

Introduction

In this assignment we consider the boundary value problem

$$\begin{aligned} -\frac{d}{dx} \left(\kappa(x) \frac{dp}{dx}(x) \right) &= f(x) \quad \text{for } x \in (0, 1) \\ p(0) &= p(1) = 0. \end{aligned} \tag{PDE}$$

The concept of a weak solution to this equation will be introduced, and an existence and uniqueness result will be shown. We will then consider how the equation may be implemented numerically, using a finite element method, and study some properties of this discretization.

This equation has a natural extension to more general domains $D \subseteq \mathbb{R}^3$, given by

$$\begin{aligned} -\nabla \cdot (\boldsymbol{\kappa}(\mathbf{x}) \nabla p(\mathbf{x})) &= f(\mathbf{x}) \quad \text{for } \mathbf{x} \in D \\ p(\mathbf{x}) &= 0 \quad \text{for } \mathbf{x} \in \partial D \end{aligned}$$

where now $\boldsymbol{\kappa} : D \rightarrow \mathbb{R}^{3 \times 3}$ is matrix-valued. The theory we look at generalizes to this case directly, and it is in this case where the weak solutions and finite element methods become key tools for understanding. The equation itself arises, for example, when considering groundwater flow, with $\boldsymbol{\kappa}$ representing the permeability of the subsurface and p the pressure of a fluid in the subsurface. Matrix-valued $\boldsymbol{\kappa}$ is considered to account for the permeability of a medium being dependent on direction. The equation is arrived at by combining Darcy's law with a conservation law. The equation also arises in electrodynamics: if an electric field E is conservative, then the equation follows from the steady state Maxwell equations, with $\boldsymbol{\kappa}$ representing the background conductivity and p representing the electric potential corresponding to E .

For exposition we work in this assignment only with the one-dimensional case $D = (0, 1)$ given by (PDE). This allows for the general theory and ideas to be introduced, without having to deal with certain technicalities. Moreover, computations in this case are much less expensive than in higher dimensions.

Throughout this assignment we will make the following assumptions about (PDE):

Assumptions 1. *The diffusion coefficient $\kappa : (0, 1) \rightarrow \mathbb{R}$ and source term $f : (0, 1) \rightarrow \mathbb{R}$ are such that:*

- (i) *there exist $\kappa_0, \kappa_1 > 0$ with $\kappa_0 \leq \kappa(x) \leq \kappa_1$ for all $x \in (0, 1)$;*
- (ii) *$f \in L^2(0, 1)$.*

Problem 1. Weak Solutions (30 points)

- (a) Denote by $C_c^\infty(0, 1)$ the set of compactly supported smooth functions on the open interval $(0, 1)$. That is,

$$C_c^\infty(0, 1) = \left\{ \varphi \in C^\infty(0, 1) \mid \begin{array}{l} \text{there exists a closed bounded set } K \subset (0, 1) \\ \text{with } \varphi(x) = 0 \text{ for all } x \notin K \end{array} \right\}.$$

(i) Let $\varphi \in C_c^\infty(0, 1)$. Show that

$$\lim_{x \downarrow 0} \varphi(x) = \lim_{x \uparrow 1} \varphi(x) = 0.$$

(ii) Let p be a solution to (PDE). Show that

$$\int_0^1 \kappa(x) \frac{dp}{dx}(x) \frac{d\varphi}{dx}(x) dx = \int_0^1 f(x) \varphi(x) dx \quad \text{for any } \varphi \in C_c^\infty(0, 1).$$

Do not solve (PDE) explicitly.

(b) Define the H_0^1 norm and inner product on $C_c^\infty(0, 1)$ by

$$\|\varphi\|_{H_0^1}^2 = \int_0^1 \left| \frac{d\varphi}{dx}(x) \right|^2 dx, \quad \langle \varphi, \psi \rangle_{H_0^1} = \int_0^1 \frac{d\varphi}{dx}(x) \frac{d\psi}{dx}(x) dx.$$

Define the space $H_0^1(0, 1)$ to be the closure of $C_c^\infty(0, 1)$ with respect to the H_0^1 norm, that is,

$$H_0^1(0, 1) = \{v : (0, 1) \rightarrow \mathbb{R} \mid \text{there exists } \{\varphi_n\} \subset C_c^\infty(0, 1) \text{ with } \|\varphi_n - v\|_{H_0^1} \rightarrow 0\}.$$

We will assume in the remainder of what follows that any $v \in H_0^1(0, 1)$ is once differentiable at almost every point in $(0, 1)$; in particular the derivative of v is well-defined whenever it appears under an integral sign. Additionally we will assume that $H_0^1(0, 1)$ is a Hilbert space when equipped with the inner product $\langle \cdot, \cdot \rangle_{H_0^1}$ defined above.

