

# MPFA schemes for Richards' equation

Truls Moholt

*Master thesis in Applied and Computational Mathematics,  
Institute of Mathematics,  
University of Bergen,  
Autumn 2021*



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
<b>2</b>	<b>Flow in Porous Media</b>	<b>9</b>
2.1	The Representative Elementary Volume . . . . .	9
2.2	Darcy's Law . . . . .	10
2.3	Mass Conservation . . . . .	12
2.4	Two-phase Flow and Richards' Equation . . . . .	12
<b>3</b>	<b>Numerical Approximation Techniques</b>	<b>17</b>
3.1	The Finite Element Method . . . . .	17
3.1.1	Function Spaces . . . . .	18
3.1.2	The Variational Problem . . . . .	21
3.1.3	Existence and Uniqueness . . . . .	23
3.1.4	Galerkin FEM . . . . .	26
3.1.5	Implementation . . . . .	29
3.1.6	Convergence . . . . .	31
3.2	The Finite Volume Method . . . . .	34
3.2.1	Two Point flux Approximation . . . . .	36
3.2.2	MPFA-O Method . . . . .	38
3.2.3	MPFA L-Method . . . . .	42
3.3	Linearization . . . . .	49
<b>4</b>	<b>Convergence of the MPFA-L Method</b>	<b>53</b>
4.1	Modified MPFA-L Method . . . . .	54
4.2	Modified Finite Element Method . . . . .	56
4.3	Convergence Rate Estimates . . . . .	64
<b>5</b>	<b>Convergence of the MPFA-L Method for Richards' Equation</b>	<b>69</b>
<b>6</b>	<b>Numerical Results</b>	<b>75</b>
6.1	Computer Code . . . . .	75
6.2	Elliptic Equation . . . . .	77

6.3	Richards' Equation . . . . .	87
6.3.1	Constant Hydraulic Conductivity . . . . .	87
6.3.2	Non-Linear Hydraulic Conductivity . . . . .	89
<b>References</b>		<b>92</b>

**Abstract**

**Acknowledgements**



# Chapter 1

## Introduction

Understanding porous media and how fluid flows through it has many useful applications, for example predicting the spread of some contaminant in an aquifer. Other examples include CO<sub>2</sub>-storage, geothermal energy and brain modelling. One way of understanding the processes involved in porous media flow is to describe it using partial differential equations, which may be time dependent, non-linear, degenerate and almost always impossible to solve with pen and paper. We therefore want numerical algorithms that solve these PDE's approximately, and at the same time respect important properties of the equations/problems we consider.

In this thesis we focus numerical techniques to solve Richards' equation, which models groundwater flow. It is time dependent, parabolic, with two non-linearities and possibly degenerate. Because it is parabolic we discretize in time with an implicit method, and get a non-linear elliptic PDE for each time step. This is then linearized with a robust linearization scheme, leading to a sequence of linear elliptic PDE's for each time step. Next, we solve these elliptic PDE's with a spatial discretization, which will be the main focus of this thesis.





# Chapter 2

## Flow in Porous Media

In this chapter we introduce the basic concepts of flow in porous media, briefly covering the modeling choices and physics that leads to Richards' equation. The theory in this chapter is to a large extent adapted from [10] and UIB's Porous media course.

### 2.1 The Representative Elementary Volume

A porous medium consists of a solid matrix and some void filled with fluid of one or more phases. In single phase flow, all the pores are filled with one fluid, in two-phase flow however, we have fluid-fluid interfaces in the void. In porous media research, one has come to the realization that the solid matrix is too complex to model. Instead, one takes averages of variables over a reasonable length scale, i.e., the *representative elementary volume* (REV). An important characterization of a porous medium is the *porosity*  $\phi$ , which is defined as

$$\phi := \frac{\text{volume of voids in REV}}{\text{volume of REV}}. \quad (2.1)$$

Another important quantity is the *saturation*  $S_\alpha$  of phase  $\alpha$ , this is defined

$$S_\alpha := \frac{\text{volume of } \alpha \text{ in REV}}{\text{volume of voids in REV}}. \quad (2.2)$$

In single phase flow, the saturation is irrelevant as the saturation is always one. Also note that the volumetric content of phase  $\alpha$  in the REV,  $\theta_\alpha$ , is given by  $\theta_\alpha = S_\alpha \phi$ .

## 2.2 Darcy's Law

In 1856, Henri Darcy performed a famous experiment where he studied the flow of water through sand. To understand his experiment we must first define some variables for measuring water content. First, we assume that the external gravitational force on some fluid is balanced by the pressure gradient force, also known as *hydrostatic equilibrium*. Then the pressure at height  $z$  above datum developed by a water column of height  $h$  above datum is given by

$$p_{abs}(z) = p_{atm} + \rho g(h - z).$$

here  $\rho$  is the density and  $g$  is the gravitational acceleration. If we define the *gauge pressure*  $p$  by  $p := p_{abs} - p_{atm}$  we get an expression for  $p$ :

$$p = \rho g(h - z).$$

This can be rearranged to give an expression for the height, which we from now on refer to as *hydraulic head*:

$$h = \frac{p}{\rho g} + z. \quad (2.3)$$

A *manometer* is a tube with one end in the reservoir and one in open atmosphere, the water level in this tube is then  $\frac{p}{\rho g}$ . The volumetric flow of water is denoted by  $q_d$ . Darcy's experiment is shown in figure 2.1, where water is poured through a cylinder filled with sand. The cylinder has length  $L$  and has cross sectional area  $A$ . His observations are given by the equation called Darcy's law:

$$q_d = -\kappa \frac{A(h_2 - h_1)}{L}. \quad (2.4)$$

Where  $\kappa$  is a positive coefficient of proportionality. Let  $q$  denote the volumetric flow-rate per area:

$$q := \frac{q_d}{A} = -\kappa \frac{h_2 - h_1}{L},$$

we will refer to this as the *flux* of hydraulic potential. We can now state the differential version of Darcy's law. Taking the limit as  $L \rightarrow 0$  we get

$$\mathbf{q} = -\kappa \nabla h. \quad (2.5)$$

We call  $\kappa$  the *hydraulic conductivity* and note that it in general is a rank two tensor, a matrix. The *hydraulic conductivity* also has the property that it is *symmetric*: This is because there are, at every point in the reservoir, two orthogonal directions; one with maximum, and one with minimum hydraulic conductivity. Thus, the matrix,  $\kappa$ , is diagonalizable by a orthogonal matrix.

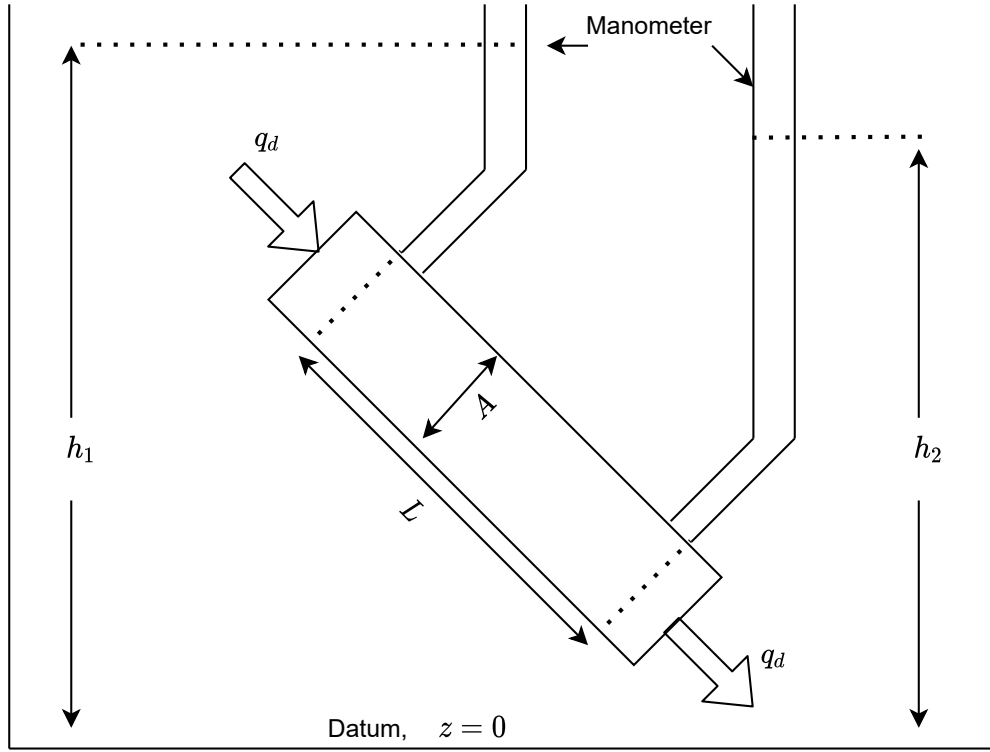


Figure 2.1: The Darcy experiment

The conductivity matrix,  $\boldsymbol{\kappa}$ , is also *positive definite*, this is because there is never flux towards higher pressure. With further experiments, similar to the one already described, we can understand what makes up  $\boldsymbol{\kappa}$ . Dimensionality analysis shows that it is a function of viscosity  $\mu$ , density of the fluid  $\rho$ , gravity  $g$  and *permeability*  $\mathbf{k}$ ,

$$\boldsymbol{\kappa} = \frac{\mathbf{k}\rho g}{\mu}. \quad (2.6)$$

The *permeability*, which is a property of the soil in the reservoir, is also a rank two tensor which is symmetric positive definite and it is in general a function of spatial coordinates, i.e., heterogeneous.

If we define the *pressure head*  $\psi$  as  $\psi := \frac{p}{\rho g}$ , we can combine (2.3), (2.5) and (2.6) to get another variant of Darcy's law;

$$\mathbf{q} = -\frac{\mathbf{k}\rho g}{\mu} \nabla(\psi + z) \quad (2.7)$$

which will be useful later.

## 2.3 Mass Conservation

Darcy's law is not enough if we want to determine the pressure or flow in a reservoir, but we can use the principle of *mass conservation* to add one more equation. The idea is that for every enclosed region in the reservoir, the change of mass inside the region is balanced by the mass flux into the region and the production of mass inside the region.

We end up with the mass balance equation, let  $\Omega$  be our domain, then:

$$\int_{\omega} \frac{\partial(\rho\phi)}{\partial t} dV = - \int_{\partial\omega} \mathbf{n} \cdot \rho \mathbf{q} dS + \int_{\omega} f dV \quad \forall \omega \subseteq \Omega \text{ with } \omega \text{ being a volume,}$$

where  $\mathbf{n}$  is an outward pointing normal vector to  $\omega$  and  $f$  corresponds to sources and/or a sinks. We can use the divergence theorem on the surface integral to get

$$\int_{\omega} \frac{\partial(\rho\phi)}{\partial t} + \nabla \cdot (\rho \mathbf{q}) - f dV = 0.$$

Since this is true for all enclosed regions  $\omega \subset \Omega$ , it also holds for the expressions inside the integral yielding the mass conservation PDE

$$\frac{\partial(\rho\phi)}{\partial t} + \nabla \cdot (\rho \mathbf{q}) = f.$$

This, together with Darcy's law (2.5) and appropriate boundary and initial conditions close the system

$$\left\{ \begin{array}{ll} \mathbf{q} = -\kappa \nabla h, & \mathbf{x} \in \Omega, \quad t > 0 \\ \frac{\partial(\rho\phi)}{\partial t} + \nabla \cdot (\rho \mathbf{q}) = f(\mathbf{x}, t), & \mathbf{x} \in \Omega, \quad t > 0 \\ h(\mathbf{x}, t) = g(\mathbf{x}, t), & \mathbf{x} \in \partial\Omega, \quad t > 0 \\ h(\mathbf{x}, t) = f(\mathbf{x}), & \mathbf{x} \in \Omega, \quad t = 0 \end{array} \right. \quad (2.8)$$

Now we have a model for single-phase flow. As it is stated now, it is a linear parabolic equation, but for incompressible fluid and matrix it becomes an elliptic equation. One often writes the density as a function of pressure, it then becomes non-linear. See chapter two of [10] for a more detailed discussion of (2.8) and modelling options.

## 2.4 Two-phase Flow and Richards' Equation

We restrict our discussion to two phases for simplicity, but the theory can be extended to more phases. In two-phase systems one has a *wetting phase* and a

*non-wetting phase*, denoted by the subscripts  $w$  and  $n$ , respectively.

When we introduce more phases, we continue with the equations we already introduced, i.e., we assume that Darcy's law (2.7) holds for both phases. Let the subscript  $\alpha$  denote the phase, then we have Darcy's law for each phase

$$\mathbf{q}_\alpha = \frac{\mathbf{k}_{r,\alpha} \mathbf{k} \rho g}{\mu} \nabla(\psi_\alpha + z), \quad (2.9)$$

where the coefficient  $\mathbf{k}_{r,\alpha}$  is known as *relative permeability* and it has to be deduced from experimental observation.

We also assume conservation of mass for each phase:

$$\frac{\partial(S_\alpha \rho_\alpha \phi)}{\partial t} + \nabla \cdot (\rho \mathbf{q}_\alpha) = f_\alpha. \quad (2.10)$$

Here, we assume that there is no mass transfer between the phases. If we combine equations (2.9) and (2.10), they give us 2 equations, but we have four unknowns  $\psi_w$ ,  $\psi_n$ ,  $S_w$  and  $S_n$ . We, therefore, introduce the algebraic relation

$$S_w + S_n = 1$$

and the physical relation

$$p_n - p_w = p_c \quad (2.11)$$

where  $p_c$  is called *capillary pressure*. As with the relative permeability,  $p_c$  also needs to be determined experimentally. With initial and boundary conditions we again have a closed system.

A common simplification is to assume that the capillary pressure and the relative permeability are functions of the saturation, and that the relative permeability is isotropic (a scalar).

Another simplification that is used, especially in groundwater hydrology, is that the non-wetting phase (air) always have  $p_n = p_{atm} = 0$ . For this assumption to hold it is important that the air always is connected to the surface. Now, equation (2.11) simplifies to

$$-p_w = -\psi_w \rho g = p_c(S_w).$$

Experiments show that the capillary pressure is a monotone decreasing function of saturation, therefore we can invert it. Equation (2.11) now becomes:

$$p_c^{-1}(\psi_w \rho g) = S_w.$$

Finally, we can multiply the above equation by the porosity to get an expression for the *water content*  $\theta_w$ :

$$\theta_w = \theta_w(\psi_w) = \phi p_c^{-1}(\psi_w \rho g).$$

Combining this with the two-phase Darcy law (2.9) and mass balance (2.10) we get **Richards' equation**

$$\frac{\partial \theta(\psi)}{\partial t} - \nabla \cdot (\boldsymbol{\kappa}(\theta(\psi))(\nabla \psi + e_z)) = f \quad (2.12)$$

where  $\theta = \theta_w$ . Note that the density is eliminated, this is because it is assumed to be constant for water. The hydraulic conductivity is parametrized as a function of water content through experiments and can be written  $\frac{k_{r,\alpha} k_{\rho g}}{\mu} = \boldsymbol{\kappa}(\theta)$ .

Richards' equation contains two non-linearities,  $\theta$  and  $\boldsymbol{\kappa}$ , which make the analysis and numerical simulation more challenging as we will see. They may also cause the equation to degenerate, i.e., the parabolic equation may "collapse" into an elliptic PDE (see figure 2.2 ) or even an ODE when the saturation is so low that there is no flow.

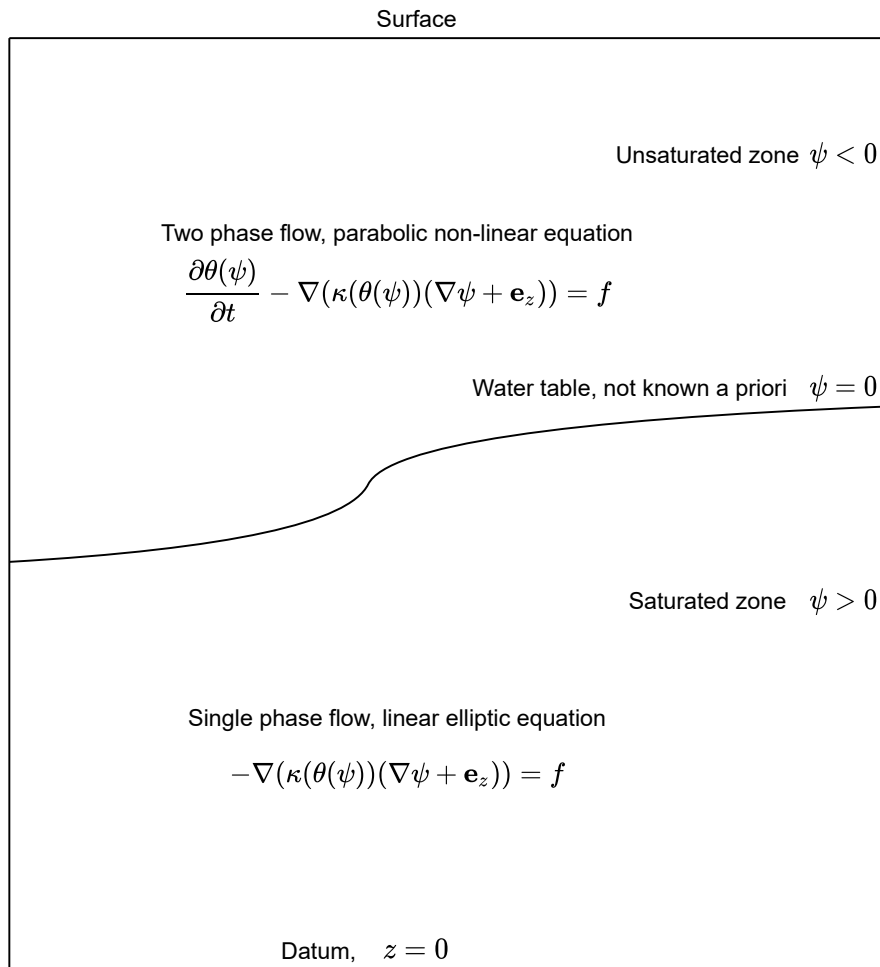


Figure 2.2: A sketch of the degeneracy of Richards' equation





# Chapter 3

## Numerical Approximation Techniques

In this chapter, we first discuss two important frameworks for spatial discretization of PDE's, followed by a brief introduction of time discretization, and at the end, an introduction to linearization. The focus will be on two dimensional elliptic and parabolic equations, but the concepts covered can easily be generalized to three dimensions. After reading this chapter, the reader hopefully has some idea of how to implement a few different methods for solving the Poisson equation, the heat equation and maybe even Richards' equation, and some of their properties.

### 3.1 The Finite Element Method

The finite element method was first developed in the 1940s by Richard Courant for problems in solid mechanics. As computers became better in the 1960s the method increased in popularity [14]. Today there are several general purpose finite element programs being used for a wide range of problems.

In this section we will introduce the finite element method and state the most important results about stability and convergence. We will concentrate on solving the Poisson equation. Let  $\Omega \subset \mathbb{R}^n$  be some open and bounded domain. Find  $u$  such that:

$$\begin{aligned} -\nabla \cdot \mathbf{K} \nabla u(x) &= f(x) \quad x \in \Omega \\ u(x) &= 0 \quad x \in \partial\Omega. \end{aligned} \tag{3.1}$$

For this equation to be well defined we require that  $u$  has double derivatives in  $\Omega$ , but it is easy to come across physical examples where this does not make sense. This is some of the motivation for the Poisson equation in the *variational formulation*. Another motivation is that it allows for a general framework for

computing the solution, as we will soon see. But first, let us introduce some spaces of functions and their properties.

### 3.1.1 Function Spaces

When discussing PDE's and the numerical schemes to solve them it is important to have a precise notion of what kind of functions we are looking for and their properties. The function spaces discussed here are all normed vector spaces. From now on we assume that  $\Omega \subset \mathbb{R}^d$  is a bounded domain.

**Definition 1** (Lebesgue spaces,  $L^p(\Omega)$ ). *For  $p \in [1, \infty)$  let  $L^p(\Omega)$  be the space of functions where  $\|u\|_p = (\int_{\Omega} u^p dx)^{1/p} < \infty$ .*

**Remark 1.** *Note that the  $L^p(\Omega)$  norm induces equivalence relations on the set of functions. Two functions in  $L^p(\Omega)$  are equal if they only differ on a set of measure zero.*

An important concept when discussing normed vector spaces are that they intuitively do not have any points missing, this is formally defined as spaces where every Cauchy sequence converges. This is known as *complete* normed vector spaces or *Banach spaces*.

**Theorem 3.1.1** (Riesz-Fischer Theorem [6] chapter 8). *Each  $L^p(\Omega)$  space is a Banach space.*

**Remark 2.** *The space  $L^2(\Omega)$  is a inner product space, with inner product*

$$\langle u, v \rangle_{L^2} = \int_{\Omega} uv \, dx.$$

*Banach spaces with an inner product, that induces the norm*

$$\langle u, u \rangle^{\frac{1}{2}} = \|u\|,$$

*are called **Hilbert spaces**.*

Before we continue the study of function spaces we develop some convenient notation for derivatives.

**Definition 2** (multi-index notation). *Let  $\bar{\alpha}$  be an ordered  $n$ -tuple. We call this a multi-index and denote the length  $|\bar{\alpha}| = \sum_{i=1}^n \alpha_i$ . For  $\phi \in C^\infty(\Omega)$  we define  $D^{\bar{\alpha}}\phi = (\frac{\partial}{\partial x_1})^{\alpha_1} (\frac{\partial}{\partial x_2})^{\alpha_2} \dots (\frac{\partial}{\partial x_n})^{\alpha_n} \phi$*

We would also like a more general notion of derivative than the one presented in a basic calculus book.

**Definition 3** (weak derivative). Let  $L_{loc}^1(\Omega) = \{ f \in L^1(K) : \forall K \subset \Omega \text{ where } K \text{ is compact} \}$ . Let  $f \in L_{loc}^1(\Omega)$ . If there exists  $g \in L_{loc}^1(\Omega)$  such that

$$\int_{\Omega} g \phi dx = (-1)^{|\bar{\alpha}|} \int_{\Omega} f D^{\bar{\alpha}} \phi dx \quad \forall \phi \in C^{\infty} \quad (3.2)$$

with  $\phi = 0$  on  $\partial\Omega$  we say that  $g$  is the weak derivative of  $f$  and denote it by  $D_w^{\bar{\alpha}} f$ .

We can now define a class of subspaces of the  $L^p$  spaces known as the **Sobolev spaces**.

**Definition 4** (Sobolev space). Let  $k$  be a non-negative integer, define the Sobolev norm as

$$\|u\|_{W^{k,p}(\Omega)} := \left( \sum_{|\bar{\alpha}| \leq k} \|D_w^{\bar{\alpha}} u\|_{L^p(\Omega)}^p \right)^{1/p}.$$

We then define the Sobolev spaces as

$$W^{k,p}(\Omega) = \{ f \in L_{loc}^1(\Omega) : \|f\|_{W^{k,p}} < \infty \}.$$

**Theorem 3.1.2.** The Sobolev spaces  $W^{k,p}(\Omega)$  are Banach spaces

*Proof.* Let  $\{u_i\}_{i=0}^{\infty} \subseteq W^{k,p}(\Omega)$  be a Cauchy sequence. This implies that for all  $\bar{\alpha}$ ,  $|\bar{\alpha}| \leq k$  we have a Cauchy sequence in  $L^p(\Omega)$ :

$$\begin{aligned} \|u_j - u_i\|_{W^{k,p}} &= \left( \sum_{|\bar{\alpha}| \leq k} \|D_w^{\bar{\alpha}} u_j - D_w^{\bar{\alpha}} u_i\|_{L^p(\Omega)}^p \right)^{1/p} < \epsilon \quad \forall i, j \geq N \\ \implies \|D_w^{\bar{\alpha}} u_j - D_w^{\bar{\alpha}} u_i\|_{L^p(\Omega)} &< \epsilon. \end{aligned}$$

By (3.1.1), every  $L^p(\Omega)$  space is a Banach space. Therefore, for each  $|\bar{\alpha}| \leq k$ ,  $D_w^{\bar{\alpha}} u_i$  converges to some limit,  $u_{\bar{\alpha}} \in L^p(\Omega)$ , as  $i \rightarrow \infty$ . In particular  $u_i \rightarrow u$  in  $L^p(\Omega)$ , so the limit in the  $\|\cdot\|_{W^{k,p}(\Omega)}$  norm,  $u$ , is well defined. Now we need to show that  $\{u_{\bar{\alpha}}\}_{\bar{\alpha}}$  are in fact the weak derivatives of  $u$ , i.e.,  $D_w^{\bar{\alpha}} u = u_{\bar{\alpha}}$ . In other words, that the limit of  $u_i$  in the  $\|\cdot\|_{W^{k,p}(\Omega)}$  norm,  $u$ , is in fact in  $W^{k,p}(\Omega)$ . By the definition of weak derivative we have:

$$\int_{\Omega} D_w^{\bar{\alpha}} u_i \phi dx = (-1)^{|\bar{\alpha}|} \int_{\Omega} u_i D^{\bar{\alpha}} \phi dx.$$

Let  $1 = \frac{1}{q} + \frac{1}{p}$ , applying Hölder's inequality on both sides we get the two inequalities:

$$\begin{aligned} \int_{\Omega} (D_w^{\bar{\alpha}} u_i - u_{\bar{\alpha}}) \phi dx &\leq \|D_w^{\bar{\alpha}} u_i - u_{\bar{\alpha}}\|_{L_p} \|\phi\|_{L_q} \\ \int_{\Omega} (u_i - u) D^{\bar{\alpha}} \phi dx &\leq \|u_i - u\|_{L_p} \|D^{\bar{\alpha}} \phi\|_{L_q}. \end{aligned}$$

Taking the limit, the right hand side goes to zero, and by the fact that we can move the limit out of the integral:

$$\begin{aligned}\lim_{i \rightarrow \infty} \int_{\Omega} D_w^{\bar{\alpha}} u_i \phi dx &= \int_{\Omega} u_{\bar{\alpha}} \phi dx \\ \lim_{i \rightarrow \infty} \int_{\Omega} u_i D^{\bar{\alpha}} \phi dx &= \int_{\Omega} u D^{\bar{\alpha}} \phi dx\end{aligned}$$

Now we can put the two equations together with the definition of the weak derivative:

$$\int_{\Omega} u_{\bar{\alpha}} \phi dx = \lim_{i \rightarrow \infty} \int_{\Omega} D_w^{\bar{\alpha}} u_i \phi dx = \lim_{i \rightarrow \infty} (-1)^{|\bar{\alpha}|} \int_{\Omega} u_i D^{\bar{\alpha}} \phi dx = \int_{\Omega} u D^{\bar{\alpha}} \phi dx.$$

We have shown  $D_w^{\bar{\alpha}} u = u_{\bar{\alpha}}$ , and therefore, that  $u \in W^{k,p}(\Omega)$ .  $\square$

**Definition 5.** We rename the  $L^2(\Omega)$  based Sobolev spaces as

$$H^k(\Omega) = W^{k,2}(\Omega),$$

with the norm of  $H^k(\Omega)$  being written in the more compact form  $\|\cdot\|_{\Omega,k}$  or just  $\|\cdot\|_k$ , and the inner product defined as follows:

$$\langle u, v \rangle_k = \sum_{|\bar{\alpha}| \leq k} \int_{\Omega} D_w^{\bar{\alpha}} u, D_w^{\bar{\alpha}} v dx.$$

In Sobolev spaces it is not obvious that a function is well defined on a lower dimensional subset of  $\Omega$ , because two functions may map elements of this zero measure subset to different values and still be of the same equivalence class. This is important to settle if we want to solve boundary value problems. The following results holds for general  $L^p(\Omega)$  based Sobolev spaces, but we will only state them for the Hilbert space  $H^1(\Omega)$ .

