# POLITECNICO DI TORINO

**Master's degree course in Mechanical Engineering**

Master's Degree Thesis

# Study of Finite Element Analysis Using the Shifted Boundary Method



**Supervisor**
Prof. Claudio Canuto

*Supervisor's signature*

**Candidate**
Yamal Abou Jokh

*Candidate's signature*

Anno Accademico 2019-2020

# Summary

Generally speaking, a differential problem has 3 main ingredients:

1. A differential equation to solve, with a physical quantity which we wish to find.

2. A physical domain where the differential equation is applied, including the boundary of such domain.

3. And boundary conditions.

Over the course of history, several methods were developed to solve differential problems. As of today, mainly three of those methods prevail: Finite Difference Method (FDM), Finite Volume Method (FVM), and Finite Element Method (FEM). This is because of the pros and cons which each of them has.

In the case of the **FDM** it is very easy to use and program, but it is restricted to specific physical domains, and some conditions of continuity and smoothness are required. The **FVM** is particularly useful where conservation laws hold, but it has issues with continuity and it has the disadvantage to smooth large gradients. On the other hand, the **FEM** is more complex algebraically and more difficult to program, but it lacks the restrictions of the FDM and the disadvantages of the FVM.

Since the previous century, the Finite Element Method (FEM) has continuously found applications to solve multiple types of physical problems, however as its range of application has increased, many Engineers, Mathematicians, Physicists, and Programmers have realized the difficulties that it encompasses.

From the 3 ingredients of a differential problem, the steps of a typical FEM are:

- Discretize the physical domain into small pieces, known as "finite elements". Across the finite elements, there will be points where it is wished to know the value of the physical quantity, those points are called "nodes".

- Then apply a suitable test function on every node, and assume that the physical quantity is a linear combination of those test functions. Test functions are functions which approximate the physical quantity over the physical domain, as well as minimize the error in the approximation. This step comes from a variational point of view.

- According to the differential equation to solve, following the algebra will result into a system of algebraic equations. Applying boundary conditions and solving the system will give an approximated numerical solution of the physical quantity on every node.

- Finally, post-processing the results.

The reasons why the FEM gained increasing attention as to become a standard to solve a wide range of situations in nowadays were two: Because it can practically solve

any differential problem, and the implementation of computers which allowed to obtain and test solutions for such problems. The latter reason being the starting point for which the scientific community paid attention to such method, even though this also represents a drawback since computers only have a finite computational capacity.

In a broad sense, the FEM in its most basic definition is more expensive in terms of computation cost compared to other methods, such as FDM and FVM. Even if the FEM can solve any differential problem, it is useless if it can not do so in an acceptable amount of time, for a certain precision. From this idea, it is logical that we wish to decrease the computational cost of the FEM.

The **total** computational cost of a FEM is directly influenced by its first step, meaning the way the physical domain is discretized, and this has implications on all of the remaining steps of the method. This is so since, when there are different shapes and sizes of finite elements along the physical domain, unique calculations have to be done for every unique finite element.

This cost could be diminished by using a constant shape of finite element across the physical domain, but this is not always possible since the physical boundary **imposes a restriction** on the shape of the finite elements. To illustrate this, imagine trying to divide a circle into small squares, no matter how small the squares are, they will never be able to equal the boundary of the circle, so special elements will be needed to fill the circle completely.

This is an obstacle only because there is information about boundary conditions at the physical boundary, but could be avoided if equivalent boundary conditions were translated into the boundaries of the finite elements. That is the fundamental idea of **Shifted Boundary Method**, which is to translate boundary conditions from the physical boundary to the boundary of the finite elements, thus diminishing the computational cost of the FEM while maintaining the same level of approximation. References to this can be found in [1, 2, 3, 4, 5, 6].

This idea also breaks the restriction of the shape of the finite elements near to the physical boundary, since there would be no need to approximate the physical boundary at all. This is indeed a powerful advantage since we could use **any** shape of finite element to fill **any** physical domain.

This advantage can be further extended if a **constant shape of finite element** is used. In this case, we propose rectangular finite elements with bilinear shape functions, instead of the typical triangular finite elements with linear shape functions. This choice will give a set of computations which are the same for each finite element, a desirable characteristic since it makes the method easier to write algebraically and to program, such as the FDM. References to this procedure can be found in [7, 8].

These ideas can be used for the case in which the physical domain changes in time, as we shall see in the development of this thesis, where we model **the heat equation**. References to this can be found in [9, 10].

The implementation of the Shifted Boundary Method for rectangular elements and its application in time-changing domains are among the original contributions of this thesis, with respect to the existing literature on the subject cited above.

# Preliminaries: Notation

Before proceeding, it is important to define a useful notation to avoid misconceptions. We define 3 important notations:

- The dot product between 2 vectors $\vec{v_1}$ and $\vec{v_2}$, will be denoted by $\langle \vec{v_1}, \vec{v_2} \rangle$. On this paper, we treat the operator $\nabla$ as a vector, therefor the divergence of vector $\vec{v_1}$ is written as $div(\vec{v_1}) = \langle \nabla, \vec{v_1} \rangle$.

- The multiplicative sign used is a simple dot, like $\cdot$

- In almost all integrals, the integrand is enclosed into brackets like $\int [I(x,y)] \, dxdy$, however it might also be enclosed into parenthesis like $\int (I) \, dxdy$. Exceptions are unless the integrand is easily distinguishable.

# Contents

# Chapter 1

# A simple beginning: Poisson's equation

We start with a simple problem. We wish to solve the Poisson equation over a rectangular region $\Omega$, with Dirichlet boundary conditions:

$$- \langle \nabla, \nabla u \rangle = f \ in \ \Omega \tag{1.1}$$

$$u = g \ on \ \partial\Omega$$

where $u$ is a certain scalar physical quantity, $g$ are the Dirichlet boundary conditions of the differential problem, the region $\Omega$ is defined as a rectangle of base $B$ and height $H$, so $\Omega = [B \times H]$, and $\partial\Omega$ is the boundary of $\Omega$. To start solving our problem, we use a Finite Element Analysis, for this we multiply the Poisson equation by an arbitrary test function $v$ and integrate over $\Omega$. Performing the algebra we get to the following equivalent problem:

$$\int_\Omega \langle \nabla u, \nabla v \rangle \, d\Omega - \int_{\partial\Omega} \left[ v \cdot \frac{du}{dn} \right] ds = \int_\Omega (f \cdot v) \, d\Omega \tag{1.2}$$

$$u = g \ in \ \partial\Omega$$

where $n$ is the outward unit vector normal to boundary $\partial\Omega$. As an additional condition, we add that $v = 0$ on $\partial\Omega$, so $\int_{\partial\Omega} \left[ v \cdot \frac{du}{dn} \right] ds = 0$. This last condition holds since we are imposing Dirichlet boundary conditions along $\partial\Omega$.

Performing the steps of the FEA:

- We choose to discretize region $\Omega$ into rectangular finite elements of constant base $h_x$ and constant height $h_y$, nodes will be placed on the corners of the finite elements, therefore each finite element has 4 nodes.

- We assign a unique number to each node on region $\Omega$ (that way, each node is uniquely identified).

- The test functions $v(x, y)$ are defined as to be equal to 1 on an internal node, and decrease to 0 in a bilinear form as we move away from it towards its consecutive most adjacent nodes, afterwards it remains zero everywhere else on domain $\Omega$.
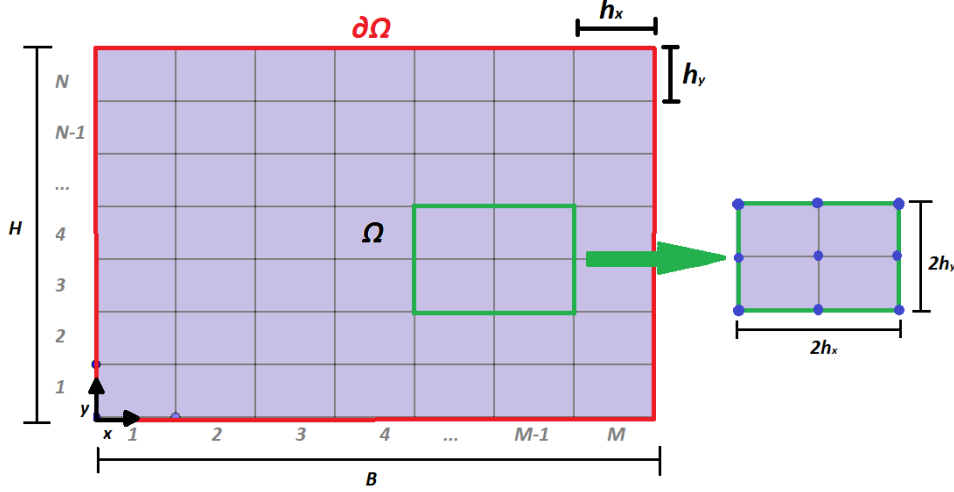
A representation of the discretization is shown next:



Figure 1.1: Discretization of region $\Omega$ and a patch of finite elements

where we can see how region $\Omega$ is discretized into a total of $M \cdot N$ finite elements, as well as a patch of finite elements (meaning the rectangle of size $2h_x \times 2h_y$) with its nodes (represented by the 9 blue dots on the corners of every rectangle of size $h_x \times h_y$), and a coordinate system useful to define all properties and points. The origin of the coordinate system is set at the bottom left corner of region $\Omega$. Notice how we have a total of $(M - 1) \cdot (N - 1)$ inner nodes (where $u$ is unknown) and $2 \cdot (M + N)$ boundary nodes (where we have boundary conditions).