Remark 1. *You will have encountered the space $H_0^1(0, 1)$ in a problem set, though it was defined differently. The two definitions can be shown to be equivalent, but in this assignment we will work only with the definition above.*

(i) Prove the Poincaré inequality: for any $v \in H_0^1(0, 1)$,

$$\int_0^1 |v(x)|^2 dx \leq \int_0^1 \left| \frac{dv}{dx}(x) \right|^2 dx.$$

Deduce that $H_0^1(0, 1)$ is continuously embedded in $L^2(0, 1)$.

(ii) Let Assumptions 1 hold, and let p be a solution to (PDE). Using Problem 1(a)(ii) and the Poincaré inequality, show that

$$\int_0^1 \kappa(x) \frac{dp}{dx}(x) \frac{dv}{dx}(x) dx = \int_0^1 f(x) v(x) dx \quad \text{for all } v \in H_0^1(0, 1). \quad (\text{WPDE})$$

In what follows, we will call any function $p \in H_0^1(0, 1)$ that satisfies (WPDE) a *weak solution* to (PDE), and any function that solves (PDE) in the classical sense a *strong solution* to (PDE). Note that if it exists, a strong solution is a weak solution.

(c) Denote $V = H_0^1(0, 1)$. Define $B : V \times V \rightarrow \mathbb{R}$, $g : V \rightarrow \mathbb{R}$ by

$$B(u, v) = \int_0^1 \kappa(x) \frac{du}{dx}(x) \frac{dv}{dx}(x) dx, \quad g(v) = \int_0^1 f(x) v(x) dx. \quad (1)$$

Then we may rewrite (WPDE) in the form

$$B(p, v) = g(v) \quad \text{for all } v \in V.$$

Using the Lax-Milgram lemma, prove that there exists a unique weak solution $p \in V$ to (PDE).

Problem 2. Finite Element Approximations (30 points)

If we wish to solve (WPDE) numerically, we will first need to make a finite dimensional approximation to the system, which will then provide a finite dimensional approximation to the solution. In this problem we consider a class of approximations known as finite element approximations.

Let $V^h \subseteq V$ be a finite dimensional subspace of V , with basis $\{\varphi_j^h\}_{j=1}^{N_h}$. The scalar parameter $h > 0$ will be related to the dimension of V^h , with $N_h \uparrow \infty$ as $h \downarrow 0$. It is this space V^h in which we will look for an approximate solution.

- (a) (i) Equip V^h with the inner product $\langle \cdot, \cdot \rangle_{H_0^1}$. Show that V^h is a Hilbert space.
(ii) Let $B : V^h \times V^h \rightarrow \mathbb{R}$, $g : V^h \rightarrow \mathbb{R}$ denote the restrictions of B, g to V^h , where B, g are defined by (1). Using the Lax-Milgram lemma, show that there exists a unique solution $p^h \in V^h$ to the problem

$$B(p^h, v^h) = g(v^h) \quad \text{for all } v^h \in V^h. \quad (\text{WPDE-}h)$$

- (iii) Let $p^h \in V^h$ be the unique solution to (WPDE- h). Since $\varphi_1^h, \dots, \varphi_{N_h}^h$ form a basis for V^h , there exist scalars $P_1^h, \dots, P_{N_h}^h \in \mathbb{R}$ such that

$$p^h = \sum_{j=1}^{N_h} P_j^h \varphi_j^h.$$

Find expressions for the entries of a matrix $A^h \in \mathbb{R}^{N_h \times N_h}$ and a vector $F^h \in \mathbb{R}^{N_h}$ such that the coefficients $P^h = (P_1^h, \dots, P_{N_h}^h)^\top \in \mathbb{R}^{N_h}$ solve

