_{H^1(\Omega)} && (a \text{ cts}) \\
 &\leq CDh \|u - u_h\|_{H^1(\Omega)} |w|_{H^2(\Omega)} && (\text{interp}) \\
 &\leq C^2 D^2 \alpha^{-1} h^2 |u|_{H^2(\Omega)} |w|_{H^2(\Omega)} && (\text{err est}) \\
 &\leq C^2 D^2 \alpha^{-1} ch^2 |u|_{H^2(\Omega)} \|u - u_h\|_{L^2(\Omega)} && (\text{ell reg})
 \end{aligned}$$

and hence there exists a constant  $C'$  such that

$$\|u - u_h\|_{L^2(\Omega)} \leq C' h^2 |u|_{H^2(\Omega)}$$

as required.

## Section 5

# Interpolation error estimates for other elements

Basic gist: if  $\mathcal{P}_p(K) \subset \mathcal{V}$  and  $\mathcal{P}_{p+1}(K) \not\subset \mathcal{V}$ , expect

$$\|u - \mathcal{I}_h u\|_{H^s(\Omega)} \leq ch^{p+1-s} |u|_{H^{p+1}(\Omega)}$$

for sufficiently regular  $u$ .

Basic gist: if  $\mathcal{P}_p(K) \subset \mathcal{V}$  and  $\mathcal{P}_{p+1}(K) \not\subset \mathcal{V}$ , expect

$$\|u - \mathcal{I}_h u\|_{H^s(\Omega)} \leq ch^{p+1-s} |u|_{H^{p+1}(\Omega)}$$

for sufficiently regular  $u$ .

This is true even for quad/hex elements: since  $\mathcal{P}_1(K) \subsetneq \mathcal{Q}_1(K)$  and  $\mathcal{P}_2(K) \not\subset \mathcal{Q}_1(K)$ , we only get first order in the  $H^1(\Omega)$  norm with CG<sub>1</sub> on quads.

Basic gist: if  $\mathcal{P}_p(K) \subset \mathcal{V}$  and  $\mathcal{P}_{p+1}(K) \not\subset \mathcal{V}$ , expect

$$\|u - \mathcal{I}_h u\|_{H^s(\Omega)} \leq ch^{p+1-s} |u|_{H^{p+1}(\Omega)}$$

for sufficiently regular  $u$ .

This is true even for quad/hex elements: since  $\mathcal{P}_1(K) \subsetneq \mathcal{Q}_1(K)$  and  $\mathcal{P}_2(K) \not\subset \mathcal{Q}_1(K)$ , we only get first order in the  $H^1(\Omega)$  norm with CG<sub>1</sub> on quads.

Let's see how this works with the Argyris element. We have (under the same assumptions)

$$\|u - \mathcal{I}_h u\|_{H^2(\Omega)} \leq Dh^4 |u|_{H^6(\Omega)},$$

$$\|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq Dh^5 |u|_{H^6(\Omega)},$$

$$\|u - \mathcal{I}_h u\|_{H^0(\Omega)} \leq Dh^6 |u|_{H^6(\Omega)}.$$

It can also be of interest to consider what interpolation does to other quantities. For example, suppose we are solving the time-dependent Maxwell's equations, which involve Gauss' law:

$$\nabla \cdot B = 0$$

for the magnetic field  $B \in H(\text{div})$ .



It can also be of interest to consider what interpolation does to other quantities. For example, suppose we are solving the time-dependent Maxwell's equations, which involve Gauss' law:

$$\nabla \cdot B = 0$$

for the magnetic field  $B \in H(\text{div})$ .

If we start with initial data  $B_0$  for the magnetic field that satisfies this, we need to make sure that

$$\nabla \cdot (\mathcal{I}_h B_0) = 0$$

also.

It can also be of interest to consider what interpolation does to other quantities. For example, suppose we are solving the time-dependent Maxwell's equations, which involve Gauss' law:

$$\nabla \cdot B = 0$$

for the magnetic field  $B \in H(\text{div})$ .

If we start with initial data  $B_0$  for the magnetic field that satisfies this, we need to make sure that

$$\nabla \cdot (\mathcal{I}_h B_0) = 0$$

also.

For  $H(\text{div})$ -conforming elements such as the BDM family, we have such a result:

$$\|\nabla \cdot (u - \mathcal{I}_h u)\|_{L^2(\Omega)} \leq Dh^s |\nabla \cdot u|_{H^s(\Omega)}.$$

If  $\nabla \cdot u = 0$ , then this forces  $\nabla \cdot \mathcal{I}_h u = 0$  also.

## C6.4 Finite Element Methods for PDEs

### Lecture 12: Nonlinear problems

Patrick E. Farrell

University of Oxford

So far we have treated *linear, coercive* problems. In these last lectures we will relax both of these assumptions.

To discretise with finite elements, we convert linear PDEs into linear variational problems. It is not a surprise that we will convert nonlinear PDEs into nonlinear variational problems!

Some samples.

Linear PDE	Nonlinear PDE
Stokes	Navier–Stokes

Some samples.

Linear PDE	Nonlinear PDE
Stokes Linear elasticity	Navier–Stokes Hyperelasticity

Some samples.

Linear PDE	Nonlinear PDE
Stokes	Navier–Stokes
Linear elasticity	Hyperelasticity
Maxwell	Magnetohydrodynamics

Some samples.

Linear PDE	Nonlinear PDE
Stokes	Navier–Stokes
Linear elasticity	Hyperelasticity
Maxwell	Magnetohydrodynamics
Schrödinger	Gross–Pitaevskii



Some samples.

Linear PDE	Nonlinear PDE
Stokes	Navier–Stokes
Linear elasticity	Hyperelasticity
Maxwell	Magnetohydrodynamics
Schrödinger	Gross–Pitaevskii
Linearised gravity	Einstein field equations

### Example (Bratu–Gelfand equation for $\lambda \in \mathbb{R}$ )

$$u''(x) + \lambda e^u = 0, \quad u(0) = 0 = u(1).$$

Example (Bratu–Gelfand equation for  $\lambda \in \mathbb{R}$ )

$$u''(x) + \lambda e^u = 0, \quad u(0) = 0 = u(1).$$

We multiply by a test function and integrate by parts: find  $u \in V$  such that

$$-\int_0^1 u'(x)v'(x) \, dx + \int_0^1 \lambda e^u v \, dx = 0 \text{ for all } v \in V,$$

and by inspection we take  $V = H_0^1(0, 1)$ .

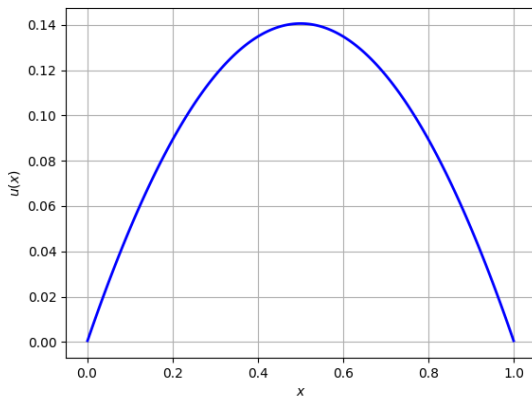
```
from firedrake import *

mesh = UnitIntervalMesh(100)
x = SpatialCoordinate(mesh)[0]
V = FunctionSpace(mesh, "CG", 5)

u = Function(V)
v = TestFunction(V)
lmbda = Constant(1)

G = -inner(grad(u), grad(v))*dx + inner(exp(u), v)*dx
bc = DirichletBC(V, 0, "on_boundary")

u.interpolate(-x*(x-1))
solve(G == 0, u, bc)
```



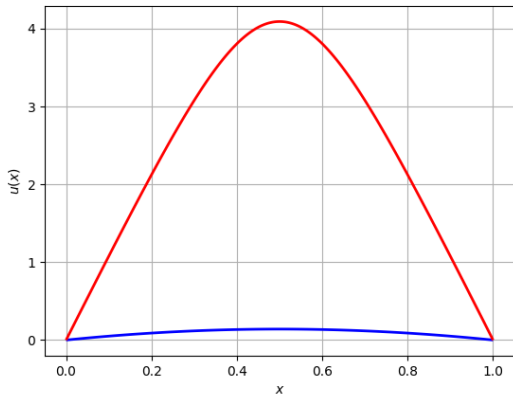
If we change

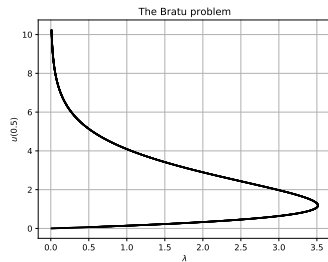
```
u.interpolate(-x*(x-1))  
solve(G == 0, u, bc)
```

to

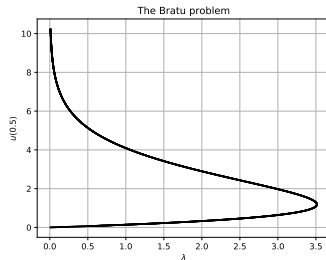
```
u.interpolate(-16*x*(x-1))  
solve(G == 0, u, bc)
```

we get ...









## Solutions

$$\text{number of solutions} = \begin{cases} 2 & \lambda \in (0, \lambda^*), \\ 1 & \lambda \in \{0, \lambda^*\}, \\ 0 & \text{otherwise,} \end{cases}$$

with

$$\lambda^* = 8 \left( \min_{x>0} \frac{x}{\cosh x} \right)^2 \approx 3.5138307.$$

## Section 2

# Variational formulation of nonlinear problems

Our basic abstraction for linear problems was:

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ .

Our basic abstraction for linear problems was:

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

Our abstraction for nonlinear problems will be:

$$\text{find } u \in V \text{ such that } G(u; v) = 0 \text{ for all } v \in V,$$

where  $G : V \times V \rightarrow \mathbb{R}$ .

Our basic abstraction for linear problems was:

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

Our abstraction for nonlinear problems will be:

$$\text{find } u \in V \text{ such that } G(u; v) = 0 \text{ for all } v \in V,$$

where  $G : V \times V \rightarrow \mathbb{R}$ .

We use  $G(u; v)$  to record that  $G$  is nonlinear in  $u$  but always linear in  $v$ .

It's often useful to reformulate the variational statement as an equation, as we did in the linear case.

It's often useful to reformulate the variational statement as an equation, as we did in the linear case.

Define  $H : V \rightarrow V^*$  via

$$(H(u))(v) = \langle H(u), v \rangle = G(u; v).$$

It's often useful to reformulate the variational statement as an equation, as we did in the linear case.

Define  $H : V \rightarrow V^*$  via

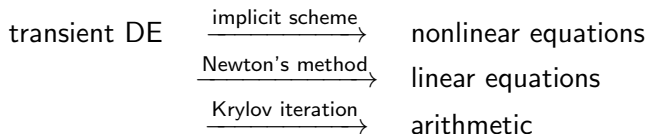
$$(H(u))(v) = \langle H(u), v \rangle = G(u; v).$$

Solutions of our nonlinear variational problem are exactly roots of  $H$ , i.e. we seek  $u \in V$  such that

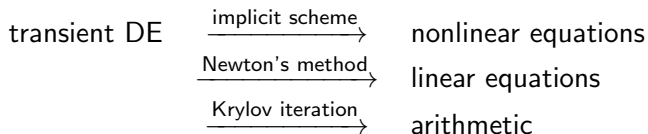
$$H(u) = 0.$$



Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have

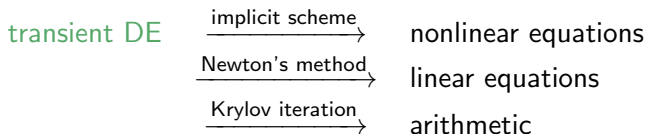


Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have



Now consider solving a stiff time-dependent PDE. We must *also* discretise in space. We could spatially discretise at any level:

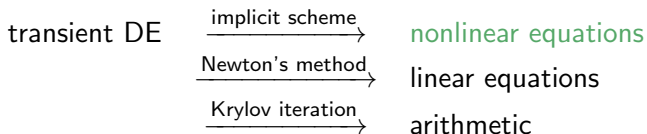
Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have



Now consider solving a stiff time-dependent PDE. We must *also* discretise in space. We could spatially discretise at any level:

- ▶ the time-dependent PDE, yielding a very large ODE system;

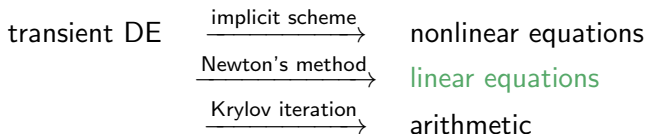
Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have



Now consider solving a stiff time-dependent PDE. We must *also* discretise in space. We could spatially discretise at any level:

- ▶ the time-dependent PDE, yielding a very large ODE system;
- ▶ the nonlinear PDEs arising from time discretisation;

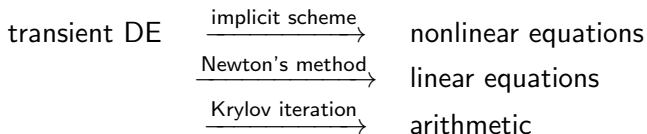
Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have



Now consider solving a stiff time-dependent PDE. We must *also* discretise in space. We could spatially discretise at any level:

- ▶ the time-dependent PDE, yielding a very large ODE system;
- ▶ the nonlinear PDEs arising from time discretisation;
- ▶ the linear PDEs arising from the nonlinear solver.

Algorithms in numerical analysis break down a problem into a sequence of simpler ones. For example, for an ODE, we have



Now consider solving a stiff time-dependent PDE. We must *also* discretise in space. We could spatially discretise at any level:

- ▶ the time-dependent PDE, yielding a very large ODE system;
- ▶ the nonlinear PDEs arising from time discretisation;
- ▶ the linear PDEs arising from the nonlinear solver.

These operations sometimes commute, and sometimes do not.

For a nonlinear PDE, we must

- ▶ devise a scheme to solve the nonlinear problem;
- ▶ discretise in space to make the problem finite-dimensional.

For a nonlinear PDE, we must

- ▶ devise a scheme to solve the nonlinear problem;
- ▶ discretise in space to make the problem finite-dimensional.

We can do these in either order.



## Section 3

Discretise, then solve

We consider the Galerkin approximation:

find  $u_h \in V_h$  such that  $G(u_h; v_h) = 0$  for all  $v_h \in V_h$ .

We consider the Galerkin approximation:

find  $u_h \in V_h$  such that  $G(u_h; v_h) = 0$  for all  $v_h \in V_h$ .

A primary goal of an analysis is bounding the error  $\|u - u_h\|_V$ .

We consider the Galerkin approximation:

$$\text{find } u_h \in V_h \text{ such that } G(u_h; v_h) = 0 \text{ for all } v_h \in V_h.$$

A primary goal of an analysis is bounding the error  $\|u - u_h\|_V$ .

But *which*  $u$ , and which  $u_h$ ? How do we pair the continuous and discrete solutions across mesh levels? How do we know that the discrete problem supports the right number of solutions? How do we know there are no spurious solutions?

We consider the Galerkin approximation:

$$\text{find } u_h \in V_h \text{ such that } G(u_h; v_h) = 0 \text{ for all } v_h \in V_h.$$

A primary goal of an analysis is bounding the error  $\|u - u_h\|_V$ .

But *which*  $u$ , and which  $u_h$ ? How do we pair the continuous and discrete solutions across mesh levels? How do we know that the discrete problem supports the right number of solutions? How do we know there are no spurious solutions?

These are difficult questions; possible to address, but we will sidestep them.

## Section 4

Solve, then discretise

We will develop an algorithm for dealing with the nonlinearity *on the infinite-dimensional level*, the *Newton–Kantorovich* iteration.

We will develop an algorithm for dealing with the nonlinearity *on the infinite-dimensional level*, the *Newton–Kantorovich* iteration.

This will achieve the solution of the nonlinear problem by solving a sequence of linear problems, each of which we then discretise with the finite element method.



We will develop an algorithm for dealing with the nonlinearity *on the infinite-dimensional level*, the *Newton–Kantorovich* iteration.

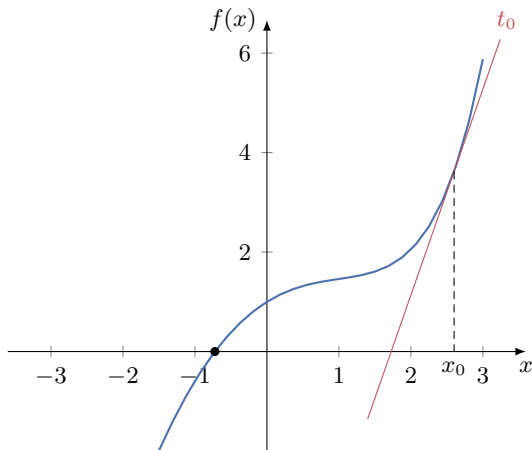
This will achieve the solution of the nonlinear problem by solving a sequence of linear problems, each of which we then discretise with the finite element method.

First, let's recall Newton's method in  $\mathbb{R}$  and  $\mathbb{R}^N$ .

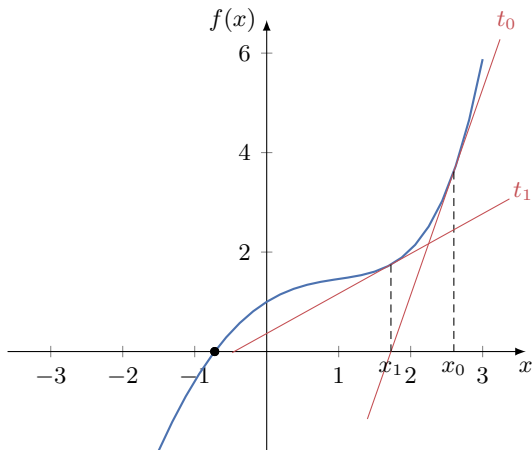
## Section 5

# Newton's method in $\mathbb{R}^N$

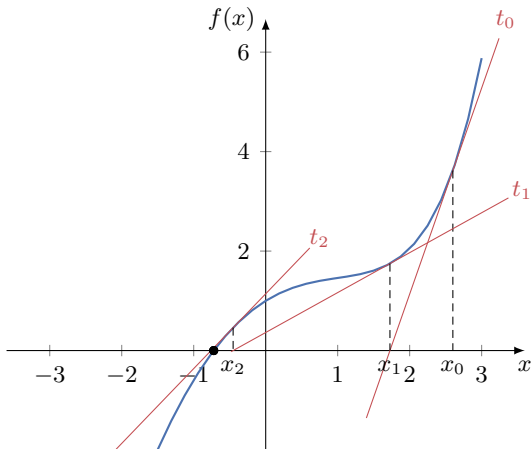
Consider  $f(x) = e^{9x/10} - x^2$  with  $x_0 = 2.6$ .