Now we assume that variable $u$ is approximated by a linear combination of the test functions $v$, let's call it $\widetilde{u}$, therefore $u \approx \widetilde{u} = \sum u_i \cdot v_i$, where $u_i$ is the numerical value of function $\widetilde{u}$ at node $i$. From here equation (1.2) becomes:

$$\int_\Omega \left\langle \nabla \left( \sum u_i \cdot v_i \right), \nabla v_j \right\rangle d\Omega = \int_\Omega \left( f \cdot v_j \right) d\Omega \tag{1.3}$$

where $v_j$ is a test function applied at node $j$. Following equation (1.3), we apply a test function $v_j$ at every unknown node on region $\Omega$ in order to build a closed system of equations. To do so, it is necessary to compute the integrals in equation (1.3) accordingly.

Equation (1.3) is based on the global numbering of the nodes, however to easily compute the above integrals, we focus on each finite element individually, use a local numbering of the nodes on the finite elements, and use a local reference frame set at its bottom left corner. **The idea is to easily get general results based on local properties, and extrapolate them to the global problem in order to build the system of equations.**

For the local numbering, our reference for a finite element $R = [h_x \times h_y]$ will be the following:
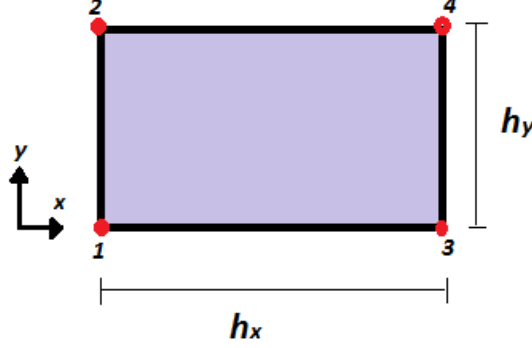


Figure 1.2: Reference for a finite element

where we locally numbered the nodes to easily identify them.

Based on the above local numbering and local reference frame, we can define bilinear functions $v_a$ applied at a local node $a$ as:

$$v_1 = \left(1 - \frac{x}{h_x}\right) \cdot \left(1 - \frac{y}{h_y}\right) \qquad v_2 = \left(1 - \frac{x}{h_x}\right) \cdot \frac{y}{h_y}$$

$$v_3 = \frac{x}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \qquad v_4 = \frac{x}{h_x} \cdot \frac{y}{h_y}$$

where $0 \leq x \leq h_x$ and $0 \leq y \leq h_y$.

# Concerning the term $\displaystyle\int_\Omega \langle \nabla v_i, \nabla v_j \rangle \, d\Omega$

For now, let us focus on the left hand side (LHS) of equation (1.3), which can be rewritten as:

$$\sum u_i \int_\Omega \langle \nabla v_i, \nabla v_j \rangle \, d\Omega$$

where both $i$ and $j$ correspond to nodes on the domain $\Omega$. The above represents a Matrix-vector product of the form $A \cdot \vec{u}$, where matrix $A$ has entries $a_{ji} = \int_\Omega \langle \nabla v_i, \nabla v_j \rangle \, d\Omega$, and $\vec{u}$ is a column vector containing the different values $u_i$.

For some given integers $i$ and $j$, because of the shape we chose for the test functions, most of the entries $\int_\Omega \langle \nabla v_i, \nabla v_j \rangle \, d\Omega$ are zero, specifically when nodes $i$ and $j$ do not belong to the same finite element. This is so since the gradients of test functions of further away nodes are zero at the region of integration (they might touch at their ends but still, the integral would be zero since they only touch along a line and not a finite area). The

geometrical interpretation of this is that the only non-zero entries correspond to nodes that are adjacent to each other.

To compute the above integral, we use local numbering of nodes $a$ and $b$. Therefore, for some given integers $a$ and $b$, we compute $\int_R \langle \nabla v_b, \nabla v_a \rangle \, d\Omega$.

As a final remark, it is convenient to express the following results in terms of a coefficient $\beta$ defined as $\beta = \frac{h_y}{h_x} + \frac{h_x}{h_y}$.

**For $a = 1$ and $b = 4$:**

$$
\begin{aligned}
\int_R \langle \nabla v_1, \nabla v_4 \rangle d\Omega &= \int_0^{h_y} \int_0^{h_x} \left[ \left( \begin{array}{c} -\frac{1}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \\ \left(1 - \frac{x}{h_x}\right) \cdot \left(-\frac{1}{h_y}\right) \end{array} \right) \cdot \left( \begin{array}{c} \frac{1}{h_x} \cdot \frac{y}{h_y} \\ \frac{x}{h_x} \cdot \frac{1}{h_y} \end{array} \right) \right] dx dy \\
&= \int_0^{h_y} \int_0^{h_x} \left[ -\left(\frac{1}{h_x}\right)^2 \cdot \left(1 - \frac{y}{h_y}\right) \frac{y}{h_y} - \left(1 - \frac{x}{h_x}\right) \frac{x}{h_x} \cdot \left(\frac{1}{h_y}\right)^2 \right] dx dy \\
&= -\frac{1}{6} \left( \frac{h_y}{h_x} + \frac{h_x}{h_y} \right) = -\frac{\beta}{6}
\end{aligned}
$$

This is the same result for whenever nodes $a$ and $b$ are in opposite corners along a diagonal of the finite element.

**For $a = 1$ and $b = 3$:**

$$
\begin{aligned}
\int_R \langle \nabla v_1, \nabla v_3 \rangle d\Omega &= \int_0^{h_y} \int_0^{h_x} \left( \begin{array}{c} -\frac{1}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \\ \left(1 - \frac{x}{h_x}\right) \cdot \left(-\frac{1}{h_y}\right) \end{array} \right) \cdot \left( \begin{array}{c} \frac{1}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \\ \frac{x}{h_x} \cdot \left(-\frac{1}{h_y}\right) \end{array} \right) dx dy \\
&= \int_0^{h_y} \int_0^{h_x} \left[ -\left(\frac{1}{h_x}\right)^2 \cdot \left(1 - \frac{y}{h_y}\right)^2 + \left(1 - \frac{x}{h_x}\right) \frac{x}{h_x} \cdot \left(\frac{1}{h_y}\right)^2 \right] dx dy \\
&= \frac{1}{6} \frac{h_x}{h_y} - \frac{1}{3} \frac{h_y}{h_x} = \frac{\beta}{6} - \frac{1}{2} \frac{h_y}{h_x}
\end{aligned}
$$

This is the same result for whenever nodes $a$ and $b$ are beside each other horizontally.

**For $a = 1$ and $b = 2$:**

We can use our previous result, which by symmetry gives:

$$
\int_R \langle \nabla v_1, \nabla v_2 \rangle d\Omega = \frac{1}{6} \frac{h_y}{h_x} - \frac{1}{3} \frac{h_x}{h_y} = \frac{\beta}{6} - \frac{1}{2} \frac{h_x}{h_y}
$$

This is the same result for whenever nodes $a$ and $b$ are one above the other vertically.

**For $a = 1$ and $b = 1$:**

$$\int_R \langle \nabla v_1, \nabla v_1 \rangle \, d\Omega = \int_0^{h_y} \int_0^{h_x} \left( \begin{array}{c} -\frac{1}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \\ \left(1 - \frac{x}{h_x}\right) \cdot \left(-\frac{1}{h_y}\right) \end{array} \right) \cdot \left( \begin{array}{c} -\frac{1}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) \\ \left(1 - \frac{x}{h_x}\right) \cdot \left(-\frac{1}{h_y}\right) \end{array} \right) dxdy$$

$$= \int_0^{h_y} \int_0^{h_x} \left[ \left(\frac{1}{h_x}\right)^2 \cdot \left(1 - \frac{y}{h_y}\right)^2 + \left(\frac{1}{h_y}\right)^2 \cdot \left(1 - \frac{x}{h_x}\right)^2 \right] dxdy$$

$$= \frac{1}{3} \left( \frac{h_y}{h_x} + \frac{h_x}{h_y} \right) = \frac{\beta}{3}$$

This is the same result for whenever $a = b$.