$$A^h P^h = F^h. \quad (\text{FEM})$$

- (iv) Define the inner product $\langle \cdot, \cdot \rangle_B$ on V by $\langle u, v \rangle_B = B(u, v)$. Show that the induced norm $\|\cdot\|_B$ is equivalent to $\|\cdot\|_{H_0^1}$.
(v) Let $p \in V, p^h \in V^h$ be the unique solutions to (WPDE) and (WPDE- h) respectively. Establish the *Galerkin orthogonality*:

$$\langle p - p^h, v^h \rangle_B = 0 \quad \text{for all } v^h \in V^h.$$

Using this, show that

$$\|p - p^h\|_B \leq \|p - v^h\|_B \quad \text{for all } v^h \in V^h$$

and hence

$$\|p - p^h\|_{H_0^1} \leq \frac{\kappa_1}{\kappa_0} \cdot \inf_{v^h \in V^h} \|p - v^h\|_{H_0^1}.$$

The first inequality above is the *Galerkin optimality property*, namely that p^h is optimal over all possible approximations in V^h with respect to the norm induced by B .

- (b) We now look at a particular choice of approximation space V^h . Let $h \in (0, 1/2)$ and set $N_h = \lfloor 1/h \rfloor - 1 \in \mathbb{N}$. Define the mesh

$$\mathbf{x}^h = \{0, h, 2h, 3h, \dots, N_h h, (N_h + 1)h\} \subset [0, 1].$$

We will denote $x_j^h = jh$ for each $j = 0, \dots, N_h + 1$. The points x_j^h will be referred to as *nodes* and the intervals (x_j^h, x_{j+1}^h) as *elements*.

Define the set of functions $\{\varphi_j^h\} \subseteq H_0^1(0, 1)$ by

$$\varphi_j^h(x) = \begin{cases} \frac{x - x_{j-1}^h}{x_j^h - x_{j-1}^h} & x \in (x_{j-1}^h, x_j^h] \\ \frac{x_{j+1}^h - x}{x_{j+1}^h - x_j^h} & x \in (x_j^h, x_{j+1}^h] \\ 0 & x \notin (x_{j-1}^h, x_{j+1}^h] \end{cases}$$

for each $j = 1, \dots, N_h$. Define the finite element space

$$V^h = \text{span}\{\varphi_1^h, \dots, \varphi_{N_h}^h\} \subset H_0^1(0, 1).$$

- (i) Draw the graphs of each $\{\varphi_j^h\}$ for $h = 1/5$.
- (ii) Show that $\varphi_1^h, \dots, \varphi_{N_h}^h$ form a basis for V^h . Are they an orthonormal basis?
- (iii) Given a function $b \in V$, describe the function $b^h \in V^h$ given by

$$b^h(x) = \sum_{j=1}^{N_h} b(x_j^h) \varphi_j^h(x).$$

- (iv) Observe that φ_j^h and φ_k^h have disjoint supports if $|j - k| > 1$. Given that this is the case, what structure should the matrix A^h in (FEM) have?
- (v) Implement the system (FEM) in MATLAB, with the terms κ and f provided as function handles. That is, given $h > 0$ and function handles for κ, f , construct A^h and F^h , and return the vector of coefficients P^h . It will be beneficial to implement P^h as a sparse matrix, using the `sparse` function in MATLAB.

For integration over elements, use the `integral` function in MATLAB, supplying relative and absolute tolerances of 10^{-14} . We choose these tolerances close to machine precision so that the errors in the integration will be dominated by the errors that arise from the finite element approximation.

Consider the case

$$\kappa(x) = e^x, \quad f(x) = 4e^x(2x + 1).$$

It can be shown that the exact solution is given by the quadratic

$$p(x) = 1 - (2x - 1)^2.$$

Test your implementation in this case by verifying that

$$P_j^h \approx p(x_j^h), \quad j = 1, \dots, N_h$$

for $h = 2^{-10}$. Construct p^h on the mesh \mathbf{x}^h using the coefficients P^h , and produce a plot on logarithmic axes of the error function $e^h(x) = |p^h(x) - p(x)|$.

Remark 2. *The definitions we give of the basis $\{\varphi_j^h\}$ can be simplified for the mesh \mathbf{x}^h that we work with. Note however that the definition makes sense for more general, non-uniform meshes, and the above still applies directly in those cases. Such meshes may be more appropriate to use if either κ or f have strong local behavior that needs to be captured.*

Problem 3. Solution via Quadrature (20 points)

As we are in one dimension, we can solve (PDE) directly by hand. This allows us to implement a reference solution against which the finite element approximations can be compared.