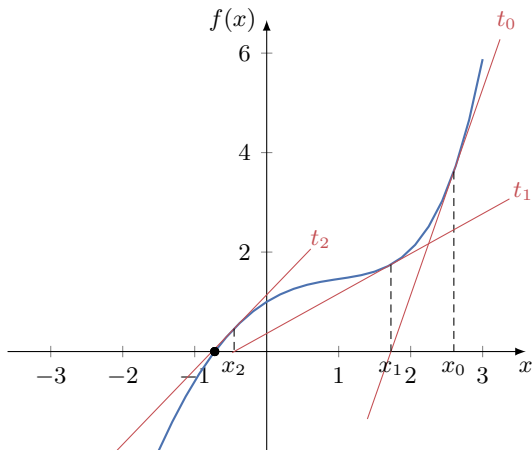
Consider  $f(x) = e^{9x/10} - x^2$  with  $x_0 = 2.6$ .



Consider  $f(x) = e^{9x/10} - x^2$  with  $x_0 = 2.6$ .



Consider  $f(x) = e^{9x/10} - x^2$  with  $x_0 = 2.6$ .



solve  $f'(x_k)\delta x_k = -f(x_k)$ ; update  $x_{k+1} = x_k + \delta x_k$ .

## Termination

The algorithm terminates if  $f(x_k) = 0$ , as desired.

## Invertibility

We require  $f'(x_k)$  to be invertible at every iteration.

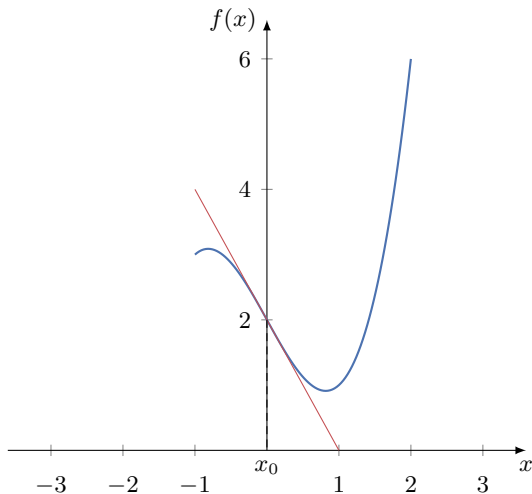
## Poor global convergence

The initial guess matters. With poor initial guesses, Newton's method may diverge to infinity, or get stuck in a cycle.

## Good local convergence

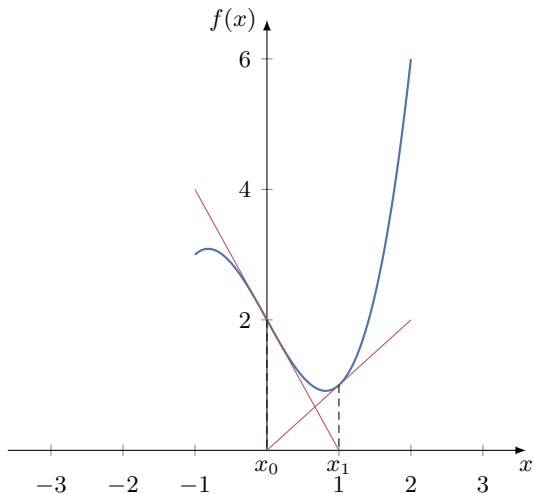
If  $f$  is smooth, the solution is isolated, and the guess close, Newton converges quadratically.

Consider  $f(x) = x^3 - 2x + 2$  with  $x_0 = 0$ .





Consider  $f(x) = x^3 - 2x + 2$  with  $x_0 = 0$ .



This geometric reasoning is hard to generalise to higher dimensions. Let's look at a derivation that *does* extend.

This geometric reasoning is hard to generalise to higher dimensions. Let's look at a derivation that *does* extend.

Consider the Taylor expansion of  $f$  around  $x_k$ :

$$f(x_k + \delta x_k) = f(x_k) + f'(x_k)\delta x_k + \mathcal{O}(\delta x_k^2).$$

This geometric reasoning is hard to generalise to higher dimensions. Let's look at a derivation that *does* extend.

Consider the Taylor expansion of  $f$  around  $x_k$ :

$$f(x_k + \delta x_k) = f(x_k) + f'(x_k)\delta x_k + \mathcal{O}(\delta x_k^2).$$

Linearise the model by ignoring higher-order terms:

$$f(x_k + \delta x) \approx f(x_k) + f'(x_k)\delta x_k$$

and find  $\delta x$  such that  $f(x_k + \delta x) \approx 0$ :

$$0 = f(x_k) + f'(x_k)\delta x_k.$$

This does extend to an  $F \in C^1(\mathbb{R}^N; \mathbb{R}^N)$ . Newton's method is to

$$\text{solve } DF(x_k)\delta x_k = -F(x_k); \text{ update } x_{k+1} = x_k + \delta x_k,$$

where  $DF$  is the Jacobian (Fréchet derivative) of  $F$ .

This does extend to an  $F \in C^1(\mathbb{R}^N; \mathbb{R}^N)$ . Newton's method is to

$$\text{solve } DF(x_k)\delta x_k = -F(x_k); \text{ update } x_{k+1} = x_k + \delta x_k,$$

where  $DF$  is the Jacobian (Fréchet derivative) of  $F$ .

All the previous remarks apply, plus

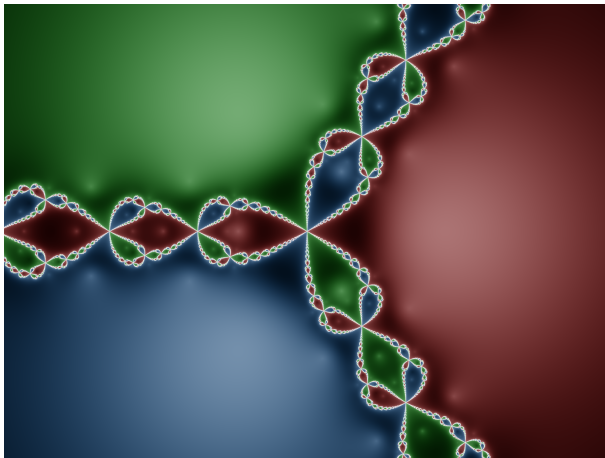
### Affine covariance

Given any nonsingular  $A \in \mathbb{R}^{N \times N}$ , Newton's method applied to  $AF$  yields *the exact same sequence of iterates* as applied to  $F$ , starting from the same initial guess.

We can visualise the erratic global convergence with a *Newton fractal*.

$$f : \mathbb{C} \rightarrow \mathbb{C}$$

$$f(z) = z^3 - 1.$$

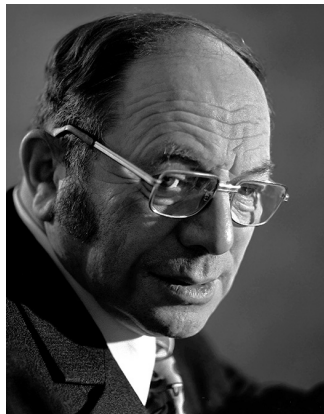


## Section 6

# Newton's method in Banach spaces



- ▶ Invented linear programming (via industrial consultancy!).
- ▶ Instrumental in saving over a million lives during the siege of Leningrad.
- ▶ Involved in the Soviet nuclear bomb project.
- ▶ Nearly sent to the gulag for “shadow prices”.
- ▶ Pseudo-Nobel prize in Economics (1975).



Leonid Kantorovich (1912–1986).

The generalisation of Newton's method to Banach spaces is called the *Newton–Kantorovich* algorithm.

The generalisation of Newton's method to Banach spaces is called the *Newton–Kantorovich* algorithm.

Kantorovich's theorem (1948) is a triumph, fundamental both to nonlinear analysis and applied mathematics. It *does not assume the existence of a solution*: given certain conditions on the residual and initial guess, it *proves* the existence and local uniqueness of a solution.

The generalisation of Newton's method to Banach spaces is called the *Newton–Kantorovich* algorithm.

Kantorovich's theorem (1948) is a triumph, fundamental both to nonlinear analysis and applied mathematics. It *does not assume the existence of a solution*: given certain conditions on the residual and initial guess, it *proves* the existence and local uniqueness of a solution.

With a good initial guess, and great cleverness, it is possible to devise *computer-assisted proofs* of the existence of solutions to infinite-dimensional nonlinear problems.

## Theorem (Kantorovich (1948))

Let  $X$  and  $Y$  be two Banach spaces. Let  $\Omega$  be an open subset of  $X$ , the set where the residual is defined. Let  $H \in C^1(\Omega, Y)$  be the residual of our nonlinear problem, and let  $x_0 \in \Omega$  be an initial guess such that the Fréchet derivative  $H'(x_0)$  is invertible (hence  $H'(x_0) \in L(X; Y)$  and  $H'(x_0)^{-1} \in L(Y; X)$ ). Let  $B(x_0, r)$  denote the open ball of radius  $r$  centred at  $x_0$ .

Assume that there exists a constant  $r > 0$  such that

- (1)  $\overline{B(x_0, r)} \subset \Omega$ ,
- (2)  $\|H'(x_0)^{-1}H(x_0)\|_X \leq \frac{r}{2}$ ,
- (3) For all  $\tilde{x}, x \in B(x_0, r)$ ,

$$\|H'(x_0)^{-1} (H'(\tilde{x}) - H'(x))\|_{L(X; X)} \leq \frac{1}{r} \|\tilde{x} - x\|_X.$$

## Theorem (Kantorovich (1948))

Then

- (1)  $H'(x) \in L(X;Y)$  is invertible at each  $x \in B(x_0, r)$ .
- (2) The Newton sequence  $(x_k)_{k=0}^{\infty}$  defined by

$$x_{k+1} = x_k - H'(x_k)^{-1}H(x_k)$$

is such that  $x_k \in B(x_0, r)$  for all  $k \geq 0$  and converges to a root  $x^* \in \overline{B(x_0, r)}$  of  $H$ .

- (3) For each  $k \geq 0$ ,

$$\|x^* - x_k\|_X \leq \frac{r}{2^k}.$$

- (4) The root  $x^*$  is the locally unique, i.e.  $x^*$  is the only root of  $H$  in the ball  $\overline{B(x_0, r)}$ .

The Bratu–Gelfand equation

$$u''(x) + \lambda e^u = 0, \quad u(0) = 0 = u(1).$$

has variational formulation:  $u \in H_0^1(0, 1)$  such that

$$G(u; v) = - \int_0^1 u'(x) v'(x) \, dx + \int_0^1 \lambda e^u v \, dx = 0$$

for all  $v \in H_0^1(0, 1)$ .

Unwinding the statement of the Newton step, the update  $\delta u_k$  solves

$$G_u(u_k; v, \delta u_k) = -G(u_k; v) \text{ for all } v \in V.$$

Here

$$G(u; v) = - \int_0^1 u'(x)v'(x) \, dx + \int_0^1 \lambda e^u v \, dx$$

so the Newton equation becomes : find  $\delta u_k \in H_0^1(0, 1)$  such that

$$- \int_0^1 \delta u'_k(x)v'(x) \, dx + \int_0^1 \lambda e^{u_k} \delta u_k v \, dx = \int_0^1 u'_k(x)v'(x) \, dx - \int_0^1 \lambda e^{u_k} v \, dx$$

for all  $v \in H_0^1(0, 1)$ .



Unwinding the statement of the Newton step, the update  $\delta u_k$  solves

$$G_u(u_k; v, \delta u_k) = -G(u_k; v) \text{ for all } v \in V.$$

Here

$$G(u; v) = - \int_0^1 u'(x)v'(x) \, dx + \int_0^1 \lambda e^u v \, dx$$

so the Newton equation becomes : find  $\delta u_k \in H_0^1(0, 1)$  such that

$$- \int_0^1 \delta u'_k(x)v'(x) \, dx + \int_0^1 \lambda e^{u_k} \delta u_k v \, dx = \int_0^1 u'_k(x)v'(x) \, dx - \int_0^1 \lambda e^{u_k} v \, dx$$

for all  $v \in H_0^1(0, 1)$ .

This is then discretised with the finite element method.

Many questions remain.

Are the linear equations well-posed? In general they are not coercive. We need a more general theory of well-posedness.

Many questions remain.

Are the linear equations well-posed? In general they are not coercive. We need a more general theory of well-posedness.

We will never compute the exact  $\delta u_k$ ; we can only compute a finite-dimensional approximation. One therefore develops *inexact* Newton–Kantorovich methods.

Many questions remain.

Are the linear equations well-posed? In general they are not coercive. We need a more general theory of well-posedness.

We will never compute the exact  $\delta u_k$ ; we can only compute a finite-dimensional approximation. One therefore develops *inexact* Newton–Kantorovich methods.

How can we control the approximation error of  $\delta u_k$ ? *Adaptive* discretisations.

Many questions remain.

Are the linear equations well-posed? In general they are not coercive. We need a more general theory of well-posedness.

We will never compute the exact  $\delta u_k$ ; we can only compute a finite-dimensional approximation. One therefore develops *inexact* Newton–Kantorovich methods.

How can we control the approximation error of  $\delta u_k$ ? *Adaptive* discretisations.

How can the algorithm be globalised? Line searches, trust regions, continuation in parameter, continuation in mesh.

## C6.4 Finite Element Methods for PDEs

### Lecture 13: Noncoercive problems

Patrick E. Farrell

University of Oxford

So far we have treated *coercive* problems. This means that the bilinear form  $a(u, v)$  in the linear variational problem we are trying to solve

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

satisfies

So far we have treated *coercive* problems. This means that the bilinear form  $a(u, v)$  in the linear variational problem we are trying to solve

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

satisfies

$$a(u, u) \geq \alpha \|u\|_V^2$$

for some  $\alpha > 0$ .



So far we have treated *coercive* problems. This means that the bilinear form  $a(u, v)$  in the linear variational problem we are trying to solve

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

satisfies

$$a(u, u) \geq \alpha \|u\|_V^2$$

for some  $\alpha > 0$ .

Recall that the best constant  $\alpha$  satisfying the definition is given by

$$\alpha := \inf_{\substack{u \in V \\ u \neq 0}} \frac{a(u, u)}{\|u\|_V^2}.$$

We now consider *noncoercive* problems, one for which no such  $\alpha > 0$  exists. We will develop more general (necessary and sufficient) criteria for well-posedness of the linear variational problem, the so-called inf-sup or Babuška conditions.

We now consider *noncoercive* problems, one for which no such  $\alpha > 0$  exists. We will develop more general (necessary and sufficient) criteria for well-posedness of the linear variational problem, the so-called inf-sup or Babuška conditions.

For coercive problems, well-posedness is inherited for  $V_h \subset V$ . This is not true for noncoercive problems. Well-posedness is **not inherited for arbitrary**  $V_h \subset V$ . One must prove the stability of each candidate discretisation individually.

## Mixed Poisson (lecture 5)

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot v u - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} f w \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

Here  $a(0, u; 0, u) = 0$  for all  $u \in L^2(\Omega)$ .

## Mixed Poisson (lecture 5)

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot v u - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} f w \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

Here  $a(0, u; 0, u) = 0$  for all  $u \in L^2(\Omega)$ .

## Mixed linear elasticity (lecture 7)

Find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx = \int_{\Omega} f \cdot v \, dx,$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

Here  $a(0, p; 0, p) \leq 0$  for all  $p \in L^2(\Omega)$ .

Not all examples have multiple variables (are *mixed* problems). For example, the (“bad”) Helmholtz equation

$$\begin{aligned} -\nabla^2 u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega, \end{aligned}$$

is well-posed if  $k^2$  is not an eigenvalue of the Dirichlet Laplacian, but is not coercive for  $k$  large enough. For  $k^2$  to be an eigenvalue of the Dirichlet Laplacian, it means that there exists  $u \neq 0$  such that

$$\begin{aligned} -\nabla^2 u &= k^2 u && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega, \end{aligned}$$

i.e.  $-\nabla^2 - k^2 I$  has a nontrivial kernel.

## Section 2

# Babuška's theorem

We will first state the theorem, so we know where we are going, and then we will try to build intuition about where the conditions come from.



## Theorem (Babuška's theorem)

*Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that*

## Theorem (Babuška's theorem)

*Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that*

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

## Theorem (Babuška's theorem)

Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that

(1)  $|a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2}$  for all  $u \in V_1, v \in V_2$ ;

(2)  $\gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$

## Theorem (Babuška's theorem)

Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

# Theorem (Babuška's theorem)

Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

Then for all  $F \in V_2^*$  there exists exactly one element  $u \in V_1$  such that

$$a(u, v) = F(v) \text{ for all } v \in V_2.$$

## Theorem (Babuška's theorem)

Let  $V_1$  and  $V_2$  be two Hilbert spaces. Let  $a : V_1 \times V_2 \rightarrow \mathbb{R}$  be a bilinear form for which there exist constants  $C < \infty$ ,  $\gamma > 0$ ,  $\gamma' > 0$  such that

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

Then for all  $F \in V_2^*$  there exists exactly one element  $u \in V_1$  such that

$$a(u, v) = F(v) \text{ for all } v \in V_2.$$

Furthermore the problem is stable in that

$$\|u\|_{V_1} \leq \frac{\|F\|_{V_2^*}}{\gamma}.$$

As a first example of how to manipulate inf-sup conditions, let's show that a coercive problem satisfies the inf-sup conditions.