**Building local matrix $A$**

With this information, we can build the local matrix $A$, which is given by:

$$A_{local} = \frac{\beta}{6} \cdot \begin{pmatrix} 2 & 1 & 1 & -1 \\ 1 & 2 & -1 & 1 \\ 1 & -1 & 2 & 1 \\ -1 & 1 & 1 & 2 \end{pmatrix} - \frac{1}{2} \cdot \begin{pmatrix} 0 & \frac{h_x}{h_y} & \frac{h_y}{h_x} & 0 \\ \frac{h_x}{h_y} & 0 & 0 & \frac{h_y}{h_x} \\ \frac{h_y}{h_x} & 0 & 0 & \frac{h_x}{h_y} \\ 0 & \frac{h_y}{h_x} & \frac{h_x}{h_y} & 0 \end{pmatrix}$$

# Concerning the term $\displaystyle\int_\Omega (f \cdot v_j) \, d\Omega$:

Let's assume that $f$ is approximated by a linear combination of the test functions $v$, let's call it $\tilde{f}$, therefore $f \approx \tilde{f} = \sum f_i \cdot v_i$, where $f_i$ is the numerical value of function $\tilde{f}$ at node $i$. Correspondingly we get:

$$\int_\Omega (f \cdot v_j) \, d\Omega \approx \int_\Omega \left( \tilde{f} \cdot v_j \right) d\Omega = \sum f_i \int_\Omega (v_i \cdot v_j) \, d\Omega$$

The above represents a Matrix-vector product of the form $K \cdot \vec{f}$, where the matrix $K$ has entries $k_{ji} = \int_\Omega (v_i \cdot v_j) \, d\Omega$ (thus, it is the mass matrix), and $\vec{f}$ is a vector containing the different values $f_i$.

For some given integers $i$ and $j$, because of the shape we chose for the test functions, most of the entries $\int_\Omega (v_i \cdot v_j) \, d\Omega$ are zero, specifically when nodes $i$ and $j$ do not belong to the same finite element (for the same reason as with matrix $A$). Performing the integrals for some local nodes $a$ and $b$, we get:

**10**

**For $a = 1$ and $b = 4$:**

$$\int_R (v_1 \cdot v_4) \, d\Omega = \int_0^{h_y} \int_0^{h_x} \left[ \left( 1 - \frac{x}{h_x} \right) \cdot \left( 1 - \frac{y}{h_y} \right) \cdot \frac{x}{h_x} \cdot \frac{y}{h_y} \right] dxdy$$
$$= \frac{h_x \cdot h_y}{36}$$

This is the same result for whenever nodes $a$ and $b$ are in opposite corners along a diagonal.

**For $a = 1$ and $b = 3$:**

$$\int_R (v_1 \cdot v_3) \, d\Omega = \int_0^{h_y} \int_0^{h_x} \left[ \left( 1 - \frac{x}{h_x} \right) \cdot \left( 1 - \frac{y}{h_y} \right) \cdot \frac{x}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) \right] dxdy$$
$$= \frac{h_x \cdot h_y}{18}$$

This is the same result for whenever nodes $a$ and $b$ are one beside the other horizontally.

**For $a = 1$ and $b = 2$:**

By symmetry, we get:

$$\int_R (v_1 \cdot v_2) \, d\Omega = \frac{h_x \cdot h_y}{18}$$

This is the same result for whenever nodes $a$ and $b$ are one above the other vertically.

**For $a = 1$ and $b = 1$:**

$$\int_R (v_1 \cdot v_1) \, d\Omega = \int_0^{h_y} \int_0^{h_x} \left[ \left( 1 - \frac{x}{h_x} \right)^2 \cdot \left( 1 - \frac{y}{h_y} \right)^2 \right] dxdy$$
$$= \frac{h_x \cdot h_y}{9}$$

This is the same result for whenever $a = b$.

**Building the local matrix $K$**

With this information, we can build the local matrix $K$, which is given by:

$$K_{local} = \frac{h_x \cdot h_y}{36} \cdot \begin{pmatrix} 4 & 2 & 2 & 1 \\ 2 & 4 & 1 & 2 \\ 2 & 1 & 4 & 2 \\ 1 & 2 & 2 & 4 \end{pmatrix}$$

# Assembly of the system of equations

Having all necessary information, we can start building the system of equations which will have a total of $(M-1)\cdot(N-1)$ equalities.

For this particular case, the equalities obtained for every test function $v_j$ follow the same pattern. To properly write such pattern, we take as model the following patch of finite elements of size $2h_x \times 2h_y$, where we have labeled the nodes to identify them:
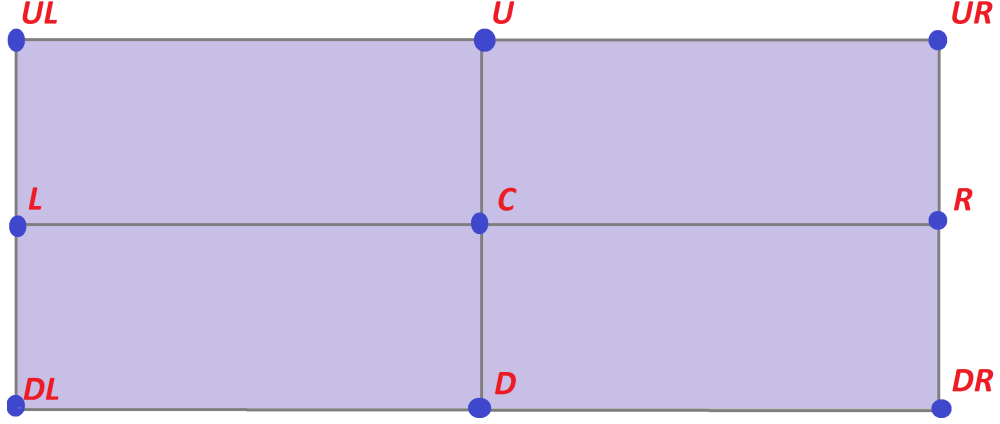


Figure 1.3: Reference for a patch of finite elements used for the assembly of system of equations

Hence, we can write for all inner nodes on region $\Omega$ that:

$$\frac{4\cdot\beta}{3}u_C - \frac{\beta}{6}\left(u_{DL} + u_{DR} + u_{UL} + u_{UR}\right) + \left(\frac{\beta}{3} - \frac{h_x}{h_y}\right)(u_U + u_D) + \left(\frac{\beta}{3} - \frac{h_y}{h_x}\right)(u_L + u_R)$$
$$= \frac{h_x \cdot h_y}{36} \cdot [f_{DL} + f_{DR} + f_{UL} + f_{UR} + 4\cdot(f_U + f_D + f_L + f_R) + 16\cdot f_C]$$

This forms one single equation of our system. It is important to mention that, depending on the test function $v_j$, on the LHS of the above equality we might have boundary nodes whose contribution has to be passed to the Right Hand Side (RHS) in order to close the system. In other words, this means that $u_i = g_i$ if node $i$ belongs to $\partial\Omega$, where $g_i$ is the Dirichlet boundary condition at node $i$.

Finally, solving the system will give a numerical approximation of variable $u$ at the unknown nodes.

# Practical implementation of the method

Now let us look at some examples. To test our model, we will do the following:

1. Set $u$ to any known function and apply it on some selected rectangular region $\Omega$. From that we will compute the right hand side $f$ and the Dirichlet boundary conditions $g$.

2. Then we will perform a discretization of region $\Omega$ and apply our finite element model. Solving the system of equations we are able to get the numerical approximation $\widetilde{u}$ of $u$.

3. With the actual function $u$ and its numerical approximation $\widetilde{u}$, we can compute the error of $\widetilde{u}$ based on some parameters.

What we are mostly interested with this procedure, is checking how does the error decrease when the discretization is enriched, or as how it's commonly called, the mesh is refined. To study this, since we have 2 discretization steps (meaning $h_x$ and $h_y$), we define a new equivalent length, which we will call $\bar{h}$:

$$\bar{h} = \sqrt{\frac{h_x^2 + h_y^2}{2}}$$

We expect that the error (computed for some given norm) follows the relationship:

$$error \approx C_{space} \cdot \bar{h}^q$$

where $C_{space}$ is a certain scalar constant and $q$ is the rate of decay of the error. We will compute the error in our approximation under 3 norms: the 1-norm, the 2-norm, and the infinity norm. The computation of them is given by the following formulas:

$$e_p = \left( \frac{\int_\Omega |u - \widetilde{u}|^p \, d\Omega}{\|\Omega\|} \right)^{1/p}$$

where for a finite element $R$, using the local numbering, we have that:

$$\int_R |u - \widetilde{u}|^p \approx \left( |u_1 - \widetilde{u}_1|^p + |u_2 - \widetilde{u}_2|^p + |u_3 - \widetilde{u}_3|^p + |u_4 - \widetilde{u}_4|^p \right) \cdot \frac{h_x \cdot h_y}{4}$$

The 1-norm corresponds to the case $p = 1$, the 2-norm to the case $p = 2$, and the infinity norm to the limit when $p \to \infty$, in which:

$$e_\infty = lim_{p \to \infty} \left( \frac{\int_\Omega |u - \widetilde{u}|^p \, d\Omega}{\|\Omega\|} \right)^{\frac{1}{p}} = max_\Omega |u - \widetilde{u}|$$

We will also be interested in computing the error for the gradient of $u$.

# Example

**Problem 1:**

- Consider region $\Omega$ as a rectangle of size $[5 \times 4]$ whose bottom left corner coincides with the origin of a coordinate system, and that the physical quantity under investigation is given by $u = 2 + sin\left(\frac{2 \cdot \pi}{5} \cdot x\right) \cdot sin\left(\frac{3 \cdot \pi}{4} \cdot y\right)$.

From function $u$ we can compute $f$ according to equation (1.1), which gives:

$$f = \left[ \left( \frac{2 \cdot \pi}{5} \right)^2 + \left( \frac{3 \cdot \pi}{4} \right)^2 \right] \cdot sin \left( \frac{2 \cdot \pi}{5} \cdot x \right) \cdot sin \left( \frac{3 \cdot \pi}{4} \cdot y \right)$$

On the other hand, $g = 2$ on $\partial \Omega$.