**Definition 6.** We denote by  $H_0^k(\Omega)$  the closure of  $C_c^\infty(\Omega)$  in  $H^k(\Omega)$ , where  $C_c^\infty(\Omega)$  is the space of infinitely differentiable functions with compact support.

**Theorem 3.1.3** (Trace theorem, (Evans [7], chapter 5)). Assume  $U$  is bounded and  $\partial U$  is  $C^1$ . Then there exists a bounded, linear operator

$$T : H^1(U) \rightarrow L^2(\partial U)$$

Such that

1.  $Tu = u|_{\partial U}$  if  $u \in H^1 \cap C(\bar{U})$

$$2. \|Tu\|_{L^p(\partial U)} \leq \|u\|_{H^1(U)}$$

We call  $Tu$  the trace of  $u$ . Note that the theorem does not state that  $T$  is surjective.

**Theorem 3.1.4.** (*Trace-zero functions in  $W^{1,p}$ , (Evans [7], chapter 5)*) Suppose  $U$  is as in the previous theorem and  $u \in W^{1,p}(U)$ , then

$$u \in H_0^1 \Leftrightarrow Tu = 0 \text{ on } \partial U$$

**Remark 3.** We often denote the image of  $T$  as:

$$H^{\frac{1}{2}}(\Omega) = T(H^1(\Omega))$$

And define the norm

$$\|f\|_{H^{\frac{1}{2}}(\Omega)} = \inf_{w \in H^1(\Omega), Tw=f} \|w\|_1$$

Now we have the theory we need to study elliptic boundary value problems and their weak solutions.

### 3.1.2 The Variational Problem

We obtain the **variational formulation** of (3.1) by multiplying (3.1) by with a function  $v$  in a suitable space  $V$  called the *test space*, integrating over  $\Omega$  and using integration by parts/divergence theorem,

$$-\int_{\Omega} v \nabla \cdot \mathbf{K} \nabla u \, dx = -\int_{\partial\Omega} v \mathbf{K} \nabla u \cdot \mathbf{n} \, dx + \int_{\Omega} (\nabla v)^T \mathbf{K} \nabla u \, dx = \int_{\Omega} v f \, dx.$$

If we choose  $v$  such that  $v = 0$  on  $\partial\Omega$ , then the integral over the boundary vanishes. The new formulation reads: Find  $u$  such that

$$\int_{\Omega} (\nabla v)^T \mathbf{K} \nabla u \, dx = \int_{\Omega} v f \, dx \quad \forall v \in V. \quad (3.3)$$

A good choice of the test space  $V$  is  $V = H_0^1(\Omega)$ . We also choose this as the solution space. We see that if  $u$  is a solution to (3.1), it also solves (3.3). But a solution to (3.3) does not necessarily solve (3.1), that is why it is also called the *weak formulation*.

The variational problems that we will look at, will all have the form: Find  $u$  such that

$$a(u, v) = b(v) \quad \forall v \in V, \quad (3.4)$$

where  $a(\cdot, \cdot)$  is a *bilinear form* on  $V$  and  $b(\cdot)$  is a *linear functional* on  $V$ . To be precise we define a famous concept from functional analysis:

**Definition 7** (dual space). *Let  $V$  be a normed vector space, then we define its dual space as the space of functions from  $V$  to  $\mathbb{R}$  that are linear and continuous, also called linear functionals. We denote it by  $V'$ . This is a normed vector space with the norm:*

$$\|v\|_{V'} = \sup_{u \in V} \{|v(u)| : \|u\|_V = 1\}.$$

In general, a variational formulation can be seen as finding the element in a Banach space that is mapped to an element in its dual space by some map.

### Boundary Conditions

Let  $\partial\Omega = \Gamma_D \cup \Gamma_N$  with  $\Gamma_D \cap \Gamma_N = \emptyset$ , then (3.1) with more complicated boundary conditions can be written:

$$\begin{aligned} -\nabla \cdot \mathbf{K} \nabla \hat{u}(x) &= f(x) & x \in \Omega \\ \hat{u}(x) &= g_D & x \in \Gamma_D \\ \mathbf{K} \nabla \hat{u}(x) &= g_N & x \in \Gamma_N \end{aligned} \quad (3.5)$$

To make a variational formulation of (3.5) we first define the test space:

$$V = \{v \in H^1(\Omega) : T(v) = 0 \text{ on } \Gamma_D\}.$$

Next, we define the bilinear form:

$$a(u, v) := \int_{\Omega} \nabla u \mathbf{K} \nabla v \, dx. \quad (3.6)$$

Further, assume there exists an element  $w$  of  $H^1(\Omega)$  that are mapped by the trace operator such that Dirichlet boundary conditions are met:  $T(w) = g_D$ . Let  $\hat{u} = u + w$ , now we can use integration by parts to as before:

$$a(u + w, v) = \int_{\Omega} (\nabla u + \nabla w)^T \mathbf{K} \nabla v \, dx = \int_{\Omega} f v \, dx - \int_{\partial\Omega} \mathbf{K} \nabla(u + w) \cdot \mathbf{n} v \, dx. \quad (3.7)$$

Using the linearity of  $a(\cdot, \cdot)$  and inserting boundary conditions we get:

$$a(u, v) = b(v) = \int_{\Omega} f v \, dx - \int_{\Omega} (\nabla w)^T \mathbf{K} \nabla v \, dx - \int_{\Gamma_N} g_N v \, dx. \quad (3.8)$$

Hence both Dirichlet and Neumann boundary conditions are incorporated into the right hand side. For homogeneous Dirichlet boundary conditions, the second term on the right hand side of (3.8) vanishes.

### 3.1.3 Existence and Uniqueness

We still need to show that (3.8) has an unique solution. First, we define some important properties that a variational problem should have in order to have a unique solution. Let  $(V, \|\cdot\|_V)$  be a Hilbert space.

**Definition 8.** Let  $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  be a bi linear form. We say that:

- $a(\cdot, \cdot)$  is **coercive with respect to**  $V$ , or **elliptic** if there exists a constant  $C_c \in \mathbb{R}$  such that  $C_c \|u\|_V^2 \leq a(u, u) \forall u \in V$
- $a(\cdot, \cdot)$  is **bounded** or **continuous** if there exists a constant  $C_B$  such that  $|a(u, v)| \leq C_B \|u\|_V \|v\|_V \forall u, v \in V$

In order to prove existence and uniqueness, we must first state some important results about the underlying space  $V$ . The following theory can be found in its entirety in the first four chapters of Cheney [6]

**Theorem 3.1.5.** If  $Y$  is a closed subspace of the Hilbert space  $X$ , then

$$X = Y \oplus Y^\perp,$$

where  $Y^\perp = \{x \in X : \langle x, y \rangle = 0 \forall y \in Y\}$  is orthogonal complement. In other words, an element in  $X$  can always be written as the sum of an element  $Y$  and an element in  $Y^\perp$ .

**Theorem 3.1.6** (Riesz representation theorem). Every continuous linear functional,  $\phi(x)$ , defined on a Hilbert space  $X$  can be written  $\phi(x) = \langle x, v \rangle$  by a uniquely determined  $v \in X$ .

*Proof.* Let  $\phi \in X'$ , define  $Y = \{x \in X : \phi(x) = 0\}$ . Take a non-zero element in the orthogonal complement  $u \in Y^\perp$  such that  $\phi(u) = 1$ , (if this does not exist then  $X = Y$  and  $\phi(x) = \langle x, 0 \rangle$ , this is ensured by theorem 3.1.5). Now we can write every vector in  $X$  as a linear combination of a vector in  $Y$  and the vector  $u$ .  $x = x - \phi(x)u + \phi(x)u$  for any  $x \in X$ . Using this, we can find an expression for the inner product of  $x$  with a scaled version of  $u$

$$\left\langle x, \frac{u}{\|u\|^2} \right\rangle = \left\langle x - \phi(x)u, \frac{u}{\|u\|^2} \right\rangle + \left\langle \phi(x)u, \frac{u}{\|u\|^2} \right\rangle. \text{ The first part of the sum vanishes as } x - \phi(x)u \in Y. \text{ So we end up with}$$

$$\left\langle x, \frac{u}{\|u\|} \right\rangle = \phi(x) \frac{\langle u, u \rangle}{\|u\|^2} = \phi(x).$$

□

**Theorem 3.1.7** (Banach fixed point theorem). Let  $X$  be a Banach space and  $F : X \rightarrow X$  an operator where  $\|Fx - Fy\|_X \leq \theta \|x - y\|_X$  for some  $\theta \in (0, 1)$ , we call this a **contraction**.

Then for all  $x \in X$  the sequence  $[x, Fx, F^2x, \dots]$  converges to a point  $x^* \in X$  called the fixed point of  $F$ .

See page 177 of [6] for a proof.

**Theorem 3.1.8** (Lax-Milgram). *Suppose that  $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  is a bilinear, bounded and coercive form and that  $b(\cdot) : V \rightarrow \mathbb{R}$  is a bounded, linear functional. Then the variational problem has a unique solution  $u \in V$ , such that*

$$a(u, v) = b(v) \quad (3.9)$$

for all  $v \in V$ .

**Remark 4.** *If  $a(\cdot, \cdot)$  also is symmetric, it defines an inner product on  $V$  giving a complete space. We can then use Riesz representation theorem 3.1.6 to show that it has an unique solution.*

*proof of Lax Milgram theorem 3.1.8.*

For each  $w$  denote the map  $a(w, v) = a_w(v)$ , this is a linear continuous functional, and follows from the assumptions on  $a$ . By Riesz representation theorem 3.1.6  $a_w(\cdot)$  uniquely determines an element  $Aw \in V$  such that  $a_w(v) = \langle Aw, v \rangle$ . The map

$$\begin{aligned} A : V &\rightarrow V \\ w &\mapsto Aw, \end{aligned}$$

- is linear:  $\langle A(x + y), v \rangle = a_{x+y}(v) = a(x + y, v) = a_x(v) + a_y(v) = \langle Ax, v \rangle + \langle Ay, v \rangle$ . Since this holds for all  $v \in V$ , we have  $A(x + y) = Ax + Ay$ .
- is bounded:  $\|Ax\| = \|a_x\| = \sup \{a(x, v) : \|v\| = 1\} \leq C_B \|x\|$ .

We can also use Riesz representation theorem on the right hand side:  $b(\cdot) = \langle f, \cdot \rangle$ . Now we have a reformulation of (3.9): Find  $u$  such that

$$Au = f. \quad (3.10)$$

Now we need to show that (3.10) has an unique solution, and for that we need the Banach fixed point theorem. Let  $\epsilon > 0$ , we define the operator

$$\begin{aligned} T : V &\rightarrow V \\ u &\mapsto u - \epsilon(Au - f). \end{aligned}$$

If  $T$  has a fixed point  $u^*$ , then  $u^* - \epsilon(Au^* - f) = u^* \Rightarrow Au^* = f$  and we have solved (3.10) and proved the theorem. We just need to show that  $T$  is a contraction. First, let  $u_1, u_2 \in V$ , and subtract what they are mapped to by  $T$ , then we get

$$\|Tu_1 - Tu_2\|^2 = \|u - \epsilon(Au)\|^2,$$



where  $u = u_1 - u_2$ , we used the linearity of  $A$ ,

$$= \|u\|^2 - 2\epsilon \langle u, Au \rangle + \epsilon^2 \langle Au, Au \rangle.$$

Now, we can use that  $a(u, u) = \langle Au, u \rangle$ , and that  $\langle Au, Au \rangle = a_u(Au) = a(u, Au)$ :

$$\|Tu_1 - Tu_2\|^2 = \|u\|^2 - 2\epsilon a(u, u) + \epsilon^2 a(u, Au).$$

Next, we use the coercivity and boundedness of  $a(\cdot, \cdot)$ . We also use the boundedness of  $A$

$$\|Tu_1 - Tu_2\|^2 \leq \|u\|^2 - 2\epsilon C_c \|u\|^2 + \epsilon^2 C_B^2 \|u\|^2.$$

This leads to the inequality

$$\|Tu_1 - Tu_2\|^2 \leq \|u_1 - u_2\|^2 (1 - 2\epsilon + \epsilon^2).$$

We choose  $\epsilon$  such that  $T$  becomes a contraction:  $\epsilon < \frac{2C_c}{C_B^2} \Rightarrow (1 - 2\epsilon + \epsilon^2) < 1$ . By the Banach fixed point theorem we have existence and uniqueness of a solution.  $\square$

**Remark 5.** *The solution,  $u$ , to our variational problem depends on the data  $b(\cdot)$ . To see this we use the coercivity:*

$$\|u\|^2 \leq \frac{a(u, u)}{C_c} = \frac{b(u)}{C_c}.$$

Now, we have proved that (3.4) has a unique solution for suitable  $a$  and  $b$ . The variational form of Poisson equation (3.3) satisfies this:

**Example 1** (Well posedness of variational form of Poisson equation). *Let  $a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx$ . Then we have that:*

- $a(\cdot, \cdot)$  is **Coercive** with respect to  $\|\cdot\|_{H_0^1}$ :

$$\begin{aligned} \|u\|_{H_0^1}^2 &= \|u\|_{L^2}^2 + \sum_{|\bar{\alpha}|=1} \|D^{\bar{\alpha}} u\|_{L^2}^2 \\ &= \|u\|_{L^2}^2 + a(u, u) \\ &\leq (C_{\Omega} + 1)a(u, u), \end{aligned}$$

where we used the **Poincaré inequality** in the last step.

- $a(\cdot, \cdot)$  is **Bounded** with respect to  $\|\cdot\|_{H_0^1}$ :

$$\begin{aligned} |a(u, v)| &\leq \left| \int_{\Omega} \nabla u \cdot \nabla v \, dx \right| \leq \int_{\Omega} |\nabla u \cdot \nabla v| \, dx \\ \int_{\Omega} \left| \sum_{|\bar{\alpha}|=1} D^{\bar{\alpha}} u D^{\bar{\alpha}} v \right| \, dx &= \sum_{|\bar{\alpha}|=1} \|D^{\bar{\alpha}} u D^{\bar{\alpha}} v\|_{L^1} \leq \sum_{|\bar{\alpha}|=1} \|D^{\bar{\alpha}} u\|_{L^2} \|D^{\bar{\alpha}} v\|_{L^2} \\ &\leq \|u\|_{H_0^1} \|v\|_{H_0^1}, \end{aligned}$$

where we used the **Cauchy-Schwarz inequality** on the second line.

- $b(\cdot)$  is in the dual space of  $H_0^1$  if  $f \in L^2(\Omega)$ :

$$\begin{aligned} |b(v)| &= \left| \int_{\Omega} f v dx \right| \leq \|f\|_{L^2} \|v\|_{L^2} \\ \Rightarrow \|b\|_{H_0^{1'}} &= \sup \left\{ \frac{|b(v)|}{\|v\|} \right\} \leq \|f\|_{L^2} \end{aligned}$$

Hence, (3.3) is well posed and we get a solution  $u \in H_0^1(\Omega)$ .

### 3.1.4 Galerkin FEM

Now we want to discretize the variational equation (3.4), we do this by replacing the test space  $V$  by a finite dimensional subspace  $V_h$ . This is called the *Galerkin method*. The discretization now reads: Find  $u \in V_h$  such that

$$a(u, v_h) = b(v_h) \quad (3.11)$$

for all  $v_h$  in  $V_h$ . Since  $a(\cdot, \cdot)$  is bilinear and  $b(\cdot)$  is linear, it is easy to see that if (3.11) holds for the basis functions of  $V_h$ , it holds for all elements in  $V_h$ . In the *finite element method*, the finite dimensional subspace are determined by the *triangulation*. In this thesis, we only consider problems in two spatial dimensions, so let  $\Omega \subset \mathbb{R}^2$ .

**Definition 9** (two dimensional triangulation, page 56 of Knabner [8]). *Let  $\tau_h$  be a partition  $\Omega$  into closed triangles  $K$  including the boundary  $\partial\Omega$ , with the following properties*

**(T1)**  $\overline{\Omega} = \bigcup_{K \in \tau_h} K$ .

**(T2)** For  $K, K' \in \tau_h$ ,  $K \neq K'$

$$\text{int}(K) \cap \text{int}(K') = \emptyset,$$

where  $\text{int}(K)$  denotes the interior of  $K$ .

**(T3)** If  $K \neq K'$ , but  $K \cap K' \neq \emptyset$ , then  $K \cap K'$  is either a point or a common edge of  $K$  and  $K'$ .

The above definition sets some rules on how we can divide our domain into triangles, often called elements. Now that we have a triangulation, we now define our finite dimensional subspace,  $V_h$ .

**Definition 10** (Linear ansatz space). *Let  $\mathcal{P}_1(K)$  be the space of linear polynomials in two variables on  $K \subset \mathbb{R}^2$ , then we define the ansatz space*

$$V_h := \left\{ u_h \in C(\overline{\Omega}) : u_h|_K \in \mathcal{P}_1(K) \ \forall K \in \tau_h, u|_{\Gamma_D} = 0 \right\},$$

*of piecewise linear functions on each  $K$ .*

**Remark 6.** *Our local ansatz space  $P_K = \{v|_K : v \in V_h\}$  is such that  $P_K = \mathcal{P}_1(K) \subset H^1(K) \cap C(K)$ . This together with **(T3)**, which ensures continuity between elements, makes  $V_h$  a conformal finite element method, i.e.,  $V_h \subset V = H_0^1(\Omega)$*

**Remark 7** (Nodes). *We will refer to the corners of the triangles in  $\tau_h$  as nodes. For more advanced element types one can have nodes also on the edges or interiors of the triangles.*

**Remark 8.** *In general, finite elements are defined by an element  $K(\in \tau_h)$ , the local ansatz space  $P_K$  and degrees of freedom  $\Sigma_K$ . In all Lagrange finite element methods  $\Sigma_K$  corresponds the evaluation of functions in  $P_K$  at the nodes of the element.*

A choice of basis for  $V_h$  could then be the hat functions. Let  $\phi_i$  be the basis function corresponding to the node  $x_i$ , it is defined by:

$$\phi_i(x_j) = \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}, \quad \phi_i \in V_h.$$

There are no basis functions defined for the nodes at the Dirichlet boundary.

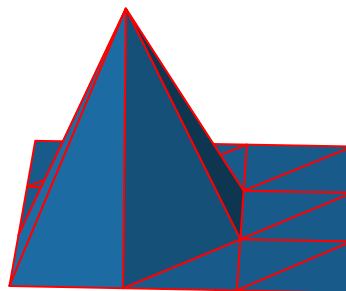


Figure 3.1: A hat function.

Now, we demonstrate how the method works in practice. We seek a solution  $u_h \in V_h$ . Write this in terms of the basis functions:  $u_h = \sum_{i=1}^n \hat{u}_i \phi_i$ . Equation (3.11) can then be written as an equation with a solution vector with real coefficients:

$$\text{Find } \hat{\mathbf{u}}_h = \begin{pmatrix} \hat{u}_1 \\ \vdots \\ \hat{u}_n \end{pmatrix} \in \mathbb{R}^n \text{ such that } \sum_{i=1}^n \hat{u}_i a(\phi_i, \phi_j) = b(\phi_j). \quad (3.12)$$

So we get a system of linear equations  $\mathbf{A}\hat{\mathbf{u}}_h = \mathbf{b}$ , where  $\mathbf{A}_{i,j} = a(\phi_i, \phi_j)$  and  $\mathbf{b}_j = b(\phi_j)$ . The matrix,  $\mathbf{A}$ , has as many rows and columns as there are nodes (the Dirichlet nodes can be removed, depending on implementation). If we solve (3.3), our variational problem, and also matrix, will be symmetric. The matrix is then often called a *stiffness matrix*. These names originated from mechanics and structural analysis, where the solution represents displacement and the force function represents load. The stiffness matrix is also sparse, which is a very important property when designing algorithms to solve it.

With the setup described in this subsection, the degrees of freedom are the same as the dimension of  $V_h$ . If we in definition 10 instead had chosen a space of quadratic polynomials on each element, we had gained three degrees of freedom on each element. In this thesis we focus on linear finite elements because we do not gain anything from increasing regularity, as the solutions to problems in porous

media flow are not expected to be very regular. Also, the finite volume methods we will discuss later, in particular MPFA-L-Method, are not higher order methods.

### 3.1.5 Implementation

Here we explain the most important parts of the algorithm for discretizing elliptic PDE's with linear Lagrange finite elements. We consider the homogenous elliptic model problem (3.3) in two dimensions, with  $\mathbf{K} = \mathbf{I}$  and zero Dirichlet boundary conditions. The procedure goes as follows:

1. Make a triangulation of the domain. This can be done in a number of different ways, see chapter 4 of Knabner [8]. If we have  $N$  nodes, our triangulation would be stored as a  $N \times 2$  array of floats, being the coordinates of the nodes. And a  $E \times 3$  array of ints being the elements, where each entry is the index of a coordinate in the coordinate matrix,  $E$  is the number of elements.
2. Allocate space for the  $N \times N$  stiffness matrix  $\mathbf{A}$  and the  $N \times 1$  source vector  $\mathbf{b}$ .
3. Define the basis functions on a reference element, this is also called the shape functions, see figure 3.2 and (3.13). Also, compute the gradients of the shape functions.

$$\begin{aligned}
 N_1(x, y) &= 1 - x - y \\
 N_2(x, y) &= x \\
 N_3(x, y) &= y
 \end{aligned} \tag{3.13}$$

Figure 3.2: The map  $F$  from element  $K$  to the reference element  $\hat{K}$ .

4. Loop through the elements. For each element  $K$  compute the affine linear map that maps it to the reference element. That means we want to find  $B \in \mathbb{R}^{2 \times 2}$  and  $d \in \mathbb{R}^2$  such that

$$\begin{aligned}
 F : K &\rightarrow \hat{K} \\
 x &\mapsto Bx + d.
 \end{aligned} \tag{3.14}$$

To achieve this we set up a system of equations inspired by figure 3.2

$$\begin{pmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ x_3 & y_3 & 1 \end{pmatrix} \begin{pmatrix} b_{1,1} & b_{2,1} \\ b_{1,2} & b_{2,2} \\ d_1 & d_2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}. \quad (3.15)$$

So for each element we solve (3.15) for  $B$  and  $d$ , that means computing an inverse of a three by three matrix and a matrix product. Note that this only needs to be done once per element and could be done in a preprocessing step.

Now that we have  $T$ , we do the following on the element:

- (a) Use the map and the shape functions to evaluate  $a(\phi_i, \phi_j)|_K$  for  $1 \leq i, j \leq 3$ . Note that for  $u : K \rightarrow \mathbb{R}$  we get by the chain rule:

$$\nabla_{\hat{x}}^T u(F^{-1}(\hat{x})) = \nabla_x^T u(F^{-1}(\hat{x})) \nabla_{\hat{x}}^T F^{-1}(\hat{x}) = \nabla_x^T u(F^{-1}(\hat{x})) B^{-1}. \quad (3.16)$$

This gives an expression for the derivative on an element expressed as a derivative in the reference element coordinate system:

$$\nabla_x u(F^{-1}(\hat{x})) = B^T \nabla_{\hat{x}} u(F^{-1}(\hat{x})). \quad (3.17)$$

Now we can compute the product of the gradients of the basis functions on an element:

$$\begin{aligned} a(\phi_i, \phi_j)|_K &= \int_K (\nabla \phi_i)^T \nabla \phi_j dx \\ &= \int_{\hat{K}} (\nabla_x \phi_i(F^{-1}(\hat{x})))^T \nabla_x \phi_j(F^{-1}(\hat{x})) |\text{Det}(J(F^{-1}))| d\hat{x} \\ &= \int_{\hat{K}} (B^T \nabla_{\hat{x}} \phi_i(F^{-1}(\hat{x})))^T B^T B \nabla_{\hat{x}} \phi_j(F^{-1}(\hat{x})) |\text{Det}(B^{-1})| d\hat{x} \\ &= \int_{\hat{K}} (\nabla_{\hat{x}} N_i(\hat{x}))^T B^T B \nabla_{\hat{x}} N_j(\hat{x}) |\text{Det}(B^{-1})| d\hat{x} \\ &= \frac{1}{2} (\nabla_{\hat{x}} N_i(\hat{x}))^T B^T B \nabla_{\hat{x}} N_j(\hat{x}) \frac{1}{|\text{Det}(B)|} \end{aligned} \quad (3.18)$$

So for each element we evaluate the last line of (3.18) for all (9) combinations of  $i$  and  $j$  on the element and add this to  $\mathbf{A}_{i,j}$ . This approach is called *element-based assembling*, and  $\mathbf{A}_{i,j} = \sum_{K \in \mathcal{N}(i)} a(\phi_i, \phi_j)|_K$ , where  $\mathcal{N}(i)$  is the set of all elements that contain node  $i$ .