As a first example of how to manipulate inf-sup conditions, let's show that a coercive problem satisfies the inf-sup conditions. Suppose  $a(u, v)$  satisfies

$$\alpha \|u\|_V^2 \leq a(u, u) \quad \text{for all } u \in V.$$

Dividing both sides of the inequality by  $\|u\|_V$  for  $u \neq 0$ , we have

$$\alpha \|u\|_V \leq \frac{a(u, u)}{\|u\|_V}$$



As a first example of how to manipulate inf-sup conditions, let's show that a coercive problem satisfies the inf-sup conditions. Suppose  $a(u, v)$  satisfies

$$\alpha \|u\|_V^2 \leq a(u, u) \quad \text{for all } u \in V.$$

Dividing both sides of the inequality by  $\|u\|_V$  for  $u \neq 0$ , we have

$$\begin{aligned} \alpha \|u\|_V &\leq \frac{a(u, u)}{\|u\|_V} \\ &\leq \sup_{\substack{v \in V \\ v \neq 0}} \frac{a(u, v)}{\|v\|_V}. \end{aligned}$$

As a first example of how to manipulate inf-sup conditions, let's show that a coercive problem satisfies the inf-sup conditions. Suppose  $a(u, v)$  satisfies

$$\alpha \|u\|_V^2 \leq a(u, u) \quad \text{for all } u \in V.$$

Dividing both sides of the inequality by  $\|u\|_V$  for  $u \neq 0$ , we have

$$\begin{aligned} \alpha \|u\|_V &\leq \frac{a(u, u)}{\|u\|_V} \\ &\leq \sup_{\substack{v \in V \\ v \neq 0}} \frac{a(u, v)}{\|v\|_V}. \end{aligned}$$

Infimising over  $u \neq 0$ , we have

$$0 < \alpha \leq \inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{a(u, v)}{\|u\|_V \|v\|_V}.$$

So the coercivity constant  $\alpha$  works for  $\gamma$  (and  $\gamma'$ ).

## Section 3

# Understanding the inf-sup conditions

Remember that we often rewrite

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

as

$$\text{find } u \in V \text{ such that } Au = F,$$

where  $A : V \rightarrow V^*$ ,  $(Au)(v) := a(u, v)$ .

Remember that we often rewrite

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

as

$$\text{find } u \in V \text{ such that } Au = F,$$

where  $A : V \rightarrow V^*$ ,  $(Au)(v) := a(u, v)$ .

The inf-sup conditions are just variational ways of stating facts about  $A$ .

Remember that we often rewrite

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

as

$$\text{find } u \in V \text{ such that } Au = F,$$

where  $A : V \rightarrow V^*$ ,  $(Au)(v) := a(u, v)$ .

The inf-sup conditions are just variational ways of stating facts about  $A$ .

The key to unpacking the inf-sup conditions is to recall the norm on the dual of a Hilbert space:

$$\|j\|_{V^*} = \sup_{\|u\|_V=1} |j(u)| = \sup_{\substack{u \in V \\ u \neq 0}} \frac{|j(u)|}{\|u\|_V}.$$

Remember that we often rewrite

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V$$

as

$$\text{find } u \in V \text{ such that } Au = F,$$

where  $A : V \rightarrow V^*$ ,  $(Au)(v) := a(u, v)$ .

The inf-sup conditions are just variational ways of stating facts about  $A$ .

The key to unpacking the inf-sup conditions is to recall the norm on the dual of a Hilbert space:

$$\|j\|_{V^*} = \sup_{\|u\|_V=1} |j(u)| = \sup_{\substack{u \in V \\ u \neq 0}} \frac{|j(u)|}{\|u\|_V}.$$

In our explanations we'll always take  $V_1 = V_2$ , but the theorem applies more generally.

Our exposition of the inf-sup conditions will follow Brezzi and Bathe. This starts off gently, looking at the stability of a linear system of equations to perturbation.

Consider the  $N$ -dimensional linear system

$$\text{find } x \in \mathbb{R}^N \text{ such that } Mx = b,$$

arising from a Galerkin discretisation of

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

We know from linear algebra that the linear system  $Mx = b$  has a unique solution for all  $b \in \mathbb{R}^N$  if and only if the associated homogeneous problem  $Mx = 0$  has only one solution,  $x = 0$ . Let us suppose for now that this is the case.



Now suppose that a perturbation  $\delta b$  is applied to the right hand side, causing a change  $\delta x$  in the solution.

Now suppose that a perturbation  $\delta b$  is applied to the right hand side, causing a change  $\delta x$  in the solution.

For the system to be stable, we hope that a small change in  $b$  will induce only a small change in  $x$ . To measure this, we introduce norms.

Now suppose that a perturbation  $\delta b$  is applied to the right hand side, causing a change  $\delta x$  in the solution.

For the system to be stable, we hope that a small change in  $b$  will induce only a small change in  $x$ . To measure this, we introduce norms.

Since we know  $x$  and  $b$  represent the coefficients of (approximations to)  $u \in V$  and  $F \in V^*$  respectively, let's use the infinite-dimensional norms. We therefore equip  $x$  and  $b$  with the norms  $\|\cdot\|_V$  and  $\|\cdot\|_{V^*}$  respectively.

## Definition (Stability constant)

The stability constant of  $M$  with respect to the norms  $\|\cdot\|_V$ ,  $\|\cdot\|_{V^*}$  is the smallest possible constant  $S$  such that

$$\frac{\|\delta x\|_V}{\|x\|_V} \leq S \frac{\|\delta b\|_{V^*}}{\|b\|_{V^*}}$$

for all vectors  $x$  and  $\delta x$  in  $\mathbb{R}^N$  such that  $Mx = b$  and  $M\delta x = \delta b$ .

Such a stability constant always exists if the matrix is invertible. However, if we consider a *sequence* of linear systems with increasing dimension  $N$  (corresponding to a finer and finer discretisation) it might be the case that the associated constants ( $S$ ) depend on  $N$  and become infinitely large as  $N \rightarrow \infty$ . We thus say that a sequence of linear systems is stable with respect to the norms  $\|\cdot\|_V$ ,  $\|\cdot\|_{V^*}$  if the sequence of stability constants is bounded.

We can use matrix (operator) norms to clarify the nature of the stability constant. Define

$$\|M\| = \sup_{\substack{y \in V \\ y \neq 0}} \frac{\|My\|_{V^*}}{\|y\|_V},$$

where we denote the input space by  $V$ .

We can use matrix (operator) norms to clarify the nature of the stability constant. Define

$$\|M\| = \sup_{\substack{y \in V \\ y \neq 0}} \frac{\|My\|_{V^*}}{\|y\|_V},$$

where we denote the input space by  $V$ . Choosing  $y = x$ ,  $My = b$ , we have

$$\|M\| \geq \frac{\|b\|_{V^*}}{\|x\|_V}.$$

This implies that

$$\|M\| \frac{\|x\|_V}{\|b\|_{V^*}} \geq 1.$$

Now let us consider the inverse norm. We have

$$\|M^{-1}\| = \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}},$$

where we denote the output space by  $V^*$ .

Now let us consider the inverse norm. We have

$$\|M^{-1}\| = \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}},$$

where we denote the output space by  $V^*$ . Choosing  $z = \delta b$ ,  $M^{-1}z = \delta x$ , we have

$$\|M^{-1}\| \geq \frac{\|\delta x\|_V}{\|\delta b\|_{V^*}}.$$



Now let us consider the inverse norm. We have

$$\|M^{-1}\| = \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}},$$

where we denote the output space by  $V^*$ . Choosing  $z = \delta b$ ,  $M^{-1}z = \delta x$ , we have

$$\|M^{-1}\| \geq \frac{\|\delta x\|_V}{\|\delta b\|_{V^*}}.$$

Multiplying the upper bound by a quantity greater than one will not change the inequality, so

$$\frac{\|\delta x\|_V}{\|\delta b\|_{V^*}} \leq \|M^{-1}\| \|M\| \frac{\|x\|_V}{\|b\|_{V^*}}$$

which implies

$$\frac{\|\delta x\|_V}{\|x\|_V} \leq \|M\| \|M^{-1}\| \frac{\|\delta b\|_{V^*}}{\|b\|_{V^*}}.$$

Now let us consider the inverse norm. We have

$$\|M^{-1}\| = \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}},$$

where we denote the output space by  $V^*$ . Choosing  $z = \delta b$ ,  $M^{-1}z = \delta x$ , we have

$$\|M^{-1}\| \geq \frac{\|\delta x\|_V}{\|\delta b\|_{V^*}}.$$

Multiplying the upper bound by a quantity greater than one will not change the inequality, so

$$\frac{\|\delta x\|_V}{\|\delta b\|_{V^*}} \leq \|M^{-1}\| \|M\| \frac{\|x\|_V}{\|b\|_{V^*}}$$

which implies

$$\frac{\|\delta x\|_V}{\|x\|_V} \leq \|M\| \|M^{-1}\| \frac{\|\delta b\|_{V^*}}{\|b\|_{V^*}}.$$

Since  $x$  and  $\delta x$  are arbitrary,

$$S = \|M\| \|M^{-1}\|.$$

This is exactly the condition number of the matrix. In a linear algebra course this is usually specialised to the Euclidean case

$$\|\cdot\|_V = \|\cdot\|_{V^*} = \|\cdot\|_{\ell_2},$$

where the condition number resolves to the ratio of the largest and smallest singular values.

This is exactly the condition number of the matrix. In a linear algebra course this is usually specialised to the Euclidean case

$$\|\cdot\|_V = \|\cdot\|_{V^*} = \|\cdot\|_{\ell_2},$$

where the condition number resolves to the ratio of the largest and smallest singular values.

For the stability of our problem to be bounded, we will require both  $\|M\|$  and  $\|M^{-1}\|$  to be bounded from above. We will now consider them in turn.

Assume  $a$  is bounded: that is, there exists a constant  $C$  such that

$$|a(u, v)| \leq C \|u\|_V \|v\|_V,$$

and since our matrix encodes the action of the bilinear form, we have

$$|y^T Mx| \leq C \|y\|_V \|x\|_V.$$

The forward operator norm  $\|M\|$  is exactly the continuity constant  $C$  in disguise. To see this, expand the definitions:

$$\|M\| = \sup_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V}$$

Assume  $a$  is bounded: that is, there exists a constant  $C$  such that

$$|a(u, v)| \leq C \|u\|_V \|v\|_V,$$

and since our matrix encodes the action of the bilinear form, we have

$$|y^T Mx| \leq C \|y\|_V \|x\|_V.$$

The forward operator norm  $\|M\|$  is exactly the continuity constant  $C$  in disguise. To see this, expand the definitions:

$$\begin{aligned} \|M\| &= \sup_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V} \\ &= \sup_{\substack{x \in V \\ x \neq 0}} \left\{ \frac{1}{\|x\|_V} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \right\} \end{aligned}$$

Assume  $a$  is bounded: that is, there exists a constant  $C$  such that

$$|a(u, v)| \leq C \|u\|_V \|v\|_V,$$

and since our matrix encodes the action of the bilinear form, we have

$$|y^T Mx| \leq C \|y\|_V \|x\|_V.$$

The forward operator norm  $\|M\|$  is exactly the continuity constant  $C$  in disguise. To see this, expand the definitions:

$$\begin{aligned} \|M\| &= \sup_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V} \\ &= \sup_{\substack{x \in V \\ x \neq 0}} \left\{ \frac{1}{\|x\|_V} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \right\} \\ &= \sup_{\substack{x, y \in V \\ x, y \neq 0}} \frac{|y^T Mx|}{\|x\|_V \|y\|_V} = C. \end{aligned}$$

If  $M$  arises from a conforming discretisation of a continuous bilinear form, then  $\|M\|$  is bounded above by the continuity constant of the form. Thus, the sequence of problems will be stable if and only if the inverse operator norm  $\|M^{-1}\|$  is bounded above.



## Lemma (Characterising the inverse operator norm)

*Let  $M \in \mathbb{R}^{N \times N}$  be nonsingular, and let  $\|\cdot\|_V$  be the norm for its input space, and let the associated dual norm be used for its output space. Then*

$$\|M^{-1}\|^{-1} = \inf_{\substack{x \in V \\ x \neq 0}} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T M x|}{\|y\|_V \|x\|_V}.$$

Proof.

$$\|M^{-1}\|^{-1} = \left( \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}} \right)^{-1}$$

Proof.

$$\begin{aligned}\|M^{-1}\|^{-1} &= \left( \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}} \right)^{-1} \\ &= \inf_{\substack{z \in V^* \\ z \neq 0}} \frac{\|z\|_{V^*}}{\|M^{-1}z\|_V}\end{aligned}$$

## Proof.

$$\begin{aligned}\|M^{-1}\|^{-1} &= \left( \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}} \right)^{-1} \\ &= \inf_{\substack{z \in V^* \\ z \neq 0}} \frac{\|z\|_{V^*}}{\|M^{-1}z\|_V} \\ &= \inf_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V} \quad (\text{set } z = Mx)\end{aligned}$$

## Proof.

$$\begin{aligned}
\|M^{-1}\|^{-1} &= \left( \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}} \right)^{-1} \\
&= \inf_{\substack{z \in V^* \\ z \neq 0}} \frac{\|z\|_{V^*}}{\|M^{-1}z\|_V} \\
&= \inf_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V} \quad (\text{set } z = Mx) \\
&= \inf_{\substack{x \in V \\ x \neq 0}} \left\{ \frac{1}{\|x\|_V} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \right\} (\text{defn of dual norm})
\end{aligned}$$

## Proof.

$$\begin{aligned}
\|M^{-1}\|^{-1} &= \left( \sup_{\substack{z \in V^* \\ z \neq 0}} \frac{\|M^{-1}z\|_V}{\|z\|_{V^*}} \right)^{-1} \\
&= \inf_{\substack{z \in V^* \\ z \neq 0}} \frac{\|z\|_{V^*}}{\|M^{-1}z\|_V} \\
&= \inf_{\substack{x \in V \\ x \neq 0}} \frac{\|Mx\|_{V^*}}{\|x\|_V} \quad (\text{set } z = Mx) \\
&= \inf_{\substack{x \in V \\ x \neq 0}} \left\{ \frac{1}{\|x\|_V} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \right\} \quad (\text{defn of dual norm}) \\
&= \inf_{\substack{x \in V \\ x \neq 0}} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V \|x\|_V}.
\end{aligned}$$

The first step in using inf-sup conditions is to unpack the inf. So rewrite

$$0 < \gamma \leq \inf_{\substack{x \in V \\ x \neq 0}} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V \|x\|_V}$$

as

$$\gamma \|x\|_V \leq \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \quad \text{for all } x \in V.$$

The first step in using inf-sup conditions is to unpack the inf. So rewrite

$$0 < \gamma \leq \inf_{\substack{x \in V \\ x \neq 0}} \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V \|x\|_V}$$

as

$$\gamma \|x\|_V \leq \sup_{\substack{y \in V \\ y \neq 0}} \frac{|y^T Mx|}{\|y\|_V} \quad \text{for all } x \in V.$$

Note that we can get rid of the absolute values:

$$\gamma \|x\|_V \leq \sup_{\substack{y \in V \\ y \neq 0}} \frac{y^T Mx}{\|y\|_V} \quad \text{for all } x \in V,$$

since if the supremum is reached with a negative value, negating the sequence of  $y$ 's will make the fraction positive (so it wasn't the supremum).



Thus, for the sequence of problems to be stable, we need that  $\|M^{-1}\|^{-1}$  to be bounded below, and so we require a constant  $\gamma \in \mathbb{R}$  such that

$$\inf_{\substack{x \in V \\ x \neq 0}} \sup_{\substack{y \in V \\ y \neq 0}} \frac{y^T M x}{\|y\|_V \|x\|_V} \geq \gamma > 0.$$

Since the matrix  $M$  encodes the bilinear form, the analogous condition for the infinite-dimensional problem is: there exists  $\gamma \in \mathbb{R}$  such that

$$\inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{a(u, v)}{\|u\|_V \|v\|_V} \geq \gamma > 0.$$

This is the famous “inf-sup” condition of Babuška. It says that (if it exists) the inverse of  $A : V \rightarrow V^*$  is bounded, where  $(Au)(v) := a(u, v)$ .

## Section 4

# The inf-sup condition and the kernel

Suppose  $A : V \rightarrow V^*$  satisfies an inf-sup condition, i.e.

$$\inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|u\|_V \|v\|_V} \geq \gamma > 0.$$

Then

$$\sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|v\|_V} \geq \gamma \|u\|_V \text{ for all } u \in V.$$

Suppose  $A : V \rightarrow V^*$  satisfies an inf-sup condition, i.e.

$$\inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|u\|_V \|v\|_V} \geq \gamma > 0.$$

Then

$$\sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|v\|_V} \geq \gamma \|u\|_V \text{ for all } u \in V.$$

The quantity on the left is just the dual norm  $\|Au\|_{V^*}$ , so we have

$$\|Au\|_{V^*} \geq \gamma \|u\|_V \geq 0 \text{ for all } u \in V.$$

Suppose  $A : V \rightarrow V^*$  satisfies an inf-sup condition, i.e.

$$\inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|u\|_V \|v\|_V} \geq \gamma > 0.$$

Then

$$\sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|v\|_V} \geq \gamma \|u\|_V \text{ for all } u \in V.$$

The quantity on the left is just the dual norm  $\|Au\|_{V^*}$ , so we have

$$\|Au\|_{V^*} \geq \gamma \|u\|_V \geq 0 \text{ for all } u \in V.$$

What is the kernel of  $A$ ? If  $Au = 0$ , then  $\|Au\|_{V^*} = 0$  and hence  $u = 0$ .

Suppose  $A : V \rightarrow V^*$  satisfies an inf-sup condition, i.e.

$$\inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|u\|_V \|v\|_V} \geq \gamma > 0.$$

Then

$$\sup_{\substack{v \in V \\ v \neq 0}} \frac{|Au(v)|}{\|v\|_V} \geq \gamma \|u\|_V \text{ for all } u \in V.$$

The quantity on the left is just the dual norm  $\|Au\|_{V^*}$ , so we have

$$\|Au\|_{V^*} \geq \gamma \|u\|_V \geq 0 \text{ for all } u \in V.$$

What is the kernel of  $A$ ? If  $Au = 0$ , then  $\|Au\|_{V^*} = 0$  and hence  $u = 0$ .