Performing a series of discretizations and applying our FEA, we can compute the errors and check how they decrease with respect to the discretization step. By doing so, we get the following results:



Figure 1.4: Errors of $\widetilde{u}$ vs $h$ for problem 1



Figure 1.5: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 1

Figure 1.6: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 1

Notice how the rate of decay is reported for every case at the upper left corner of each graph. We can appreciate quadratic convergence for function $u$ and linear convergence for $\nabla u$.

**14**

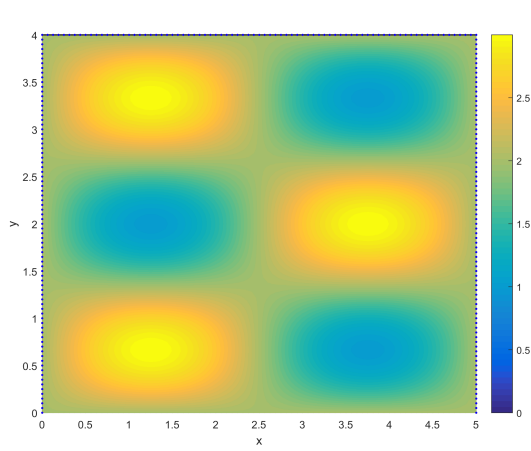As a plus, we also show a graph of $\widetilde{u}$ and its absolute error:



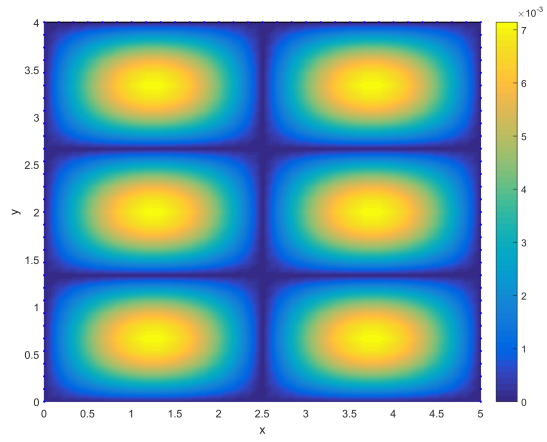Figure 1.7: Graph of $\widetilde{u}$ for problem 1



Figure 1.8: Graph of the absolute error $|u - \widetilde{u}|$ for problem 1

# Chapter 2

# Applying Dirichlet boundary conditions through Nitsche's method

So far, we managed to compute all necessary information on equation (1.3) to properly apply our FEA, however we did it by imposing boundary conditions in the "strong" form. The "strong" form means that the discrete solution takes the Dirichlet data at each node on $\partial\Omega$, meaning that the contribution of boundary nodes in the system of equations had to be passed towards the RHS in order to close the system.

However, to apply shifted boundary conditions, it is necessary to apply boundary conditions in another way, which we call the "weak" form. To do so, we use Nitsche's method which is a formulation that enforces the boundary conditions through a penalization technique. The formulation is the following, from equation (1.2) we get that:

$$\int_\Omega \langle \nabla u, \nabla v \rangle \, d\Omega - \int_{\partial\Omega} \left( v \cdot \frac{du}{dn} \right) ds - \int_{\partial\Omega} \left[ (u-g) \cdot \frac{dv}{dn} \right] ds$$
$$+ \frac{\gamma}{h} \int_{\partial\Omega} \left[ (u-g) \cdot v \right] ds = \int_\Omega (f \cdot v) \, d\Omega \quad (2.1)$$

where $\gamma$ is a penalty constant at least larger than 1, and $g$ are the Dirichlet boundary conditions. Equation (2.1) is correct for the exact solution $u$ of Problem (1.1) (since $u = g$ on $\partial\Omega$), however not for its numerical approximation $\widetilde{u}$. The idea is that by applying a penalization, equation (2.1) will force $\widetilde{u}$ to approach $g$ on $\partial\Omega$. Let us rewrite equation (2.1) as:

$$\int_\Omega \langle \nabla u, \nabla v \rangle \, d\Omega - \int_{\partial\Omega} \left( v \frac{du}{dn} \right) ds - \int_{\partial\Omega} \left( u \frac{dv}{dn} \right) ds + \frac{\gamma}{h} \int_{\partial\Omega} (uv) \, ds =$$
$$\int_\Omega (f \cdot v) \, d\Omega - \int_{\partial\Omega} \left( g \frac{dv}{dn} \right) ds + \frac{\gamma}{h} \int_{\partial\Omega} (gv) \, ds$$

By assuming that $u \approx \widetilde{u} = \sum u_i \cdot v_i$, that $g \approx \widetilde{g} = \sum g_i \cdot v_i$ where $\widetilde{g}$ is the numerical approximation of the boundary conditions, applying a test function $v_j$ to each node on region $\Omega$ now including those on $\partial\Omega$, and rearranging the LHS as to get a matrix vector product, we get to a system of equations of the shape:

$$\left[ A - \int_{\partial\Omega} \left[ v_j \cdot \frac{dv_i}{dn} \right] ds - \int_{\partial\Omega} \left[ v_i \cdot \frac{dv_j}{dn} \right] ds + \frac{\gamma}{h} \int_{\partial\Omega} (v_i \cdot v_j)\, ds \right] \cdot \vec{u} =$$

$$K \cdot \vec{f} + \left[ \frac{\gamma}{h} \int_{\partial\Omega} (v_i \cdot v_j)\, ds - \int_{\partial\Omega} \left[ v_i \cdot \frac{dv_j}{dn} \right] ds \right] \cdot \vec{g} \quad (2.2)$$

where the entries of matrices $A$ and $K$ were computed previously, and $\vec{g}$ is a column vector containing the different values of $g_i$. It is important to notice that in this new case, matrices $A$ and $K$ are larger than before, since now we will also have equations for all nodes on $\partial\Omega$. Nitsche's formulation also has the advantage that it makes the coefficient matrix symmetric and diagonally dominant.

## Concerning the term $\displaystyle\int_{\partial\Omega} \left[ v_j \cdot \frac{dv_i}{dn} \right] ds$

The term $\int_{\partial\Omega} \left[ v_j \cdot \frac{dv_i}{dn} \right] ds$ represents a matrix, which we will call $C$, with entries $c_{ji} = \int_{\partial\Omega} \left( v_j \cdot \frac{dv_i}{dn} \right) ds$.

For some given global nodes $i$ and $j$, because of the shape we chose for the test functions, most of the coefficients $\int_{\partial\Omega} \left( v_j \cdot \frac{dv_i}{dn} \right) ds$ are zero, specifically when nodes $j$ do not belong to the boundary $\partial\Omega$ and when nodes $i$ do not belong to finite elements being on the boundary $\partial\Omega$. In order to compute the above coefficients, we need to differentiate between the different boundaries of $\partial\Omega$, meaning if it's a left, right, upper or lower boundary. ***Following our reference finite element, we will compute the above integrals using a local numbering of the nodes***, therefore:

**For the left boundary we have that:**

$$\int_{\partial R} \left( v_1 \cdot \frac{dv_1}{dn} \right) ds = \int_0^{h_y} \left( 1 - \frac{y}{h_y} \right)^2 \cdot \frac{1}{h_x} dy = \frac{1}{3} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_1 \cdot \frac{dv_3}{dn} \right) ds = -\frac{1}{3} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_1 \cdot \frac{dv_2}{dn} \right) ds = \int_0^{h_y} \left( 1 - \frac{y}{h_y} \right) \cdot \frac{y}{h_y} \cdot \frac{1}{h_x} dy = \frac{1}{6} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_1 \cdot \frac{dv_4}{dn} \right) ds = -\frac{1}{6} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_2 \cdot \frac{dv_1}{dn} \right) ds = \int_0^{h_y} \frac{y}{h_y} \cdot \left( 1 - \frac{y}{h_y} \right) \cdot \frac{1}{h_x} dy = \frac{1}{6} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_2 \cdot \frac{dv_3}{dn} \right) ds = -\frac{1}{6} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_2 \cdot \frac{dv_2}{dn} \right) ds = \int_0^{h_y} \left( \frac{y}{h_y} \right)^2 \cdot \frac{1}{h_x} dy = \frac{1}{3} \cdot \frac{h_y}{h_x}$$

$$\int_{\partial R} \left( v_1 \cdot \frac{dv_4}{dn} \right) ds = -\frac{1}{3} \cdot \frac{h_y}{h_x}$$

With this information, we can build the local matrix $C$ for the left boundary:

$$C_{local}^{left} = \frac{1}{6} \frac{h_y}{h_x} \cdot \begin{pmatrix} 2 & 1 & -2 & -1 \\ 1 & 2 & -1 & -2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

**For the right boundary we have that:**

By symmetry, we can get the local $C$ matrix for the right boundary:

$$C_{local}^{right} = \frac{1}{6} \frac{h_y}{h_x} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -2 & -1 & 2 & 1 \\ -1 & -2 & 1 & 2 \end{pmatrix}$$

**For the upper boundary we have that:**

In a similar way, we can also use symmetry to get the local $C$ matrix for the upper boundary:

$$C_{local}^{upper} = \frac{1}{6} \frac{h_x}{h_y} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ -2 & 2 & -1 & 1 \\ 0 & 0 & 0 & 0 \\ -1 & 1 & -2 & 2 \end{pmatrix}$$

**For the lower boundary we have that:**

Also, symmetry works for the local matrix $C$ for the lower boundary:

$$C_{local}^{lower} = \frac{1}{6} \frac{h_x}{h_y} \cdot \begin{pmatrix} 2 & -2 & 1 & -1 \\ 0 & 0 & 0 & 0 \\ 1 & -1 & 2 & -2 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

# Concerning the term $\int_{\partial\Omega} \left[ v_i \cdot \frac{dv_j}{dn} \right] ds$

In our previous step, we computed matrix $C$ which corresponds to $\int_{\partial\Omega} \left[ v_j \cdot \frac{dv_i}{dn} \right] ds$. The term $\int_{\partial\Omega} \left[ v_i \cdot \frac{dv_j}{dn} \right] ds$ has nodes $i$ and $j$ interchanged with respect to matrix $C$, therefore this means that:

$$\int_{\partial\Omega} \left[ v_i \cdot \frac{dv_j}{dn} \right] ds = C^T$$

where $^T$ represents the transpose of a matrix.