- (b) In almost the same way we compute  $b(\phi_i)|_K$  and add this to  $\mathbf{b}_i$ . As in

(3.18), we compute the integral on the reference element:

$$\begin{aligned} b(\phi_i)|_K &= \int_{\hat{K}} f(F^{-1}(\hat{x})) \phi_i(F^{-1}(\hat{x})) \frac{1}{\text{Det}(B)} d\hat{x} \\ &= \int_{\hat{K}} \hat{f}(\hat{x}) N_i(F^{-1}(\hat{x})) \frac{1}{\text{Det}(B)} d\hat{x} \\ &\approx \frac{1}{\text{Det}(B)} \sum_k \omega_k \hat{f}(\hat{p}_k) N_i(\hat{p}_k) \end{aligned} \quad (3.19)$$

Where  $\hat{f} := f(F^{-1}(\hat{x}))$  and  $\{(\omega_k, \hat{p}_k)\}_k$  defines a *quadrature rule*. This can be chosen in different ways, for higher order finite elements this may even affect the convergence behaviour. But for linear Lagrange elements, the trapezoidal rule works fine, i.e., using three points per element with appropriate weights.

5. Loop through the Dirichlet boundary nodes  $x_j$  at the boundary and set  $\mathbf{A}_{j,i} = \delta_{ij}$ ,  $b_j = 0$ . This fixes the value of  $u$  at the Dirichlet boundary to zero.

**Remark 9.** *If we have inhomogeneous Dirichlet boundary conditions this is in practice done the same way as in the homogenous case, eliminating the degrees of freedom on the boundary. For Neumann conditions one has to evaluate integrals along the boundary as in (3.8), using one-dimensional elements.*

### 3.1.6 Convergence

In this subsection, we review the most important concepts regarding the convergence of FEM. For a detailed discussion see [8]. The starting point of convergence rate estimates for the finite element method already described are **C  a's lemma**:

**Theorem 3.1.9** (C  a's lemma). *Let  $u$  solve the variational problem (3.4) and  $u_h$  solve the corresponding Galerkin approximation (3.11), where the bilinear form  $a(\cdot, \cdot)$  is bounded and coercive. Then we have:*

$$\|u - u_h\|_V \leq \frac{C_b}{C_c} \min \{\|u - v_h\| : v_h \in V_h\}. \quad (3.20)$$

*Proof.* By the coercivity and linearity of  $a(\cdot, \cdot)$  we have:

$$C_c \|u - u_h\|_V^2 \leq a(u - u_h, u - u_h) = a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h).$$

The last term equals zero, since both  $u$  and  $u_h$  solves the variational problem in  $V_h$ :  $v_h - u_h = v \in V_h$  and  $a(u - u_h, v) = a(u, v) - a(u_h, v) = b(v) - b(v) = 0$ , this is called *Galerkin orthogonality*. Then, we use the boundedness of  $a(\cdot, \cdot)$ :

$$C_c \|u - u_h\|_V^2 \leq a(u - u_h, u - u_h) \leq C_b \|u - u_h\|_V \|u - v_h\|_V.$$

We divide by  $C_c$  and  $\|u - u_h\|_V$  and take the infimum over  $v_h \in V_h$ :

$$\|u - u_h\|_V \leq \frac{C_b}{C_c} \inf \{ \|u - v_h\|_V : v_h \in V_h \}.$$

By (Cheney [6], page 64, theorem 2), as  $V_h$  is closed and convex subspace of a Hilbert space, there exist an unique element of  $V_h$  closest to  $u$  and minimum is attained.  $\square$

Hence the solution to Galerkin problem is the best in the subspace  $V_h$  up to a constant. We can therefore study convergence rate estimates for a suitable comparison element in  $V_h$ . In one dimension it is easy to picture what this comparison element might be, see figure 3.3. A direct proof with techniques from calculus is possible in this case.

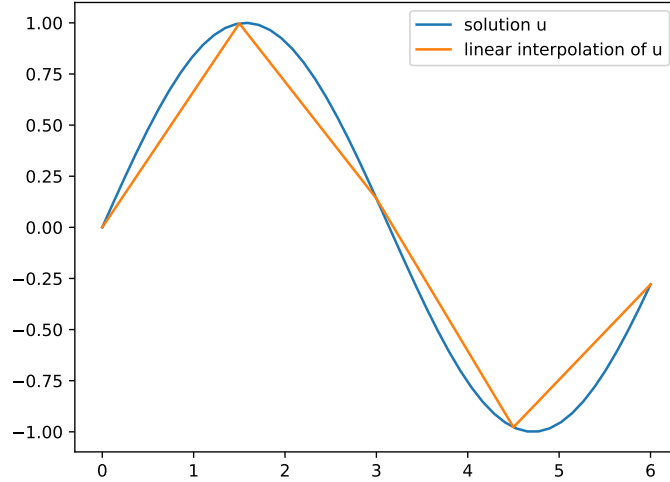


Figure 3.3: The unique linear interpolation of a function in one dimension.

The idea for more dimensions are the same, to be precise we define the interpolation operator.

**Definition 11** (Global interpolation operator).

$$\begin{aligned} I_h : C(\bar{\Omega}) &\rightarrow V_h \\ v &\mapsto \sum_i v(n_i) \phi_i \end{aligned} \tag{3.21}$$

Where  $\{n_i\}_i$  are the nodes and  $\{\phi_i\}_i$  the corresponding basis functions.



**Remark 10.** *The global interpolator operator (3.21) maps from continuous functions, so we need to make sure our solution is continuous. By the Sobolev embedding theorem (Evans [7], page 286), we are okay if our space dimension is such that  $\Omega \subset \mathbb{R}^d$  for  $d \leq 3$ , and  $u \in H^k(\Omega)$  for  $k \geq 2$ .*

We remind the reader of the notation  $\|\cdot\|_1 = \|\cdot\|_{H^1(\Omega)}$ , and similarly for the semi-norms. In the setting of the model problem (3.3), we hope to reach an estimate

$$\|u - u_h\|_1 \leq C \|u - I_h(u)\|_1 \leq C^* h^k |u|_{k+1}, \quad (3.22)$$

where  $h$  is the maximum diameter of the elements in the triangulation, and  $k$  is the polynomial degree on the ansatz space. This bound is indeed attainable if we make sure the triangles in our triangulation have maximum angle less than  $\pi$ . In chapter 3.4 of Knabner [8], there is a detailed proof of (3.22).

Note that this means that our linear finite element method has a linear convergence in the  $\|\cdot\|_1$  norm, if our variational problem admits a solution with sufficient regularity. We tie these observations together in a theorem:

**Theorem 3.1.10** (energy norm estimate, Knabner [8] page 144). *Consider a finite element discretization as described by (3.12) in  $\mathbb{R}^d$  for  $d \leq 3$  on a family of triangulations with an uniform upper bound on the maximal angle. Suppose we have a linear ansatz space as in 10, then*

$$\|u - u_h\|_1 \leq Ch |u|_2. \quad (3.23)$$

The above is called an energy norm estimate due to the equivalence of  $\|\cdot\|_1^2$  and  $a(\cdot, \cdot)$  in case of a symmetric bi linear form, in structural mechanics  $a(\cdot, \cdot)$  corresponds to the potential energy.

Often we are happy with a convergence rate estimate in the  $\|\cdot\|_0$  norm, which does not measure an error in the approximation of the derivative. We then expect a better convergence rate, as can be shown by the *Aubin-Nitsche technique*. We consider the dual problem of our variational problem (3.3):  $a(v, u_f) = \langle f, v \rangle_0$ , and assume some uniqueness and stability of the solution  $u_f$  of this.

**Theorem 3.1.11** ( $L^2$  estimate). *Suppose the situation of theorem 3.1.10 and assume there exists a unique solution to the adjoint problem with  $|u_f| \leq C \|f\|_0$ , then there exist a constant  $C^*$  such that:*

$$\|u - u_h\|_0 \leq C^* h \|u - u_h\|_1. \quad (3.24)$$

See [8] for a proof. When it comes to the assumption on the dual problem, this is satisfied for our elliptic model problem 3.1. If we put the last two theorems together we obtain quadratic convergence in the  $L^2$  norm.

**Remark 11.** *In this chapter we have only discussed the convergence behaviour of the solution to the Galerkin problem (3.11). In practice, one often only solves this approximately. For example the term  $b(v_h) = \int_{\Omega} f v_h dx$  is impossible to evaluate exactly depending on  $f$ . We will later see error estimates with this taken into account.*

## 3.2 The Finite Volume Method

Finite volume methods are designed such that the conservation law we solve hold everywhere in the domain. Consider our elliptic model problem (3.1): Find  $u$  such that

$$\begin{aligned} -\nabla \cdot \mathbf{K} \nabla u(x) &= f(x) \quad x \in \Omega \\ u(x) &= 0 \quad x \in \partial\Omega. \end{aligned} \quad (3.25)$$

First we divide our domain  $\Omega$  into convex quadrilaterals (control volumes, cells),  $\{\Omega_i\}_i$ . Then we integrate our equation over  $\Omega_i$  and apply the divergence theorem:

$$\int_{\Omega_i} -\nabla \cdot \mathbf{K} \nabla u dx = - \int_{\partial\Omega_i} \mathbf{K} \nabla u \cdot \hat{\mathbf{n}} ds = \int_{\Omega_i} f dx. \quad (3.26)$$

The above equation equates the fluxes through the boundary of a control volume, with the source or sinks inside the control volume. The finite volume methods are discrete versions of this. Let  $E_{i,j}$  be the edge between cells  $i$  and  $j$ . Then the main idea is to approximate the flux through  $E_{i,j}$ , from cell  $i$  to cell  $j$ ,

$$q_{E_{i,j}} = - \int_{E_{i,j}} \mathbf{K} \nabla u \cdot \hat{\mathbf{n}} ds \quad (3.27)$$

by a linear combination of  $u_i$  at neighbouring cell centers (geometric center of a cell)

$$q_{E_{i,j}} \approx \tilde{q}_{E_{i,j}} = \sum_k t_{i,j}^k u^k. \quad (3.28)$$

Where the *transmissibility*  $t_{i,j}^k$  has the property  $\sum_k t_{i,j}^k = 0$ . Note that with this notation, we have  $q_{E_{i,j}} = -q_{E_{j,i}}$ .

We also approximate the integral on the right side,  $\int_{\Omega_i} f dx$ , with some quadrature rule. In porous media flow, the space discretization used, usually has a truncation error of at most second order. This is because the solution has low regularity due to heterogeneous permeability. The upshot is that we use the midpoint rule for evaluating the right hand side, as this also has a second order truncation error.

Hence we evaluate  $f$  at the cell center and multiply by the area of  $\Omega_i$ . We then end up with a system of equations

$$\sum_{j \in \mathcal{S}_i} \tilde{q}_{E_{i,j}} = |\Omega_i| f(x_i), \quad (3.29)$$

where  $\mathcal{S}_i$  is the set of indices of neighbouring cells. The system of equations (3.29) ensures local mass conservation. It can also be written in matrix form as:

$$\mathbf{A}^V \tilde{\mathbf{u}}_h = \mathbf{f}. \quad (3.30)$$

We will discuss different ways of constructing the transmissibility coefficients, as they result in very different discretizations.

The motivation for using finite volume methods for problems in porous media, for example Richards' equation, is that the flux appears explicitly in our discretization. If one, for example, wants to simulate the spread of some contaminant by groundwater flow, one can easily obtain a local mass conservative flux field using the finite volume method. This flux field can then be used in the desired transport equation.

We always make sure that our control volumes are inside our domain, with the boundary of our domain aligning the edges of our grid. To set our boundary conditions we thus need to specify the flux across the boundary. If we have Neumann boundary conditions, this is straightforward. One way of implementing zero Neumann boundary conditions is to make a strip of cells outside our domain with zero permeability.

Dirichlet boundary conditions is not always so natural for the equations we consider, that is knowing the pressure or the saturation on a thin line in a two dimensional domain. This is reflected in the finite volume framework, where a common approach is to make ghost cells outside our domain where the potential is known, see figure 3.4. We can make these ghost cells as small as we like, and this approach is easy to implement. On the other hand, if we insist we only know the potential at  $\partial\Omega$ , there exists techniques for determining the flux across the boundary as well. In section 4.1 demonstrate an approach for the MPFA-L-method where this is achieved.

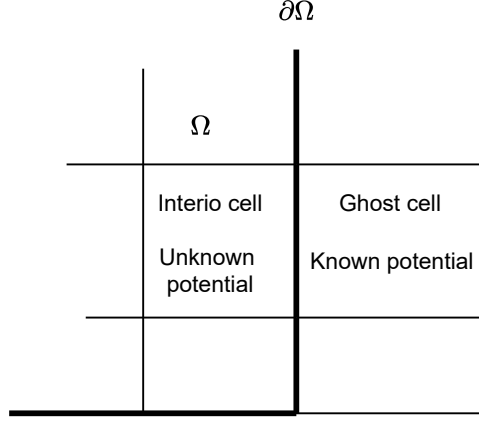


Figure 3.4: Ghost cell outside the boundary.

We will focus on discretizing the interior of the domain in the following sections.

### 3.2.1 Two Point flux Approximation

The simplest way of constructing  $t_{i,j}^k$  is also the most popular in the industry. As the name suggests, we only use the function value at two points,  $x_0$  and  $x_1$ , which are the cell center of two neighbouring cells, to compute the numerical flux  $\tilde{q}_{E_{0,1}}$  between them.

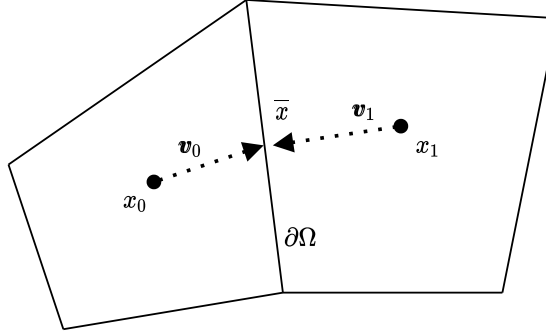


Figure 3.5: The two point flux approximation (TPFA) setup.

Let  $\mathbf{v}_1$  be the vector from cell center  $x_1$  to the midpoint of the edge between the cells,  $\bar{x}$ . Then we approximate the flux out of cell  $x_0$  into cell  $x_1$  by:

$$\tilde{q}_{E_{0,1},0} = -\mathbf{n}_0^T \mathbf{K}_0 \frac{\mathbf{v}_0}{\|\mathbf{v}_0\|} (u(\bar{x}) - u(x_0)) \quad (3.31)$$

or as

$$\tilde{q}_{E_{0,1},1} = -\mathbf{n}_1^T \mathbf{K}_1 \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|} (u(x_1) - u(\bar{x})) \quad (3.32)$$

where  $\mathbf{n}_i$  is the normal vector pointing out of cell  $i$  with length equal to  $\partial\Omega$ . In figure 3.5 we have  $\mathbf{n}_1 = -\mathbf{n}_0$ , and they are in general not aligned with  $\mathbf{v}_1$  or  $\mathbf{v}_0$ . Because we require flux continuity we have that

$$\tilde{q}_{E_{0,1},0} = \tilde{q}_{E_{0,1},1} = t^0 u(x_0) + t^1 u(x_1) \quad (3.33)$$

where, as before,  $t^0 + t^1 = 0 \Rightarrow t^0 = -t^1$ , and the subscript on  $t$  is dropped for readability. We now have three equations and three unknowns,  $u(\bar{x})$ ,  $t^0$  and  $t^1$ . To simplify, we introduce the quantity  $\mathcal{T}_i := \mathbf{n}_i^T \mathbf{K}_i \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}$  to represent the cell *transmissivity*. So first we solve for  $u(\bar{x})$ :

$$\mathcal{T}_0(u(\bar{x}) - u(x_0)) = \mathcal{T}_1(u(x_1) - u(\bar{x})) \Rightarrow u(\bar{x}) = \frac{\mathcal{T}_0 u(x_0) + \mathcal{T}_1 u(x_1)}{\mathcal{T}_0 + \mathcal{T}_1}. \quad (3.34)$$

Next, we insert this into the expression for  $\tilde{q}_0$ :

$$\begin{aligned} \tilde{q}_{E_{0,1},0} &= -\mathcal{T}_0(u(\bar{x}) - u(x_0)) \\ &= -\mathcal{T}_0 \left( \frac{\mathcal{T}_0 u(x_0) + \mathcal{T}_1 u(x_1)}{\mathcal{T}_0 + \mathcal{T}_1} - u(x_0) \right) \\ &= -\mathcal{T}_0 \left( \frac{\mathcal{T}_0 u(x_0) + \mathcal{T}_1 u(x_1) - u(x_0)\mathcal{T}_0 - u(x_0)\mathcal{T}_1}{\mathcal{T}_0 + \mathcal{T}_1} \right) \\ &= -\mathcal{T}_0 \left( \frac{\mathcal{T}_1 u(x_1) - u(x_0)\mathcal{T}_1}{\mathcal{T}_0 + \mathcal{T}_1} \right) \\ &= \frac{u(x_0) - u(x_1)}{\frac{1}{\mathcal{T}_1} + \frac{1}{\mathcal{T}_0}}. \end{aligned} \quad (3.35)$$

Now, we have solved the equations for the transmissivity coefficients:

$$\begin{aligned} \tilde{q}_{E_{0,1},0} &= t^0 u(x_0) + t^1 u(x_1) \\ \frac{u(x_0) - u(x_1)}{\frac{1}{\mathcal{T}_1} + \frac{1}{\mathcal{T}_0}} &= t^0 u(x_0) + t^1 u(x_1) \\ \Rightarrow t^0 &= \frac{1}{\frac{1}{\mathcal{T}_1} + \frac{1}{\mathcal{T}_0}}. \end{aligned} \quad (3.36)$$

Hence, the transmissibility is the *harmonic mean* of the transmissivities. This kind of mean appears naturally when one wants to find the permeability of flow through layers of different permeability.

One way of looking at this discretization, is that we assume the potential to be a linear function of one variable, with its gradient pointing in the  $\mathbf{v}_i$  direction between the cell center and the edge in figure 3.5. So for each edge, we have two linear functions on each side, which gives us four degrees of freedom. Two of

them are used to respect the cell center potential values, the other two are used on potential and flux continuity across the edge. With these assumptions, expressions (3.31) and (3.32) are exact. And we only have to solve for the transmissibility coefficients.

Two point flux approximation has the advantage of being fast to assemble and simple to code. It yields a pleasant five point stencil for two dimensional problems. However, there is one big disadvantage with two point flux approximation: Computing the flux with only two points is not consistent when the grid is not aligned with the principal directions of  $\mathbf{K}$ . If our grid is aligned with  $\mathbf{K}$ , we have that

$$\mathbf{n}_2 \cdot \mathbf{K} \mathbf{n}_1 = 0 \quad (3.37)$$

for a uniform parallelogram mesh with the normal vectors  $\mathbf{n}_1$  and  $\mathbf{n}_2$ . We then call the grid **K-orthogonal**. In the setting of figure 3.5, our grid would not be K-orthogonal as the control volumes are not parallelograms. All meshes with orthogonal control volumes are K-orthogonal if the permeability is isotropic. In figure 6.6 we observe the failed convergence of the TPFA-method for a parallelogram mesh.

### 3.2.2 MPFA-O Method

The O-method is a multi-point flux approximation method, these types of methods were developed to make control volume methods converge for grids that are not K-orthogonal. It is described in detail in [1], we only give a brief introduction. Consider the control volumes in 3.6.

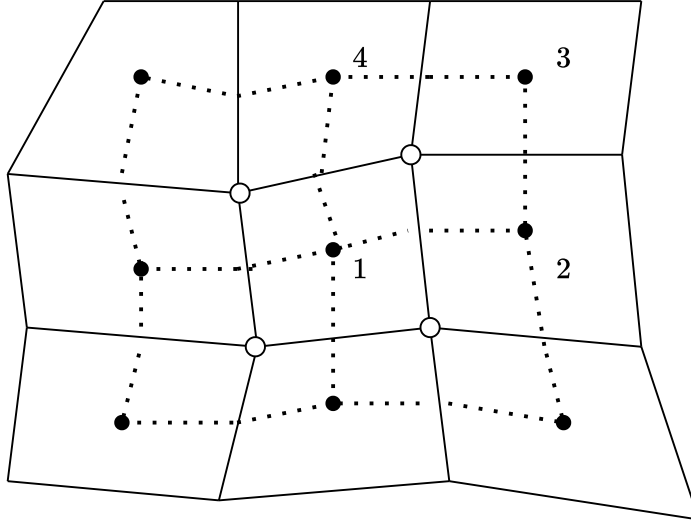


Figure 3.6: The solid lines are the control volumes, the dashed lines are the dual mesh connecting the cell centers, going through the midpoints of each edge. The solid circles are cell centers, the white circles are grid points.

For each grid point, that means where four control volumes intersect, we consider an interaction region. This is the polygon drawn by the dual mesh around the grid point.

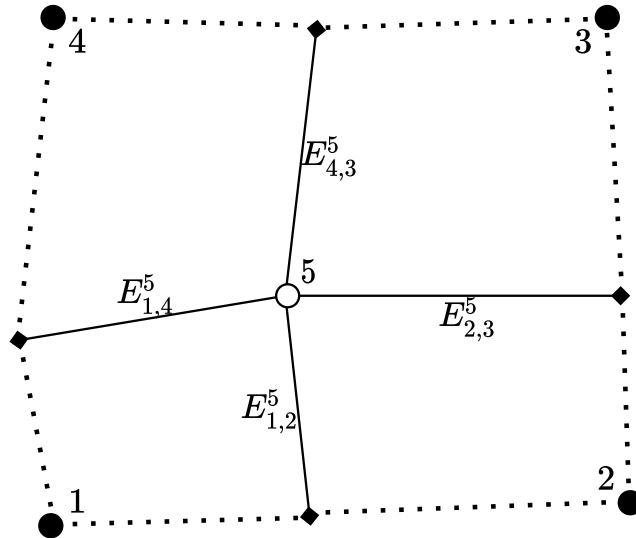


Figure 3.7: The four subcells in the interaction region corresponding to cells 1, 2, 3, 4 and grid point 5. Here,  $\mathcal{R}_5 = \{1, 2, 3, 4\}$ .

In each interaction region there are four half edges. Our goal is to obtain an

expression

$$\tilde{q}_{E_{i,j}^n} = \sum_{k \in \mathcal{R}_n} t_{i,j}^{k,n} u^k \approx \int_{E_{i,j}} \hat{\mathbf{n}}_j^T \mathbf{K} \nabla u \, ds \quad i, j \in \mathcal{R}_n \quad (3.38)$$

for the flux through each half edge  $E_{i,j}^n$  in the interaction region corresponding to grid point  $n$  (figure 3.7). Where  $\mathcal{R}_n$  is the index set of the four cells neighbouring grid point  $n$ .

We assume for now that the potential is linear in each of the four sub cells in the interaction region, figure 3.7. This gives  $4 \cdot 3 = 12$  degrees of freedom. The linear potential must of course equal the cell center values of the potential in the cell centres, this removes four degrees of freedom. We also require flux continuity on the four half edges in the interaction region, this removes an additional four degrees of freedom. The last four degrees of freedom are spent on potential continuity of the midpoints of the edges.

By these assumptions on flux and potential continuity, the linear potential in each sub cell is well defined given values at the cell center. We can now use this to compute the four by four matrix of transmissibility coefficients for each of the four half edges. In the situation of figure 3.7 and equation (3.38) it would look like

$$\mathbf{T}^5 = \begin{bmatrix} t_{1,2}^{1,5} & t_{1,2}^{2,5} & t_{1,2}^{3,5} & t_{1,2}^{4,5} \\ t_{1,5}^{1,5} & t_{1,5}^{2,5} & t_{1,5}^{3,5} & t_{1,5}^{4,5} \\ t_{2,3}^{1,5} & t_{2,3}^{2,5} & t_{2,3}^{3,5} & t_{2,3}^{4,5} \\ t_{4,3}^{1,5} & t_{4,3}^{2,5} & t_{4,3}^{3,5} & t_{4,3}^{4,5} \\ t_{1,4}^{1,5} & t_{1,4}^{2,5} & t_{1,4}^{3,5} & t_{1,4}^{4,5} \end{bmatrix}, \quad (3.39)$$

where each row corresponds to the flux over some half edge, and each column corresponds to some cell center. Computing (3.39) involves inverting a four by four matrix with coefficients depending on the mesh and permeability, see [1] for details. Finally, we assemble the system of equations (3.29) with the transmissibility coefficients. Note that we write the flux over the  $j$ th edge of cell  $i$ ,  $\tilde{q}_{i,j}$  as the flux over the two half edges:

$$\sum_{j \in \mathcal{S}_i} (\tilde{q}_{\hat{E}_{i,j}^1} + \tilde{q}_{\hat{E}_{i,j}^2}) = |\Omega_i| f(x_i) \quad (3.40)$$

Where  $\hat{E}_{i,j}^1 = E_{i,j}^n$  for some  $n$ , see (3.38). Hence, computing the transmissibility coefficients and assembling them into the discretization matrix, requires two different indexing systems.

Next, we see that the interaction regions of the two half edges sharing same edge overlaps, so we get a six point flux stencil. In other words, for each  $j$  in (3.40), the union of the two interaction regions used to compute  $\tilde{q}_{E_{i,j}^1}$  and  $\tilde{q}_{E_{i,j}^2}$  consists of six points. Taking the union of the four flux stencils connected to a cell, we



observe that the O-method yields a nine point stencil

$$\sum_{k \in \mathcal{M}_i} \hat{t}^k u^k = |\Omega_i| f(x_i),$$

where  $\mathcal{M}_i$  is the set of nine indices corresponding to cell  $i$  and its eight neighbour cells.

The O-method is consistent for non K-orthogonal grids, and reduces to two point flux approximation when the grid is K-orthogonal. This happens because the systems of equations to be solved for the transmissibility coefficients in each interaction region, becomes diagonal. This is because  $\mathbf{n}^T \mathbf{K} \nabla u$  can be expressed as two points when  $u$  is a linear function given by three points which forms two K-orthogonal vectors.

In [12], Nordbotten and Keilegavlen describes a framework of MPFA methods where the O-method is a special case. They consider the problem of finding the four linear potential functions in each interaction region that minimizes the discontinuity across the edges. The discontinuity should be minimized given that the functions respect cell center potential values, that the flux models the constitute law and flux continuity. The O-method is then defined for some special cost function measuring the discontinuity. Other methods, with the potential continuity at other places than the edge midpoint, are also common.

With our implementation of MPFA-O method, one needs for each interaction region to assemble four, four by four, matrices. Compute the inverse of one of them, and do two matrix multiplications and one subtraction. All of this could be done in parallel. However, for my implementation, it slows matrix assembly down a lot compared to two point flux approximation. Another drawback of the O-method is the *monotonicity* properties: One can risk having positive entries off the main diagonal of the discretization matrix for difficult meshes. This may lead to oscillations in the solution and violation of the minimum principle. For two point flux approximation we avoid this issue altogether, as the signs of the five point stencil always are one plus and four negatives. Even for the linear finite element method, this issue is avoided if one imposes some maximum angle condition, see [Knabner,[8]] page 175. When solving for example the Richards' equation, violating the minimum principle can lead to air bubbles being formed spontaneously in the saturated region. For a discussion on monotonicity see [11].