So the inf-sup condition guarantees the injectivity of  $A$ .

## Section 5

### Necessity of the inf-sup condition

If the operator equation

$$Au = F$$

is well-posed,  $A^{-1}$  exists.



If the operator equation

$$Au = F$$

is well-posed,  $A^{-1}$  exists.

If the problem is stable,  $\|A^{-1}\| < \gamma^{-1}$  for some  $\gamma < \infty$ , so

$$0 < \gamma = \|A^{-1}\|^{-1} = \inf_{\substack{u \in V \\ u \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{a(u, v)}{\|u\|_V \|v\|_V}.$$

## Section 6

### The transpose condition

Conditions for Babuška's theorem:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

Conditions for Babuška's theorem:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

For square linear systems in finite dimensions, injectivity is equivalent to surjectivity is equivalent to bijectivity. So if you satisfy an inf-sup condition on a finite-dimensional space  $V_1 = V_2$ , one inf-sup condition is enough.

Conditions for Babuška's theorem:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

For square linear systems in finite dimensions, injectivity is equivalent to surjectivity is equivalent to bijectivity. So if you satisfy an inf-sup condition on a finite-dimensional space  $V_1 = V_2$ , one inf-sup condition is enough.

For  $V_1 \neq V_2$  or infinite-dimensional problems we need (3).

Consider a linear system with a rectangular matrix  $M \in \mathbb{R}^{P \times N}$ :

$$Mx = b.$$

Consider a linear system with a rectangular matrix  $M \in \mathbb{R}^{P \times N}$ :

$$Mx = b.$$

We know that this cannot be well-posed for arbitrary  $b \in \mathbb{R}^P$ , so we want the Babuška conditions to fail.

Consider a linear system with a rectangular matrix  $M \in \mathbb{R}^{P \times N}$ :

$$Mx = b.$$

We know that this cannot be well-posed for arbitrary  $b \in \mathbb{R}^P$ , so we want the Babuška conditions to fail.

If  $M$  is underdetermined (fat,  $P < N$ ), rank-nullity guarantees there exist  $x \in \mathbb{R}^N$  such that  $Mx = 0$ , so inf-sup (2) fails.



Consider a linear system with a rectangular matrix  $M \in \mathbb{R}^{P \times N}$ :

$$Mx = b.$$

We know that this cannot be well-posed for arbitrary  $b \in \mathbb{R}^P$ , so we want the Babuška conditions to fail.

If  $M$  is underdetermined (fat,  $P < N$ ), rank-nullity guarantees there exist  $x \in \mathbb{R}^N$  such that  $Mx = 0$ , so inf-sup (2) fails.

What if  $M$  is overdetermined (tall and skinny,  $P > N$ )?

Take

$$M = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}.$$

Take

$$M = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}.$$

The first step is to unpack. We want to show that there exists  $\gamma$  such that

$$\gamma \|x\| \leq \sup_{\substack{y \in \mathbb{R}^3 \\ y \neq 0}} \frac{y^T Mx}{\|y\|} \text{ for all } x \in \mathbb{R}^2.$$

For  $x = [x_1, x_2]^T$ , take  $\tilde{y} = [x_1, x_2, x_1 + x_2]^T = Mx$  to get

$$\sup_{\substack{y \in \mathbb{R}^3 \\ y \neq 0}} \frac{y^T Mx}{\|y\|} \geq \frac{\tilde{y}^T \tilde{y}}{\|\tilde{y}\|} = \|\tilde{y}\| \geq \|x\|$$

so we can take  $\gamma = 1$ .

So the associated bilinear form satisfies the inf-sup condition!

So the associated bilinear form satisfies the inf-sup condition!

We thus need a third condition: that the *nullspace of the transpose* is trivial. That is, we also require that  $y^T M = 0 \implies y = 0$ . Recall that the fundamental theorem of linear algebra tells us that

$$\text{range}(M) = \text{kernel}(M^T)^\perp,$$

that is, the range of  $M$  is the orthogonal complement of the nullspace of  $M^T$ . Thus, in order for the operator to be surjective (and have a complete range), we must therefore require the nullspace of  $M^T$  to be trivial. This is the condition that fails in this case; for example, choose  $y^T = [1, 0, -1]$ .

Let's now distinguish  $u \in V_1$  from  $v \in V_2$ . We've already introduced

$$A : V_1 \rightarrow V_2^*, \quad (Au)(v) := a(u, v).$$

Let's now distinguish  $u \in V_1$  from  $v \in V_2$ . We've already introduced

$$A : V_1 \rightarrow V_2^*, \quad (Au)(v) := a(u, v).$$

Now consider

$$A^* : V_2 \rightarrow V_1^*, \quad (A^*v)(u) := a(u, v).$$

Let's now distinguish  $u \in V_1$  from  $v \in V_2$ . We've already introduced

$$A : V_1 \rightarrow V_2^*, \quad (Au)(v) := a(u, v).$$

Now consider

$$A^* : V_2 \rightarrow V_1^*, \quad (A^*v)(u) := a(u, v).$$

These really are adjoints:

$$\langle Au, v \rangle = a(u, v) = \langle A^*v, u \rangle.$$



Let's now distinguish  $u \in V_1$  from  $v \in V_2$ . We've already introduced

$$A : V_1 \rightarrow V_2^*, \quad (Au)(v) := a(u, v).$$

Now consider

$$A^* : V_2 \rightarrow V_1^*, \quad (A^*v)(u) := a(u, v).$$

These really are adjoints:

$$\langle Au, v \rangle = a(u, v) = \langle A^*v, u \rangle.$$

The *closed range theorem* of functional analysis guarantees that if  $\ker(A^*)$  is trivial, then  $A$  is surjective. This gives the “transpose” inf-sup condition:

$$0 < \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

## Section 7

### Review

The conditions again:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

The conditions again:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

To summarise: (1) is the boundedness of the operator. (2) is the injectivity of  $A$ . (3) is the surjectivity of  $A$ .

The conditions again:

$$(1) \quad |a(u, v)| \leq C \|u\|_{V_1} \|v\|_{V_2} \text{ for all } u \in V_1, v \in V_2;$$

$$(2) \quad \gamma \leq \inf_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}};$$

$$(3) \quad \gamma' \leq \inf_{\substack{v \in V_2 \\ v \neq 0}} \sup_{\substack{u \in V_1 \\ u \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}}.$$

To summarise: (1) is the boundedness of the operator. (2) is the injectivity of  $A$ . (3) is the surjectivity of  $A$ .

Note that for symmetry we could rewrite (1) as

$$\sup_{\substack{u \in V_1 \\ u \neq 0}} \sup_{\substack{v \in V_2 \\ v \neq 0}} \frac{a(u, v)}{\|u\|_{V_1} \|v\|_{V_2}} \leq C.$$

so we have a “sup-sup” condition and two “inf-sup” conditions.

## Section 8

# Discretisation and quasioptimality

Start with

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ ,

and take the Galerkin approximation over closed  $V_h \subset V$ :

find  $u_h \in V_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in V_h$ .

Start with

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ ,

and take the Galerkin approximation over closed  $V_h \subset V$ :

find  $u_h \in V_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in V_h$ .

Note that Galerkin orthogonality still holds.



Start with

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ ,

and take the Galerkin approximation over closed  $V_h \subset V$ :

find  $u_h \in V_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in V_h$ .

Note that Galerkin orthogonality still holds.

Is the discrete problem well-posed?

Let's check the Babuška conditions.

Satisfaction of (1) is inherited. What about (2)? (We don't need to check (3) in this case! The discrete system is square and finite-dimensional, so (2)  $\iff$  (3) by rank-nullity.) That is, does there exist  $\tilde{\gamma}$  such that

$$\inf_{\substack{u_h \in V_h \\ u_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{a(u_h, v_h)}{\|u_h\|_V \|v_h\|_V} \geq \tilde{\gamma} > 0,$$

with  $\tilde{\gamma}$  independent of the mesh size  $h$ ?

Let's check the Babuška conditions.

Satisfaction of (1) is inherited. What about (2)? (We don't need to check (3) in this case! The discrete system is square and finite-dimensional, so (2)  $\iff$  (3) by rank-nullity.) That is, does there exist  $\tilde{\gamma}$  such that

$$\inf_{\substack{u_h \in V_h \\ u_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{a(u_h, v_h)}{\|u_h\|_V \|v_h\|_V} \geq \tilde{\gamma} > 0,$$

with  $\tilde{\gamma}$  independent of the mesh size  $h$ ?

**No!** Examples later.

## Theorem

*Assume we have a well-posed discretisation of a well-posed problem. Then*

$$\|u - u_h\|_V \leq \left(1 + \frac{C}{\tilde{\gamma}}\right) \inf_{v_h \in V_h} \|u - v_h\|_V.$$

## Proof.

For every  $v_h \in V_h$ , we have

$$\tilde{\gamma} \|v_h - u_h\|_V \leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u_h, w_h)}{\|w_h\|_V} \quad (\text{discrete inf-sup})$$

## Proof.

For every  $v_h \in V_h$ , we have

$$\tilde{\gamma} \|v_h - u_h\|_V \leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u_h, w_h)}{\|w_h\|_V} \quad (\text{discrete inf-sup})$$

$$= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h) + a(u - u_h, w_h)}{\|w_h\|_V} \quad (\text{bilinearity of } a)$$

## Proof.

For every  $v_h \in V_h$ , we have

$$\tilde{\gamma} \|v_h - u_h\|_V \leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u_h, w_h)}{\|w_h\|_V} \quad (\text{discrete inf-sup})$$

$$= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h) + a(u - u_h, w_h)}{\|w_h\|_V} \quad (\text{bilinearity of } a)$$

$$= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h)}{\|w_h\|_V} \quad (\text{Galerkin orth.})$$

## Proof.

For every  $v_h \in V_h$ , we have

$$\begin{aligned}
 \tilde{\gamma} \|v_h - u_h\|_V &\leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u_h, w_h)}{\|w_h\|_V} && \text{(discrete inf-sup)} \\
 &= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h) + a(u - u_h, w_h)}{\|w_h\|_V} && \text{(bilinearity of } a) \\
 &= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h)}{\|w_h\|_V} && \text{(Galerkin orth.)} \\
 &\leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{C \|v_h - u\|_V \|w_h\|_V}{\|w_h\|_V} && \text{(bddness of } a)
 \end{aligned}$$



## Proof.

For every  $v_h \in V_h$ , we have

$$\begin{aligned}
 \tilde{\gamma} \|v_h - u_h\|_V &\leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u_h, w_h)}{\|w_h\|_V} && \text{(discrete inf-sup)} \\
 &= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h) + a(u - u_h, w_h)}{\|w_h\|_V} && \text{(bilinearity of } a) \\
 &= \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{a(v_h - u, w_h)}{\|w_h\|_V} && \text{(Galerkin orth.)} \\
 &\leq \sup_{\substack{w_h \in V_h \\ w_h \neq 0}} \frac{C \|v_h - u\|_V \|w_h\|_V}{\|w_h\|_V} && \text{(bddness of } a) \\
 &= C \|v_h - u\|_V.
 \end{aligned}$$

## Proof.

Now apply the triangle inequality to  $\|u - u_h\|_V$ :

$$\|u - u_h\|_V \leq \|u - v_h\|_V + \|v_h - u_h\|_V$$



## Proof.

Now apply the triangle inequality to  $\|u - u_h\|_V$ :

$$\begin{aligned}\|u - u_h\|_V &\leq \|u - v_h\|_V + \|v_h - u_h\|_V \\ &\leq \|u - v_h\|_V + \frac{C}{\tilde{\gamma}} \|u - v_h\|_V\end{aligned}$$



## Proof.

Now apply the triangle inequality to  $\|u - u_h\|_V$ :

$$\begin{aligned}\|u - u_h\|_V &\leq \|u - v_h\|_V + \|v_h - u_h\|_V \\ &\leq \|u - v_h\|_V + \frac{C}{\tilde{\gamma}} \|u - v_h\|_V \\ &= \left(1 + \frac{C}{\tilde{\gamma}}\right) \|u - v_h\|_V.\end{aligned}$$



## Proof.

Now apply the triangle inequality to  $\|u - u_h\|_V$ :

$$\begin{aligned}\|u - u_h\|_V &\leq \|u - v_h\|_V + \|v_h - u_h\|_V \\ &\leq \|u - v_h\|_V + \frac{C}{\tilde{\gamma}} \|u - v_h\|_V \\ &= \left(1 + \frac{C}{\tilde{\gamma}}\right) \|u - v_h\|_V.\end{aligned}$$



As before, we can combine this with an approximation result and a regularity result to derive error estimates for finite element discretisations.

## C6.4 Finite Element Methods for PDEs

### Lecture 14: Saddle point problems

Patrick E. Farrell

University of Oxford

We have now seen the general necessary and sufficient Babuška conditions for the well-posedness of

find  $u \in V$  such that  $a(u, v) = F(v)$  for all  $v \in V$ .

We have now seen the general necessary and sufficient Babuška conditions for the well-posedness of

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

Many noncoercive problems arise via mixed formulations (solving for more than one variable), and in this lecture we will rephrase the well-posedness conditions for saddle point problems: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q) \end{aligned}$$

for all  $(v, q) \in V \times Q$ .



We have now seen the general necessary and sufficient Babuška conditions for the well-posedness of

$$\text{find } u \in V \text{ such that } a(u, v) = F(v) \text{ for all } v \in V.$$

Many noncoercive problems arise via mixed formulations (solving for more than one variable), and in this lecture we will rephrase the well-posedness conditions for saddle point problems: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q) \end{aligned}$$

for all  $(v, q) \in V \times Q$ .

These are the *Brezzi* conditions. The Brezzi conditions are easier to understand and verify than the Babuška conditions if you have a saddle point problem.

Note that the problem: find  $(u, p) \in V \times Q$  such that

$$a(u, v) + b(v, p) = F(v)$$

$$b(u, q) = G(q)$$

for all  $(v, q) \in V \times Q$

Note that the problem: find  $(u, p) \in V \times Q$  such that

$$a(u, v) + b(v, p) = F(v)$$

$$b(u, q) = G(q)$$

for all  $(v, q) \in V \times Q$

is equivalent to

Note that the problem: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned}a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q)\end{aligned}$$

for all  $(v, q) \in V \times Q$

is equivalent to

find  $(u, p) \in V \times Q$  such that

$$a(u, v) + b(v, p) + b(u, q) = F(v) + G(q)$$

for all  $(v, q) \in V \times Q$ . (Set  $v = 0$  and vary  $q \in Q$ , set  $q = 0$  and vary  $v \in V$ .)

We've already seen one example:

## Mixed Poisson (lecture 5)

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot vu - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} fw \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

We've already seen one example:

## Mixed Poisson (lecture 5)

Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot vu - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} fw \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

Here

$$a(\sigma, v) = \int_{\Omega} \sigma \cdot v \, dx, \quad b(v, u) = - \int_{\Omega} \nabla \cdot vu \, dx.$$

We've already seen one system that is *not* an example:

## Mixed linear elasticity (lecture 7)

Find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx = \int_{\Omega} f \cdot v \, dx,$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

We've already seen one system that is *not* an example:

## Mixed linear elasticity (lecture 7)

Find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx = \int_{\Omega} f \cdot v \, dx,$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

We can restructure this as

$$a(u, v) + b(v, p) = F(v)$$

$$b(u, q) + c(p, q) = G(q)$$

$$a(u, v) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx, \quad b(v, p) = \int_{\Omega} \nabla \cdot v p \, dx, \quad c(p, q) = -\frac{1}{\lambda} \int_{\Omega} p q \, dx.$$



We've already seen one system that is *not* an example:

## Mixed linear elasticity (lecture 7)

Find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx = \int_{\Omega} f \cdot v \, dx,$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

We can restructure this as

$$a(u, v) + b(v, p) = F(v)$$

$$b(u, q) + c(p, q) = G(q)$$

$$a(u, v) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx, \quad b(v, p) = \int_{\Omega} \nabla \cdot v p \, dx, \quad c(p, q) = -\frac{1}{\lambda} \int_{\Omega} p q \, dx.$$

We've already seen one system that is *not* an example:

## Mixed linear elasticity (lecture 7)

Find  $(u, p) \in H_0^1(\Omega; \mathbb{R}^n) \times L^2(\Omega)$  such that

$$\int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx + \int_{\Omega} p \nabla \cdot v \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx + \int_{\Omega} q \nabla \cdot u \, dx = \int_{\Omega} f \cdot v \, dx,$$

for all  $(v, q) \in H_0^1(\Omega; \mathbb{R}^2) \times L^2(\Omega)$ .

We can restructure this as

$$a(u, v) + b(v, p) = F(v)$$

$$b(u, q) + c(p, q) = G(q)$$

$$a(u, v) = \int_{\Omega} 2\mu \varepsilon(u) : \varepsilon(v) \, dx, \quad b(v, p) = \int_{\Omega} \nabla \cdot v p \, dx, \quad c(p, q) = -\frac{1}{\lambda} \int_{\Omega} p q \, dx.$$

This *is* a saddle point problem for  $\lambda \rightarrow \infty$ .

Let's consider one more example.

The Stokes equations are an elementary model in fluid mechanics. They describe the motion of a steady, incompressible, viscous, Newtonian, isothermal, slow-moving fluid.

$$\begin{aligned} -\nabla^2 u + \nabla p &= f && \text{in } \Omega, \\ \nabla \cdot u &= 0 && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Let's consider one more example.

The Stokes equations are an elementary model in fluid mechanics. They describe the motion of a steady, incompressible, viscous, Newtonian, isothermal, slow-moving fluid.

$$\begin{aligned} -\nabla^2 u + \nabla p &= f && \text{in } \Omega, \\ \nabla \cdot u &= 0 && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Here  $u : \Omega \rightarrow \mathbb{R}^n$  is the flow velocity and  $p : \Omega \rightarrow \mathbb{R}$  is the pressure.