# Concerning the term $\frac{\gamma}{h} \int_{\partial\Omega} (v_i \cdot v_j) \, ds$

For this term we just have to evaluate the integral of the product of the test functions over $\partial\Omega$, it will be zero for all test functions $v_i$ or $v_j$ associated with nodes that do not belong to $\partial\Omega$. It is to be mentioned that coefficient $h$ is the discretization step aligned with the direction of the normal vector $n$, so for the upper and lower boundary we have that $h = h_y$, and for the left and right boundary we have that $h = h_x$. The terms $\frac{\gamma}{h} \int_{\partial\Omega} (v_i \cdot v_j) \, ds$ are the entries $d_{ji}$ of a matrix which we will call $D$. ***To compute matrix $D$ we follow our reference finite element, and will compute the above integrals using a local numbering of the nodes***

To ease the computations, we follow the same procedure as with matrix $C$, first we compute the entries for a given boundary (either the left, right, lower or upper boundary) and use symmetry to compute all the other local matrices. For example, for the left boundary we have that:

$$\frac{\gamma}{h} \int_{\partial R} (v_1 \cdot v_1) \, ds = \frac{\gamma}{h_x} \int_0^{h_y} \left( 1 - \frac{y}{h_y} \right)^2 dy \quad = \frac{\gamma}{3} \cdot \frac{h_y}{h_x}$$

$$\frac{\gamma}{h} \int_{\partial R} (v_2 \cdot v_2) \, ds = \frac{\gamma}{h_x} \int_0^{h_y} \left( \frac{y}{h_y} \right)^2 dy \quad\quad = \frac{\gamma}{3} \cdot \frac{h_y}{h_x}$$

$$\frac{\gamma}{h} \int_{\partial R} (v_1 \cdot v_2) \, ds = \frac{\gamma}{h_x} \int_0^{h_y} \left( 1 - \frac{y}{h_y} \right) \cdot \frac{y}{h_y} dy = \frac{\gamma}{6} \cdot \frac{h_y}{h_x}$$

From here, we can write the local $D$ matrix for the left boundary:

$$D_{local}^{left} = \frac{\gamma}{6} \frac{h_y}{h_x} \cdot \begin{pmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Using symmetry, we can write the remaining local $D$ matrices for each boundary:

$$D_{local}^{right} = \frac{\gamma}{6}\frac{h_y}{h_x} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{pmatrix}$$

$$D_{local}^{upper} = \frac{\gamma}{6}\frac{h_x}{h_y} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 2 \end{pmatrix}$$

$$D_{local}^{lower} = \frac{\gamma}{6}\frac{h_x}{h_y} \cdot \begin{pmatrix} 2 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

# Example

**Problem 2:**

- Consider the same rectangle of size $[5 \times 4]$ as in the previous example, with a new physical quantity under investigation given by $u = \left(\frac{x+1}{5}\right)^{2.5} + \left(\frac{y+1}{4}\right)^{3.2}$. Now we are interested in checking how do the errors behave with respect to the penalization parameter $\gamma$. To study this, we choose a fixed discretization ($h_x = 0.1$ and $h_y = 0.08$) and compute the errors for several values of $\gamma$. By doing so we get that:



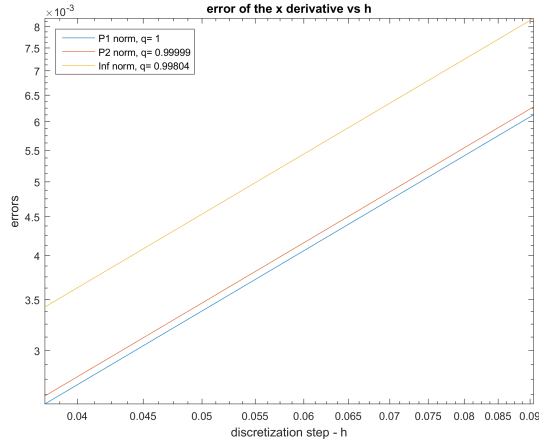Figure 2.1: Errors of $\widetilde{u}$ vs $\gamma$ for problem 2

Figure 2.2: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $\gamma$ for problem 2

Figure 2.3: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $\gamma$ for problem 2

Similar results are obtained with different values of the discretization parameters. Notice how as $\gamma$ increases, the errors decrease (as expected, since BC's are imposed "stronger"). What happens with Nitsche's method is that as $\gamma \longrightarrow \infty$, the numerical solution using Nitsche's method converges towards the numerical solution imposing boundary conditions in the "strong" form.

In this new case, we are also interested in checking how do the errors behave with respect to the discretization step while keeping $\gamma$ fixed. By performing several discretizations and computing the errors for $\gamma = 1000$ we get that:



Figure 2.4: Errors of $\widetilde{u}$ vs $h$ for problem 2

Figure 2.5: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 2
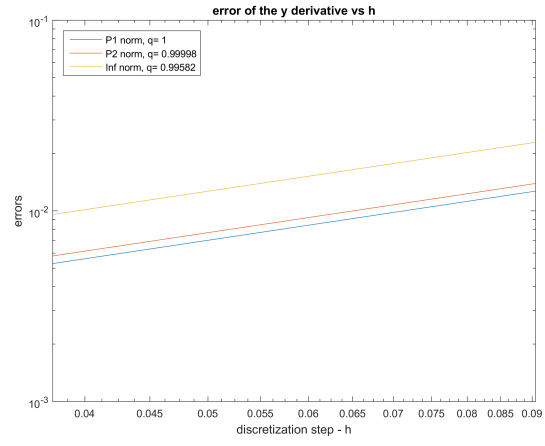


Figure 2.6: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 2

In our results, we can appreciate quadratic convergence for the error $\widetilde{u}$ and linear convergence for $\nabla \widetilde{u}$, as in Chapter 1.

# Chapter 3

# Applying shifted boundary conditions

So far we have applied non-shifted boundary conditions since our rectangular finite elements filled completely region $\Omega$, meaning that the boundary $\partial\Omega$ coincided perfectly with the boundaries of the finite elements. Now, we are interested in the case in which region $\Omega$ is not rectangular, therefore the rectangular finite elements contained in $\Omega$ will not fill region $\Omega$ completely.

To properly apply boundary conditions in this new case, we can use our previous results to shift the boundary conditions from the actual physical boundary $\partial\Omega$ towards the boundary of the region $\widetilde{\Omega}$ formed by the finite elements, which we will call the "numerical boundary" and denote by $\partial\widetilde{\Omega}$. A representation of this is shown in the following picture:
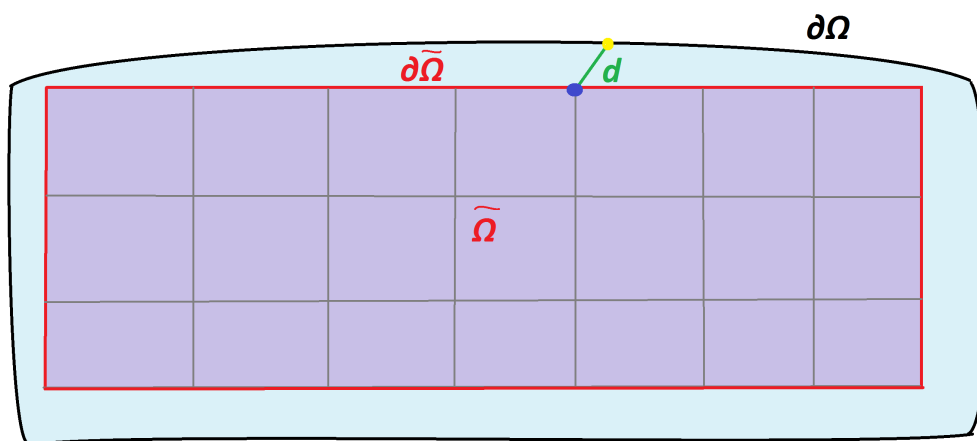


Figure 3.1: Comparison between region $\Omega$ and $\widetilde{\Omega}$

Here, the black curve represents the physical boundary $\partial\Omega$ (where we have information about boundary conditions), and the red curve represents the numerical boundary $\partial\widetilde{\Omega}$ (where we wish to have information about boundary conditions).