### 3.2.3 MPFA L-Method

The L-method is the Ferrari of discretization techniques for porous media flow problems, while conformal finite elements is the Volvo.

---

*Professor Jan Martin Nordbotten*

Like the O-method, the L-method is also a multi-point flux approximation method. It was introduced in [2], where the authors demonstrate improved monotonicity properties with numerical experiments. This method is similar to the O-method, in that it goes through the half edges and uses information from the same interaction regions. But instead of using four points for the flux across each half edge, we use three, with two half edges between them.

As in the O-method, we assume linear potential in each cell, this gives us  $3 \cdot 3 = 9$  degrees of freedom. Three are eliminated because we respect the cell center value of the potential, this leaves six degrees of freedom. We use two, one at each edge, for flux continuity. The last four are used for potential continuity at the two edges.

We have two choices of flux stencil for each half edge, see figure 3.8. We compute the transmissibility coefficients for both, then we choose the one "best" aligned with the flow: Let  $t_1^i$  be the  $i$ th transmissibility coefficient of  $T_1$ , then

$$\begin{aligned}
 &\text{if } |t_1^1| < |t_2^2| \\
 &\text{choose } T_1 \text{ else} \\
 &\text{choose } T_2.
 \end{aligned}
 \tag{3.41}$$

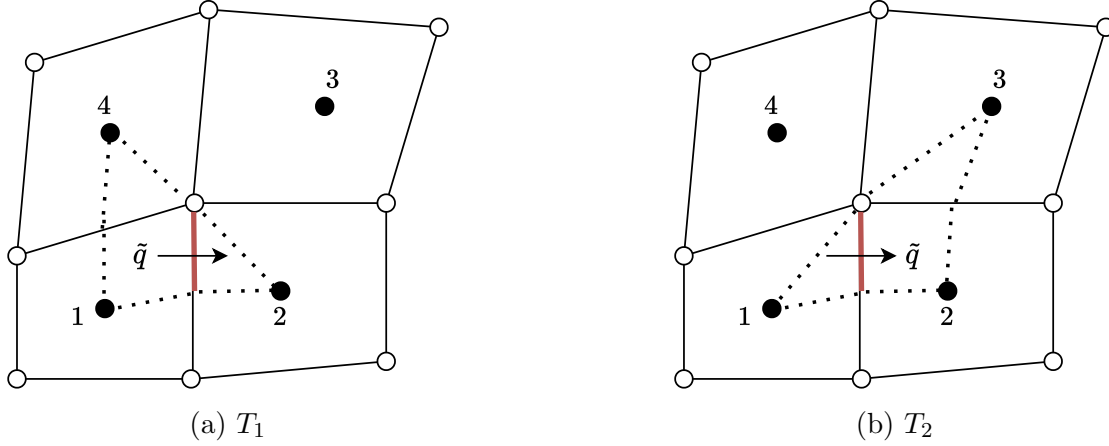


Figure 3.8: The two choices of which cell centers to use for computing the flux over the half edge in red. We call the hexagons spanned  $T_1$  and  $T_2$  for L-triangles, as they consists of three cell centers.

A cheap intuition behind (3.41) is that if  $|t_1^1| < |t_2^2|$ , it is more likely that  $\text{sgn}(t_1^1) = \text{sgn}(t_1^4)$  and if not,  $\text{sgn}(t_2^2) = \text{sgn}(t_2^3)$  is more likely. This is due to the fact that  $\sum t^i = 0$ . Choosing L-triangle as in (3.41) increases the chances that we get the same sign of  $t^i$  on the same side of the half edge, thus increasing the chance that we get a monotone discretization. See [4] for a more detailed geometric intuition of choosing L-triangle in the case of homogenous permeability.

To compute transmissibility coefficients in a given L-triangle, we use the assumptions on flux and potential continuity, to construct a linear system. The coefficients depends on mesh and permeability in the three cells. As with the O-method, we end up with a system assembled from the fluxes over the half edges:

$$\sum_{j=1}^4 (\tilde{q}_{i,j}^1 + \tilde{q}_{i,j}^2) = |\Omega_i| f(x_i) \quad (3.42)$$

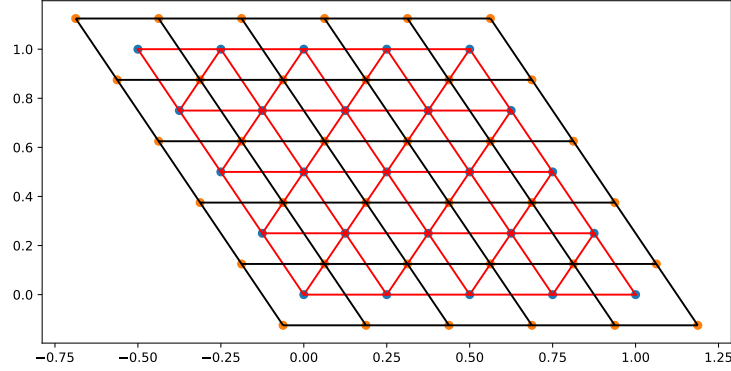
$$\sum_{j=1}^4 \left( \sum_{k=1}^3 t_{i,j}^{k,1} u^k + \sum_{k=1}^3 t_{i,j}^{k,2} u^k \right) = |\Omega_i| f(x_i).$$

But the flux stencil across each edge is possibly smaller, often just four points.

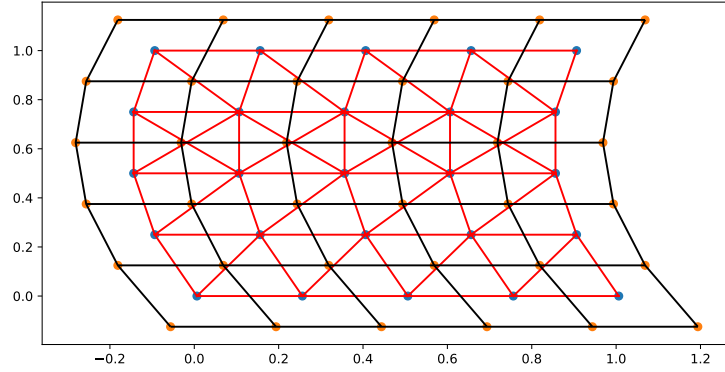
**Remark 12.** *In the L-method, we need to construct and solve a matrix equation twice for each half edge to compute the transmissibility coefficients, as there are always two choices. In contrast, the O-method only needs this done once for each grid point, and its four half edges.*

In figure 3.9 we see the criterion in practice for a homogenous medium: In figure 3.9a all L-triangles are used by two half-edges, and they are chosen in the

same way throughout the domain. In figure 3.9b there are some triangles that overlap, this is due to the fact that some L-triangles are used by only one half edge.



(a) Parallelogram grid, all triangles are chose similarly.



(b) Complicated grid, note that some of the L-triangles overlap.

Figure 3.9: Examples of L-triangles(in red) in a domain with homogenous permeability tensor.

The observation in figure 3.9a can be stated as a theorem:

**Theorem 3.2.1** (Cao, Y., Helmig, R. and Wohlmuth, B.I. (2009),[5]). *For homogeneous media and uniform parallelogram grids, the MPFA L-method has a seven-point cell stencil for the discretization of each interior cell, i.e., the discretization of each cell is a seven point stencil including the center cell and the six closest potential cells, as shown in 3.9a.*

In case of parallelogram grid with heterogeneous permeability, it may also happen that one gets overlapping L-triangles. This is the case even if the permeability only changes as a scalar in the domain. In figure 3.10 the L-triangles are shown for a random, scalar permeability. Let  $K_{m,n}$  be the permeability of the  $m$ th cell in  $y$  direction and  $n$ th cell in  $x$  direction. Then the random permeability used in figure 3.10 given by

$$K_{n,m} = (e^{\hat{x}} - 1)^2 \quad (3.43)$$

where  $\hat{x}$  is a random sample drawn from a uniform distribution over  $[0, 1)$ . We see that two of the L-triangles overlap. This is due to some combination of permeability at four neighbouring cells. Also note that the permeability is not so low that it causes numerical rounding errors, as  $\min_{m,n} K_{m,n} = 0.0017$  in figure 3.10.

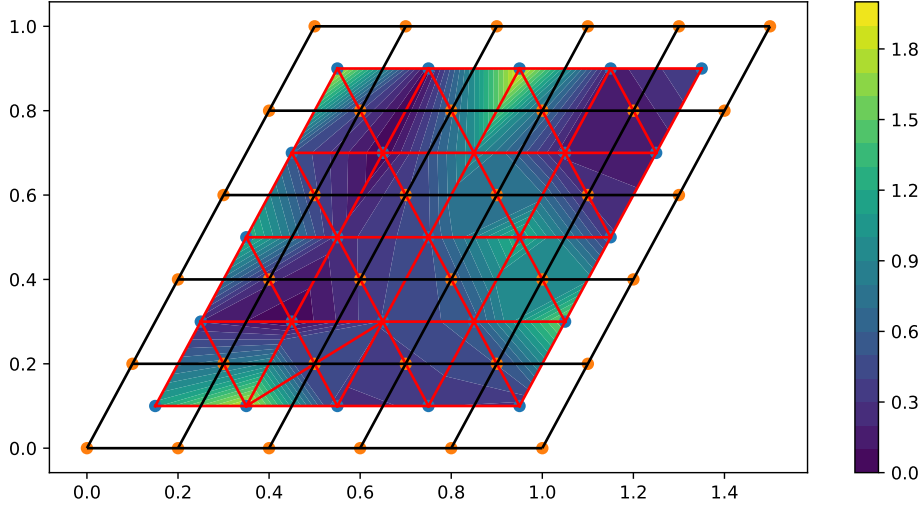


Figure 3.10: L-triangles on a random permeability.

For homogenous media the L-method becomes easier to simpler. We continue with a useful theorem which we will use later:

**Lemma 3.2.2** (Cao, Y., Helmig, R. and Wohlmuth, B.I. (2009),[4]). *Assume that the permeability  $\mathbf{K}$  is homogeneous on  $\Omega$ , then the flux through each half edge  $e$ , computed by the L-method, can be written as*

$$\tilde{q}_e = -\mathbf{K} \nabla u \cdot \mathbf{n}_e \quad (3.44)$$

Where  $\mathbf{n}_e$  is the scaled normal vector to the half edge  $e$ , having the same length as  $e$ .  $u$  is a linear scalar field uniquely given by the potential values at the three cell centers chosen by the L-method.

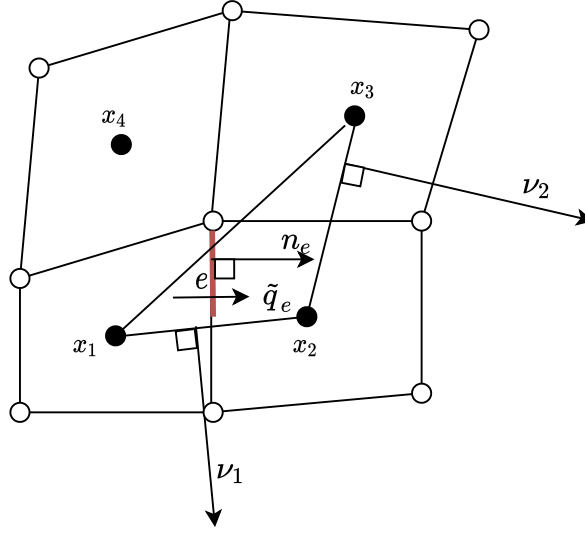


Figure 3.11: Simplified L-triangle, the original L-triangle is shown in figure 3.8b or 3.12. The vector  $\nu_1$  is perpendicular to the edge between  $x_1$  and  $x_2$ , with the same length as the edge it is perpendicular to. Same for  $\nu_2$ , with  $x_2$  and  $x_3$ .

Moreover, the gradient  $\nabla u$ , is given by:

$$\nabla u = -\frac{1}{2F}[(u_1 - u_2)\nu_2 + (u_3 - u_2)\nu_1], \quad (3.45)$$

where  $F$  is the area of the simplified L-triangle with corners  $x_1$ ,  $x_2$  and  $x_4$ , see figure 3.11. An expression like (3.45) can be obtained for the other choice of L-triangle as well.

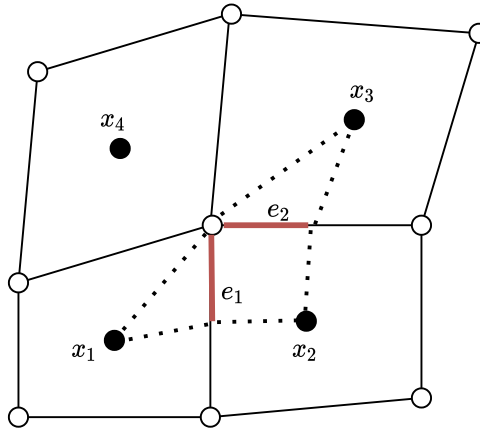


Figure 3.12: Original L-triangle with notations in proof.

*Proof.* It is enough to check that the jump  $[\nabla u]$  is zero on  $e_1$  and  $e_2$  on the original L-triangle in figure 3.12. Let  $\mathbf{t}_{e_1}$  and  $\mathbf{n}_{e_1}$  be the tangent and normal vector to  $e_1$ . Since we require potential continuity on each half edge, we get:

$$[\nabla u \cdot \mathbf{t}_{e_1}] = 0. \quad (3.46)$$

Using the fact that  $\mathbf{K}$  is symmetric and homogenous, we obtain:

$$[\mathbf{K} \nabla u \cdot \mathbf{n}_{e_1}] = [\nabla u \cdot \mathbf{K}^T \mathbf{n}_{e_1}] = [\nabla u \cdot \mathbf{K} \mathbf{n}_{e_1}] = 0. \quad (3.47)$$

Where we used flux continuity across each half edge in the last equality. Since  $\mathbf{K}$  is positive definite, we have that  $\mathbf{K} \mathbf{n}_{e_1}$  and  $\mathbf{t}_{e_1}$  are independent, thus  $[\nabla u] = 0$  on  $e_1$ . Same arguments holds for  $e_2$ . Hence  $\nabla u$  is constant on the original L-triangle and the desired result follows.  $\square$

**Remark 13.** *The above lemma suggests that we can obtain the transmissibility coefficients without solving a system of equations for each half edge. This simplifies implementation, but it's only possible for homogenous media.*

To conclude; the L-method is the most sophisticated method. It has the best monotonicity properties, it is consistent for non K orthogonal grids, but it is more complicated than the O-method.

## Time Discretization

We start by considering the most famous parabolic equation, namely the heat equation. Let  $u = u(x, t)$ , given appropriate boundary and initial conditions, find  $u$  such that:

$$\begin{cases} \partial_t u - \nabla \cdot \mathbf{K} \nabla u = f, & \text{in } \Omega \times (0, T] \\ u = 0, & \text{on } \partial\Gamma_D \times (0, T] \\ -\mathbf{K} \nabla u = g_N, & \text{on } \partial\Gamma_N \times (0, T] \\ u = u_0, & \text{on } \Omega \times \{t = 0\} \end{cases} \quad (3.48)$$

The well-posedness of (3.48) is discussed in chapter seven of [7], it requires a more detailed discussion of Sobolev spaces and Bochner spaces, i.e., spaces containing functions from the real numbers to some Sobolev space.

We expect low regularity in time, so there is not much to be gained by using a higher order discretization in time. The two choices we have left is the forward Euler (explicit) and the backward Euler (implicit). The obvious choice is backward Euler, as it is stable for large time-step sizes. This can be understood intuitively by considering the parabolic nature of the equation, the signals propagate through

the domain instantaneously. A careful analysis of time discretizations of parabolic equations is done in ([8], chapter 7). There, it is shown that explicit schemes only are stable for time-step sizes proportional to the square of the diameter of the space discretization, whereas fully implicit schemes are stable for all time-step sizes.

Let  $\{t_n\}_n$  be a sequence of  $N + 1$  uniformly distributed numbers from 0 to  $T$  and let  $\tau = \frac{T}{N}$  be the time-step size. Then we state the semi-discrete version of (3.48) by exchanging the time derivative by a difference quotient  $\partial_t u = \frac{u^n - u^{n-1}}{\tau}$ . We end up with: Given  $u^{n-1}$  and  $f^n$ , find  $u^n$  such that

$$\begin{aligned} u^n - \tau \nabla \cdot \mathbf{K} \nabla u^n &= \tau f^n + u^{n-1} & x \in \Omega \\ u^n &= 0 & x \in \partial\Gamma_D \\ \mathbf{K} \nabla u &= g_N & x \in \partial\Gamma_N \\ u^0 &= u_0 & x \in \Omega. \end{aligned} \quad (3.49)$$

The above equation shows that this time discretization is implicit, i.e., we cannot solve (3.49) for  $u^n$  with simple algebraic manipulation. Now, we have an elliptic problem (3.49) for each time-step. This has almost the same structure as the elliptic model problem (3.1) we solved in the previous chapters, the difference being the  $u^n$  term.

## Finite element approach

We are now ready to fit this problem into our finite element framework from chapter 2. The variational formulation of (3.49) is achieved as before by multiplying by test functions in  $H_0^1(\Omega)$ : Given  $u^{n-1} \in V$ ,  $f^n \in V'$ , find  $u^n \in V$  such that

$$\langle u^n, v \rangle_0 + \tau \langle \mathbf{K} \nabla u^n, \nabla v \rangle_0 = \tau \langle f^n, v \rangle_0 + \langle u^{n-1}, v \rangle_0 \quad (3.50)$$

for all  $v$  in  $V$ . If we exchange  $V$  with a finite dimensional subspace  $V_h$ , and write  $u_h^n = \sum_{i=1}^d \hat{u}_i^n \phi_i$ , as in the Galerkin FEM section 3.1.4, we end up with the system: Find  $\hat{\mathbf{u}}^n \in \mathbb{R}^d$  such that

$$(\mathbf{B} + \tau \mathbf{A}) \hat{\mathbf{u}}^n = \tau \mathbf{f}^n + \mathbf{B} \hat{\mathbf{u}}^{n-1}, \quad (3.51)$$

where the *stiffness matrix*,  $\mathbf{A}$ , is as before, that is  $\mathbf{A}_{i,j} = \langle \mathbf{K} \nabla \phi_i, \nabla \phi_j \rangle$ . The matrix  $\mathbf{B}$  is often called the *mass matrix* and is defined as  $\mathbf{B}_{i,j} = \int_{\Omega} \phi_i \phi_j dx$ .

## Finite volume approach

As before, we divide our domain  $\Omega$  into  $d$  control volumes  $\{\Omega_i\}_i$ . Either, one can write the heat equation (3.48) in conservation form on each control volume

$$\partial_t \int_{\Omega_i} u \, dx - \int_{\partial\Omega_i} \mathbf{K} \nabla u \cdot \hat{\mathbf{n}} \, dx = \int_{\Omega_i} f \, dx, \quad (3.52)$$



and discretize the first term with backward Euler, or one can make sure the semi-discrete heat equation (3.48) holds for each control volume and use the divergence theorem. Both ways, we end up with

$$\int_{\Omega_i} u^n dx - \tau \int_{\partial\Omega_i} \mathbf{K} \nabla u^n \cdot \hat{\mathbf{n}} dx = \tau \int_{\Omega_i} f^n dx + \int_{\Omega_i} u^{n-1} dx, \quad (3.53)$$

if we, as discussed earlier, use the midpoint rule to evaluate the integrals, we get

$$\int_{\Omega_i} u^n(x_i) dx - \tau \int_{\partial\Omega_i} \mathbf{K} \nabla u^n \cdot \hat{\mathbf{n}} dx = \tau \int_{\Omega_i} f^n(x_i) dx + \int_{\Omega_i} u^{n-1}(x_i) dx. \quad (3.54)$$

As in the previous section we end up with a system of equations, where superscript  $V$  is just to distinct between FVM and FEM. Find  $\tilde{\mathbf{u}} \in \mathbb{R}^d$ , such that

$$(\mathbf{B}^V + \tau \mathbf{A}^V) \tilde{\mathbf{u}}^n = \tau \mathbf{f}^n + \mathbf{B}^V \tilde{\mathbf{u}}^{n-1} \quad (3.55)$$

The matrix  $\mathbf{A}^V$  is as in chapter 3, with the fluxes through the edges of cell  $i$  described by the  $j$ th row of  $\mathbf{A}^V$ . The matrix  $\mathbf{B}^V$  is diagonal with the entry  $i$  being the volumes of the volume of cell  $i$ . That is, for two dimensional problems, the entries of  $\mathbf{B}^V$  are the areas of the control volumes.

If  $\mathbf{A} = \mathbf{A}^V$ , then the discretization of the constitutive law is the same for both the finite volume and the finite element method. As we will see later, this is challenging.

### 3.3 Linearization

We have seen that the heat equation leads to a sequence of linear systems. In the same way, we expect that our non-linear Richards' equation (2.12) leads to a system of non-linear equations. We start by discussing this in a general setting: Find  $x \in U$  such that

$$\mathbf{f}(\mathbf{x}) = \mathbf{0} \text{ where } \mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^n. \quad (3.56)$$

The solution of (3.56) is called a *root*, it is almost always found using an iterative method.

A common iterative scheme to solve (3.56) is the *Newton's method*. Let  $D\mathbf{f}(\mathbf{x}_{j-1})^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be the Jacobian of  $\mathbf{f}(\mathbf{x}_{j-1})$ , then the newton iteration is given by:

$$\mathbf{x}_j = \mathbf{x}_{j-1} - D\mathbf{f}(\mathbf{x}_{j-1})^{-1} \mathbf{f}(\mathbf{x}_{j-1}). \quad (3.57)$$

In one dimension a convergence proof is easily obtained by techniques from calculus, the following theorem is found in slightly more detail in (Cheney[6], chapter 3):

**Theorem 3.3.1.** *Let  $f'' < 2$  with  $f(\bar{x}) = 0$  and  $f'(x) > \delta \forall x \in B_\epsilon(\bar{x})$ , then the Newton method is locally quadratic convergent: For  $x_0 \in B_\epsilon(\bar{x})$  we have*

$$|x_{j+1} - \bar{x}| \leq \frac{1}{\delta} |x_j - \bar{x}|^2 < |x_j - \bar{x}|. \quad (3.58)$$

*Proof.* Define  $e_j := x_j - \bar{x}$ . Then we have by Taylor expansion

$$0 = f(\bar{x}) = f(x_j - e_j) = f(x_j) - f'(x_j)e_j + \frac{f''(\psi)e_j^2}{2}. \quad (3.59)$$

For some  $\psi$  between  $x_j$  and  $\bar{x}$ . Further, we get by the definition of the newton method:

$$\begin{aligned} e_{j+1} = x_{j+1} - \bar{x} &= x_j - \frac{f(x_j)}{f'(x_j)} - \bar{x} \\ &= e_j - \frac{f(x_j)}{f'(x_j)} \\ &= \frac{e_j f'(x_j) - f(x_j)}{f'(x_j)} \end{aligned} \quad (3.60)$$

By the Taylor expansion around  $x_j$ , (3.59), we get

$$f'(x_j) = \frac{f(x_j)}{e_j} + \frac{f''(\psi)e_j}{2}. \quad (3.61)$$

Inserting this into (3.60), we get the equality

$$e_{j+1} = \frac{e_j^2 f''(\psi)}{2f'(x_j)}. \quad (3.62)$$

The assumptions on  $f'$  and  $f''$  combined with  $|e_0| < \delta$  give us the estimate

$$|e_1| \leq \frac{2}{2\delta} |e_0|^2 < |e_0| \quad (3.63)$$

The above equation implies  $x_1 \in B_\epsilon(\bar{x})$ , and by induction we get:

$$|e_{j+1}| < |e_j|, \quad (3.64)$$

and the quadratic convergence

$$|e_{j+1}| \leq \frac{1}{\delta} |e_j|^2 \quad (3.65)$$

□

For a similar result in more dimensions see (Knabner [8], chapter 8). One apparent drawback of this method is that it is only locally convergent, i.e., one needs to start the iteration in a neighbourhood of the root where the Jacobian is well defined. In practice one often solves the system

$$D\mathbf{f}(\mathbf{x}_{j-1})\boldsymbol{\delta}_j = -\mathbf{f}(\mathbf{x}_{j-1}), \quad (3.66)$$

and then update the current iterate:  $\mathbf{x}_j = \mathbf{x}_{j-1} + \boldsymbol{\delta}_j$ . Typically, the matrix  $D\mathbf{f}(\mathbf{x}_{j-1})$ , needs to be computed and assembled for every iteration. This may be computationally expensive. So Newton's method may be slow despite its quadratic convergence, if it even converges.

A simpler approach is to exchange the Jacobian with a diagonal matrix  $L\mathbf{I}$  such that

$$L\boldsymbol{\delta}_j = -\mathbf{f}(\mathbf{x}_{j-1}). \quad (3.67)$$

This is called the *L-scheme*, and will be the method we use for linearization in this thesis. In one dimension it is easy to prove convergence:

**Theorem 3.3.2.** *Let  $f \in C(\mathbb{R})$  and  $L > \sup_{x \in \mathbb{R}} f'(x)$ , then the L-scheme converges linearly for all  $x_0 \in \mathbb{R}$ .*

*Proof.* Define  $e_j := x_j - \bar{x}$ , then we get

$$e_{j+1} = x_j - \frac{f(x_j)}{L} - \bar{x} = e_j - \frac{f(x_j)}{L}. \quad (3.68)$$

We use the same trick as before with the Taylor expansion around the root:

$$0 = f(\bar{x}) = f(x_j - e_j) = f(e_j) - f'(\psi)e_j \Rightarrow e_j = \frac{f(x_j)}{f'(\psi)}. \quad (3.69)$$

Using this and the assumption on  $L$ , we get the estimate:

$$|e_{j+1}| = |e_j(1 - \frac{f'(\psi)f(x_j)}{f(x_j)L})| \leq |e_j||1 - \frac{f'(\psi)}{L}| < |e_j|. \quad (3.70)$$

□

In practice we need to stop the linearization scheme at some point, and we decide on a **stopping criterion**. A common choice is, and the one we will use, is

$$|x_j - x_{j-1}| < TOL + TOL|x_{j-1}|, \quad (3.71)$$

Where  $TOL$  is some constant chosen to be close to machine epsilon. See (Storvik, [15]) for a discussion of the L-scheme and how to choose the  $L$  and  $TOL$  parameters in a smart way. One can also study other linearization approaches with different properties. In (List and Radu, [9]) the authors compare different iterative linearization methods for the Richards equation and propose a method that combines the Newton method and the L-scheme with desirable convergence rate and robustness.