Multiply the momentum equation by a vector-valued test function  $v \in V$ , and the continuity equation by a scalar-valued test function  $q \in Q$ :

$$\int_{\Omega} -\nabla \cdot \nabla u \cdot v \, dx + \int_{\Omega} \nabla p \cdot v \, dx = \int_{\Omega} f \cdot v \, dx,$$
$$\int_{\Omega} q \nabla \cdot u \, dx = 0.$$

Multiply the momentum equation by a vector-valued test function  $v \in V$ , and the continuity equation by a scalar-valued test function  $q \in Q$ :

$$\begin{aligned}\int_{\Omega} -\nabla \cdot \nabla u \cdot v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

Integrate the vector Laplacian by parts:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx - \int_{\partial\Omega} n \cdot \nabla u \cdot v \, ds + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

Multiply the momentum equation by a vector-valued test function  $v \in V$ , and the continuity equation by a scalar-valued test function  $q \in Q$ :

$$\begin{aligned}\int_{\Omega} -\nabla \cdot \nabla u \cdot v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

Integrate the vector Laplacian by parts:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx - \int_{\partial\Omega} n \cdot \nabla u \cdot v \, ds + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

We have nowhere to weakly enforce  $u = 0$ , so take  $V = H_0^1(\Omega; \mathbb{R}^n)$ .

Multiply the momentum equation by a vector-valued test function  $v \in V$ , and the continuity equation by a scalar-valued test function  $q \in Q$ :

$$\begin{aligned} \int_{\Omega} -\nabla \cdot \nabla u \cdot v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0. \end{aligned}$$

Integrate the vector Laplacian by parts:

$$\begin{aligned} \int_{\Omega} \nabla u : \nabla v \, dx - \int_{\partial\Omega} n \cdot \nabla u \cdot v \, ds + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0. \end{aligned}$$

We have nowhere to weakly enforce  $u = 0$ , so take  $V = H_0^1(\Omega; \mathbb{R}^n)$ .



Multiply the momentum equation by a vector-valued test function  $v \in V$ , and the continuity equation by a scalar-valued test function  $q \in Q$ :

$$\begin{aligned}\int_{\Omega} -\nabla \cdot \nabla u \cdot v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

Integrate the vector Laplacian by parts:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0.\end{aligned}$$

We have nowhere to weakly enforce  $u = 0$ , so take  $V = H_0^1(\Omega; \mathbb{R}^n)$ .

The formulation

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

requires  $u \in H_0^1(\Omega; \mathbb{R}^n)$  and  $p \in H^1(\Omega)$ .

The formulation

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

requires  $u \in H_0^1(\Omega; \mathbb{R}^n)$  and  $p \in H^1(\Omega)$ .

We can weaken the regularity requirement to  $p \in L^2(\Omega)$  by integrating by parts, and then negating the second equation for symmetry:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

The formulation

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

requires  $u \in H_0^1(\Omega; \mathbb{R}^n)$  and  $p \in H^1(\Omega)$ .

We can weaken the regularity requirement to  $p \in L^2(\Omega)$  by integrating by parts, and then negating the second equation for symmetry:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ - \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

The formulation

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx + \int_{\Omega} \nabla p \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

requires  $u \in H_0^1(\Omega; \mathbb{R}^n)$  and  $p \in H^1(\Omega)$ .

We can weaken the regularity requirement to  $p \in L^2(\Omega)$  by integrating by parts, and then negating the second equation for symmetry:

$$\begin{aligned}\int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx, \\ - \int_{\Omega} q \nabla \cdot u \, dx &= 0,\end{aligned}$$

Here

$$a(u, v) = \int_{\Omega} \nabla u : \nabla v \, dx, \quad b(v, p) = - \int_{\Omega} p \nabla \cdot v \, dx.$$

In the strong form of the problem,  $p$  only appears via  $\nabla p$ .

In the strong form of the problem,  $p$  only appears via  $\nabla p$ .

So if  $(u, p)$  is a solution, so is  $(u, p + c)$  for  $c \in \mathbb{R}$ . We can see this variationally:

$$\int_{\Omega} (p + c) \nabla \cdot v \, dx = \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\Omega} \nabla \cdot v \, dx$$

In the strong form of the problem,  $p$  only appears via  $\nabla p$ .

So if  $(u, p)$  is a solution, so is  $(u, p + c)$  for  $c \in \mathbb{R}$ . We can see this variationally:

$$\begin{aligned}\int_{\Omega} (p + c) \nabla \cdot v \, dx &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\Omega} \nabla \cdot v \, dx \\ &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\partial\Omega} v \cdot n \, ds\end{aligned}$$



In the strong form of the problem,  $p$  only appears via  $\nabla p$ .

So if  $(u, p)$  is a solution, so is  $(u, p + c)$  for  $c \in \mathbb{R}$ . We can see this variationally:

$$\begin{aligned}\int_{\Omega} (p + c) \nabla \cdot v \, dx &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\Omega} \nabla \cdot v \, dx \\ &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\partial\Omega} v \cdot n \, ds \\ &= \int_{\Omega} p \nabla \cdot v \, dx.\end{aligned}$$

In the strong form of the problem,  $p$  only appears via  $\nabla p$ .

So if  $(u, p)$  is a solution, so is  $(u, p + c)$  for  $c \in \mathbb{R}$ . We can see this variationally:

$$\begin{aligned} \int_{\Omega} (p + c) \nabla \cdot v \, dx &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\Omega} \nabla \cdot v \, dx \\ &= \int_{\Omega} p \nabla \cdot v \, dx + c \int_{\partial\Omega} v \cdot n \, ds \\ &= \int_{\Omega} p \nabla \cdot v \, dx. \end{aligned}$$

To fix a unique pressure we choose

$$Q = L_0^2(\Omega) = \{q \in L^2(\Omega) : \int_{\Omega} q \, dx = 0\}.$$

## Section 2

# Energy minimisation

Many weak formulations arise from energy minimisation.

Many weak formulations arise from energy minimisation.

Many problems of saddle-point form arise from *constrained* minimisation.

Consider

$$u = \operatorname{argmin}_{v \in H_0^1(\Omega; \mathbb{R}^n)} \frac{1}{2} \int_{\Omega} \nabla v : \nabla v \, dx - \int_{\Omega} f \cdot v \, dx,$$

subject to  $\nabla \cdot v = 0.$

Consider

$$u = \operatorname{argmin}_{v \in H_0^1(\Omega; \mathbb{R}^n)} \frac{1}{2} \int_{\Omega} \nabla v : \nabla v \, dx - \int_{\Omega} f \cdot v \, dx,$$

subject to  $\nabla \cdot v = 0.$

We introduce a Lagrange multiplier  $p$  and write the Lagrangian  $L : H_0^1(\Omega; \mathbb{R}^n) \times L_0^2(\Omega) \rightarrow \mathbb{R}$ :

$$L(u, p) = \frac{1}{2} \int_{\Omega} \nabla u : \nabla u \, dx - \int_{\Omega} f \cdot u \, dx - \int_{\Omega} p \nabla \cdot u \, dx.$$

$$L(u, p) = \frac{1}{2} \int_{\Omega} \nabla u : \nabla u \, dx - \int_{\Omega} p \nabla \cdot u \, dx - \int_{\Omega} f \cdot u \, dx.$$



$$L(u, p) = \frac{1}{2} \int_{\Omega} \nabla u : \nabla u \, dx - \int_{\Omega} p \nabla \cdot u \, dx - \int_{\Omega} f \cdot u \, dx.$$

Calculating the Euler–Lagrange equations, we have

$$L_u(u, p; v) = \int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx - \int_{\Omega} f \cdot v \, dx = 0,$$

$$L(u, p) = \frac{1}{2} \int_{\Omega} \nabla u : \nabla u \, dx - \int_{\Omega} p \nabla \cdot u \, dx - \int_{\Omega} f \cdot u \, dx.$$

Calculating the Euler–Lagrange equations, we have

$$L_u(u, p; v) = \int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx - \int_{\Omega} f \cdot v \, dx = 0,$$

$$L_p(u, p; q) = - \int_{\Omega} q \nabla \cdot u \, dx = 0,$$

the Stokes equations in weak form.

$$L(u, p) = \frac{1}{2} \int_{\Omega} \nabla u : \nabla u \, dx - \int_{\Omega} p \nabla \cdot u \, dx - \int_{\Omega} f \cdot u \, dx.$$

Calculating the Euler–Lagrange equations, we have

$$L_u(u, p; v) = \int_{\Omega} \nabla u : \nabla v \, dx - \int_{\Omega} p \nabla \cdot v \, dx - \int_{\Omega} f \cdot v \, dx = 0,$$

$$L_p(u, p; q) = - \int_{\Omega} q \nabla \cdot u \, dx = 0,$$

the Stokes equations in weak form.

In general constrained optimisation problems give you saddle point problems, because the constraint equation does not involve the Lagrange multiplier.

## Section 3

### Prelude: Orthogonal decompositions in Hilbert spaces

A very useful fact: Hilbert spaces can be cleanly separated into any closed subspace and its orthogonal complement.

A very useful fact: Hilbert spaces can be cleanly separated into any closed subspace and its orthogonal complement.

### Theorem (Orthogonal decomposition of a Hilbert space)

*Let  $H$  be a Hilbert space, and suppose  $K \subset H$  is a closed subspace of  $H$ . Then its orthogonal complement*

$$K^\perp := \{v \in H : v \perp k \text{ for all } k \in K\}$$

*is also a closed subspace, and*

$$H = K \oplus K^\perp,$$

*which means that every  $v \in H$  can be uniquely written as*

$$v = v^K + v^\perp,$$

*with  $v^K \in K$  and  $v^\perp \in K^\perp$ .*

## Section 4

### Saddle point systems in finite dimensions: homogeneous case

Consider the following  $N \times N$  linear system:

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ 0 \end{pmatrix},$$

where  $A \in \mathbb{R}^{NA \times NA}$ ,  $B \in \mathbb{R}^{NB \times NA}$ ,  $NA + NB = N$ ,  $NA > NB$ .



Consider the following  $N \times N$  linear system:

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ 0 \end{pmatrix},$$

where  $A \in \mathbb{R}^{NA \times NA}$ ,  $B \in \mathbb{R}^{NB \times NA}$ ,  $NA + NB = N$ ,  $NA > NB$ .

For Stokes,  $A$  is the vector Laplacian,  $B^T$  is the pressure gradient,  $B$  is the divergence.

Consider the following  $N \times N$  linear system:

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ 0 \end{pmatrix},$$

where  $A \in \mathbb{R}^{NA \times NA}$ ,  $B \in \mathbb{R}^{NB \times NA}$ ,  $NA + NB = N$ ,  $NA > NB$ .

For Stokes,  $A$  is the vector Laplacian,  $B^T$  is the pressure gradient,  $B$  is the divergence.

Since the second equation says that  $Bu = 0$ , we know that

$$u \in K := \text{kernel}(B) = \{v \in V : Bv = 0\}.$$

Consider the following  $N \times N$  linear system:

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ 0 \end{pmatrix},$$

where  $A \in \mathbb{R}^{NA \times NA}$ ,  $B \in \mathbb{R}^{NB \times NA}$ ,  $NA + NB = N$ ,  $NA > NB$ .

For Stokes,  $A$  is the vector Laplacian,  $B^T$  is the pressure gradient,  $B$  is the divergence.

Since the second equation says that  $Bu = 0$ , we know that

$$u \in K := \text{kernel}(B) = \{v \in V : Bv = 0\}.$$

If  $K = \{0\}$ , then the equation reduces to  $B^T p = f$ . Let's suppose  $\dim(K) > 0$ .

Let's test the first equation with  $v \in K$ :

$$v^T Au + v^T B^T p = v^T f,$$

Let's test the first equation with  $v \in K$ :

$$v^T Au + v^T B^T p = v^T f,$$

and since  $v^T B^T p = p^T(Bv) = 0$ , we derive the problem

find  $u \in K$  such that  $v^T Au = v^T f$  for all  $v \in K$ .

Let's test the first equation with  $v \in K$ :

$$v^T Au + v^T B^T p = v^T f,$$

and since  $v^T B^T p = p^T(Bv) = 0$ , we derive the problem

$$\text{find } u \in K \text{ such that } v^T Au = v^T f \text{ for all } v \in K.$$

This is a variational problem posed on a closed subspace of a Hilbert space, and we know how to ensure its well-posedness:

Let's test the first equation with  $v \in K$ :

$$v^T Au + v^T B^T p = v^T f,$$

and since  $v^T B^T p = p^T(Bv) = 0$ , we derive the problem

$$\text{find } u \in K \text{ such that } v^T Au = v^T f \text{ for all } v \in K.$$

This is a variational problem posed on a closed subspace of a Hilbert space, and we know how to ensure its well-posedness:

In the simpler case, assume the form is coercive on the kernel, i.e.

$$v^T Av \geq \alpha \|v\|^2 \text{ for all } v \in K.$$

In general, we assume the Babuška conditions hold.

Suppose we have solved the problem on the kernel for  $u$ , and let us see what we have achieved.



Suppose we have solved the problem on the kernel for  $u$ , and let us see what we have achieved.

Using the orthogonal decomposition, we know that we can write

$$f = f^K + f^\perp,$$

where  $f^K \in K$  and  $f^\perp \perp K$ .

Suppose we have solved the problem on the kernel for  $u$ , and let us see what we have achieved.

Using the orthogonal decomposition, we know that we can write

$$f = f^K + f^\perp,$$

where  $f^K \in K$  and  $f^\perp \perp K$ .

Testing with  $v \in K$  yields

$$v^T f = v^T f^K + v^T f^\perp = v^T f^K,$$

and so  $Au = f^K$ , and  $f - Au = f^\perp$ .

Having solved for the variable  $u$ , we must complete the solution of the problem by computing the unique  $p \in Q$  such that

$$B^T p = f - Au = f^\perp.$$

Having solved for the variable  $u$ , we must complete the solution of the problem by computing the unique  $p \in Q$  such that

$$B^T p = f - Au = f^\perp.$$

By the fundamental theorem of linear algebra,

$$\text{range}(B^T) = \text{kernel}(B)^\perp = K^\perp,$$

and since  $f^\perp \in \text{range}(B^T)$ , there exists at least one  $p \in Q$  such that  $B^T p = f - Au$ .

Having solved for the variable  $u$ , we must complete the solution of the problem by computing the unique  $p \in Q$  such that

$$B^T p = f - Au = f^\perp.$$

By the fundamental theorem of linear algebra,

$$\text{range}(B^T) = \text{kernel}(B)^\perp = K^\perp,$$

and since  $f^\perp \in \text{range}(B^T)$ , there exists at least one  $p \in Q$  such that  $B^T p = f - Au$ .

However, we must ensure that there is only one such  $p$ ; that is, we must ensure that  $B^T$  is injective.

To ensure injectivity, we wish to ensure that  $B^T p = 0 \implies p = 0$ .

To ensure injectivity, we wish to ensure that  $B^T p = 0 \implies p = 0$ .

We saw in the previous lecture that the inf-sup condition expresses this variationally: there exists  $\gamma \in \mathbb{R}$  such that

$$0 < \gamma \leq \inf_{\substack{q \in Q \\ q \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{v^T B^T q}{\|q\| \|v\|}.$$

To ensure injectivity, we wish to ensure that  $B^T p = 0 \implies p = 0$ .

We saw in the previous lecture that the inf-sup condition expresses this variationally: there exists  $\gamma \in \mathbb{R}$  such that

$$0 < \gamma \leq \inf_{\substack{q \in Q \\ q \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{v^T B^T q}{\|q\| \|v\|}.$$

If we assume this holds, then the operator  $B^T : Q \rightarrow K^\perp$  is a bijection, and we can solve for  $p \in Q$  uniquely.



## Section 5

Saddle point systems in finite dimensions: the inhomogeneous case

Now consider the modified problem

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix},$$

the inhomogeneous case.

Now consider the modified problem

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix},$$

the inhomogeneous case.

Again, define  $K$  to be the kernel of  $B$

$$K = \text{kernel}(B)$$

and write

$$u = u^K + u^\perp.$$

Suppose that we change our basis so that we may write

$$u = \begin{pmatrix} u^K \\ u^\perp \end{pmatrix}, \quad A = \begin{pmatrix} A^{KK} & A^{K\perp} \\ A^{\perp K} & A^{\perp\perp} \end{pmatrix}.$$

Such a change of basis is always possible.

Suppose that we change our basis so that we may write

$$u = \begin{pmatrix} u^K \\ u^\perp \end{pmatrix}, \quad A = \begin{pmatrix} A^{KK} & A^{K\perp} \\ A^{\perp K} & A^{\perp\perp} \end{pmatrix}.$$

Such a change of basis is always possible.

We can therefore write

$$\begin{aligned} A^{KK}u^K + A^{K\perp}u^\perp &= f^K, \\ A^{\perp K}u^K + A^{\perp\perp}u^\perp + B^T p &= f^\perp, \\ Bu^\perp &= g. \end{aligned}$$

There is no  $B^T p$  term in the first equation because its range is  $K^\perp$  and so it can only contribute to the second equation after our change of basis.

Can we solve  $Bu^\perp = g$  for  $u^\perp$ ? Yes! The inf-sup condition guarantees  $B^T : Q \rightarrow K^\perp$  is a bijection, and so  $B : K^\perp \rightarrow Q$  is also a bijection. We solve this for  $u^\perp$ .

Can we solve  $Bu^\perp = g$  for  $u^\perp$ ? Yes! The inf-sup condition guarantees  $B^T : Q \rightarrow K^\perp$  is a bijection, and so  $B : K^\perp \rightarrow Q$  is also a bijection. We solve this for  $u^\perp$ .

Testing the equation

$$A^{KK}u^K = f^K - A^{K\perp}u^\perp$$

with  $v \in K$  yields a linear variational problem over  $K$

find  $u^K \in K$  such that  $v^T A^{KK}u^K = v^T f^K - v^T A^{K\perp}u^\perp$  for all  $v \in K$

as before. We solve the problem on the kernel for  $u^K$ .