In order to translate the boundary conditions from the physical boundary, which we have called $g$, to boundary conditions on the numerical boundary, which we will call $g_{\partial\widetilde{\Omega}}$, we perform a first order Taylor approximation. This is expressed as:

$$g = g_{\partial\widetilde{\Omega}} + \langle \nabla u, d \rangle + r$$

where $d$ is a function which gives the distance vector between two points in boundaries $\partial\Omega$ and $\partial\widetilde{\Omega}$, and $r$ is the remainder. We can appreciate this in figure (3.1), where the blue dot represents a point in $\partial\widetilde{\Omega}$, the yellow dot represents a point in $\partial\Omega$, and $d$ is the vector that connects these two points.

The idea is that for every node on boundary $\partial\widetilde{\Omega}$, we find a suitable point on the physical boundary $\partial\Omega$ and shift the boundary condition of that point into the node. By doing so, new coefficients will be added into the system of equations, which will depend on the behaviour of function $d$, after neglecting the remainder. There is not a single way of doing this, since a node on boundary $\partial\widetilde{\Omega}$ can take information from any point on boundary $\partial\Omega$, however because we are doing a first order Taylor approximation it is best to take information from points which are the closest to the node.
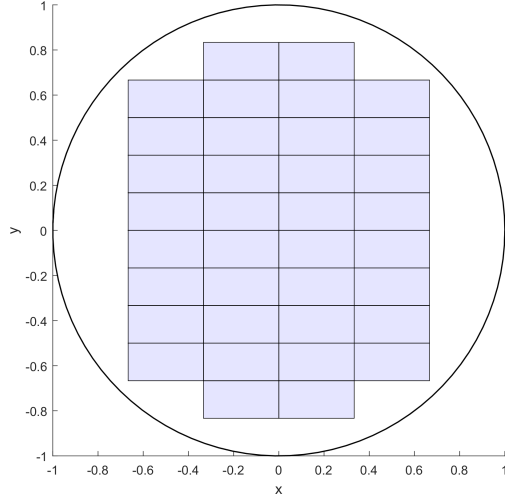
Therefore we define function $d$ such that it connects a node on $\partial\widetilde{\Omega}$ to its **closest points on the physical boundary** and give the distance vector between them. Since in our formulation we will have a set of line integrals over $\partial\widetilde{\Omega}$ which involve function $d$, it is upon us to find such function. A way to do this is to find the values of function $d$ at the different nodes belonging to $\partial\widetilde{\Omega}$, and then interpolating between them with linear functions.

On the following parts, we will assume that we already have the different values of function $d$ for every node $i$ on boundary $\partial\widetilde{\Omega}$, which we will call $d_i$. We will also assume that vectors $d_i$ have components $d = (d_{ix}, d_{iy})$. Interpolating with linear functions between the different $d_i$ vectors, in order to approximate the behaviour of function $d$ along boundary $\partial\widetilde{\Omega}$, we get that $d \approx \sum d_i \cdot v_i$.

It is important to emphasise that function $d$ is what allows to shift the boundary conditions, and the better the choice of it will give better results. In practice, the computation of our previously explained function $d$ is not difficult, since it is a geometrical problem and it can be resumed in 3 simple steps:
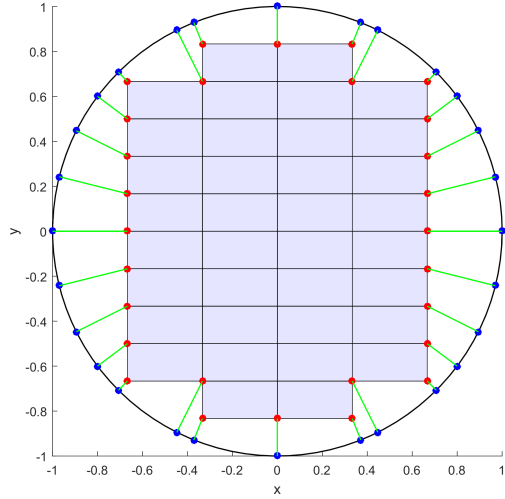
1. First, start with region $\Omega$ and choose a certain discretization. With such discretization, fill region $\Omega$ with as many finite elements as possible without crossing boundary $\partial\Omega$. As an example, refer to image (3.2) where region $\Omega$ is a circle, and the finite elements are the blue rectangles.

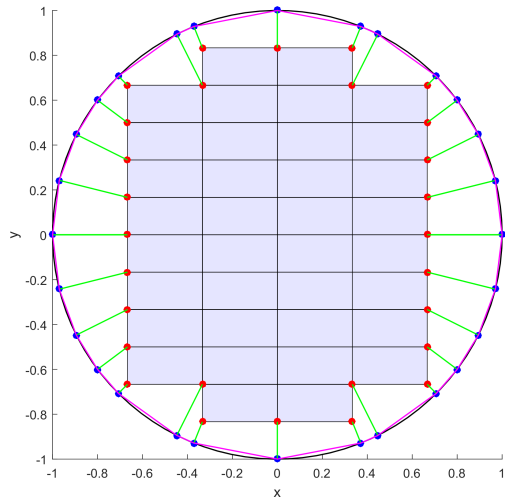Figure 3.2: Region $\Omega$ and a mesh of finite elements

2. Then identify the nodes belonging in $\partial\widetilde{\Omega}$ and compute their closest points on the physical boundary. This is seen on image (3.3) where the nodes of $\partial\widetilde{\Omega}$ are the red dots, and their closest points on $\partial\Omega$ are the blue dots, notice how each pair of blue and red dots are connected with a green line. These green lines are the different $d_i$ vectors.

Figure 3.3: Nodes on $\partial\widetilde{\Omega}$ (red) with their closest points on $\partial\Omega$ (blue)

3. Finally, we interpolate between the different $d_i$ vectors with linear functions, which approximates the behaviour of function $d$ across boundary $\partial\widetilde{\Omega}$. Since test functions $v_i$ become linear at $\partial\widetilde{\Omega}$, this means that $d \approx \sum d_i \cdot v_i$. This is visible on image (3.4), where the pink lines represent what function $d$ believes what the physical boundary is. The approximation of function $d$ towards $\partial\Omega$ becomes better as more finite elements are added into the discretization.

Figure 3.4: Behaviour of function $d$ along $\partial\widetilde{\Omega}$

**27**

It is important to notice that this whole procedure is to shift boundary conditions from the blue points into the red nodes. Now that we have completely explained how to obtain function $d$, we continue by writing equation (2.1) in the region $\widetilde{\Omega}$ rather than $\Omega$, and by substituting $g$ with $g_{\partial\widetilde{\Omega}}$, which gives:

$$
\int_{\widetilde{\Omega}} \langle \nabla u, \nabla v \rangle \, d\Omega - \int_{\partial\widetilde{\Omega}} \left( v \cdot \frac{du}{dn} \right) ds - \int_{\partial\widetilde{\Omega}} \left[ (u - g_{\partial\widetilde{\Omega}}) \cdot \frac{dv}{dn} \right] ds + \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} \left[ (u - g_{\partial\widetilde{\Omega}}) \cdot v \right] ds
$$
$$
= \int_{\widetilde{\Omega}} (f \cdot v) \, d\Omega
$$

Substituting for $g_{\partial\widetilde{\Omega}}$ as $g_{\partial\widetilde{\Omega}} = g - \langle \nabla u, d \rangle$, gives:

$$
\int_{\widetilde{\Omega}} \langle \nabla u, \nabla v \rangle \, d\Omega - \int_{\partial\widetilde{\Omega}} \left[ v \cdot \frac{du}{dn} \right] ds - \int_{\partial\widetilde{\Omega}} \left[ (u - g + \langle \nabla u, d \rangle) \cdot \frac{dv}{dn} \right] ds
$$
$$
+ \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} \left[ (u - g + \langle \nabla u, d \rangle) \cdot v \right] ds = \int_{\widetilde{\Omega}} (f \cdot v) \, d\Omega \quad (3.1)
$$

Approximating $u$, $g$, and $f$ as a linear combination of the test functions, and doing the same steps of the FEA as done in previous chapters, we only have 2 new terms appearing into equation (3.1) in comparison to our previous cases, since equation (3.1) can be rewritten as:

$$
\left( A - C - C^T + D \right) \cdot \vec{u} - \int_{\partial\widetilde{\Omega}} \left[ \langle \nabla \widetilde{u}, d \rangle \cdot \frac{dv_j}{dn} \right] ds + \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} \left[ \langle \nabla \widetilde{u}, d \rangle \cdot v_j \right] ds
$$
$$
= K \cdot \vec{f} + \left( D - C^T \right) \cdot \vec{g}
$$

Therefore we need to compute their contribution.

## Concerning the term $\int_{\partial\widetilde{\Omega}} \left[ \langle \nabla \widetilde{u}, d \rangle \cdot \frac{dv_j}{dn} \right] ds$:

The above term represents a Matrix-vector product of the form $E \cdot \vec{u}$, where matrix $E$ has entries $e_{ji} = \int_{\partial\widetilde{\Omega}} \left[ \langle \nabla v_i, d \rangle \cdot \frac{dv_j}{dn} \right] ds$.