# Chapter 4

## Convergence of the MPFA-L Method

In this chapter we show equivalence between a modified MPFA-L method and a modified Lagrange finite element method, for linear time dependent problems discretized in time with backward Euler (3.49). That is, we prove equivalence between the two discretizations of the equation: Let  $x \in \Omega \subset \mathbb{R}^2$ , find  $u(x)$  such that

$$\begin{cases} u - \nabla \cdot \mathbf{K} \nabla u = f, & \text{in } \Omega \\ u = 0, & \text{on } \partial\Gamma_D \\ -\mathbf{K} \nabla u = g_N, & \text{on } \partial\Gamma_N, \end{cases} \quad (4.1)$$

where  $\mathbf{K}$  is homogeneous, in addition to being symmetric positive definite. Once equivalence is obtained, we prove convergence for the finite element method using techniques from section 3.1.6.

After reading this chapter, the reader should be convinced that the finite element method covered in section 3.1 is almost the same as the L-method. Moreover, that the L-method can be used as a locally mass conservative flux recovery algorithm on the modified finite element solution. See section 6.2 for a comparison of the MPFA-L method and normal linear Lagrange finite element method.

We saw in the section about the MPFA-L method that the interaction regions (L-triangles) may form a triangulation of our domain. With this observation in mind, modifications are made to both methods so that we obtain equivalence. This entire chapter is adapted from (Cao, Y., Helmig, R. and Wohlmuth, B.I. (2009),[5]), where, convergence is proved for the Poisson equation,i.e., without the  $u$  term.

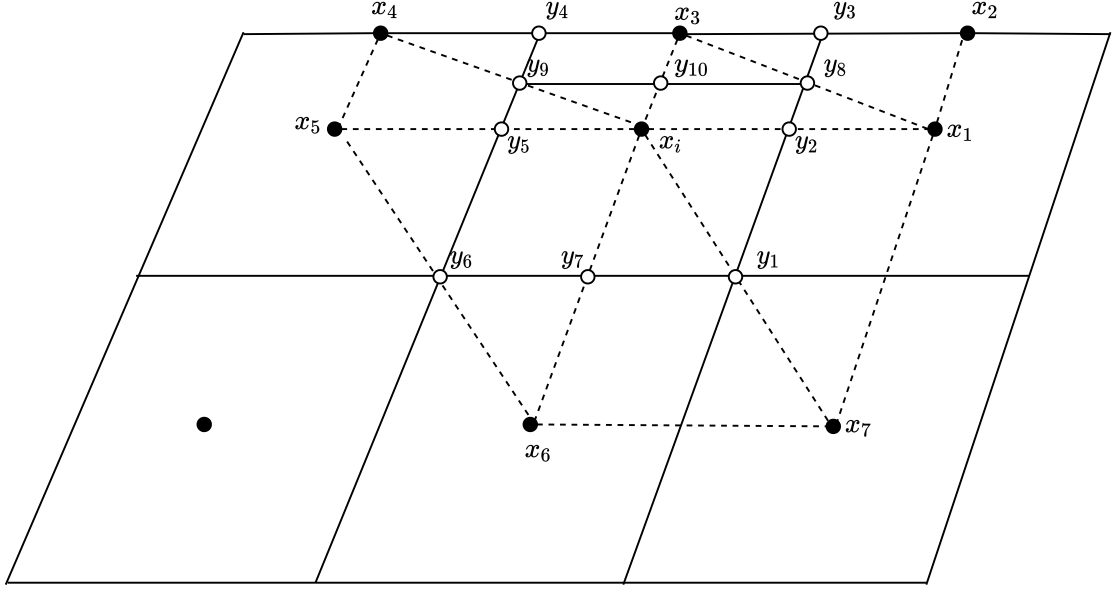


Figure 4.1: Control volumes in solid lines and interaction regions in dashed lines at the boundary.

## 4.1 Modified MPFA-L Method

First of all, we assume that we have a uniform parallelogram grid, as in 3.9a. As we saw in the previous chapter, one gets with the finite volume method the following relation for all control volumes  $\Omega_i$ :

$$\int_{\Omega_i} u \, dx - \int_{\partial\Omega_i} \mathbf{K} \nabla u \cdot \hat{\mathbf{n}} \, dx = \int_{\Omega_i} f \, dx. \quad (4.2)$$

The MPFA-L method deals with the second term, approximating the constitutive law. The other two terms are common to all control volume methods solving time dependent problems or (4.1).

On the interior control volumes, we use the original MPFA-L method already covered. On the Neumann boundaries we need a modification. This is to be expected, as control volume methods handle flux at the boundary in a very simple way; specifying it the same way we deal with with the source term, adding it to the load vector. In finite element methods however, we have degrees of freedom on the boundary, one dimensional elements. We will also make a special treatment of the Dirichlet boundary, in a way that is equivalent to the finite element method. In [5] they claim that this is a very natural way of dealing with the Dirichlet boundary conditions, and a good practical alternative to other ways of enforcing Dirichlet boundaries in the MPFA-L method.

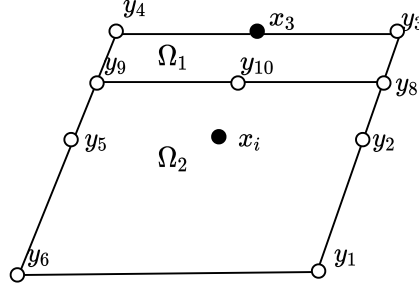


Figure 4.2: Control volume along top boundary.

Consider the control volume  $y_1y_6y_4y_3$ . For the **Neumann** boundary conditions, we split the control volume into two,  $y_1y_6y_9y_8$  as  $\Omega_2$  and  $y_8y_9y_4y_3$  as  $\Omega_1$ , see figure 4.2 or 4.1. We therefore get one equation each for  $u_3$  and  $u_i$  as the potential at  $x_3$  and  $x_i$ . For the fluxes on  $\Omega_2$  we have six interaction triangles and a normal seven point stencil. For the  $\Omega_1$  we compute the flux through  $\overline{y_3y_8}$  using  $\triangle x_1x_3x_2$ , the flux through  $\overline{y_8y_{10}}$  using  $\triangle x_1x_ix_3$ , for  $\overline{y_{10}y_9}$  and  $\overline{y_9y_4}$  the L triangle  $\triangle x_ix_4x_3$  is used. Finally the Neumann boundary condition is used at the the edge  $\overline{y_4x_3}$  and  $\overline{x_3y_3}$ . We are not able to eliminate the unknown value at  $x_3$  and it remains a degree of freedom, which makes sense if we want equivalence with finite element method.

In the case of **Dirichlet** boundary conditions, we compute the fluxes into  $y_1y_6y_4y_3$  using seven L-triangles, as can be seen in figure 4.1. The flux over the edge  $\overline{y_3y_1}$  are computed as the sum of the flux over  $\overline{y_3y_8}$ ,  $\overline{y_8y_2}$  and  $\overline{y_2y_1}$  using the L-triangles  $\triangle x_1x_3x_2$ ,  $\triangle x_1x_ix_3$  and  $\triangle x_1x_7x_i$  respectively. Similarly for the edge  $\overline{y_6y_4}$ . For  $\overline{y_1y_6}$  we only use the two big L-triangles at the bottom,  $\triangle x_ix_7x_6$  and  $\triangle x_ix_6x_5$ .

The flux over  $\overline{y_4y_3}$ , at the boundary, we compute by balancing with the other fluxes out of the small control volume  $\Omega_1$ , see figure 4.3. Let  $\tilde{q}_{\overline{y_iy_j}}$  be the flux through edge  $\overline{y_iy_j}$ , out of the volume  $\Omega_1$ . Then we get the expression for the flux through the Dirichlet boundary:

$$\tilde{q}_{\overline{y_3y_4}} = -(\tilde{q}_{\overline{y_3y_8}} + \tilde{q}_{\overline{y_{10}y_8}} + \tilde{q}_{\overline{y_9y_{10}}} + \tilde{q}_{\overline{y_4y_9}}) + \int_{\Omega_1} f \, dx. \quad (4.3)$$

The fluxes on the right hand side of (4.3) are computed as for the Neumann case.

On the **corners**, special treatment is needed. Our modified MPFA-L method is modified to become equivalent to the finite element method here. This is done by splitting the corner control volume into four smaller cells, where mass conservation does not necessarily hold, see [5] for details.

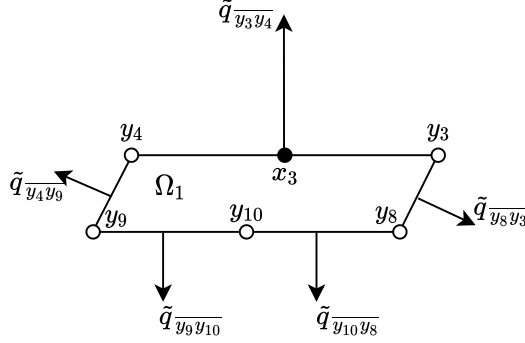


Figure 4.3: The fluxes on the Dirichlet boundary.

## 4.2 Modified Finite Element Method

In this section we introduce a finite element method for solving (4.1). By theorem 3.2.1 the L-triangles form a triangulation  $\{\tau_h\}$ , we will use linear Lagrange elements on this triangulation. The only modifications we need to make are to the mass matrix and the load vector, we let the stiffness matrix stay the same as before. That is, we do not touch the discretization of the constitutive law. We do want however, to define an interpolation operator such that the inner products that make up the mass matrix and load vector, become mass conservative in each control volume.

We need some notation so that we can distinguish between the cell centers in the interior, at cell centers along the boundary and the nodes at the boundary. In addition, corner cells need special treatment. Let  $\mathcal{N}_h^*$  be a set of indices corresponding to all interior nodes of  $\{\tau_h\}$ , which are also the cell centers of the control volume mesh. This index set contains two disjoint sets  $\mathcal{N}_h^* = \mathcal{N}_h^b \cup \mathcal{N}_h^i$ , where superscript  $i$  denotes the cell centers of the interior cells and  $b$  the boundary cells. The index set  $\mathcal{N}_h^b$  are further subdivided as we see in figure 4.4. The nodes at the boundary is indexed by the set  $\mathcal{N}_h^N \cup \mathcal{N}_h^D$ , where  $N$  and  $D$  represent Neumann and Dirichlet boundary nodes, these are further subdivided as illustrated in figure 4.4.



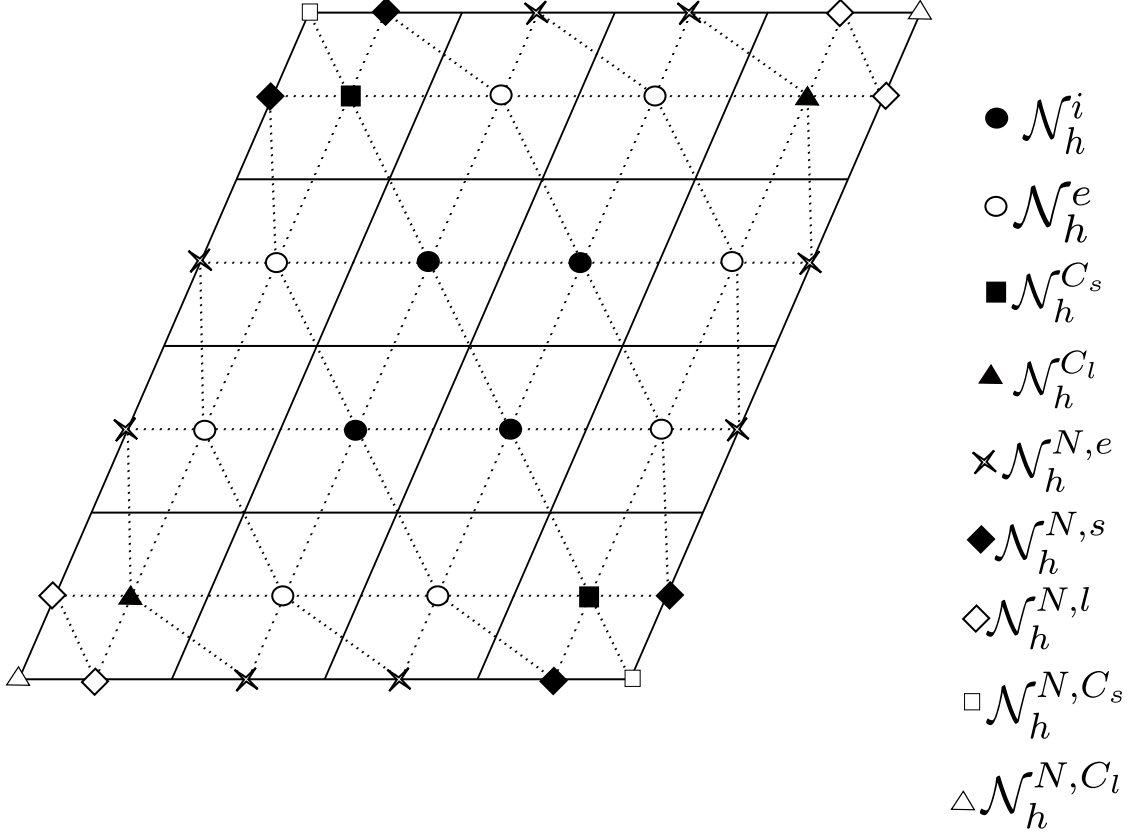


Figure 4.4: A parallelogram mesh with finite element triangles in dotted lines and control volumes in solid lines. In this case we have a pure Neumann problem.

As before we denote by  $V_h$  the linear ansatz space as in definition 10:

$$V_h = \{u_h \in C(\overline{\Omega}) : u_h|_K \in \mathcal{P}_1(K) \ \forall K \in \tau_h, u|_{\Gamma_D} = 0\}$$

similarly  $\phi_i$  is the standard nodal basis function, where  $i \in \mathcal{N}_h \setminus \mathcal{N}_h^D$ . In addition to our global interpolation operator, definition 11, we define an operator that maps functions  $v_h \in V_h$  to functions that are piecewise constant on the control volumes. This piecewise function are equal to  $v_h$  at the nodes of the triangulation. This is an example of *mass lumping*, see [3] for more examples.

**Definition 12** (Piecewise global interpolator). *Let  $\hat{I}_h$  be an operator that maps from the test space to functions that are piecewise constant on control volumes.*

$$\hat{I}_h : V_h \rightarrow L^2(\Omega)$$

And

$$\hat{I}_h v_h = \sum_{i \in \mathcal{N}_h \setminus \mathcal{N}_h^d} v_h(x_i) \hat{I}_h \phi_i(x)$$

Where

$$\hat{I}_h \phi_i(x) = \begin{cases} 1 & \text{if } x \in D_i \\ 0 & \text{otherwise} \end{cases} \quad (4.4)$$

In interior cells,  $i \in \mathcal{N}_h^i$ , we have  $D_i = \Omega_i$ , i.e., the support of (4.4) is the control volume corresponding to  $\phi_i$ . If we are close or on the boundary the situation is more complicated:

- $i \in \mathcal{N}_h^e$ : In this case the function vanishes for the quarter of the parallelogram closest to the boundary, i.e.,  $D_i = \Omega_2$  from figure 4.2
- $i \in \mathcal{N}_h^{N,e}$  In this case of the neumann boundary node  $\hat{I}_h \phi_i(x)$  vanishes outside the quarter of the control volume closest to the edge, i.e.,  $D_i = \Omega_1$  in figure 4.2
- On the corners there are special definitions, see (Cao Wolmuth [5], 2009)

Let  $\hat{I}_{\Gamma_N} = \hat{I}_h|_{\Gamma_N}$  be the trace of the piecewise interpolation operator on the neumann boundary. The finite element method we end up with reads as follows: Find  $u_h \in V_h$  such that

$$\left\langle \hat{I}_h u_h, \hat{I}_h v_h \right\rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla v_h \rangle_{0,\Omega} = \left\langle f, \hat{I}_h v_h \right\rangle_{0,\Omega} + \left\langle g, \hat{I}_{\Gamma_N} v_h \right\rangle_{0,\Gamma_N}, \quad (4.5)$$

for all  $v_h \in V_h$ . The key takeaway here is the local support of the inner products, this will make the mass matrix diagonal.

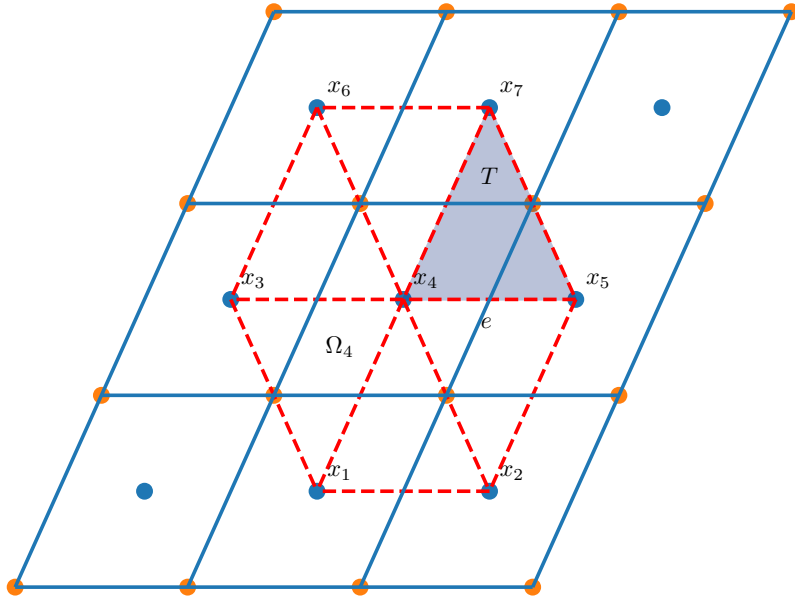


Figure 4.5: The support of  $\phi_4$ , the coloured area corresponds to one triangle (element) in the support of  $\phi_4$ .

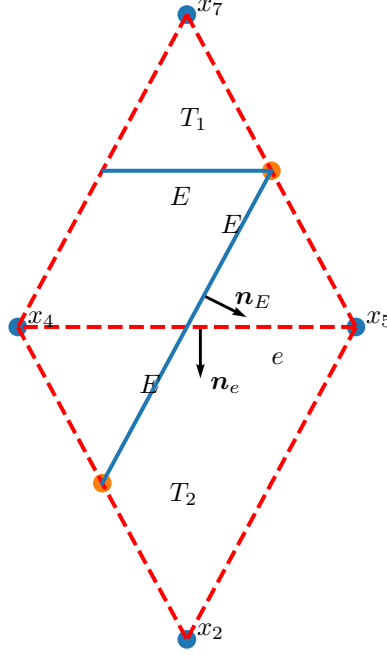


Figure 4.6: Notation in the proof

Now we can state the equivalence theorem:

**Theorem 4.2.1.** *The modified finite element method (4.5) and the modified MPFA-L method are equivalent on uniform parallelogram grid for the time discretized heat equation, i.e., (4.1), on homogeneous media.*

*Proof.* We do the proof in four steps:

1. First, we show the equivalence for the interior, so let  $\Omega_i$  be an interior control volume and  $\phi_i$  be the corresponding basis function evaluating to one at the centre of  $\Omega_i$ , where  $i \in \mathcal{N}_h^i$ . We test (4.5) with  $v_h = \phi_i$ :

$$\langle \hat{I}_h u_h, \hat{I}_h \phi_i \rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla \phi_i \rangle_{0,\Omega} = \langle f, \hat{I}_h \phi_i \rangle_{0,\Omega}. \quad (4.6)$$

Let  $T \in \tau_h \cap \text{supp}(\phi_i)$  be one of the elements in the triangulation that makes up the support of  $\phi_i$ .  $S = T \cap \Omega_i$  is a part of the control volume that lies in

some element, and  $E \subset S \cap \partial\Omega_i$  are the half edges of  $\Omega_i$ .  $e$  are the interior edges of  $\tau_h$  inside the support of  $\phi_i$ , see fig 4.6 and 4.5.  $\mathbf{n}_e$  is the unit normal on  $e$  with fixed and arbitrary orientation, and  $\mathbf{n}_E$  is the unit normal on  $E$  pointing out of  $\Omega_i$ . Let  $T_{e,0}$  and  $T_{e,1}$  be the two elements having  $e$  as a common edge, with the numbering corresponding to the orientation of  $\mathbf{n}_e$ . Since  $u_h$  and  $\phi_i$  are piecewise linear and  $\mathbf{K}$  is constant on each triangle  $T$ , we have:

$$\begin{aligned}
\langle \mathbf{K} \nabla u_h, \nabla \phi_i \rangle_0 &= \int_{\text{supp}(\phi_i)} (\mathbf{K} \nabla u_h)^T \nabla \phi_i \, dx = \sum_{T \in \text{supp}(\phi_i)} \int_T (\mathbf{K} \nabla u_h)^T \nabla \phi_i \, dx \\
&= \sum_{T \in \text{supp}(\phi_i)} \left( \int_{\partial T} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_i \, ds - \int_T \nabla \cdot \mathbf{K} \nabla u_h \phi_i \, dx \right) \\
&= \sum_{T \in \text{supp}(\phi_i)} \int_{\partial T} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_i \, ds \\
&= \sum_{e \in \text{supp}(\phi_i)} \int_e ((\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,0}} - (\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,1}}) \phi_i \, ds \\
&= \sum_{e \in \text{supp}(\phi_i)} ((\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,0}} - (\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,1}}) \frac{|e|}{2} \\
&= \sum_{S \in \text{supp}(\phi)} \int_{\partial S} (\mathbf{K} \nabla u_h)^T \mathbf{n} \, ds - \sum_{E \in \partial\Omega_i} \int_E (\mathbf{K} \nabla u_h)^T \mathbf{n}_E \, ds \\
&= \sum_{S \in \text{supp}(\phi)} \int_S \nabla \cdot \mathbf{K} \nabla u_h \, ds - \sum_{E \in \partial\Omega_i} \int_E (\mathbf{K} \nabla u_h)^T \mathbf{n}_E \, ds \\
&= - \sum_{E \in \partial\Omega_i} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E|.
\end{aligned} \tag{4.7}$$

Note that this last sum is a sum of integrals over the half edges of  $\Omega_i$ . Further, we have that

$$\langle \hat{I}_h u_h, \hat{I}_h \phi_i \rangle_0 = \int_{\Omega} \hat{I}_h u_h \hat{I}_h \phi_i \, dx = \int_{\Omega_i} u_h(x_i) \, dx \tag{4.8}$$

and

$$\langle f, \hat{I}_h \phi_i \rangle_0 = \int_{\Omega_i} f \, dx. \tag{4.9}$$

Combining equation (4.7), (4.8) and (4.9) we get that (4.6) is equivalent to:

$$\int_{\Omega_i} u_h(x_i) \, dx - \sum_{E \in \partial\Omega_i} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| = \int_{\Omega_i} f \, dx. \tag{4.10}$$

We know from theorem 3.2.2 that the flux over each half edge in the L-method is given uniquely by the potential values of the three cell centers in the L-triangle. Since the L-triangles and the elements are the same,  $\nabla u_h$  corresponds to the gradient used in the L-method, see equation (3.44). Hence, if  $\hat{u}_h$  is the solution to (4.1) with the original L-method in the interior, then  $\hat{u}_h(x_i) = u_h(x_i)$  for  $x_i \in \mathcal{N}_h^i$ .

2. For a control volume bordering the **Neumann** boundary, first let  $i \in \mathcal{N}_h^e$ , we have:

$$\left\langle \hat{I}_h u_h, \hat{I}_h \phi_i \right\rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla \phi_i \rangle_{0,\Omega} = \left\langle f, \hat{I}_h \phi_i \right\rangle_{0,\Omega}. \quad (4.11)$$

With similar computations and reasoning as for (4.7) we get:

$$\langle \mathbf{K} \nabla u_h, \nabla \phi_i \rangle_{0,\Omega} = - \sum_{E \in \partial \Omega_{i,2}} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E|, \quad (4.12)$$

where  $\Omega_{i,2}$  is as  $\Omega_2$  in figure 4.2. As  $\hat{I}_h$  is carefully defined close to the Neumann boundary, we get that (4.11) is equivalent to:

$$\int_{\Omega_{i,2}} u_h(x_i) dx - \sum_{E \in \partial \Omega_{i,2}} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| = \int_{\Omega_{i,2}} f(x_i) dx. \quad (4.13)$$

Next, let  $j \in \mathcal{N}_h^{N,e}$ , i.e., the index of a node on the boundary. Then we have

$$\left\langle \hat{I}_h u_h, \hat{I}_h \phi_j \right\rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla \phi_j \rangle_{0,\Omega} = \left\langle f, \hat{I}_h \phi_j \right\rangle_{0,\Omega} + \left\langle g, \hat{I}_{\Gamma_N} \phi_j \right\rangle_{0,\Gamma_N}. \quad (4.14)$$

Similarly as in (4.7) we have

$$\begin{aligned} \langle \mathbf{K} \nabla u_h, \nabla \phi_j \rangle_0 &= \int_{\text{supp}(\phi_j)} (\mathbf{K} \nabla u_h)^T \nabla \phi_j dx = \sum_{T \in \text{supp}(\phi_j)} \int_T (\mathbf{K} \nabla u_h)^T \nabla \phi_j dx \\ &= \sum_{T \in \text{supp}(\phi_j)} \left( \int_{\partial T} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_j ds - \int_T \nabla \cdot \mathbf{K} \nabla u_h \phi_j dx \right) \\ &= \sum_{T \in \text{supp}(\phi_j)} \int_{\partial T} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_j ds. \end{aligned} \quad (4.15)$$

But because  $\phi_j \neq 0$  on  $\text{supp}(\phi_j) \cap \Gamma_N$ , we get

$$\begin{aligned}
\sum_{T \in \text{supp}(\phi_j)} \int_{\partial T} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_j \, ds &= \sum_{e \in \text{supp}(\phi_j)} \int_e ((\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,0}} - (\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,1}}) \phi_j \, ds \\
&\quad + \int_{\Gamma_N \cap \text{supp}(\phi_j)} (\mathbf{K} \nabla u_h)^T \mathbf{n} \phi_j \, ds \\
&= \sum_{e \in \text{supp}(\phi_j)} ((\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,0}} - (\mathbf{K} \nabla u_h)^T \mathbf{n}_e|_{T_{e,1}}) \frac{|e|}{2} \, ds \\
&\quad + (\mathbf{K} \nabla u_h)^T \mathbf{n} |E_{\Gamma_N}| \, ds \\
&= - \sum_{E \in \partial \Omega_j \setminus \Gamma_N} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| \, ds
\end{aligned} \tag{4.16}$$

Combining (4.15) and (4.16) and using the definition of  $\hat{I}_h$ , definition 11, we get that (4.14) is equivalent to:

$$\int_{\Omega_{j,1}} u_h(x_i) \, dx - \sum_{E \in \partial \Omega_{j,1}} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| = \int_{\Omega_{j,1}} f(x_i) \, dx. \tag{4.17}$$

Where  $\Omega_{j,1}$  is as  $\Omega_1$  in figure 4.2. Now, (4.13) and (4.17) are exactly the L-method for the Neumann boundary, as described earlier, see figure 4.1.