- (a) Assume that $f(x) = F'(x)$ for some $F : [0, 1] \rightarrow \mathbb{R}$. By integrating (PDE) twice, show that the solution to (PDE) is given by

$$p(x) = - \int_0^x \frac{F(y)}{\kappa(y)} dy + \int_0^1 \frac{F(y)}{\kappa(y)} dy \left(\int_0^1 \frac{1}{\kappa(y)} dy \right)^{-1} \int_0^x \frac{1}{\kappa(y)} dy.$$

- (b) Fix $h = h_* := 2^{-14}$. Using the `integral` function in MATLAB, implement the solution p of (PDE) in MATLAB on the mesh \mathbf{x}^h , using the above expression. When calling `integral`, supply relative and absolute tolerances of 10^{-14} . As for the finite element case, the terms κ and f should be provided as function handles.

Test your implementation using the same κ and f as used when testing the finite element method, noting that $f(x) = F'(x)$ where $F(x) = 4e^x(2x - 1)$. Produce a plot on logarithmic axes of the error function for this approximation.

Problem 4. Convergence rates (20 points)

We now study the rates of convergence of the finite element approximations as $h \rightarrow 0$, i.e. as the number of basis elements increases. Throughout this question, we fix

$$\kappa(x) = 1.1 + \sin(25x^2), \quad f(x) = \cos(x).$$

We saw in Problem 2(v) that the error between the true solution and the finite element approximation satisfies the following bound in the $H_0^1(0, 1)$ norm:

$$\|p - p^h\|_{H_0^1} \leq \frac{\kappa_1}{\kappa_0} \cdot \inf_{v^h \in V^h} \|p - v^h\|_{H_0^1}.$$

Choosing v^h to be the piecewise linear interpolant of the true solution p at the mesh points, which belongs to V^h , it is known that there exists a constant $C(p) > 0$ such that

$$\|p - v^h\|_{H_0^1} \leq C(p)h.$$

We can hence deduce that the convergence of p^h to p as $h \rightarrow 0$ is at least linearly fast in the $H_0^1(0, 1)$ norm. We will see how this compares with numerical experiments, and also look at convergence in the $L^2(0, 1)$ norm.

- (a) Denote by p_* the solution to (PDE) using the implementation from Problem 3(b). Define $h_k := 2^{-k}$. Compute the coefficients P^{h_k} of the finite element approximation p^{h_k} for $k = 4, \dots, 12$. Construct each p^{h_k} on the mesh \mathbf{x}^{h_k} using these coefficients, then calculate and store the H_0^1 and L^2 errors,

$$e_k = \|p_* - p^{h_k}\|_{H_0^1}, \quad e'_k = \|p_* - p^{h_k}\|_{L^2},$$

using the `integral`, `gradient` and `interp1` functions in MATLAB. It may also be beneficial to run `format long` to increase the number of decimal places in the output. Present these errors in a table.

- (b) Given a norm $\|\cdot\|$ and $h, h' > 0$, we define the *experimental order of convergence* (EOC) with respect to $\|\cdot\|$ by

$$\text{EOC}(h, h') = \frac{\log(\|p - p^h\|/\|p - p^{h'}\|)}{\log(h/h')}.$$

- (i) Explain why if the convergence is of the form

$$\|p - p^h\| \approx Ch^\alpha,$$

then $\text{EOC}(h, h')$ provides an estimate for α .

- (ii) For each of the H_0^1 and L^2 norms, compute $\text{EOC}(h_k, h_{k+1})$ for $k = 4, \dots, 11$ using the stored values of e_k, e'_k , and present these values in a table. How do the experimental rates of convergence in the H_0^1 norm compare with the theoretical bound above? What do you think the theoretical rate of convergence in the L^2 norm is?

Remark 3. *Faster rates of convergence can be attained by choosing a different family of finite element basis functions $\{\varphi_j^h\}$. Rather than being piecewise linear, they can be chosen to be piecewise quadratic, piecewise cubic, and so on. These basis functions have a larger support than the piecewise linear basis functions, and so the matrix P^h becomes less sparse. Thus there is an increase in, for example, memory requirements, in exchange for the higher rate of convergence.*