For some given global nodes $i$ and $j$, because of the shape we chose for the test functions, most of the coefficients $\int_{\partial\widetilde{\Omega}} \left[ \langle \nabla v_i, d \rangle \cdot \frac{dv_j}{dn} \right] ds$ are zero, specifically when nodes $j$ do not belong to finite elements being on the boundary $\partial\widetilde{\Omega}$. In order to compute the above coefficients, we need to differentiate between the different boundaries of $\partial\widetilde{\Omega}$, meaning if it's a left, right, upper or lower boundary.

It is to be noted that, in this particular case, due to the extensive algebra that comes with substituting for $\widetilde{u} = \sum u_i \cdot v_i$, we will leave the term $\int_{\partial\widetilde{\Omega}} \left[ \langle \nabla \widetilde{u}, d \rangle \cdot \frac{dv_j}{dn} \right] ds$ as it is and work with it accordingly.

Also, we need to develop further the above term, since it can be rewritten as:

$$\int_{\partial\widetilde{\Omega}} \left[ \langle \nabla \widetilde{u}, d \rangle \cdot \frac{dv_j}{dn} \right] ds = \int_{\partial\widetilde{\Omega}} \left[ \frac{\partial\widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_j}{dn} \right] ds + \int_{\partial\widetilde{\Omega}} \left[ \frac{\partial\widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_j}{dn} \right] ds$$

For all upcoming computations, we will follow our reference finite element and use a **local numbering** of the nodes. Therefore we can write that:

$$\frac{\partial\widetilde{u}}{\partial x} = \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y}$$

$$\frac{\partial\widetilde{u}}{\partial y} = \frac{(u_2 - u_1)}{h_y} \cdot \left( 1 - \frac{x}{h_x} \right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}$$

**For the left boundary:**

- For the $x$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial\widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_1}{dn} \right] ds = \int_0^{h_y} \left[ \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y} \right] \cdot \ldots$$

$$\ldots \cdot \left[ d_{x1} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{x2} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \left( 1 - \frac{y}{h_y} \right) dy$$

$$= \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( 3 \cdot d_{x1} + d_{x2} \right) + (u_4 - u_2) \left( d_{x1} + d_{x2} \right) \right]$$

$$\int_{\partial R} \left[ \frac{\partial\widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_3}{dn} \right] ds = -\frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( 3 \cdot d_{x1} + d_{x2} \right) + (u_4 - u_2) \left( d_{x1} + d_{x2} \right) \right]$$

$$\int_{\partial R} \left[ \frac{\partial\widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_2}{dn} \right] ds = \int_0^{h_y} \left[ \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y} \right] \cdot \ldots$$

$$\ldots \left[ d_{x1} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{x2} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \frac{y}{h_y} dy$$

$$= \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( d_{x1} + d_{x2} \right) + (u_4 - u_2) \left( d_{x1} + 3 \cdot d_{x2} \right) \right]$$

$$\int_{\partial R} \left[ \frac{\partial\widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_4}{dn} \right] ds = -\frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( d_{x1} + d_{x2} \right) + (u_4 - u_2) \left( d_{x1} + 3 \cdot d_{x2} \right) \right]$$

- For the $y$ component we have that:

**29**

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_1}{dn} \right] ds = \int_0^{h_y} \frac{(u_2 - u_1)}{h_y} \cdot \left[ d_{y1} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{y2} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \left( 1 - \frac{y}{h_y} \right) dy$$

$$= \frac{(u_2 - u_1)}{6 \cdot h_x} \cdot (2 \cdot d_{y1} + d_{y2})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_3}{dn} \right] ds = - \frac{(u_2 - u_1)}{6 \cdot h_x} \cdot (2 \cdot d_{y1} + d_{y2})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_2}{dn} \right] ds = \int_0^{h_y} \frac{(u_2 - u_1)}{h_y} \cdot \left[ d_{y1} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{y2} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \cdot \frac{y}{h_y} dy$$

$$= \frac{(u_2 - u_1)}{6 \cdot h_x} \cdot (d_{y1} + 2 \cdot d_{y2})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_4}{dn} \right] ds = - \frac{(u_2 - u_1)}{6 \cdot h_x} \cdot (d_{y1} + 2 \cdot d_{y2})$$

With this information we can write the local matrix $E$ for the left boundary:

$$E_{local}^{left} = \frac{1}{12} \frac{h_y}{h_x^2} \cdot \begin{pmatrix} -(3 \cdot d_{x1} + d_{x2}) & -(d_{x1} + d_{x2}) & 3 \cdot d_{x1} + d_{x2} & d_{x1} + d_{x2} \\ -(d_{x1} + d_{x2}) & -(d_{x1} + 3 \cdot d_{x2}) & d_{x1} + d_{x2} & d_{x1} + 3 \cdot d_{x2} \\ 3 \cdot d_{x1} + d_{x2} & d_{x1} + d_{x2} & -(3 \cdot d_{x1} + d_{x2}) & -(d_{x1} + d_{x2}) \\ d_{x1} + d_{x2} & d_{x1} + 3 \cdot d_{x2} & -(d_{x1} + d_{x2}) & -(d_{x1} + 3 \cdot d_{x2}) \end{pmatrix}$$

$$\ldots + \frac{1}{6 \cdot h_x} \cdot \begin{pmatrix} -(2 \cdot d_{y1} + d_{y2}) & 2 \cdot d_{y1} + d_{y2} & 0 & 0 \\ -(d_{y1} + 2 \cdot d_{y2}) & d_{y1} + 2 \cdot d_{y2} & 0 & 0 \\ 2 \cdot d_{y1} + d_{y2} & -(2 \cdot d_{y1} + d_{y2}) & 0 & 0 \\ d_{y1} + 2 \cdot d_{y2} & -(d_{y1} + 2 \cdot d_{y2}) & 0 & 0 \end{pmatrix}$$

**For the right boundary we have that:**

- For the $x$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_3}{dn} \right] ds = \int_0^{h_y} \left[ \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y} \right] \cdot \ldots$$

$$\ldots \left[ d_{x3} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{x4} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \left( 1 - \frac{y}{h_y} \right) dy$$

$$= \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( 3 \cdot d_{x3} + d_{x4} \right) + (u_4 - u_2) \left( d_{x3} + d_{x4} \right) \right]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_1}{dn} \right] ds = - \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1) \left( 3 \cdot d_{x3} + d_{x4} \right) + (u_4 - u_2) \left( d_{x3} + d_{x4} \right) \right]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_4}{dn} \right] ds = \int_0^{h_y} \left[ \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y} \right] \cdot \dots$$

$$\dots \left[ d_{x3} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{x4} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \frac{y}{h_y} dy$$

$$= \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1)(d_{x3} + d_{x4}) + (u_4 - u_2)(d_{x3} + 3 \cdot d_{x4}) \right]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_2}{dn} \right] ds = - \frac{1}{12} \frac{h_y}{h_x^2} \cdot \left[ (u_3 - u_1)(d_{x3} + d_{x4}) + (u_4 - u_2)(d_{x3} + 3 \cdot d_{x4}) \right]$$

- For the $y$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_3}{dn} \right] ds = \int_0^{h_y} \frac{(u_4 - u_3)}{h_y} \cdot \left[ d_{y3} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{y4} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \left( 1 - \frac{y}{h_y} \right) dy$$

$$= \frac{(u_4 - u_3)}{6 \cdot h_x} \cdot (2 \cdot d_{y3} + d_{y4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_1}{dn} \right] ds = - \frac{(u_4 - u_3)}{6 \cdot h_x} \cdot (2 \cdot d_{y3} + d_{y4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_4}{dn} \right] ds = \int_0^{h_y} \frac{(u_4 - u_3)}{h_y} \cdot \left[ d_{y3} \cdot \left( 1 - \frac{y}{h_y} \right) + d_{y4} \cdot \frac{y}{h_y} \right] \cdot \frac{1}{h_x} \cdot \frac{y}{h_y} dy$$

$$= \frac{(u_4 - u_3)}{6 \cdot h_x} \cdot (d_{y3} + 2 \cdot d_{y4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_2}{dn} \right] ds = - \frac{(u_4 - u_3)}{6 \cdot h_x} \cdot (d_{y3} + 2 \cdot d_{y4})$$

With this information we can write the local matrix $E$ for the right boundary:

$$E_{local}^{right} = \frac{1}{12} \frac{h_y}{h_x^2} \cdot \begin{pmatrix} 3 \cdot d_{x3} + d_{x4} & d_{x3} + d_{x4} & -(3 \cdot d_{x3} + d_{x4}) & -(d_{x3} + d_{x4}) \\ d_{x3} + d_{x4} & d_{x3} + 3 \cdot d_{x4} & -(d_{x3} + d_{x4}) & -(d_{x3} + 3 \cdot d_{x4}) \\ -(3 \cdot d_{x3} + d_{x4}) & -(d_{x3} + d_{x4}) & 3 \cdot d_{x3} + d_{x4} & d_{x3} + d_{x4} \\ -(d_{x3} + d_{x4}) & -(d_{x3} + 3 \cdot d_{x4}) & d_{x3} + d_{x4} & d_{x3} + 3 \cdot d_{x4} \end{pmatrix}$$

$$\dots + \frac{1}{6 \cdot h_x} \cdot \begin{pmatrix} 0 & 0 & 2 \cdot d_{y3} + d_{y4} & -(2 \cdot d_{y3} + d_{y4}) \\ 0 & 0 & d_{y3} + 2 \cdot d_{y4} & -(d_{y3} + 2 \cdot d_{y4}) \\ 0 & 0 & -(2 \cdot d_{y3} + d_{y4}) & 2 \cdot d_{y3} + d_{y4} \\ 0 & 0 & -(d_{y3} + 2 \cdot d_{y4}) & d_{y3} + 2 \cdot d_{y4} \end{pmatrix}$$