3. For a control volume near the **Dirichlet** boundary, let first  $i \in \mathcal{N}_h^e$ , i.e., the cell center. Then, our modified finite element method

$$\left\langle \hat{I}_h u_h, \hat{I}_h \phi_j \right\rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla \phi_j \rangle_{0,\Omega} = \left\langle f, \hat{I}_h \phi_j \right\rangle_{0,\Omega} \tag{4.18}$$

is equivalent to

$$\int_{\Omega_{i,2}} u_h(x_i) \, dx - \sum_{E \in \partial \Omega_{i,2}} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| = \int_{\Omega_{i,2}} f \, dx, \tag{4.19}$$

with the same reasoning as in (4.7), (4.8) and (4.9). As  $\Omega_i = \Omega_{i,1} \cup \Omega_{i,2}$  and  $\Omega_{i,1} \cap \Omega_{i,2} = \emptyset$ , see figure 4.2, we have:

$$- \sum_{E \in \Omega_i \setminus \Gamma_D} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| + \sum_{E \in \Omega_{i,1} \setminus \Gamma_D} (\mathbf{K} \nabla u_h)^T \mathbf{n}_E |E| + \int_{\Omega_{i,1}} f \, dx = \int_{\Omega_i} f \, dx. \tag{4.20}$$

We recognize the second and third terms in the above equation as the flux across the Dirichlet boundary in the modified L method, see (4.3).

4. See [5] for equivalence on the corner cells.

□

### 4.3 Convergence Rate Estimates

Our modified finite element method only approximates the bi-linear and linear form, and we need to take this into account when proving a convergence rate estimate. The following lemma is an extension of C  a's lemma 3.1.9, it is useful for estimating the error when our bi-linear and linear form is not exact.

**Lemma 4.3.1** (First Lemma of Strang, page 155 [8]). *Suppose there exists some  $\alpha > 0$  such that for all  $h > 0$  and  $v_h \in V_h$*

$$\alpha \|v_h\|_1^2 \leq a_h(v_h, v_h)$$

*and let  $a$  be continuous in  $V \times V$ . Then there exist some constant  $C$  independent of  $V_h$  such that*

$$\begin{aligned} \|u - u_h\|_1 \leq C \left\{ \inf_{v_h \in V_h} \left\{ \|u - v_h\|_1 + \sup_{w_h \in V_h} \frac{|a(v_h, w_h) - a_h(v_h, w_h)|}{\|w_h\|_1} \right\} \right. \\ \left. + \sup_{w_h \in V_h} \frac{|l(w_h) - l_h(w_h)|}{\|w_h\|_1} \right\} \end{aligned} \quad (4.21)$$

From (4.5) we see that we have a bi-linear form

$$a_h(u_h, v_h) = \left\langle \hat{I}_h u_h, \hat{I}_h v_h \right\rangle_{0,\Omega} + \langle \mathbf{K} \nabla u_h, \nabla v_h \rangle_{0,\Omega}. \quad (4.22)$$

And the linear form:

$$b_h(v_h) = \left\langle F, \hat{I}_h v_h \right\rangle_{0,\Omega} + \left\langle g, \hat{I}_{\Gamma_N} v_h \right\rangle_{0,\Gamma_N}. \quad (4.23)$$

To apply the first Lemma of Strang 4.3.1, we first show that  $a_h(\cdot, \cdot)$  is coercive. We write out the Sobolev norm

$$\|u_h\|_1^2 = \langle \nabla u_h, \nabla u_h \rangle_0 + \|u_h\|_0^2. \quad (4.24)$$

Using the Poincar   inequality on the second term:

$$\begin{aligned} \|u_h\|_1^2 &\leq \langle \nabla u_h, \nabla u_h \rangle_0 + C_\Omega \langle \nabla u_h, \nabla u_h \rangle_0 \\ &\leq \left( \frac{1 + C_\Omega}{\tau \|\mathbf{K}\|} \right) \tau \langle \mathbf{K} \nabla u_h, \nabla u_h \rangle_0 \\ &\leq \left( \frac{1 + C_\Omega}{\tau \|\mathbf{K}\|} \right) \left( \tau \langle \mathbf{K} \nabla u_h, \nabla u_h \rangle_0 + \left\langle \hat{I}_h u_h, \hat{I}_h u_h \right\rangle_0 \right) \\ &= \frac{1}{\alpha} a_h(u_h, u_h), \end{aligned} \quad (4.25)$$

we obtain coercivity with  $\alpha = \tau \|\mathbf{K}\| / (1 + C_\Omega)$ , where  $C_\Omega$  is some constant depending on the domain and the boundary conditions.

Another important piece that must be in place for a convergence proof is the piecewise interpolation error:



**Lemma 4.3.2.** *For the previously defined piecewise global interpolator  $\hat{I}_h$ , definition 12, we have the estimate:*

$$\left\| \hat{I}_h u_h - u_h \right\|_{0,\Omega} \leq Ch |u_h|_{1,\Omega} \quad \forall u_h \in V_h, \quad (4.26)$$

for some constant  $C$  independent of the mesh diameter.

*Proof.*

$$\begin{aligned} \left\| \hat{I}_h u_h - u_h \right\|_0^2 &= \sum_{i \in \mathcal{N}_h^*} \left\| \hat{I}_h u_h - u_h \right\|_{0,\Omega_i}^2 \\ &= \sum_{i \in \mathcal{N}_h^*} \int_{\Omega_i} (u_h(x_i) - u_h(x))^2 dx \\ &= \sum_{i \in \mathcal{N}_h^*} \int_{\Omega_i} h^2 \left( \frac{u_h(x_i) - u_h(x)}{h} \right)^2 dx \\ &\leq \sum_{i \in \mathcal{N}_h^*} \int_{\Omega_i} h (\nabla u_h)^T \nabla u_h dx \\ &= Ch |\nabla u_h|_1^2 \end{aligned} \quad (4.27)$$

□

Still working on a short proof here.

The proof is trivial as our test space,  $V_h \subset H^1(\Omega) \cap C(\Omega)$ , consists of continuous functions. We are now ready to state the  $H^1$  error estimate for the modified finite element method and thus the MPFA-L method.

**Theorem 4.3.3.** *Let  $u$  solve (4.1) and  $u_h$  be the solution resulting from MPFA-L, then there exists a positive constant  $C$  independent of the mesh diameter,  $h$ , such that*

$$\|u - u_h\|_1 \leq Ch (\|u\|_2 + \|f\|_0 + \|g\|_{\frac{1}{2},\Gamma_N}). \quad (4.28)$$

*Proof.* The hypothesis in Strang's lemma 4.3.1 on continuity and coercivity are fulfilled. Let  $C$  be a generic positive constant. We start by controlling the second

term on the right hand side in (4.21), the truncation error in the bi-linear form:

$$\begin{aligned}
& \sup_{w_h \in V_h} \frac{|a(v_h, w_h) - a_h(v_h, w_h)|}{\|w_h\|_1} \\
&= \sup_{w_h \in V_h} \frac{|\langle v_h, w_h \rangle + \tau \langle \mathbf{K} \nabla v_h, \nabla w_h \rangle - \langle \hat{I}_h v_h, \hat{I}_h w_h \rangle - \tau \langle \mathbf{K} \nabla v_h, \nabla w_h \rangle|}{\|w_h\|_1} \\
&= \sup_{w_h \in V_h} \frac{|\langle v_h, w_h \rangle - \langle \hat{I}_h v_h, w_h \rangle + \langle \hat{I}_h v_h, w_h \rangle - \langle \hat{I}_h v_h, \hat{I}_h w_h \rangle|}{\|w_h\|_1} \\
&= \sup_{w_h \in V_h} \frac{|\langle \hat{I}_h v_h - v_h, w_h \rangle + \langle \hat{I}_h v_h, \hat{I}_h w_h - w_h \rangle|}{\|w_h\|_1}.
\end{aligned} \tag{4.29}$$

We see from the above computations, that the truncation error in the bi-linear form, only has a contribution from the *mass lumping*. By Cauchy Schwarz inequality and lemma 4.3.2 we get:

$$\begin{aligned}
& \leq \sup_{w_h \in V_h} \frac{Ch|v_h|_1 \|w_h\|_0 + \|\hat{I}_h v_h\|_0 Ch|w_h|_1}{\|w_h\|_1} \\
& \leq \sup_{w_h \in V_h} \frac{Ch|v_h|_1 \|w_h\|_0 + \|\hat{I}_h v_h\|_0 Ch|w_h|_1}{\|w_h\|_1} + \frac{Ch \|v_h\|_0 \|w_h\|_0 + \|\hat{I}_h v_h\|_0 Ch \|w_h\|_0}{\|w_h\|_1} \\
& \leq Ch \left( \|v_h\|_0 + \|\hat{I}_h v_h\|_0 \right).
\end{aligned} \tag{4.30}$$

The third term in (4.21), the linear form, can be controlled similarly:

$$\begin{aligned}
& \sup_{w_h \in V_h} \frac{l(w_h) - l_h(w_h)}{\|w_h\|_1} = \sup_{w_h \in V_h} \frac{\langle f, w_h - \hat{I}_h w_h \rangle_{0,\Omega} + \langle g, w_h - \hat{I}_{\Gamma_N} w_h \rangle_{0,\Gamma_N}}{\|w_h\|_1} \\
& \leq \sup_{w_h \in V_h} \frac{\|f\|_0 Ch \|w_h\|_1 + \|g\|_{\frac{1}{2},\Gamma_N} \|w_h - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2},\Gamma_N}}{\|w_h\|_1}.
\end{aligned} \tag{4.31}$$

Now, we want to bound  $\|w - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2},\Gamma_N}$  by  $\|w_h\|_1$ . Let  $v_h$  be a piecewise constant function on the boundary in each Neumann boundary triangle. Then we

have:

$$\begin{aligned} \|w_h - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2}, \Gamma_N} &= \sup_{0 \neq v \in H^{\frac{1}{2}}(\Omega)} \frac{\langle w_h - \hat{I}_{\Gamma_N} w_h, v \rangle_{\Gamma_N}}{\|v\|_{\frac{1}{2}, \Gamma_N}} \\ &= \sup_{0 \neq v \in H^{\frac{1}{2}}(\Omega)} \frac{\langle w_h - \hat{I}_{\Gamma_N} w_h, v - v_h \rangle_{\Gamma_N}}{\|v\|_{\frac{1}{2}, \Gamma_N}}, \end{aligned} \quad (4.32)$$

as  $\int_{\Gamma_N} (w_h - \hat{I}_{\Gamma_N} w_h) dx = 0$ . Now, we can use Cauchy Schwarz inequality:

$$\|w_h - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2}, \Gamma_N} \leq \sup_{0 \neq v \in H^{\frac{1}{2}}(\Omega)} \frac{\|w_h - \hat{I}_{\Gamma_N} w_h\|_{0, \Gamma_N} \|v - v_h\|_{0, \Gamma_N}}{\|v\|_{\frac{1}{2}, \Gamma_N}}. \quad (4.33)$$

By the inequality

$$\|v - v_h\|_{0, \Gamma_N} \leq Ch^{\frac{1}{2}} \|v\|_{\frac{1}{2}, \Gamma_N}, \quad (4.34)$$

we can bound the right hand side of (4.33):

$$\|w_h - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2}, \Gamma_N} \leq Ch^{\frac{1}{2}} \|w_h - \hat{I}_{\Gamma_N} w_h\|_{0, \Gamma_N}. \quad (4.35)$$

Using (4.34) again, we get

$$\|w_h - \hat{I}_{\Gamma_N} w_h\|_{-\frac{1}{2}, \Gamma_N} \leq Ch \|w_h\|_{\frac{1}{2}, \Gamma_N} \leq Ch \|w_h\|_1. \quad (4.36)$$

Where the last inequality is due to the definition of the  $H^{\frac{1}{2}}$  norm. Inserting this into (4.31), gives us a bound on the truncation error of our linear form:

$$\sup_{w_h \in V_h} \frac{l(w_h) - l_h(w_h)}{\|w_h\|_1} \leq Ch(\|f\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N}). \quad (4.37)$$

Hence, from (4.21), we have the error estimate:

$$\|u - u_h\|_1 \leq \inf_{v_h \in V_h} \left\{ \|u - v_h\|_1 + Ch \left( \|v_h\|_0 + \|\hat{I}_h v_h\|_0 + \|f\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N} \right) \right\}. \quad (4.38)$$

If we let  $v_h = I_h u \in V_h$ , in (4.38), where  $I_h : C(\Omega) \rightarrow V_h$  is the global interpolation operator, we get the inequality:

$$\|u - u_h\|_1 \leq \|u - I_h u\|_1 + Ch \left( \|I_h u\|_0 + \|\hat{I}_h I_h u\|_0 + \|f\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N} \right). \quad (4.39)$$

As discussed earlier, in section 3.1.6 about convergence of finite element method, we have the estimate:

$$\|u - I_h u\|_1 \leq Ch|u|_2. \quad (4.40)$$

If we insert this into (4.39), we get:

$$\|u - u_h\|_1 \leq Ch \left( \|u\|_2 + \|I_h u\|_0 + \left\| \hat{I}_h I_h u \right\|_0 + \|f\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N} \right). \quad (4.41)$$

As  $I_h u, \hat{I}_h u \rightarrow u$  as  $h \rightarrow 0$ , we can control the first three terms inside the parenthesis by the  $H^2$  norm of  $u$ , and we get the desired result:

$$\|u - u_h\|_1 \leq Ch \left( \|u\|_2 + \|f\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N} \right). \quad (4.42)$$

□

In this chapter we have introduced a way to handle Dirichlet and Neumann boundary conditions for the MPFA-L method, and showed convergence when applied to (4.1). The convergence was obtained with showing equivalence to a modified linear Lagrange finite element method.

**Remark 14.** *In our convergence rate estimate, we showed that  $\|u - u_h\|_1$  decreases proportional to the mesh diameter,  $h$ . In [5], the authors show, using the Aubin-Nitsche technique, an estimate where  $\|u - u_h\|_0$  decreases proportional to the square of the mesh diameter. We expect similar results can be shown here, as the equations are similar.*

## Chapter 5

# Convergence of the MPFA-L Method for Richards' Equation

In this chapter we use the results from chapter three to prove convergence of the Richards equation discretized in space with MPFA-L method, in time with backward Euler and linearized with the L-scheme. We start by considering the Richards' equation without the gravity term in an isotropic medium: Find  $\psi = \psi(x, t)$  such that

$$\begin{cases} \partial_t \theta(\psi) - \nabla \cdot (\kappa(\theta(\psi)) \nabla \psi) = f, & \text{in } \Omega \times (0, T] \\ \psi = 0, & \text{on } \partial\Gamma_D \times (0, T] \\ -\kappa(\theta(\psi)) \nabla \psi = g_N, & \text{on } \partial\Gamma_N \times (0, T) \\ \psi = u_0, & \text{on } \Omega \times \{t = 0\} \end{cases} \quad (5.1)$$

This equation has a non linearity in the flux,  $\mathbf{q} = -\kappa(\theta(\psi)) \nabla \psi$ , which makes it hard to apply our results, as they require a homogeneous medium. The hydraulic conductivity,  $\kappa(\theta(\psi))$ , depends on our solution which is heterogeneous, and is thus itself heterogeneous. To remedy this we use the Kirchhoff transform

$$\begin{aligned} \mathcal{K} : \mathbb{R} &\rightarrow \mathbb{R}^+ \\ \psi &\mapsto \int_0^\psi \kappa(\theta(\phi)) \, d\phi = u. \end{aligned} \quad (5.2)$$

As discussed in 2.4, the functions  $\theta(\cdot)$  and  $\kappa(\cdot)$  are continuous, monotone increasing functions, the Kirchhoff transform,  $\mathcal{K}$ , therefore has an inverse,  $\mathcal{K}^{-1}$ . We define

$$\begin{aligned} b(u) &:= \theta(\mathcal{K}^{-1}(u)) \\ k(u) &:= \kappa(\theta(\mathcal{K}^{-1}(u))). \end{aligned} \quad (5.3)$$

Further, assume that the hydraulic conductivity is bounded below and above

$$0 < \kappa_m \leq \kappa(\theta) \leq \kappa_M, \quad (5.4)$$

and that the water content is a Lipschitz continuous function of pressure

$$\sup_{\psi} |\theta(\psi)| \leq L_{\theta}. \quad (5.5)$$

Then  $b(\cdot)$  is also Lipschitz continuous

$$\sup_u |b'(u)| = \left| \theta'(\mathcal{K}^{-1}(u)) \frac{1}{\mathcal{K}'(\mathcal{K}^{-1}(u))} \right| = \left| \theta'(\mathcal{K}^{-1}(u)) \frac{1}{\kappa(\mathcal{K}^{-1}(u))} \right| \leq L_B. \quad (5.6)$$

We note that we in fact remove the non linearity in the constitutive law; by the chain rule, we get

$$\nabla u = \kappa(\theta(\psi)) \nabla \psi. \quad (5.7)$$

We can write the Richards' equation (5.1) in the transformed variable  $u$  to get: Find  $u = u(x, t)$  such that

$$\begin{cases} \partial_t b(u) - \nabla \cdot \nabla u = f, & \text{in } \Omega \times (0, T] \\ u = 0, & \text{on } \partial\Gamma_D \times (0, T] \\ -\nabla u = g_N, & \text{on } \partial\Gamma_N \times (0, T] \\ u = u_0, & \text{on } \Omega \times \{t = 0\} \end{cases} \quad (5.8)$$

We start by discretizing (5.8) with the MPFA-L method, we divide our domain into  $d$  quadrilaterals (control volumes). Writing (3.52) in vector form we find  $\tilde{u}_h \in \mathbb{R}^d$  such that:

$$\partial_t \mathbf{B}^V b(\tilde{u}_h) + \mathbf{A}^V \tilde{u}_h = \mathbf{q}^V. \quad (5.9)$$

We can then discretize in time using backward Euler. Given  $\tilde{u}_h^{n-1}, \mathbf{q}^n \in \mathbb{R}^d$  we should then find  $\tilde{u}_h^n \in \mathbb{R}^d$  such that:

$$\mathbf{B}^V b(\tilde{u}_h)^n + \tau \mathbf{A}^V \tilde{u}_h^n = \tau \mathbf{q}^{Vn} + \mathbf{B}^V b(\tilde{u}_h)^{n-1}. \quad (5.10)$$

Now we need to linearize (5.10) with the L-scheme. We see from (3.67) that the applying this linearization leads to the equation: Given  $\tilde{u}_h^{n,j-1}, \tilde{u}_h^{n-1} \in \mathbb{R}^d$  find  $\tilde{u}_h^{n,j} \in \mathbb{R}^d$  such that

$$L \mathbf{B}^V (\tilde{u}_h^{n,j} - \tilde{u}_h^{n,j-1}) + \tau \mathbf{A} \tilde{u}_h^{n,j} = -\mathbf{B}^V \theta(\tilde{u}_h^{n,j-1}) + \tau \mathbf{q}^{Vn} + \mathbf{B}^V \theta(\tilde{u}_h^{n-1}). \quad (5.11)$$

We set  $\tilde{u}_h^{n,0} = \tilde{u}_h^{n-1}$  and solve the above equation until we reach the stopping criterion (3.71) each time step. See listing 6.2 for the code. Note that (5.11) is equivalent to a finite element discretization, this is what we will use to prove convergence.

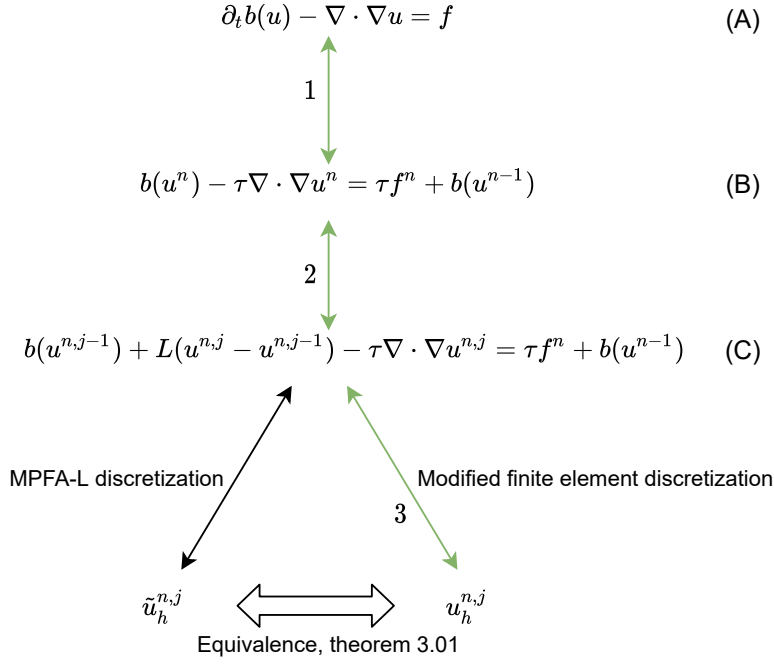


Figure 5.1

To prove convergence we list some assumptions:

**A-1** The domain  $\Omega \subset \mathbb{R}^2$  is bounded and have a Lipschitz continuous boundary, i.e., it is locally the graph of a Lipschitz continuous function.

**A-2**  $b \in C^1$  is non-decreasing and Lipschitz continuous.

**A-3**  $b(u_0)$  is essentially bounded in  $\Omega$  and  $u_0 \in L^2(\Omega)$ .

**Theorem 5.0.1.** *Assume **A 1-3**, then Richards' equation after Kirchhoff transform (5.8) discretized with backward Euler in time,  $L$ -scheme linearization and MPFA-L discretization in space (5.11) converges, and we have the estimate*

$$\|\tilde{u}_h^n - u(t^n)\|_0 \leq C(h + \tau), \quad (5.12)$$

where  $C$  is some constant independent of the maximum mesh diameter  $h$ , the time step length  $\tau$  and the minimum number of linearization iterations  $j$ .

*Proof.* As the MPFA-L method and the modified finite element method are equivalent, we prove the convergence of our finite element solution,  $u_h^{n,j}$ . The proof will be done in three steps, see figure 5.1. We have by the triangle inequality

$$\|u(t_n) - u_h^{n,j}\|_0 \leq \|u(t_n) - u^n\|_0 + \|u^n - u^{n,j}\|_0 + \|u^{n,j} - u_h^{n,j}\|_0. \quad (5.13)$$

- The third term is the error of solving the elliptic problem (C) in 5.1, or

$$u^{n,j} - \frac{\tau}{L} \nabla \cdot \nabla u^{n,j} = \frac{Lu^{n,j-1} + \tau f^n - b(u^{n,j-1}) + u^{n-1}}{L}. \quad (5.14)$$

By theorem 4.3.3 we have the error estimate

$$\|u_h^{n,j} - u^{n,j}\|_1 \leq Ch(\|u^{n,j-1}\|_2 + \left\| \frac{Lu^{n,j-1} + \tau f^n - b(u^{n,j-1}) + u^{n-1}}{L} \right\|_0 + \|g\|_{\frac{1}{2}, \Gamma_N}). \quad (5.15)$$

The expression involving  $u^{n-1}$  and  $u^{n,j-1}$  are bounded independently of the mesh, see [13] page 1459 for a stability estimate, hence  $\|u_h^{n,j} - u^{n,j}\|_0 \leq C_3 h$ .

- To bound the second term,  $\|u^{n,j} - u^n\|_0$ , we will use techniques found in (Radu, List, [9]). First, we subtract the variational form of (B) from the variational form of (C) in figure 5.1: For any  $v \in H_0^1$  and  $j > 1$

$$\langle b(u^{n,j-1}) - b(u^n), v \rangle + \tau \langle \nabla(u^{n,j} - u^n), \nabla v \rangle + L \langle u^{n,j} - u^{n,j-1}, v \rangle = 0. \quad (5.16)$$

Let  $e^{n,j} := u^{n,j} - u^n$ , then test (5.16) with  $e^{n,j}$ :

$$\langle b(u^{n,j-1}) - b(u^n), e^{n,j} \rangle + \tau \|\nabla e^{n,j}\| + L \langle u^{n,j} - u^{n,j-1}, e^{n,j} \rangle = 0. \quad (5.17)$$

Now, we use the relation  $\langle b - a, b \rangle = \frac{1}{2} \|b\|^2 + \frac{1}{2} \|b - a\|^2 - \frac{1}{2} \|a\|^2$  and some simple algebraic manipulation to obtain

$$\begin{aligned} \langle b(u^{n,j-1}) - b(u^n), e^{n,j-1} \rangle + \tau \|\nabla e^{n,j}\| + \frac{L}{2} \|e^{n,j}\|^2 + \frac{L}{2} \|e^{n,j} - e^{n,j-1}\|^2 \\ \leq \frac{L}{2} \|e^{n,j-1}\|^2 - \langle b(u^{n,j-1}) - b(u^n), e^{n,j} - e^{n,j-1} \rangle. \end{aligned} \quad (5.18)$$

Next, we use Cauchy Schwarz inequality on the first term, and Young's inequality on the last term

$$\begin{aligned} \|b(u^{n,j-1}) - b(u^n)\| \|e^{n,j-1}\| + \tau \|\nabla e^{n,j}\| + \frac{L}{2} \|e^{n,j}\|^2 + \frac{L}{2} \|e^{n,j} - e^{n,j-1}\|^2 \\ \leq \frac{L}{2} \|e^{n,j-1}\|^2 + \frac{1}{2L} \|b(u^{n,j-1}) - b(u^n)\|^2 + \frac{L}{2} \|e^{n,j} - e^{n,j-1}\|^2. \end{aligned} \quad (5.19)$$

We cancel the last term on the right side against the last term on the left side. Since  $b(\cdot)$  is Lipschitz continuous with  $\|b(x) - b(y)\| \leq L_B \|x - y\|$ , we have

$$\begin{aligned} \frac{1}{L_B} \|b(u^{n,j-1}) - b(u^n)\|^2 + \tau \|\nabla e^{n,j}\| + \frac{L}{2} \|e^{n,j}\|^2 \\ \leq \frac{L}{2} \|e^{n,j-1}\|^2 + \frac{1}{2L} \|b(u^{n,j-1}) - b(u^n)\|^2. \end{aligned} \quad (5.20)$$



Using the Poincaré inequality we obtain

$$\left(\frac{L}{2} + \frac{\tau}{C_\Omega}\right) \|e^{n,j}\|^2 \leq \frac{L}{2} \|e^{n,j-1}\|^2 + \left(\frac{1}{2L} - \frac{1}{L_B}\right) \|b(u^{n,j-1}) - b(u^n)\|^2. \quad (5.21)$$

Since  $L_B < 2L$  we reach the convergence estimate

$$\|e^{n,j}\|^2 \leq \frac{L}{L + \frac{2\tau}{C_\Omega}} \|e^{n,j-1}\|^2 \quad (5.22)$$

Hence,  $\|u^{n,j} - u^n\| = \|e^{n,j}\| \rightarrow 0$  as  $j \rightarrow \infty$ . Therefore, we choose  $j$  such that

$$\|u^{n,j} - u^n\| \leq (h + \tau), \quad \forall n \in [1, N],$$

where  $N = \lceil \frac{1}{\tau} \rceil$  is the number of time steps.