Can we solve  $Bu^\perp = g$  for  $u^\perp$ ? Yes! The inf-sup condition guarantees  $B^T : Q \rightarrow K^\perp$  is a bijection, and so  $B : K^\perp \rightarrow Q$  is also a bijection. We solve this for  $u^\perp$ .

Testing the equation

$$A^{KK}u^K = f^K - A^{K\perp}u^\perp$$

with  $v \in K$  yields a linear variational problem over  $K$

find  $u^K \in K$  such that  $v^T A^{KK}u^K = v^T f^K - v^T A^{K\perp}u^\perp$  for all  $v \in K$

as before. We solve the problem on the kernel for  $u^K$ .

We can then solve

$$B^T p = f^\perp - A^{\perp K}u^K - A^{\perp\perp}u^\perp$$

for  $p$  as before.



Can we solve  $Bu^\perp = g$  for  $u^\perp$ ? Yes! The inf-sup condition guarantees  $B^T : Q \rightarrow K^\perp$  is a bijection, and so  $B : K^\perp \rightarrow Q$  is also a bijection. We solve this for  $u^\perp$ .

Testing the equation

$$A^{KK}u^K = f^K - A^{K\perp}u^\perp$$

with  $v \in K$  yields a linear variational problem over  $K$

find  $u^K \in K$  such that  $v^T A^{KK}u^K = v^T f^K - v^T A^{K\perp}u^\perp$  for all  $v \in K$

as before. We solve the problem on the kernel for  $u^K$ .

We can then solve

$$B^T p = f^\perp - A^{\perp K}u^K - A^{\perp\perp}u^\perp$$

for  $p$  as before.

So no further assumptions are required for the inhomogeneous case.

## Section 6

# Brezzi's theorem

We now state the Brezzi conditions for the well-posedness of the abstract saddle point problem.

### Theorem (Well-posedness of saddle point problems)

*Let  $V$  and  $Q$  be Hilbert spaces. Given  $F \in V^*$  and  $G \in Q^*$ , we consider the problem: find  $(u, p) \in V \times Q$  such that*

$$\begin{aligned}a(u, v) + b(v, p) &= F(v), \\ b(u, q) &= G(q),\end{aligned}$$

*for all  $(v, q) \in V \times Q$ . Let*

$$K = \{v \in V : b(v, q) = 0 \text{ for all } q \in Q\}.$$

# Theorem (Well-posedness of saddle point problems)

*Suppose that*

- (1)  $a : V \times V \rightarrow \mathbb{R}$  and  $b : V \times Q \rightarrow \mathbb{R}$  are bounded bilinear forms;*
- (2) The variational problem*

*find  $u \in K$  such that  $a(u, v) = F(v)$  for all  $v \in K$*

*is well-posed;*

- (3)  $b$  satisfies the following inf-sup condition: there exists  $\gamma \in \mathbb{R}$  such that*

$$0 < \gamma \leq \inf_{\substack{q \in Q \\ q \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{b(v, q)}{\|v\|_V \|q\|_Q}.$$

*Then there exists a unique pair  $(u, p) \in V \times Q$  that solves the variational problem, and the solution is stable with respect to the data  $F$  and  $G$ .*

## Section 7

# Finite element discretisations of mixed problems

Take  $V_h \times Q_h \subset V \times Q$ , and consider: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$a(u_h, v_h) + b(v_h, p_h) = F(v_h),$$

$$b(u_h, q_h) = G(q_h),$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

Take  $V_h \times Q_h \subset V \times Q$ , and consider: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned}a(u_h, v_h) + b(v_h, p_h) &= F(v_h), \\ b(u_h, q_h) &= G(q_h),\end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

For this to be well-posed, Brezzi's conditions require that the LVP involving  $a$  is well-posed on the *discrete* kernel

$$K_h = \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}.$$

Take  $V_h \times Q_h \subset V \times Q$ , and consider: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned}a(u_h, v_h) + b(v_h, p_h) &= F(v_h), \\ b(u_h, q_h) &= G(q_h),\end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

For this to be well-posed, Brezzi's conditions require that the LVP involving  $a$  is well-posed on the *discrete* kernel

$$K_h = \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}.$$

Compare with

$$K \cap V_h = \{v_h \in V_h : b(v_h, q) = 0 \text{ for all } q \in Q\}.$$



In general, for  $v_h \in V_h$ , the property

$$b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h$$

will not imply

$$b(v_h, q) = 0 \text{ for all } q \in Q.$$

In general, for  $v_h \in V_h$ , the property

$$b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h$$

will not imply

$$b(v_h, q) = 0 \text{ for all } q \in Q.$$

(It will sometimes, but not always.)

In general, for  $v_h \in V_h$ , the property

$$b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h$$

will not imply

$$b(v_h, q) = 0 \text{ for all } q \in Q.$$

(It will sometimes, but not always.)

So in general  $K_h \not\subset K$ . This means that *well-posedness of  $a$  on the discrete kernel  $K_h$  does not necessarily follow automatically from well-posedness of  $a$  on the full kernel  $K$ .*

In general, for  $v_h \in V_h$ , the property

$$b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h$$

will not imply

$$b(v_h, q) = 0 \text{ for all } q \in Q.$$

(It will sometimes, but not always.)

So in general  $K_h \not\subset K$ . This means that *well-posedness of  $a$  on the discrete kernel  $K_h$  does not necessarily follow automatically from well-posedness of  $a$  on the full kernel  $K$ .*

One way to look at it: we have a *non-conforming* discretisation of the kernel problem.

That's one way a discretisation might fail. Any others?

That's one way a discretisation might fail. Any others?

Given that  $b$  satisfies the inf-sup condition over  $V$  and  $Q$ , it does *not* follow that  $b$  satisfies the inf-sup condition: there exists  $\tilde{\gamma} \in \mathbb{R}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\| \|q_h\|}.$$

That's one way a discretisation might fail. Any others?

Given that  $b$  satisfies the inf-sup condition over  $V$  and  $Q$ , it does *not* follow that  $b$  satisfies the inf-sup condition: there exists  $\tilde{\gamma} \in \mathbb{R}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\| \|q_h\|}.$$

We will see this by counterexample (later).

That's one way a discretisation might fail. Any others?

Given that  $b$  satisfies the inf-sup condition over  $V$  and  $Q$ , it does *not* follow that  $b$  satisfies the inf-sup condition: there exists  $\tilde{\gamma} \in \mathbb{R}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\| \|q_h\|}.$$

We will see this by counterexample (later).

So to analyse our discretisation error, we must *additionally* assume the Brezzi conditions hold for our discrete problem. This is a *compatibility* condition on the elements we choose for  $V_h$  and  $Q_h$ : they must work together.



## Section 8

# Quasi-optimality

## Theorem

*Consider the Galerkin approximation of our saddle point problem over  $V_h \times Q_h$ , a closed subspace of  $V \times Q$ :*

$$\begin{aligned}a(u_h, v_h) + b(v_h, p_h) &= F(v_h), \\ b(u_h, q_h) &= G(q_h).\end{aligned}$$

*Let*

$$K_h = \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}.$$

## Theorem

*In addition to the assumptions of Brezzi's theorem that guarantee well-posedness of the continuous problem, suppose that*

*(1) The variational problem*

*find  $u_h \in K_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in K_h$*

*is well-posed.*

*(2)  $b$  satisfies the following inf-sup condition over  $V_h \times Q_h$ : there exists  $\tilde{\gamma} \in \mathbb{R}$  such that*

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\|_V \|q_h\|_Q}.$$

*Then the Galerkin approximation is well-posed.*

## Theorem

*Furthermore, the approximate solutions are quasi-optimal: there exists  $c < \infty$  such that*

$$\|u - u_h\|_V + \|p - p_h\|_Q \leq c \left( \inf_{v_h \in V_h} \|u - v_h\|_V + \inf_{q_h \in Q_h} \|p - q_h\|_Q \right).$$

## Theorem

*Furthermore, the approximate solutions are quasi-optimal: there exists  $c < \infty$  such that*

$$\|u - u_h\|_V + \|p - p_h\|_Q \leq c \left( \inf_{v_h \in V_h} \|u - v_h\|_V + \inf_{q_h \in Q_h} \|p - q_h\|_Q \right).$$

This means you generally have to think about the quality of the approximation for  $u$  and  $p$  together: there's probably no point having a very high-order discretisation for  $u$  and a very low-order one for  $p$ !

## C6.4 Finite Element Methods for PDEs

### Lecture 15: The mixed Poisson equation

Patrick E. Farrell

University of Oxford

Many interesting problems are of saddle point form: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned}a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q)\end{aligned}$$

for all  $(v, q) \in V \times Q$ .

Many interesting problems are of saddle point form: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned}a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q)\end{aligned}$$

for all  $(v, q) \in V \times Q$ .

For this to be well-posed, we needed continuity of  $a$  and  $b$ , and



Many interesting problems are of saddle point form: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q) \end{aligned}$$

for all  $(v, q) \in V \times Q$ .

For this to be well-posed, we needed continuity of  $a$  and  $b$ , and

(1) The variational problem

find  $u \in K$  such that  $a(u, v) = F(v)$  for all  $v \in K$

over  $K := \{v \in V : b(v, q) = 0 \text{ for all } q \in Q\}$  is well-posed;

Many interesting problems are of saddle point form: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} a(u, v) + b(v, p) &= F(v) \\ b(u, q) &= G(q) \end{aligned}$$

for all  $(v, q) \in V \times Q$ .

For this to be well-posed, we needed continuity of  $a$  and  $b$ , and

(1) The variational problem

find  $u \in K$  such that  $a(u, v) = F(v)$  for all  $v \in K$

over  $K := \{v \in V : b(v, q) = 0 \text{ for all } q \in Q\}$  is well-posed;

(2)  $b$  satisfies the following inf-sup condition: there exists  $\gamma \in \mathbb{R}$  such that

$$0 < \gamma \leq \inf_{\substack{q \in Q \\ q \neq 0}} \sup_{\substack{v \in V \\ v \neq 0}} \frac{b(v, q)}{\|v\|_V \|q\|_Q}.$$

Consider a Galerkin approximation: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} a(u_h, v_h) + b(v_h, p_h) &= F(v_h) \\ b(u_h, q_h) &= G(q_h) \end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

Consider a Galerkin approximation: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} a(u_h, v_h) + b(v_h, p_h) &= F(v_h) \\ b(u_h, q_h) &= G(q_h) \end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

We similarly require

Consider a Galerkin approximation: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} a(u_h, v_h) + b(v_h, p_h) &= F(v_h) \\ b(u_h, q_h) &= G(q_h) \end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

We similarly require

(1) The variational problem

find  $u_h \in K_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in K_h$

over  $K_h := \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}$  is well-posed;

Consider a Galerkin approximation: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} a(u_h, v_h) + b(v_h, p_h) &= F(v_h) \\ b(u_h, q_h) &= G(q_h) \end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

We similarly require

(1) The variational problem

find  $u_h \in K_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in K_h$

over  $K_h := \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}$  is well-posed;

(2)  $V_h \times Q_h$  satisfies the following inf-sup condition: there exists  $\tilde{\gamma} \in \mathbb{R}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\|_V \|q_h\|_Q}.$$

Consider a Galerkin approximation: find  $(u_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} a(u_h, v_h) + b(v_h, p_h) &= F(v_h) \\ b(u_h, q_h) &= G(q_h) \end{aligned}$$

for all  $(v_h, q_h) \in V_h \times Q_h$ .

We similarly require

(1) The variational problem

find  $u_h \in K_h$  such that  $a(u_h, v_h) = F(v_h)$  for all  $v_h \in K_h$

over  $K_h := \{v_h \in V_h : b(v_h, q_h) = 0 \text{ for all } q_h \in Q_h\}$  is well-posed;

(2)  $V_h \times Q_h$  satisfies the following inf-sup condition: there exists  $\tilde{\gamma} \in \mathbb{R}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{v_h \in V_h \\ v_h \neq 0}} \frac{b(v_h, q_h)}{\|v_h\|_V \|q_h\|_Q}.$$

In this lecture we apply this theory to the mixed Poisson equation.

## Section 2

# Mixed Poisson in 1D



Let's consider the mixed Poisson equation in one dimension. Start with

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

Let's consider the mixed Poisson equation in one dimension. Start with

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

and introduce  $\sigma = -u'$  to get the system

$$\begin{aligned}\sigma + u' &= 0, \\ \sigma' &= f.\end{aligned}$$

Let's consider the mixed Poisson equation in one dimension. Start with

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

and introduce  $\sigma = -u'$  to get the system

$$\begin{aligned}\sigma + u' &= 0, \\ \sigma' &= f.\end{aligned}$$

Testing the equations with  $(\tau, v) \in V \times Q$ , we get

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx + \int_{\Omega} u' \tau \, dx &= 0, \\ \int_{\Omega} \sigma' v \, dx &= \int_{\Omega} f v \, dx.\end{aligned}$$

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx + \int_{\Omega} u' \tau \, dx &= 0, \\ \int_{\Omega} \sigma' v \, dx &= \int_{\Omega} f v \, dx.\end{aligned}$$

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx + \int_{\Omega} u' \tau \, dx &= 0, \\ \int_{\Omega} \sigma' v \, dx &= \int_{\Omega} f v \, dx.\end{aligned}$$

As it stands we need both  $\sigma, \tau \in H^1(\Omega)$  and  $u, v \in H^1(\Omega)$ . Let's integrate by parts to remove the derivative from  $u$  onto  $\tau$ , and negate:

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx - \int_{\Omega} u \tau' \, dx + \int_{\partial\Omega} u \tau \, ds &= 0, \\ - \int_{\Omega} \sigma' v \, dx &= - \int_{\Omega} f v \, dx.\end{aligned}$$

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx + \int_{\Omega} u' \tau \, dx &= 0, \\ \int_{\Omega} \sigma' v \, dx &= \int_{\Omega} f v \, dx.\end{aligned}$$

As it stands we need both  $\sigma, \tau \in H^1(\Omega)$  and  $u, v \in H^1(\Omega)$ . Let's integrate by parts to remove the derivative from  $u$  onto  $\tau$ , and negate:

$$\begin{aligned}\int_{\Omega} \sigma \tau \, dx - \int_{\Omega} u \tau' \, dx + \int_{\partial\Omega} u \tau \, ds &= 0, \\ - \int_{\Omega} \sigma' v \, dx &= - \int_{\Omega} f v \, dx.\end{aligned}$$

We can impose the Dirichlet BCs on  $u$  naturally by dropping the boundary term.

We thus have: find  $(\sigma, u) \in V \times Q := H^1(\Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \tau \, dx - \int_{\Omega} u \tau' \, dx - \int_{\Omega} \sigma' v \, dx = - \int_{\Omega} f v \, dx$$

for all  $(\tau, v) \in V \times Q$ .

We thus have: find  $(\sigma, u) \in V \times Q := H^1(\Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \tau \, dx - \int_{\Omega} u \tau' \, dx - \int_{\Omega} \sigma' v \, dx = - \int_{\Omega} f v \, dx$$

for all  $(\tau, v) \in V \times Q$ .

Let's think about well-posedness. Is

$$a(\sigma, \tau) = \int_{\Omega} \sigma \tau \, dx = (\sigma, \tau)_{L^2(\Omega)}$$

coercive over the kernel

$$K := \{\tau \in H^1(\Omega) : \int_{\Omega} \tau' v \, dx = 0 \text{ for all } v \in L^2(\Omega)\}?$$



We thus have: find  $(\sigma, u) \in V \times Q := H^1(\Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \tau \, dx - \int_{\Omega} u \tau' \, dx - \int_{\Omega} \sigma' v \, dx = - \int_{\Omega} f v \, dx$$

for all  $(\tau, v) \in V \times Q$ .

Let's think about well-posedness. Is

$$a(\sigma, \tau) = \int_{\Omega} \sigma \tau \, dx = (\sigma, \tau)_{L^2(\Omega)}$$

coercive over the kernel

$$K := \{\tau \in H^1(\Omega) : \int_{\Omega} \tau' v \, dx = 0 \text{ for all } v \in L^2(\Omega)\}?$$

Since  $\tau' \in L^2(\Omega)$ , choosing  $v = \tau'$  as test function yields that

$$\tau \in K \iff \tau' = 0.$$

For coercivity on the kernel, note that

$$a(\tau, \tau) = \|\tau\|_{L^2(\Omega)}^2 = \|\tau\|_{L^2(\Omega)}^2 + \|\tau'\|_{L^2(\Omega)}^2 = \|\tau\|_{H^1(\Omega)}^2,$$

so it is coercive with constant  $\alpha = 1$ . This is *only* true on the kernel, where  $\tau' = 0$ . It is *not* coercive on the whole of  $H^1(\Omega)$ !

For the inf-sup condition, we require that there exists  $\gamma$  such that

$$0 < \gamma \leq \inf_{\substack{v \in L^2(\Omega) \\ v \neq 0}} \sup_{\substack{\tau \in H^1(\Omega) \\ \tau \neq 0}} \frac{\int_{\Omega} \tau' v \, dx}{\|\tau\|_{H^1(\Omega)} \|v\|_{L^2(\Omega)}}.$$

For the inf-sup condition, we require that there exists  $\gamma$  such that

$$0 < \gamma \leq \inf_{\substack{v \in L^2(\Omega) \\ v \neq 0}} \sup_{\substack{\tau \in H^1(\Omega) \\ \tau \neq 0}} \frac{\int_{\Omega} \tau' v \, dx}{\|\tau\|_{H^1(\Omega)} \|v\|_{L^2(\Omega)}}.$$

For a given  $v \in L^2(\Omega)$ , choose

$$\tau(x) = \int_0^x v(x) \, dx$$

so that  $\tau' = v$  and  $\tau(0) = 0$ . We saw in Lecture 5, slide 7–9 that for such a function  $\|\tau\|_{L^2(\Omega)} \leq c \|\tau'\|_{L^2(\Omega)} = c \|v\|_{L^2(\Omega)}$  and hence

$$\|\tau\|_{H^1(\Omega)} \leq c \|v\|_{L^2(\Omega)}$$

for some (different)  $c$ .