- The first term  $\|u(t^n) - u^n\|$  can be bounded by the techniques used in (Radu, Pop and Knabner, [13]). We reach the estimate

$$\|u^n - u(t^n)\|_0 \leq C_1 \tau. \quad (5.23)$$

Using all of the above, we get

$$\|\tilde{u}_h^n - u(t^n)\|_0 \leq C_3 h + C_1 \tau + \tau + h \leq (\max(C_1, C_3) + 1)(h + \tau). \quad (5.24)$$

We see from the above equation, that the desired result (5.12) is reached with  $C = (\max(C_1, C_3) + 1)$   $\square$

We have showed convergence for Richards' equation after Kirchhoff transform (5.8), discretized in space by MPFA-L method, in time by backward Euler and L-scheme for linearization. In section 6.3 we do numerical results to confirm this.

**Remark 15.** *We expect that a better convergence rate estimate*

$$\|\tilde{u}_h^n - u(t^n)\|_0 \leq C(h^2 + \tau), \quad (5.25)$$

*with the square of the mesh diameter, is possible. This is because we use the  $\|\cdot\|_1$  estimate for the spatial discretization, but according to remark 14, there exists a better  $\|\cdot\|_0$  estimate we could use instead in the above proof.*



# Chapter 6

## Numerical Results

In this chapter we do several numerical experiments with the algorithms covered in this thesis. And briefly discuss the code used to do the experiments.

### 6.1 Computer Code

Most of the code used in this thesis can be found on <https://github.com/trulsmoholt/masterthesis>, and is written in python and numpy. It is primarily intended for educational purposes and for comparing convergence rates of different spatial approximation techniques. An example of a very simple use case can be seen in 6.1

```

1
2 from discretization.mesh import Mesh
3 from discretization.FVML import compute_matrix, compute_vector
4 import numpy as np
5 import math
6
7 #Function to perturb mesh from unit square. Takes a 2d numpy
  vector and returns a 2d numpy vector. This particular choice
  makes a parallelogram mesh.
8 perturbation=lambda p: np.array([p[0],0.5*p[0]+p[1]])
9
10 #Number of grid points in x and y direction
11 nx=ny=10
12
13 mesh = Mesh(nx,ny,perturbation,ghostboundary=True)
14 source = lambda x,y:math.sin(y)*math.cos(x)
15 boundary_condition = lambda x,y:0
16 tensor = np.eye(2)
17 permeability = np.ones((mesh.num_unknowns,mesh.num_unknowns))
18
19 A = np.zeros((mesh.num_unknowns,mesh.num_unknowns))#stiffness
  matrix
20 f = np.zeros(mesh.num_unknowns)#load vector
21
22 compute_matrix(mesh,A,tensor,permeability)
23 compute_vector(mesh,f,source,boundary_condition)
24
25 u = np.linalg.solve(A,f)
26 mesh.plot_vector(u)
27

```

Listing 6.1: Solving simple Poisson equation.

The code is centred around the mesh class, which contains information about how the domain is discretized into quadrilaterals and its properties. This class also contains public functions to make plots of different kinds and compute errors. The spatial numerical methods implemented are: TPFA, MPFA-L, MPFA-O and linear Lagrange finite elements. They all have the same call signature, as in 6.1 line 22 and 23. The control volume methods also has the option of taking a matrix to store the flux stencils. One can use sparse matrices instead of dense numpy matrices in 6.1, as long as the indexing signature is the same as in numpy, for example scipy has a compatible sparse matrix library.

The code also has implementations of mass matrix and the gravitation term. Also included in the github are an example of how to solve Richards' equation using L-scheme linearization and backward Euler, see 6.2 for some of the code.

```

1 u_l = u.copy() #linearization/L-scheme iterate
2 u_t = u.copy() #timestep iterate
3 F = u.copy() #source vector
4 A = np.zeros((mesh.num_unknowns, mesh.num_unknowns)) #stiffness matrix
5 B = mass_matrix(mesh)
6
7 #time iteration
8 for t in time_partition[1:]:
9     #empty source vector
10    F.fill(0)
11    #compute source vector
12    compute_vector(mesh, F, lambda x, y: f(x, y, t), lambda x, y: u_exact(x, y, t))
13    #L-scheme iteration
14    while True:
15        #compute the heterogeneous hydraulic conductivity, kappa
16        conductivity = kappa(np.reshape(u_l, (mesh.cell_centers.shape[0], mesh
17        .cell_centers.shape[1]), order='F'))
18        A.fill(0) #empty the stiffness matrix
19        compute_matrix(mesh, A, K, conductivity) #compute stiffness matrix
20        lhs = L*B*u_l + B@theta(u_t) - B@theta(u_l) + tau*F
21        u = np.linalg.solve(lhs, rhs)
22        #check if L-scheme linearization has acceptable error
23        if np.linalg.norm(u-u_l) <= TOL+TOL*np.linalg.norm(u_l):
24            #quit linearization and do another time step
25            break
26        else:
27            #update linearization iterate
28            u_l = u
29    #update time step iterate
30    u_t = u
31    #update linearization iterate
32    u_l = u

```

Listing 6.2: Linearization and time stepping of Richards' equation.

## 6.2 Elliptic Equation

The convergence tests in this section are similar to some of the tests done in chapter three of [2]. We consider the elliptic model problem (3.1), find  $u \in H_0^1(\Omega)$  such that

$$\begin{aligned}\nabla \cdot \mathbf{q} &= f \\ \mathbf{q} &= -\mathbf{K} \nabla u.\end{aligned}\tag{6.1}$$

We set the solution

$$u = \cosh(\pi x) \cos(\pi y)\tag{6.2}$$

And set  $\mathbf{K}$  to be the identity matrix. As in [2] page 1340 we define the normalized discrete  $L_2$  norms:

$$\|u - u_h\|_{0,h} = \left( \frac{1}{V} \sum_i V_i (u_{h,i} - u_i)^2 \right)^{\frac{1}{2}} \quad (6.3)$$

$$\|q - q_h\|_{0,h} = \left( \frac{1}{Q} \sum_a Q_a (q_{h,a} - q_a)^2 \right)^{\frac{1}{2}}, \quad (6.4)$$

where  $q_a := -\hat{\mathbf{n}} \cdot \mathbf{q}$  is the normal flow density over edge  $a$ , with  $\hat{\mathbf{n}}$  being unit normal to the edge and  $\mathbf{q}$  evaluated at the midpoint of the edge.  $q_{h,a}$  is the discrete flux over  $a$  defined similarly with  $\mathbf{q}_h$  being the discrete normal flow density, for a finite volume method this would be the flux across some edge divided by the edge length. For the finite element method, we use the MPFA-L flux stencil to recover the flux in the experiments where it is present. Let  $u_{h,i}$  denote the discrete potential at cell  $i$ , and  $u_i$  is the potential evaluated at cell  $i$ . For the finite element method, this would be the function value at the grid points/nodes.  $Q_a$  is the volume associated with edge  $a$ , i.e., the sum of the two volumes sharing edge  $a$ .  $V = \sum_i V_i$  and  $Q = \sum_a Q_a$ .

In the figures (5.12-5.15) in this section we see convergence rates in the  $\|\cdot\|_{0,h}$  norm for different spatial numerical methods. The y-axis is the  $\log_2$  of the error and the x-axis is  $\log_2 n$ , where  $n$  is the number of points in x and y direction, and thus proportional with the inverse of the mesh diameter. The slope of the graph in the plots are the convergence rate.

In figures 6.3 and 6.3 we see that all the methods converge with the same quadratic rate. This fits well with the fact that all the methods covered in this thesis are equivalent to the TPFA method for uniform grids.

In figures 6.5, 6.6 and 6.7 we observe that the TPFA method does not converge for parallelogram mesh. This makes sense as the grid is not K-orthogonal. The others methods still have quadratic convergence for potential and flow density.

In figures 6.8, 6.9 and 6.10 the MPFA-L, MPFA-O and FEM still converges quadratically for the potential on a rough grid, where every grid point is perturbed randomly by a factor proportional to the grid diameter. For the normal flow density however, the convergence rate drops to about one.

In figure 6.11 we introduce grids with aspect ratio, that is grids with more points in the y direction. In figure 6.12 we observe that MPFA-L, MPFA-O and FEM has a convergence rates for the potential of about 1.5 for the grid with aspect ratio 0.1. In figure 6.14 we see that the MPFA-O method fail to converge for the grid with aspect ratio 0.01. Thus the MPFA-L method wins this round of numerical experiments.

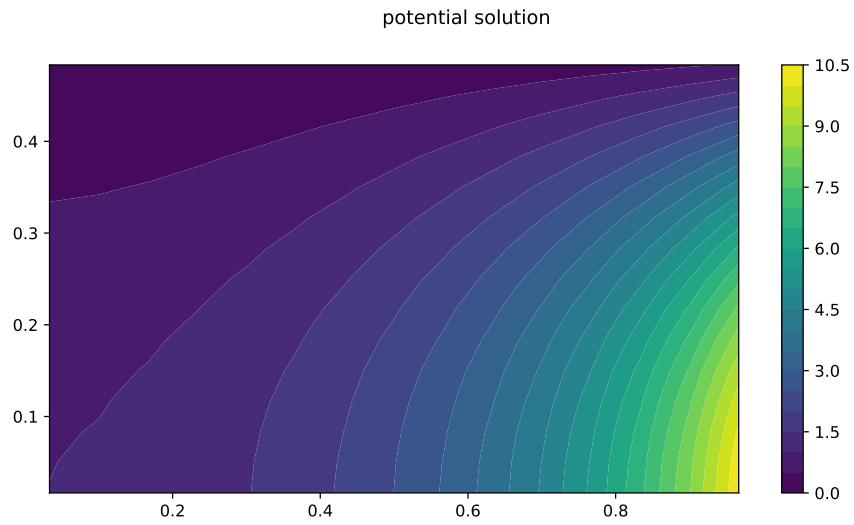


Figure 6.1: The solution (6.2) on half the unit square

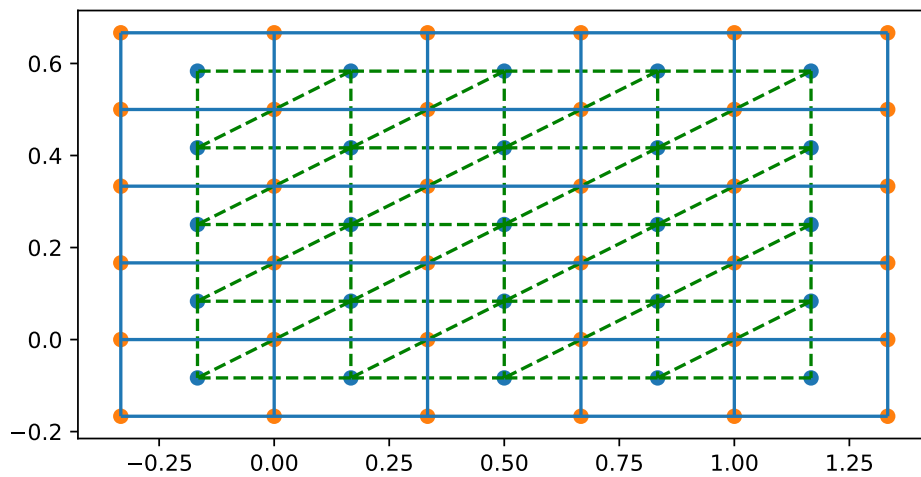


Figure 6.2: Uniform rectangular mesh on half the unit square. The triangles are used for the finite element solution and are spanned between the nodes of the cell centers of the finite volume methods. The ghost cell boundary is included, so this mesh has nine degrees of freedom, i.e., the interior cells.

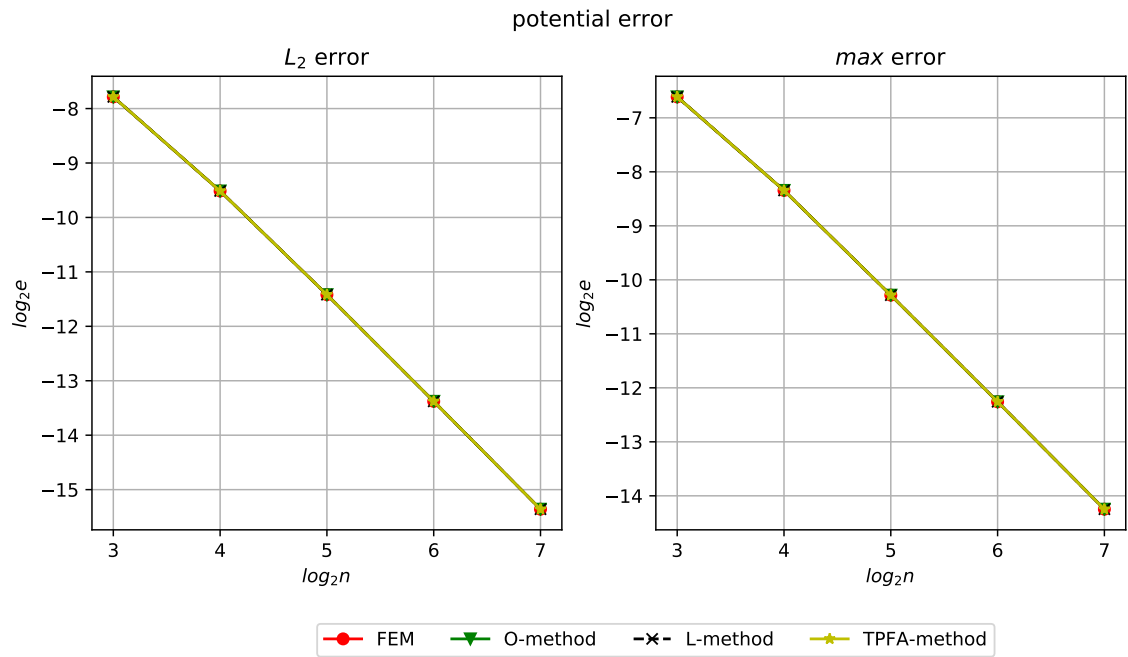


Figure 6.3: Potential error on refinements of the uniform rectangular mesh 6.2. The convergence is the same for all the schemes.



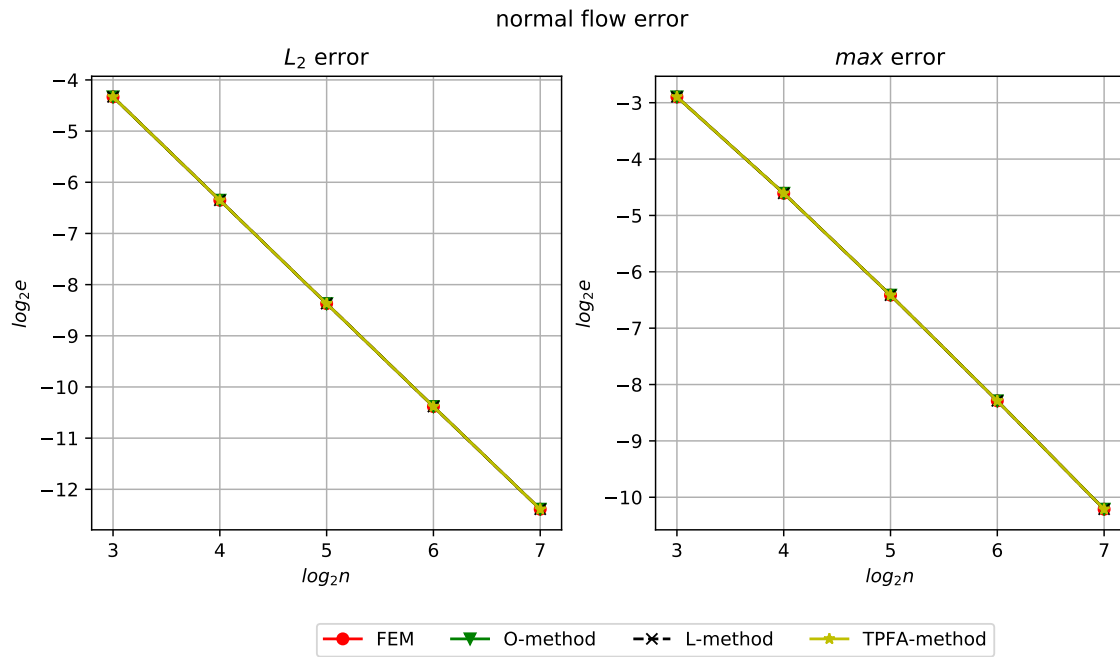


Figure 6.4: Normal flow density error on refinements of the uniform rectangular mesh 6.2. The convergence is the same for all the schemes.

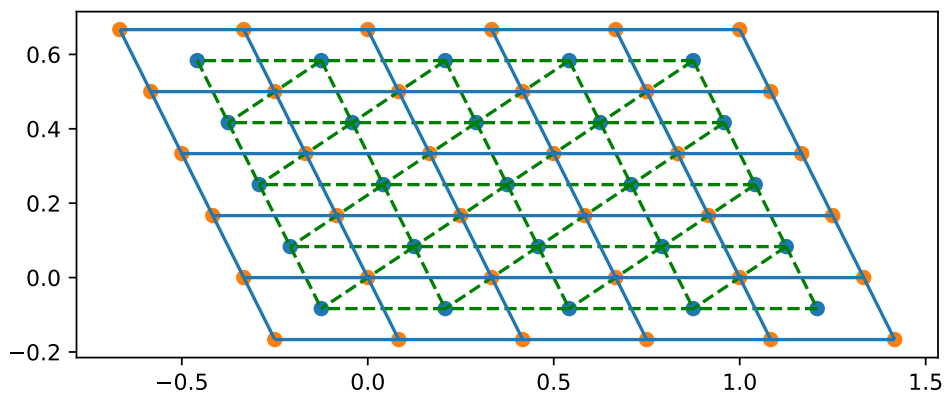


Figure 6.5: Trapezoidal mesh, now every point is transformed by  $(x, y) \mapsto (x - 0.5y, y)$

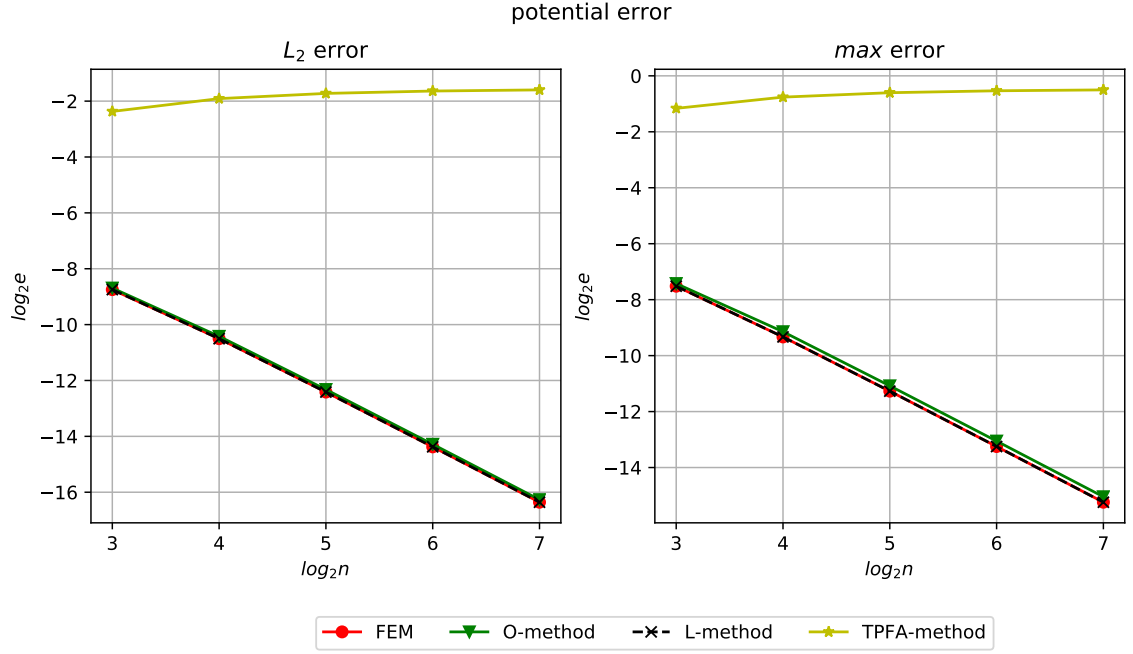


Figure 6.6: Pressure error on refinements of the mesh 6.5

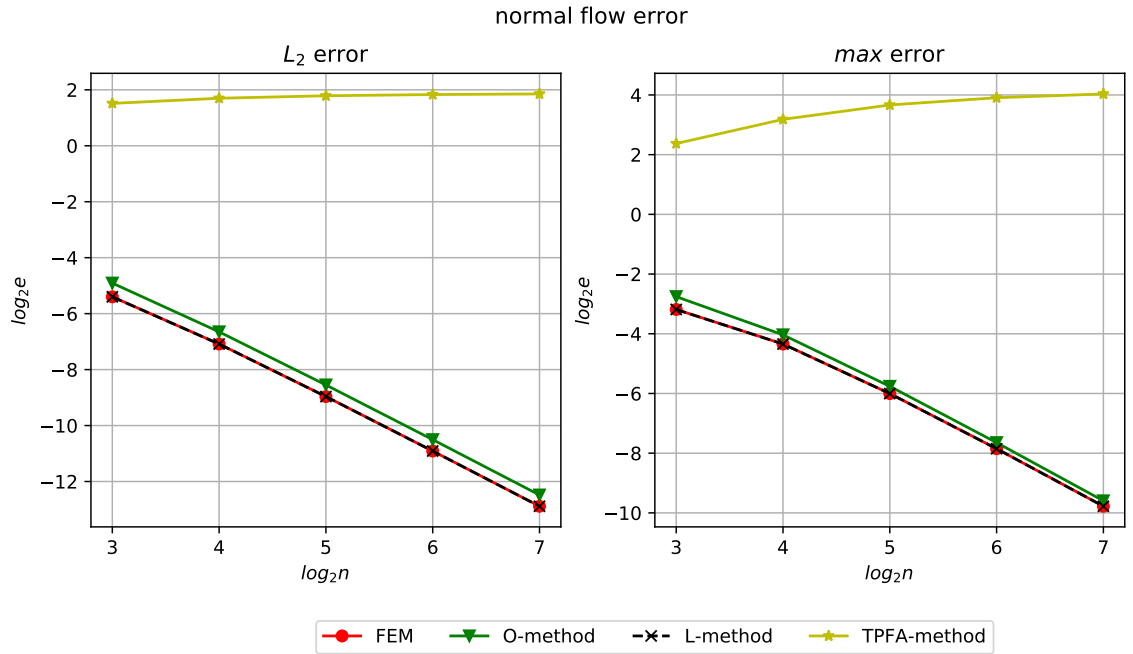


Figure 6.7: Normal flow density error on refinements of the mesh 6.5

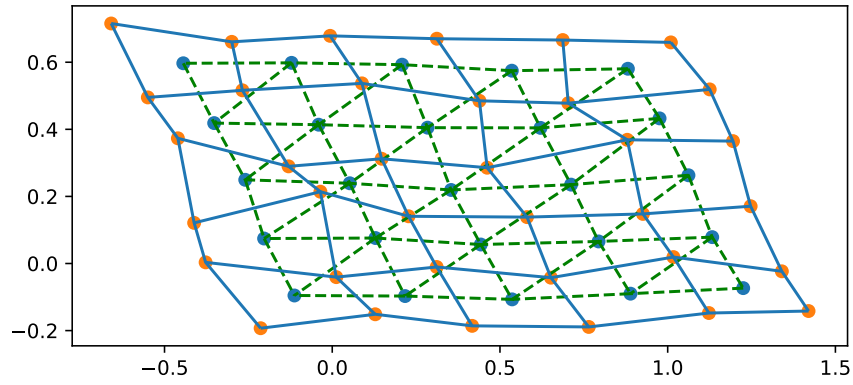


Figure 6.8: Perturbed mesh, every point in the mesh is perturbed by a random number which is  $O(\frac{h}{5})$ , in both x and y direction.