With this choice, for any  $v \in L^2(\Omega)$ ,

$$\sup_{\substack{\tau \in V \\ \tau \neq 0}} \frac{\int_{\Omega} \tau' v \, dx}{\|\tau\|_{H^1(\Omega)} \|v\|_{L^2(\Omega)}} \geq \frac{\|v\|_{L^2(\Omega)}^2}{c \|v\|_{L^2(\Omega)}^2} = \frac{1}{c} > 0$$

so the inf-sup condition holds.

With this choice, for any  $v \in L^2(\Omega)$ ,

$$\sup_{\substack{\tau \in V \\ \tau \neq 0}} \frac{\int_{\Omega} \tau' v \, dx}{\|\tau\|_{H^1(\Omega)} \|v\|_{L^2(\Omega)}} \geq \frac{\|v\|_{L^2(\Omega)}^2}{c \|v\|_{L^2(\Omega)}^2} = \frac{1}{c} > 0$$

so the inf-sup condition holds.

Applying Brezzi's theorem, we conclude that the mixed formulation is well-posed.

With this choice, for any  $v \in L^2(\Omega)$ ,

$$\sup_{\substack{\tau \in V \\ \tau \neq 0}} \frac{\int_{\Omega} \tau' v \, dx}{\|\tau\|_{H^1(\Omega)} \|v\|_{L^2(\Omega)}} \geq \frac{\|v\|_{L^2(\Omega)}^2}{c \|v\|_{L^2(\Omega)}^2} = \frac{1}{c} > 0$$

so the inf-sup condition holds.

Applying Brezzi's theorem, we conclude that the mixed formulation is well-posed.

Here we see that the inf-sup condition is really surjectivity of

$$\frac{d}{dx} : H^1(\Omega) \rightarrow L^2(\Omega),$$

plus continuity.

## Section 3

# Discretising the mixed Poisson equation in 1D



Let's consider three different discretisations for  $V_h \times Q_h$ :

(A)  $CG_1 \times CG_1$

(B)  $CG_1 \times DG_0$

(C)  $CG_2 \times DG_0$

Let's consider three different discretisations for  $V_h \times Q_h$ :

(A)  $CG_1 \times CG_1$

(B)  $CG_1 \times DG_0$

(C)  $CG_2 \times DG_0$

Try them with  $f = 8$ , so that the exact solution is  $u = -4x(x - 1)$ :

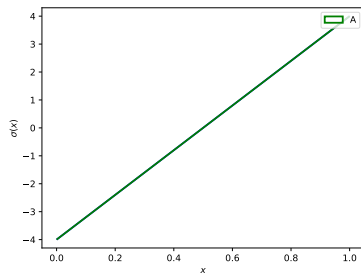
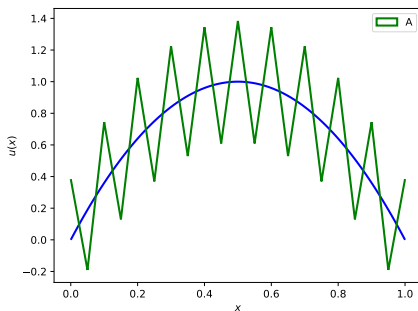
Let's consider three different discretisations for  $V_h \times Q_h$ :

(A)  $CG_1 \times CG_1$

(B)  $CG_1 \times DG_0$

(C)  $CG_2 \times DG_0$

Try them with  $f = 8$ , so that the exact solution is  $u = -4x(x - 1)$ :



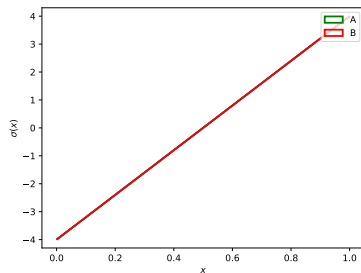
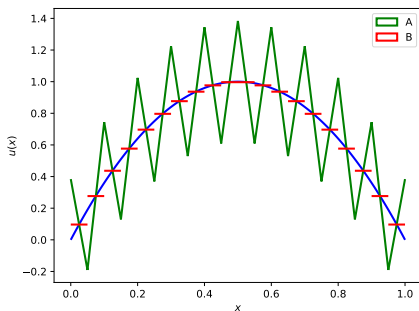
Let's consider three different discretisations for  $V_h \times Q_h$ :

(A)  $CG_1 \times CG_1$

(B)  $CG_1 \times DG_0$

(C)  $CG_2 \times DG_0$

Try them with  $f = 8$ , so that the exact solution is  $u = -4x(x - 1)$ :



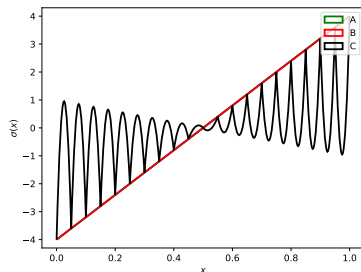
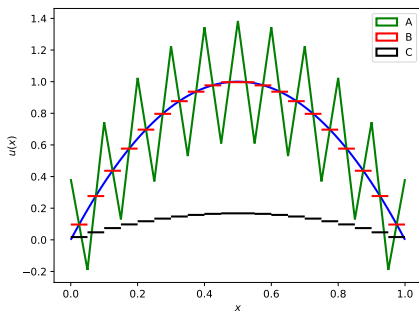
Let's consider three different discretisations for  $V_h \times Q_h$ :

(A)  $CG_1 \times CG_1$

(B)  $CG_1 \times DG_0$

(C)  $CG_2 \times DG_0$

Try them with  $f = 8$ , so that the exact solution is  $u = -4x(x - 1)$ :



## Subsection 1

Discretisation (A):  $CG_1 \times CG_1$

The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Let's choose  $v_h \in Q_h$  in the worst possible way. On each cell  $v_h$  is a linear polynomial and  $\tau'_h$  is a constant.



The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Let's choose  $v_h \in Q_h$  in the worst possible way. On each cell  $v_h$  is a linear polynomial and  $\tau'_h$  is a constant.

For integrands of polynomial degree one, midpoint quadrature is exact.

The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Let's choose  $v_h \in Q_h$  in the worst possible way. On each cell  $v_h$  is a linear polynomial and  $\tau'_h$  is a constant.

For integrands of polynomial degree one, midpoint quadrature is exact.

If we choose  $v_h$  to have value zero at each midpoint, then

$$\int_{\Omega} \tau'_h v_h \, dx = \sum_{K \in \mathcal{M}} \int_K \tau'_h v_h \, dx$$

The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Let's choose  $v_h \in Q_h$  in the worst possible way. On each cell  $v_h$  is a linear polynomial and  $\tau_h'$  is a constant.

For integrands of polynomial degree one, midpoint quadrature is exact.

If we choose  $v_h$  to have value zero at each midpoint, then

$$\begin{aligned} \int_{\Omega} \tau_h' v_h \, dx &= \sum_{K \in \mathcal{M}} \int_K \tau_h' v_h \, dx \\ &= \sum_{K \in \mathcal{M}} \tau_h' v_h(\text{midpoint}(K)) \end{aligned}$$

The discrete inf-sup condition is that there exists  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Let's choose  $v_h \in Q_h$  in the worst possible way. On each cell  $v_h$  is a linear polynomial and  $\tau_h'$  is a constant.

For integrands of polynomial degree one, midpoint quadrature is exact.

If we choose  $v_h$  to have value zero at each midpoint, then

$$\begin{aligned} \int_{\Omega} \tau_h' v_h \, dx &= \sum_{K \in \mathcal{M}} \int_K \tau_h' v_h \, dx \\ &= \sum_{K \in \mathcal{M}} \tau_h' v_h(\text{midpoint}(K)) \\ &= 0 \end{aligned}$$

for all  $\tau_h \in V_h$ .

Thus, for this discretisation,

$$0 = \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

Thus, for this discretisation,

$$0 = \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}.$$

## Conclusion

The  $CG_1 \times CG_1$  discretisation is unstable: it fails the inf-sup condition.

## Subsection 2

Discretisation (B):  $CG_1 \times DG_0$

Let's check the two discrete Brezzi conditions.



Let's check the two discrete Brezzi conditions.

The discrete kernel is

$$K_h := \{\tau_h \in V_h : \int_{\Omega} \tau_h' v_h \, dx = 0 \text{ for all } v_h \in Q_h\}.$$

Let's check the two discrete Brezzi conditions.

The discrete kernel is

$$K_h := \{\tau_h \in V_h : \int_{\Omega} \tau_h' v_h \, dx = 0 \text{ for all } v_h \in Q_h\}.$$

The derivative of a  $\text{CG}_1$  function is a  $\text{DG}_0$  function.

Let's check the two discrete Brezzi conditions.

The discrete kernel is

$$K_h := \{\tau_h \in V_h : \int_{\Omega} \tau_h' v_h \, dx = 0 \text{ for all } v_h \in Q_h\}.$$

The derivative of a  $\text{CG}_1$  function is a  $\text{DG}_0$  function.

Thus, if  $\tau_h \in K_h$ , choosing  $v_h = \tau_h'$  yields that  $\tau_h' = 0$ , so  $K_h = K \cap V_h$ .

Let's check the two discrete Brezzi conditions.

The discrete kernel is

$$K_h := \{\tau_h \in V_h : \int_{\Omega} \tau'_h v_h \, dx = 0 \text{ for all } v_h \in Q_h\}.$$

The derivative of a  $CG_1$  function is a  $DG_0$  function.

Thus, if  $\tau_h \in K_h$ , choosing  $v_h = \tau'_h$  yields that  $\tau'_h = 0$ , so  $K_h = K \cap V_h$ .

Hence, for  $\tau_h \in K_h$ ,

$$a(\tau_h, \tau_h) = \|\tau_h\|_{L^2(\Omega)}^2 = \|\tau_h\|_{L^2(\Omega)}^2 + \|\tau'_h\|_{L^2(\Omega)}^2 = \|\tau_h\|_{H^1(\Omega)}^2,$$

so the  $a(\tau_h, \tau_h)$  form is coercive over  $K_h$ .

What about the discrete inf-sup condition? Does there exist a  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}?$$

What about the discrete inf-sup condition? Does there exist a  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}?$$

In infinite dimensions we proved this by constructing a  $\tau \in V$  for any  $v \in Q$

$$\tau(x) = \int_0^x v(x) \, dx.$$

What about the discrete inf-sup condition? Does there exist a  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}?$$

In infinite dimensions we proved this by constructing a  $\tau \in V$  for any  $v \in Q$

$$\tau(x) = \int_0^x v(x) \, dx.$$

If  $v_h$  is piecewise constant, the associated  $\tau$  is piecewise linear and continuous, so  $\tau \in V_h$ ! So the inf-sup argument works in the same way.

What about the discrete inf-sup condition? Does there exist a  $\tilde{\gamma}$  such that

$$0 < \tilde{\gamma} \leq \inf_{\substack{v_h \in Q_h \\ v_h \neq 0}} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau_h' v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}?$$

In infinite dimensions we proved this by constructing a  $\tau \in V$  for any  $v \in Q$

$$\tau(x) = \int_0^x v(x) \, dx.$$

If  $v_h$  is piecewise constant, the associated  $\tau$  is piecewise linear and continuous, so  $\tau \in V_h$ ! So the inf-sup argument works in the same way.

## Conclusion

This discretisation is stable, by Brezzi's theorem.



These are the same arguments that worked in infinite dimensions. Why did they work again?

These are the same arguments that worked in infinite dimensions. Why did they work again?

$$\begin{array}{ccc} H^1 & \xrightarrow{\frac{d}{dx}} & L^2 \\ \downarrow & & \downarrow \\ V_h & \xrightarrow{\frac{d}{dx}} & Q_h \end{array}$$

These are the same arguments that worked in infinite dimensions. Why did they work again?

$$\begin{array}{ccc} H^1 & \xrightarrow{\frac{d}{dx}} & L^2 \\ \downarrow & & \downarrow \\ V_h & \xrightarrow{\frac{d}{dx}} & Q_h \end{array}$$

## Structure preservation

It worked because our choice of function spaces mimics the structure of the infinite-dimensional problem: the diagram commutes.

## Subsection 3

Discretisation (C):  $CG_2 \times DG_0$

First let's consider the inf-sup condition. Let  $V_h$  be constructed with  $\text{CG}_2$  elements, and let  $\tilde{V}_h \subsetneq V_h$  be constructed with  $\text{CG}_1$  elements. Then for any  $v_h \in Q_h$ ,

$$\sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} \geq \sup_{\substack{\tau_h \in \tilde{V}_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}}$$

First let's consider the inf-sup condition. Let  $V_h$  be constructed with  $\text{CG}_2$  elements, and let  $\tilde{V}_h \subsetneq V_h$  be constructed with  $\text{CG}_1$  elements. Then for any  $v_h \in Q_h$ ,

$$\begin{aligned} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} &\geq \sup_{\substack{\tau_h \in \tilde{V}_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} \\ &\geq \tilde{\gamma} \|v_h\|_{L^2(\Omega)}, \end{aligned}$$

since discretisation (B) is stable.

First let's consider the inf-sup condition. Let  $V_h$  be constructed with  $\text{CG}_2$  elements, and let  $\tilde{V}_h \subsetneq V_h$  be constructed with  $\text{CG}_1$  elements. Then for any  $v_h \in Q_h$ ,

$$\begin{aligned} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} &\geq \sup_{\substack{\tau_h \in \tilde{V}_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} \\ &\geq \tilde{\gamma} \|v_h\|_{L^2(\Omega)}, \end{aligned}$$

since discretisation (B) *is* stable.

In other words, enriching your discretisation of the  $V$ -space can only *improve* the inf-sup condition.

First let's consider the inf-sup condition. Let  $V_h$  be constructed with  $\text{CG}_2$  elements, and let  $\tilde{V}_h \subsetneq V_h$  be constructed with  $\text{CG}_1$  elements. Then for any  $v_h \in Q_h$ ,

$$\begin{aligned} \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} &\geq \sup_{\substack{\tau_h \in \tilde{V}_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \tau'_h v_h \, dx}{\|\tau_h\|_{H^1(\Omega)} \|v_h\|_{L^2(\Omega)}} \\ &\geq \tilde{\gamma} \|v_h\|_{L^2(\Omega)}, \end{aligned}$$

since discretisation (B) is stable.

In other words, enriching your discretisation of the  $V$ -space can only *improve* the inf-sup condition.

What about well-posedness of the LVP involving  $a$  on the kernel?



Increasing the size of  $V_h$  makes the discrete kernel  $K_h$  larger, so it is harder to satisfy coercivity.

Increasing the size of  $V_h$  makes the discrete kernel  $K_h$  larger, so it is harder to satisfy coercivity.

Consider a single mesh cell  $K = [\bar{x}, \bar{x} + h]$ . Define

$$\tau_h(x) = (x - \bar{x})(x - (\bar{x} + h))$$

on  $K$ , and zero elsewhere. We have  $\tau_h \in V_h$ .

Increasing the size of  $V_h$  makes the discrete kernel  $K_h$  larger, so it is harder to satisfy coercivity.

Consider a single mesh cell  $K = [\bar{x}, \bar{x} + h]$ . Define

$$\tau_h(x) = (x - \bar{x})(x - (\bar{x} + h))$$

on  $K$ , and zero elsewhere. We have  $\tau_h \in V_h$ .

Here the discrete kernel is:

$$K_h = \{\tau_h \in V_h : \int_{\Omega} \tau_h' v_h \, dx = 0 \text{ for all } v_h \in Q_h\}.$$

Claim:  $\tau_h \in K_h$ .

Calculating, on  $K$ ,

$$\tau_h(x) = x^2 - x(\bar{x} + h) - x\bar{x} + \bar{x}(\bar{x} + h),$$

and so

$$\tau'_h(x) = 2x - 2\bar{x} - h.$$

Calculating, on  $K$ ,

$$\tau_h(x) = x^2 - x(\bar{x} + h) - x\bar{x} + \bar{x}(\bar{x} + h),$$

and so

$$\tau'_h(x) = 2x - 2\bar{x} - h.$$

We know that for integrands of degree 1, midpoint quadrature is exact, so

$$\int_{\bar{x}}^{\bar{x}+h} \tau'_h v_h \, dx = v_h \int_{\bar{x}}^{\bar{x}+h} \tau'_h \, dx = h v_h \tau'_h(\bar{x} + h/2)$$

and thus we evaluate  $\tau'_h(\bar{x} + h/2)$ :

Calculating, on  $K$ ,

$$\tau_h(x) = x^2 - x(\bar{x} + h) - x\bar{x} + \bar{x}(\bar{x} + h),$$

and so

$$\tau'_h(x) = 2x - 2\bar{x} - h.$$

We know that for integrands of degree 1, midpoint quadrature is exact, so

$$\int_{\bar{x}}^{\bar{x}+h} \tau'_h v_h \, dx = v_h \int_{\bar{x}}^{\bar{x}+h} \tau'_h \, dx = h v_h \tau'_h(\bar{x} + h/2)$$

and thus we evaluate  $\tau'_h(\bar{x} + h/2)$ :

$$\tau'_h(\bar{x} + \frac{h}{2}) = 2\bar{x} + h - 2\bar{x} - h = 0.$$

Calculating, on  $K$ ,

$$\tau_h(x) = x^2 - x(\bar{x} + h) - x\bar{x} + \bar{x}(\bar{x} + h),$$

and so

$$\tau'_h(x) = 2x - 2\bar{x} - h.$$

We know that for integrands of degree 1, midpoint quadrature is exact, so

$$\int_{\bar{x}}^{\bar{x}+h} \tau'_h v_h \, dx = v_h \int_{\bar{x}}^{\bar{x}+h} \tau'_h \, dx = h v_h \tau'_h(\bar{x} + h/2)$$

and thus we evaluate  $\tau'_h(\bar{x} + h/2)$ :

$$\tau'_h(\bar{x} + \frac{h}{2}) = 2\bar{x} + h - 2\bar{x} - h = 0.$$

So  $\tau_h \in K_h$ .