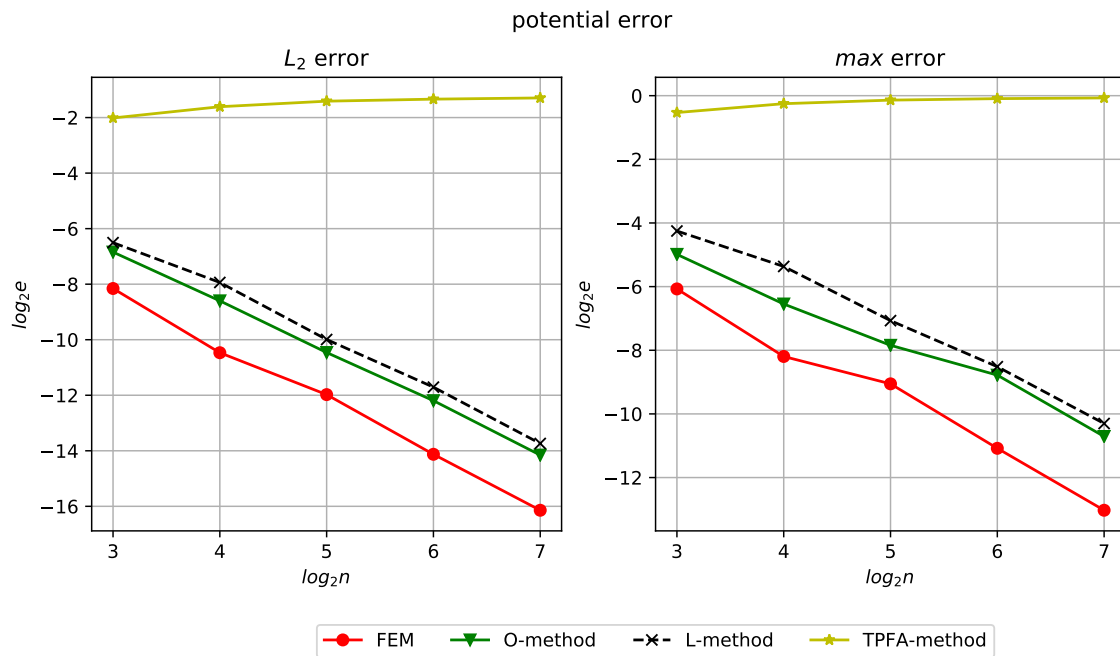


Figure 6.9: The pressure error of perturbed mesh.

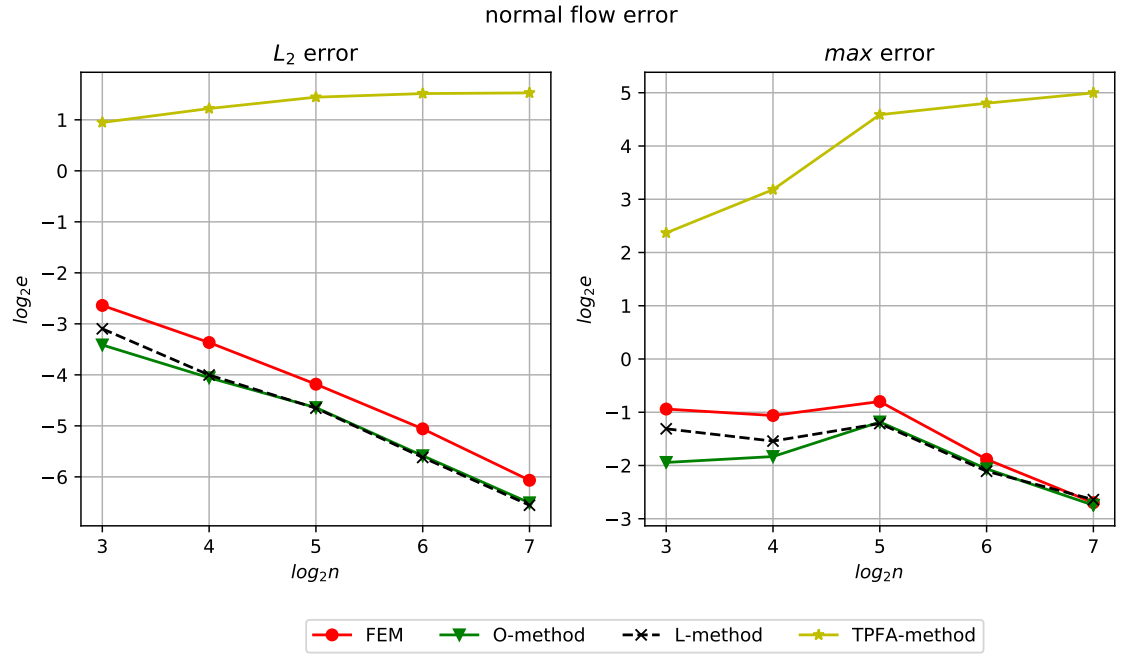


Figure 6.10: The normal flow density error of perturbed mesh

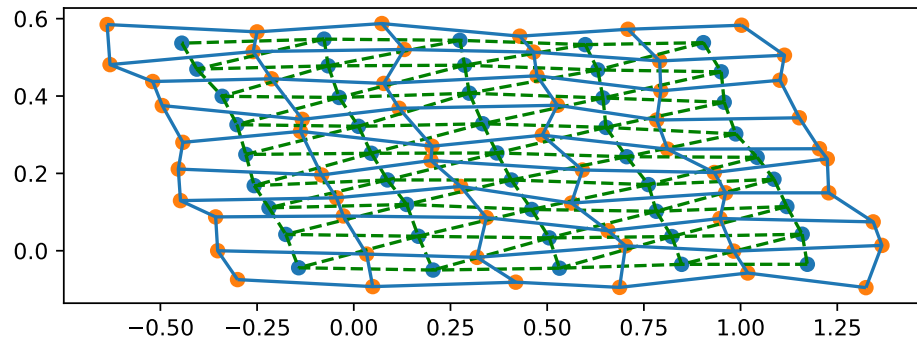


Figure 6.11: Perturbed mesh with aspect ratio 0.5, there are half as many points in the x-direction as in the y-direction.

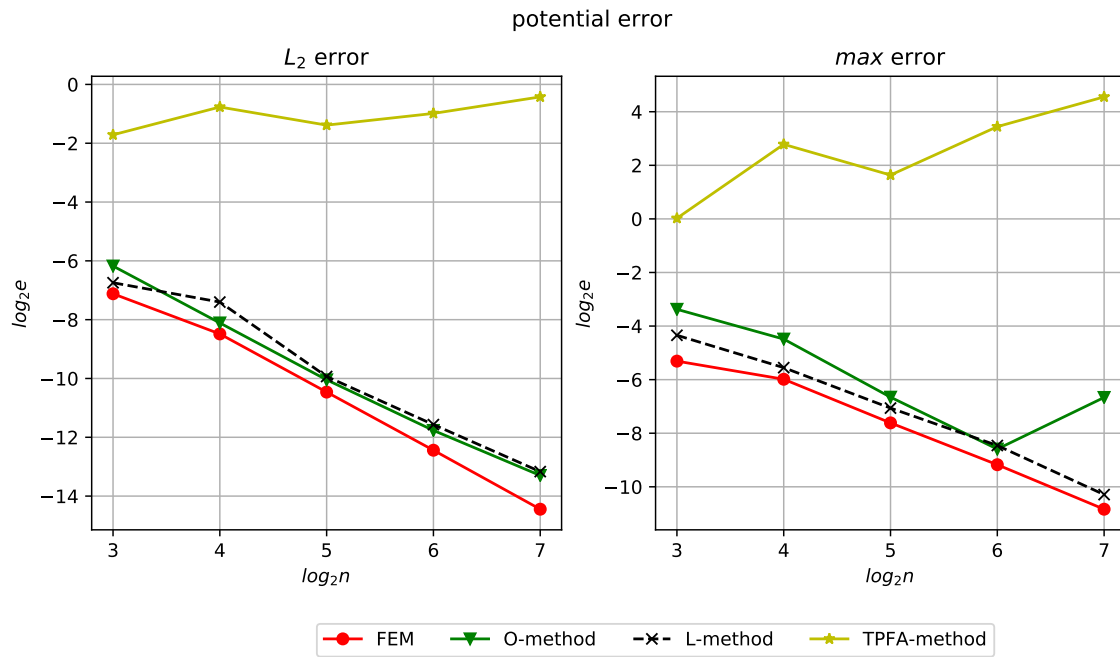


Figure 6.12: The pressure error of perturbed mesh with aspect ratio 0.1.

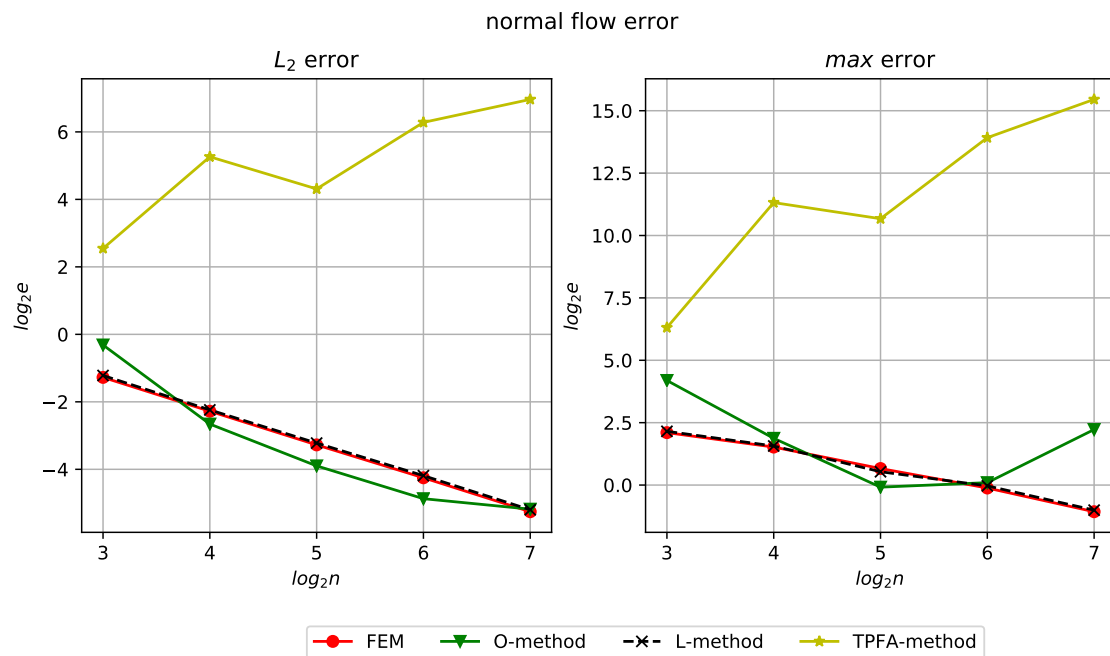


Figure 6.13: The normal flow density error of perturbed mesh with aspect ratio 0.1.

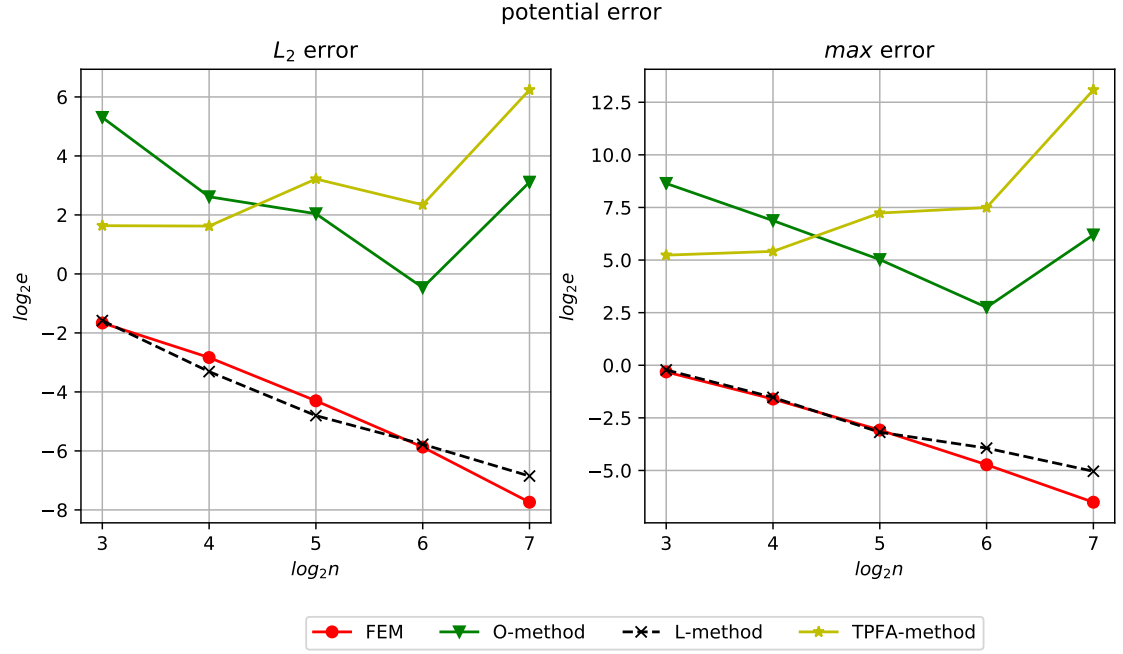


Figure 6.14: The pressure error of perturbed mesh with aspect ratio 0.01.

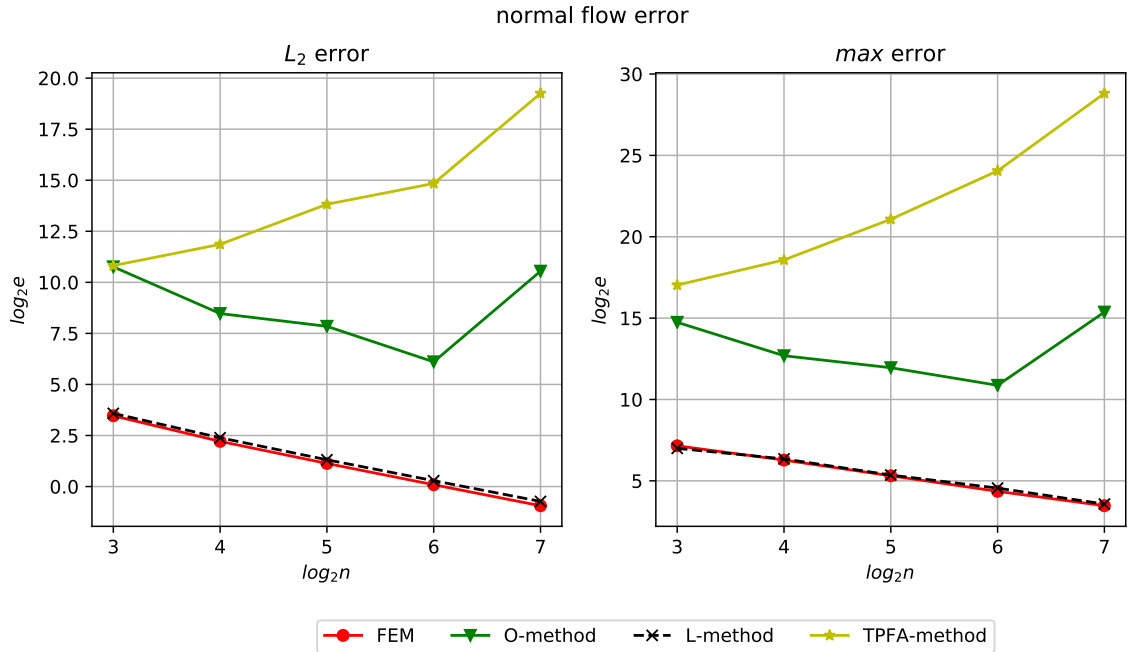


Figure 6.15: The normal flow density error of perturbed mesh with aspect ratio 0.01.

## 6.3 Richards' Equation

This section is not yet done

### 6.3.1 Constant Hydraulic Conductivity

In this section we consider numerical experiments for (5.8), with Dirichlet boundary conditions: Find  $u = u(x, t)$  such that

$$\begin{cases} \partial_t b(u) - \nabla \cdot \nabla u = f, & \text{in } \Omega \times (0, T] \\ u = u|_{\Gamma_D}, & \text{on } \partial\Gamma_D \times (0, T] \\ u = u_0, & \text{on } \Omega \times \{t = 0\} \end{cases} \quad (6.5)$$

We define

$$b(u) := \frac{1}{1 - u}, \quad (6.6)$$

and compute the source term,  $f$ , such that the solution becomes

$$u = -tx(1 - x)y(1 - y) - 1. \quad (6.7)$$

We let  $\Omega$  be the unit square perturbed by  $(x, y) \mapsto (x - 0.5y, y)$ . The L-scheme linearization has the parameters  $L := 1.5$  and error tolerance  $TOL = 5e^{-9}$ . We perform grid refinement of a parallelogram grid, see figure (6.16). In table 6.3 we observe a quadratic convergence when the time step length is equal to the square of the grid diameter. In table 6.2 we observe a linear convergence when the time step length and the grid diameter is equal.

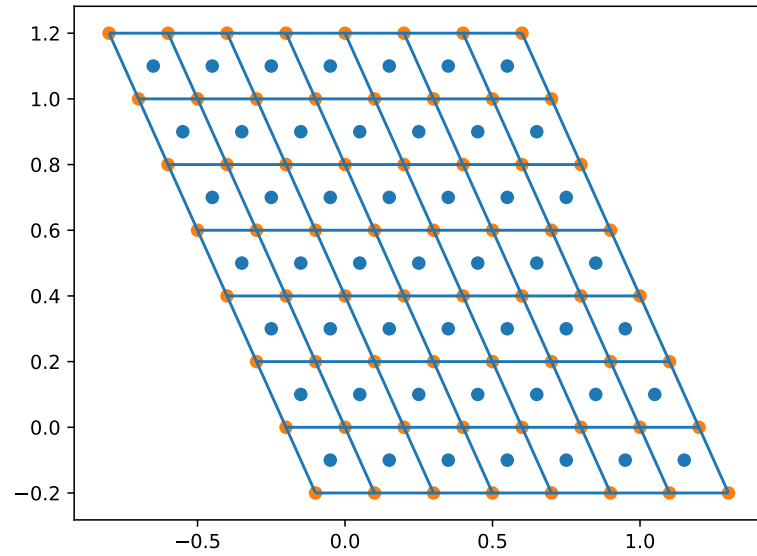


Figure 6.16: Parallelogram grid, with ghost Dirichlet boundary cells.

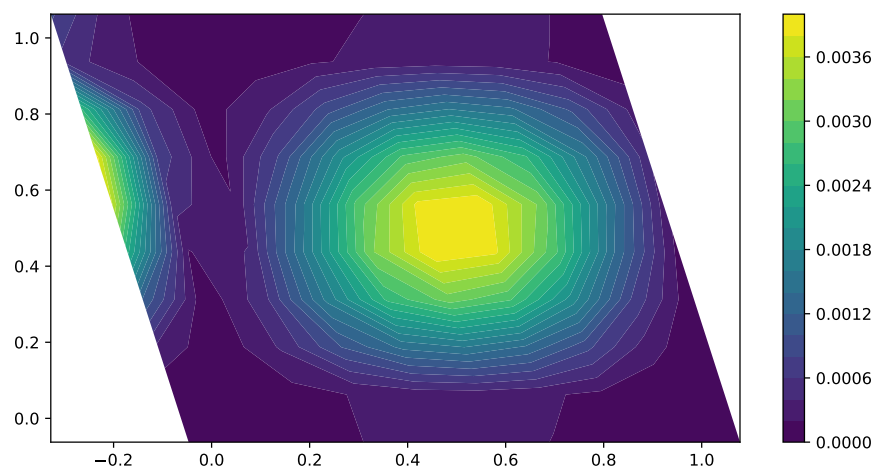


Figure 6.17: The solution of (6.5) at  $T = 1$ , with the ghost Dirichlet boundary.



number	mesh diameter, $h$	time step length, $\tau$	discrete $L_2(\Omega)$ error	improvement
1	0.45069	0.20312	0.001695	-
2	0.22535	0.05078	0.000375	4.51623
3	0.11267	0.01270	0.000087	4.31530
4	0.05634	0.00317	0.000021	4.20040

Table 6.1: Convergence table for (6.5),(6.6) and (6.7). The time step length,  $\tau$ , is set proportional to the square of the mesh diameter, that is  $\tau = h^2$ .

number	mesh diameter, $h$	time step length, $\tau$	discrete $L_2(\Omega)$ error	improvement
1	0.45069	0.45069	0.001694	-
2	0.22535	0.22535	0.000374	4.52868
3	0.11267	0.11267	0.000086	4.33993
4	0.05634	0.05634	0.000020	4.24067

Table 6.2: Convergence table for (6.5),(6.6) and (6.7). The time step length,  $\tau$ , is set proportional to the square of the mesh diameter, that is  $\tau = h$ .

### 6.3.2 Non-Linear Hydraulic Conductivity

Here, we consider Richards' equation (2.12) in pressure variable, find  $p = p(x, t)$  such that

$$\begin{cases} \partial_t \theta(p) - \nabla \cdot \kappa(\theta(p)) \nabla u = f, & \text{in } \Omega \times (0, T] \\ p = p|_{\Gamma_D}, & \text{on } \partial\Gamma_D \times (0, T] \\ p = u_0, & \text{on } \Omega \times \{t = 0\} \end{cases} \quad (6.8)$$

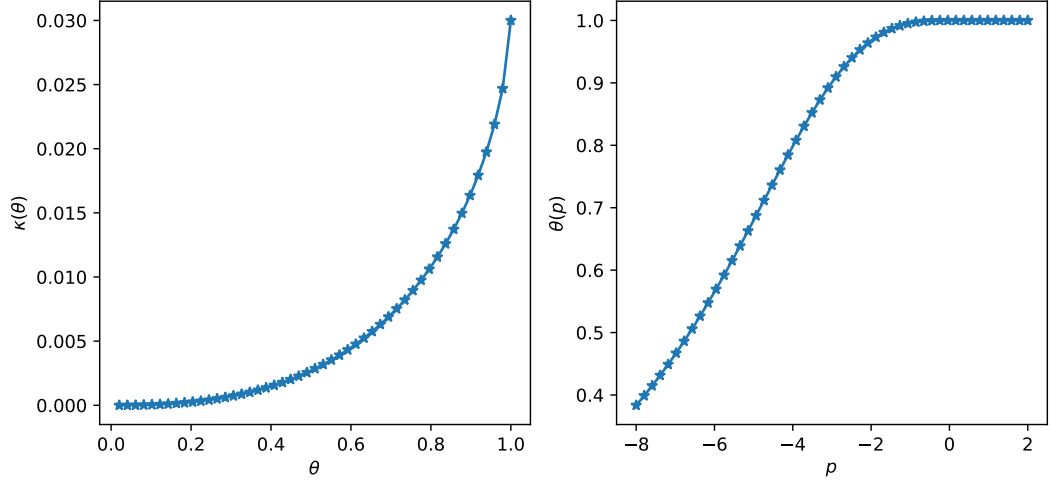
With  $\Omega$  being the perturbed unit square as before, and  $T = 1$ . We define

$$\theta(u) := \begin{cases} (1 + (-\alpha_{vG} p)^{n_{vG}})^{-\frac{n_{vG}-1}{n_{vG}}}, & p \leq 0 \\ 1, & p > 0 \end{cases} \quad (6.9)$$

and

$$\kappa(\theta) := \frac{\kappa_{abs}}{\mu} \sqrt{\theta} \left( 1 - \left( 1 - \theta^{\frac{n_{vG}}{n_{vG}-1}} \right)^{\frac{n_{vG}-1}{n_{vG}}} \right)^2 \quad (6.10)$$

and compute the source term,  $f$ , such that the solution becomes (6.7).



number	mesh diameter, $h$	time step length, $\tau$	discrete $L_2(\Omega)$ error	improvement
1	0.45069	0.20312	0.001863	-
2	0.22535	0.05078	0.000455	4.09403
3	0.11267	0.01270	0.000110	4.13154

Table 6.3: Convergence table for (6.5),(6.6) and (6.7). The time step length,  $\tau$ , is set proportional to the square of the mesh diameter, that is  $\tau = h^2$ .

number	mesh diameter, $h$	time step length, $\tau$	discrete $L_2(\Omega)$ error	improvement
1	0.45069	0.45069	0.017904	-
2	0.22535	0.22535	0.006011	2.97848
3	0.11267	0.11267	0.001682	3.57467

Table 6.4: Convergence table for (6.5),(6.6) and (6.7). The time step length,  $\tau$ , is set proportional to the square of the mesh diameter, that is  $\tau = h$ .

# Bibliography

- [1] I. AAVATSMARK, *An introduction to multipoint flux approximations for quadrilateral grids*, Computational Geosciences, 6 (2002), pp. 405–432.
- [2] I. AAVATSMARK, G. EIGESTAD, B. MALLISON, AND J. NORDBOTTEN, *A compact multipoint flux approximation method with improved robustness*, Numerical Methods for Partial Differential Equations, 24 (2008), pp. 1329–1360.
- [3] J. BARANGER, J.-F. MAITRE, AND F. OUDIN, *Connection between finite volume and mixed finite element methods*, ESAIM: Mathematical Modelling and Numerical Analysis-Modélisation Mathématique et Analyse Numérique, 30 (1996), pp. 445–465.
- [4] Y. CAO, R. HELMIG, AND B. I. WOHLMUTH, *Geometrical interpretation of the multi-point flux approximation l-method*, International Journal for Numerical Methods in Fluids, 60 (2009), pp. 1173–1199.
- [5] ———, *Convergence of the multipoint flux approximation l-method for homogeneous media on uniform grids*, Numerical Methods for Partial Differential Equations, 27 (2011), pp. 329–350.
- [6] W. CHENEY, *Analysis for Applied Mathematics*, Springer-Verlag New York Inc., 2001.
- [7] L. C. EVANS, *Partial differential equations*, American Mathematical Society, Providence, R.I., 2010.
- [8] P. KNABNER AND L. ANGERMAN, *Numerical Methods for Elliptic and Parabolic Partial Differential Equations*, vol. 44 of Texts in Applied Mathematics, Springer-Verlag New York, 2003.
- [9] F. LIST AND F. A. RADU, *A study on iterative methods for solving richards' equation*, Computational Geosciences, 20 (2016), pp. 341–353.

- [10] J. NORDBOTTEN AND M. CELIA, *Geological storage of CO<sub>2</sub>: Modeling Approaches for Large-Scale Simulation*, "John Wiley & Sons", 2011.
- [11] J. M. NORDBOTTEN, I. AAVATSMARK, AND G. T. EIGESTAD, *Monotonicity of control volume methods*, Numer. Math., 106 (2007), p. 255–288.
- [12] J. M. NORDBOTTEN AND E. KEILEGAVLEN, *An introduction to multi-point flux (mpfa) and stress (mpsa) finite volume methods for thermo-poroelasticity*, 2020.
- [13] F. A. RADU, I. S. POP, AND P. KNABNER, *Order of convergence estimates for an euler implicit, mixed finite element discretization of richards' equation*, SIAM Journal on Numerical Analysis, 42 (2004), pp. 1452–1478.
- [14] E. STEIN, *History of the Finite Element Method – Mathematics Meets Mechanics – Part I: Engineering Developments*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2014, pp. 399–442.
- [15] E. STORVIK, *On the optimization of iterative schemes for solving non-linear and/or coupled pdes*, Master's thesis, University of Bergen, 2018.