So now we must decide whether  $a(\sigma, \tau) = (\sigma, \tau)_{L^2(\Omega)}$  is coercive on  $K_h$ .



So now we must decide whether  $a(\sigma, \tau) = (\sigma, \tau)_{L^2(\Omega)}$  is coercive on  $K_h$ .

Taking  $\tau_h$  as previously constructed,

$$a(\tau_h, \tau_h) = \mathcal{O}(h^5), \quad \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^3),$$

so

$$a(\tau_h, \tau_h) / \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^2) \rightarrow 0 \text{ as } h \rightarrow 0.$$

So now we must decide whether  $a(\sigma, \tau) = (\sigma, \tau)_{L^2(\Omega)}$  is coercive on  $K_h$ .

Taking  $\tau_h$  as previously constructed,

$$a(\tau_h, \tau_h) = \mathcal{O}(h^5), \quad \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^3),$$

so

$$a(\tau_h, \tau_h) / \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^2) \rightarrow 0 \text{ as } h \rightarrow 0.$$

So there is a loss of coercivity on the kernel as the mesh is refined. Assembling on any given mesh yields an invertible linear system, but the resulting approximations do not converge to the exact solution as the mesh is refined.

So now we must decide whether  $a(\sigma, \tau) = (\sigma, \tau)_{L^2(\Omega)}$  is coercive on  $K_h$ .

Taking  $\tau_h$  as previously constructed,

$$a(\tau_h, \tau_h) = \mathcal{O}(h^5), \quad \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^3),$$

so

$$a(\tau_h, \tau_h) / \|\tau_h\|_{H^1(\Omega)}^2 = \mathcal{O}(h^2) \rightarrow 0 \text{ as } h \rightarrow 0.$$

So there is a loss of coercivity on the kernel as the mesh is refined. Assembling on any given mesh yields an invertible linear system, but the resulting approximations do not converge to the exact solution as the mesh is refined.

## Conclusion

The  $\text{CG}_2 \times \text{DG}_0$  discretisation is unstable: loss of coercivity on the kernel.

- (A)  $CG_1 \times CG_1$  ✗ (does not satisfy inf-sup)
- (B)  $CG_1 \times DG_0$  ✓
- (C)  $CG_2 \times DG_0$  ✗ (coercivity on the kernel)

## Section 4

# Higher dimensions

Let's now recall the mixed Poisson equation in higher dimensions. As in Lecture 5, introducing  $\sigma = -\nabla u$ , we have

$$\begin{aligned}\sigma + \nabla u &= 0 \text{ in } \Omega, \\ \nabla \cdot \sigma &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega.\end{aligned}$$

Let's now recall the mixed Poisson equation in higher dimensions. As in Lecture 5, introducing  $\sigma = -\nabla u$ , we have

$$\begin{aligned}\sigma + \nabla u &= 0 \text{ in } \Omega, \\ \nabla \cdot \sigma &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega.\end{aligned}$$

This yields the variational formulation: Find  $(\sigma, u) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\int_{\Omega} \sigma \cdot v \, dx - \int_{\Omega} \nabla \cdot v u - \int_{\Omega} \nabla \cdot \sigma w \, dx = - \int_{\Omega} f w \, dx$$

for all  $(v, w) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$ .

In two dimensions, the  $L^2(\Omega)$  de Rham complex is

$$H^1 \xrightarrow{\text{curl}} H(\text{div}) \xrightarrow{\text{div}} L^2$$



In two dimensions, the  $L^2(\Omega)$  de Rham complex is

$$\begin{array}{ccccc}
 H^1 & \xrightarrow{\text{curl}} & H(\text{div}) & \xrightarrow{\text{div}} & L^2 \\
 \downarrow & & \downarrow & & \downarrow \\
 \Sigma_h & \xrightarrow{\text{curl}} & V_h & \xrightarrow{\text{div}} & Q_h
 \end{array}$$

In two dimensions, the  $L^2(\Omega)$  de Rham complex is

$$\begin{array}{ccccc}
 H^1 & \xrightarrow{\text{curl}} & H(\text{div}) & \xrightarrow{\text{div}} & L^2 \\
 \downarrow & & \downarrow & & \downarrow \\
 \Sigma_h & \xrightarrow{\text{curl}} & V_h & \xrightarrow{\text{div}} & Q_h
 \end{array}$$

In two dimensions, curl of a scalar field is: take its gradient and rotate it  $90^\circ$  anticlockwise. On a simply connected domain, anything in  $\text{kernel}(\text{div})$  is in the range of curl.

In two dimensions, the  $L^2(\Omega)$  de Rham complex is

$$\begin{array}{ccccc}
 H^1 & \xrightarrow{\text{curl}} & H(\text{div}) & \xrightarrow{\text{div}} & L^2 \\
 \downarrow & & \downarrow & & \downarrow \\
 \Sigma_h & \xrightarrow{\text{curl}} & V_h & \xrightarrow{\text{div}} & Q_h
 \end{array}$$

In two dimensions, curl of a scalar field is: take its gradient and rotate it  $90^\circ$  anticlockwise. On a simply connected domain, anything in  $\text{kernel}(\text{div})$  is in the range of curl.

It would make sense to choose  $V_h \subset H(\text{div}, \Omega)$  and  $Q_h \subset L^2(\Omega)$  so that

$$\text{div}(V_h) = Q_h.$$

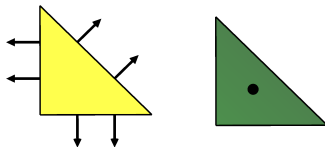
In two dimensions, the  $L^2(\Omega)$  de Rham complex is

$$\begin{array}{ccccc}
 H^1 & \xrightarrow{\text{curl}} & H(\text{div}) & \xrightarrow{\text{div}} & L^2 \\
 \downarrow & & \downarrow & & \downarrow \\
 \Sigma_h & \xrightarrow{\text{curl}} & V_h & \xrightarrow{\text{div}} & Q_h
 \end{array}$$

In two dimensions, curl of a scalar field is: take its gradient and rotate it  $90^\circ$  anticlockwise. On a simply connected domain, anything in  $\text{kernel}(\text{div})$  is in the range of curl.

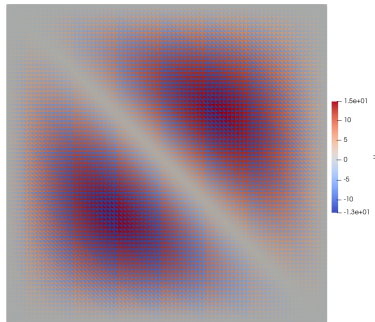
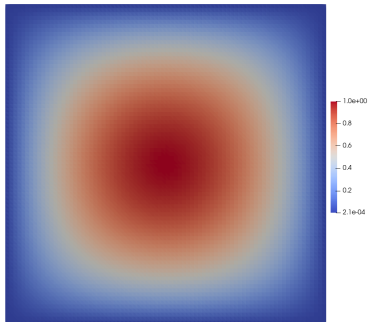
It would make sense to choose  $V_h \subset H(\text{div}, \Omega)$  and  $Q_h \subset L^2(\Omega)$  so that

$$\text{div}(V_h) = Q_h.$$

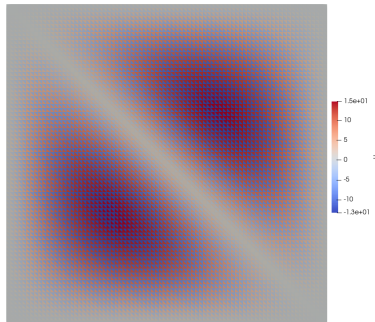
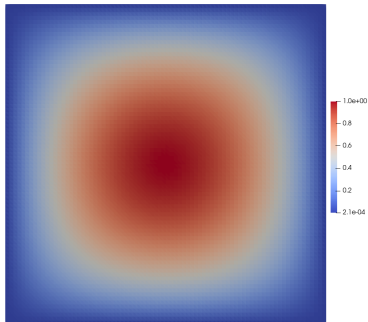


Let's try  $\text{BDM}_1 \times \text{DG}_0$ , and  $[\text{CG}_1]^2 \times \text{DG}_0$  with exact  
 $u = 4x(x-1)y(y-1)$ :

Let's try  $\text{BDM}_1 \times \text{DG}_0$ , and  $[\text{CG}_1]^2 \times \text{DG}_0$  with exact  
 $u = 4x(x-1)y(y-1)$ :



Let's try  $\text{BDM}_1 \times \text{DG}_0$ , and  $[\text{CG}_1]^2 \times \text{DG}_0$  with exact  
 $u = 4x(x-1)y(y-1)$ :



The results with Brezzi–Douglas–Marini look correct (they are). The results with continuous Lagrange elements for the flux are rubbish. The element pair is not stable; the inf-sup condition fails.

## Subsection 1

### Well-posedness of $\text{BDM}_1 \times \text{DG}_0$



Let's investigate the Brezzi conditions for the  $\text{BDM}_1 \times \text{DG}_0$  discretisation.

Let's investigate the Brezzi conditions for the  $\text{BDM}_1 \times \text{DG}_0$  discretisation.

We have

$$K_h = \{\tau_h \in V_h : \int_{\Omega} \nabla \cdot \tau_h v_h \, dx = 0 \text{ for all } v_h \in Q_h\},$$

and since  $\text{div}(V_h) = Q_h$ , we can rewrite this as

Let's investigate the Brezzi conditions for the  $\text{BDM}_1 \times \text{DG}_0$  discretisation.

We have

$$K_h = \{\tau_h \in V_h : \int_{\Omega} \nabla \cdot \tau_h v_h \, dx = 0 \text{ for all } v_h \in Q_h\},$$

and since  $\text{div}(V_h) = Q_h$ , we can rewrite this as

$$K_h = \{\tau_h \in V_h : \nabla \cdot \tau_h = 0\} = K \cap V_h.$$

Let's investigate the Brezzi conditions for the  $\text{BDM}_1 \times \text{DG}_0$  discretisation.

We have

$$K_h = \{\tau_h \in V_h : \int_{\Omega} \nabla \cdot \tau_h v_h \, dx = 0 \text{ for all } v_h \in Q_h\},$$

and since  $\text{div}(V_h) = Q_h$ , we can rewrite this as

$$K_h = \{\tau_h \in V_h : \nabla \cdot \tau_h = 0\} = K \cap V_h.$$

So in this case  $K_h \subset K$  and coercivity on the kernel is inherited with  $\alpha = 1$ .

The inf-sup condition will follow from the following results. Consider the smoother complex

The inf-sup condition will follow from the following results. Consider the smoother complex

$$\begin{array}{ccc} H^1 & \xrightarrow{\text{div}} & L^2 \\ \downarrow \pi_V & & \downarrow \pi_Q \\ V_h & \xrightarrow{\text{div}} & Q_h \end{array}$$

The inf-sup condition will follow from the following results. Consider the smoother complex

$$\begin{array}{ccc} H^1 & \xrightarrow{\text{div}} & L^2 \\ \downarrow \pi_V & & \downarrow \pi_Q \\ V_h & \xrightarrow{\text{div}} & Q_h \end{array}$$

There exist *bounded cochain projections*

$$\pi_V : H^1(\Omega) \rightarrow V_h,$$

$$\pi_Q : L^2(\Omega) \rightarrow Q_h,$$

The inf-sup condition will follow from the following results. Consider the smoother complex

$$\begin{array}{ccc} H^1 & \xrightarrow{\text{div}} & L^2 \\ \downarrow \pi_V & & \downarrow \pi_Q \\ V_h & \xrightarrow{\text{div}} & Q_h \end{array}$$

There exist *bounded cochain projections*

$$\pi_V : H^1(\Omega) \rightarrow V_h,$$

$$\pi_Q : L^2(\Omega) \rightarrow Q_h,$$

(so  $\pi_V^2 = \pi_V, \pi_Q^2 = \pi_Q$ ) such that for all  $\tau \in H^1(\Omega)$ ,

$$\|\pi_V \tau\|_{H(\text{div}, \Omega)} \leq C_V \|\tau\|_{H^1(\Omega)},$$

$$\nabla \cdot (\pi_V \tau) = \pi_Q \nabla \cdot \tau.$$



In this context we use  $H^1(\Omega)$  as the domain of  $\pi_V$  in order to ensure that the Brezzi–Douglas–Marini degrees of freedom are all defined.

In this context we use  $H^1(\Omega)$  as the domain of  $\pi_V$  in order to ensure that the Brezzi–Douglas–Marini degrees of freedom are all defined.

This is analogous to the requirement that we apply the continuous Lagrange interpolation operator to  $H^2(\Omega)$ -smooth functions.

In this context we use  $H^1(\Omega)$  as the domain of  $\pi_V$  in order to ensure that the Brezzi–Douglas–Marini degrees of freedom are all defined.

This is analogous to the requirement that we apply the continuous Lagrange interpolation operator to  $H^2(\Omega)$ -smooth functions.

The development of a bounded cochain projection onto  $\text{BDM}_k$  with domain  $H(\text{div}, \Omega)$  is much more technical.

## Theorem

*There exists  $\tilde{\gamma}$  such that for all  $v_h \in Q_h$ ,*

$$\tilde{\gamma} \|v_h\|_{L^2(\Omega)} \leq \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \nabla \cdot \tau_h v_h \, dx}{\|\tau_h\|_{H(\text{div}, \Omega)}}.$$

## Theorem

*There exists  $\tilde{\gamma}$  such that for all  $v_h \in Q_h$ ,*

$$\tilde{\gamma} \|v_h\|_{L^2(\Omega)} \leq \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \nabla \cdot \tau_h v_h \, dx}{\|\tau_h\|_{H(\text{div}, \Omega)}}.$$

## Proof.

Let's construct a  $\tau_h \in V_h$  such that  $\nabla \cdot \tau_h = v_h$  and  $\|\tau_h\|_V \leq c \|v_h\|_Q$  for some  $c$ .

## Theorem

*There exists  $\tilde{\gamma}$  such that for all  $v_h \in Q_h$ ,*

$$\tilde{\gamma} \|v_h\|_{L^2(\Omega)} \leq \sup_{\substack{\tau_h \in V_h \\ \tau_h \neq 0}} \frac{\int_{\Omega} \nabla \cdot \tau_h v_h \, dx}{\|\tau_h\|_{H(\text{div}, \Omega)}}.$$

## Proof.

Let's construct a  $\tau_h \in V_h$  such that  $\nabla \cdot \tau_h = v_h$  and  $\|\tau_h\|_V \leq c \|v_h\|_Q$  for some  $c$ .

First find  $\sigma \in H^1(\Omega)$  with  $\nabla \cdot \sigma = v_h$ ,  $\|\sigma\|_{H^1(\Omega)} \leq c \|v_h\|_{L^2(\Omega)}$ . We can do this by extending the domain to a disc (so an elliptic regularity result applies), extending  $v_h$  by zero, then solve  $-\nabla^2 u = v_h$  with zero Dirichlet BCs for  $u \in H^2$  on the extended domain. We then set  $\sigma = -\nabla u|_{\Omega}$ .

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\nabla \cdot \tau_h = \nabla \cdot \pi_V \sigma$$



Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\begin{aligned}\nabla \cdot \tau_h &= \nabla \cdot \pi_V \sigma \\ &= \pi_Q \nabla \cdot \sigma\end{aligned}$$

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\begin{aligned}\nabla \cdot \tau_h &= \nabla \cdot \pi_V \sigma \\ &= \pi_Q \nabla \cdot \sigma \\ &= \pi_Q v_h\end{aligned}$$

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\begin{aligned}\nabla \cdot \tau_h &= \nabla \cdot \pi_V \sigma \\ &= \pi_Q \nabla \cdot \sigma \\ &= \pi_Q v_h \\ &= v_h,\end{aligned}$$

and

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\begin{aligned}\nabla \cdot \tau_h &= \nabla \cdot \pi_V \sigma \\ &= \pi_Q \nabla \cdot \sigma \\ &= \pi_Q v_h \\ &= v_h,\end{aligned}$$

and

$$\|\tau_h\|_V \leq C_v \|\sigma\|_{H^1(\Omega)} \leq C_v c \|v_h\|_{L^2(\Omega)}.$$

Proof.

Now set  $\tau_h = \pi_V \sigma$ .

Using the commuting diagram property,

$$\begin{aligned}\nabla \cdot \tau_h &= \nabla \cdot \pi_V \sigma \\ &= \pi_Q \nabla \cdot \sigma \\ &= \pi_Q v_h \\ &= v_h,\end{aligned}$$

and

$$\|\tau_h\|_V \leq C_v \|\sigma\|_{H^1(\Omega)} \leq C_v c \|v_h\|_{L^2(\Omega)}.$$

Using  $\tau_h$  to give a lower-bound for the supremum concludes the argument, as before. □

This kind of analysis, using *Hilbert complexes* to understand and preserve the structure, is now fundamental to the finite element method. The theory goes by the name of the *finite element exterior calculus*.

This kind of analysis, using *Hilbert complexes* to understand and preserve the structure, is now fundamental to the finite element method. The theory goes by the name of the *finite element exterior calculus*.

We develop stable and structure-preserving discretisations of PDEs by constructing *subcomplexes* with *bounded cochain projections*.

This kind of analysis, using *Hilbert complexes* to understand and preserve the structure, is now fundamental to the finite element method. The theory goes by the name of the *finite element exterior calculus*.

We develop stable and structure-preserving discretisations of PDEs by constructing *subcomplexes* with *bounded cochain projections*.

The finite element exterior calculus connects numerical approximation of PDEs with geometry and topology (the de Rham complex encodes the topology of the domain, and so does a subcomplex!).



This kind of analysis, using *Hilbert complexes* to understand and preserve the structure, is now fundamental to the finite element method. The theory goes by the name of the *finite element exterior calculus*.

We develop stable and structure-preserving discretisations of PDEs by constructing *subcomplexes* with *bounded cochain projections*.

The finite element exterior calculus connects numerical approximation of PDEs with geometry and topology (the de Rham complex encodes the topology of the domain, and so does a subcomplex!).

It is now used to develop structure-preserving methods in fluid mechanics, elasticity, electromagnetism, general relativity, and beyond.