**For the upper boundary:**

- For the $x$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_2}{dn} \right] ds = \int_0^{h_y} \frac{(u_4 - u_2)}{h_x} \cdot \left[ d_{x2} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{x4} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \left( 1 - \frac{x}{h_x} \right) dx$$

$$= \frac{(u_4 - u_2)}{6 \cdot h_y} \cdot (2 \cdot d_{x2} + d_{x4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_1}{dn} \right] ds = - \frac{(u_4 - u_2)}{6 \cdot h_y} \cdot (2 \cdot d_{x2} + d_{x4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_4}{dn} \right] ds = \int_0^{h_y} \frac{(u_4 - u_2)}{h_x} \cdot \left[ d_{x2} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{x4} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \frac{x}{h_x} dx$$

$$= \frac{(u_4 - u_2)}{6 \cdot h_y} \cdot (d_{x2} + 2 \cdot d_{x4})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_3}{dn} \right] ds = - \frac{(u_4 - u_2)}{6 \cdot h_y} \cdot (d_{x2} + 2 \cdot d_{x4})$$

- For the $y$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_2}{dn} \right] ds = \int_0^{h_x} \left[ \frac{(u_2 - u_1)}{h_y} \cdot \left( 1 - \frac{x}{h_x} \right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x} \right] \cdot \ldots$$

$$\ldots \left[ d_{y2} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{y4} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \left( 1 - \frac{x}{h_x} \right) dx$$

$$= \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(3 \cdot d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + d_{y4})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_1}{dn} \right] ds = - \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(3 \cdot d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + d_{y4})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_4}{dn} \right] ds = \int_0^{h_x} \left[ \frac{(u_2 - u_1)}{h_y} \cdot \left( 1 - \frac{x}{h_x} \right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x} \right] \cdot \ldots$$

$$\ldots \left[ d_{y2} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{y4} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \frac{x}{h_x} dx$$

$$= \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + 3 \cdot d_{y4})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_3}{dn} \right] ds = - \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + 3 \cdot d_{y4})]$$

**32**

With this information we can write the local matrix $E$ for the upper boundary:

$$E_{local}^{upper} = \frac{1}{6 \cdot h_y} \cdot \begin{pmatrix} 0 & 2 \cdot d_{x2} + d_{x4} & 0 & -(2 \cdot d_{x2} + d_{x4}) \\ 0 & -(2 \cdot d_{x2} + d_{x4}) & 0 & 2 \cdot d_{x2} + d_{x4} \\ 0 & d_{x2} + 2 \cdot d_{x4} & 0 & -(d_{x2} + 2 \cdot d_{x4}) \\ 0 & -(d_{x2} + 2 \cdot d_{x4}) & 0 & d_{x2} + 2 \cdot d_{x4} \end{pmatrix}$$

$$\dots + \frac{1}{12} \frac{h_x}{h_y^2} \cdot \begin{pmatrix} 3 \cdot d_{y2} + d_{y4} & -(3 \cdot d_{y2} + d_{y4}) & d_{y2} + d_{y4} & -(d_{y2} + d_{y4}) \\ -(3 \cdot d_{y2} + d_{y4}) & 3 \cdot d_{y2} + d_{y4} & -(d_{y2} + d_{y4}) & d_{y2} + d_{y4} \\ d_{y2} + d_{y4} & -(d_{y2} + d_{y4}) & d_{y2} + 3 \cdot d_{y4} & -(d_{y2} + 3 \cdot d_{y4}) \\ -(d_{y2} + d_{y4}) & d_{y2} + d_{y4} & -(d_{y2} + 3 \cdot d_{y4}) & d_{y2} + 3 \cdot d_{y4} \end{pmatrix}$$

**For the lower boundary:**

- For the $x$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_1}{dn} \right] ds = \int_0^{h_y} \frac{(u_3 - u_1)}{h_x} \cdot \left[ d_{x1} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{x3} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \left( 1 - \frac{x}{h_x} \right) dx$$

$$= \frac{(u_3 - u_1)}{6 \cdot h_y} \cdot (2 \cdot d_{x1} + d_{x3})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_2}{dn} \right] ds = - \frac{(u_3 - u_1)}{6 \cdot h_y} \cdot (2 \cdot d_{x1} + d_{x3})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_3}{dn} \right] ds = \int_0^{h_y} \frac{(u_3 - u_1)}{h_x} \cdot \left[ d_{x1} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{x3} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \frac{x}{h_x} dx$$

$$= \frac{(u_3 - u_1)}{6 \cdot h_y} \cdot (d_{x1} + 2 \cdot d_{x3})$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot \frac{dv_4}{dn} \right] ds = - \frac{(u_3 - u_1)}{6 \cdot h_y} \cdot (d_{x1} + 2 \cdot d_{x3})$$

- For the $y$ component we have that:

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_1}{dn} \right] ds = \int_0^{h_x} \left[ \frac{(u_2 - u_1)}{h_y} \cdot \left( 1 - \frac{x}{h_x} \right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x} \right] \cdot \dots$$

$$\dots \left[ d_{y1} \cdot \left( 1 - \frac{x}{h_x} \right) + d_{y3} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \left( 1 - \frac{x}{h_x} \right) dx$$

$$= \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(3 \cdot d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + d_{y3})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_2}{dn} \right] ds = - \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(3 \cdot d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + d_{y3})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_3}{dn} \right] ds = \int_0^{h_x} \left[ \frac{(u_2 - u_1)}{h_y} \cdot \left(1 - \frac{x}{h_x}\right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x} \right] \cdot \ldots$$

$$\ldots \left[ d_{y1} \cdot \left(1 - \frac{x}{h_x}\right) + d_{y3} \cdot \frac{x}{h_x} \right] \cdot \frac{1}{h_y} \frac{x}{h_x} dx$$

$$= \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + 3 \cdot d_{y3})]$$

$$\int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot \frac{dv_4}{dn} \right] ds = - \frac{1}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + 3 \cdot d_{y3})]$$

With this information we can write the local matrix $E$ for the lower boundary:

$$E_{local}^{lower} = \frac{1}{6 \cdot h_y} \cdot \begin{pmatrix} -(2 \cdot d_{x1} + d_{x3}) & 0 & 2 \cdot d_{x1} + d_{x3} & 0 \\ 2 \cdot d_{x1} + d_{x3} & 0 & -(2 \cdot d_{x1} + d_{x3}) & 0 \\ -(d_{x1} + 2 \cdot d_{x3}) & 0 & d_{x1} + 2 \cdot d_{x3} & 0 \\ d_{x1} + 2 \cdot d_{x3} & 0 & -(d_{x1} + 2 \cdot d_{x3}) & 0 \end{pmatrix}$$

$$\ldots + \frac{1}{12} \frac{h_x}{h_y^2} \cdot \begin{pmatrix} -(3 \cdot d_{y1} + d_{y3}) & 3 \cdot d_{y1} + d_{y3} & -(d_{y1} + d_{y3}) & d_{y1} + d_{y3} \\ 3 \cdot d_{y1} + d_{y3} & -(3 \cdot d_{y1} + d_{y3}) & d_{y1} + d_{y3} & -(d_{y1} + d_{y3}) \\ -(d_{y1} + d_{y3}) & d_{y1} + d_{y3} & -(d_{y1} + 3 \cdot d_{y3}) & d_{y1} + 3 \cdot d_{y3} \\ d_{y1} + d_{y3} & -(d_{y1} + d_{y3}) & d_{y1} + 3 \cdot d_{y3} & -(d_{y1} + 3 \cdot d_{y3}) \end{pmatrix}$$

# Concerning the term $\frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} [\langle \nabla \widetilde{u}, d \rangle \cdot v_j] \, ds$:

The above represents a Matrix-vector product of the form $L \cdot \vec{u}$, where matrix $L$ has entries $l_{ji} = \frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} [\langle \nabla v_i, d \rangle \cdot v_j] \, ds$.

For some given global nodes $i$ and $j$, because of the shape we chose for the test functions, most of the coefficients $\frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} [\langle \nabla v_i, d \rangle \cdot v_j] \, ds$ are zero, specifically when nodes $j$ do not belong to the boundary $\partial \widetilde{\Omega}$. In order to compute the above coefficients, we need to differentiate between the different boundaries of $\partial \widetilde{\Omega}$, meaning if it's a left, right, upper or lower boundary.

Due to the extensive algebra that comes with substituting for $\widetilde{u} = \sum u_i \cdot v_i$, we will leave it as it is and work with it accordingly following a ***local numbering*** of the nodes. Also:

$$\frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} [\langle \nabla \widetilde{u}, d \rangle \cdot v_j] \, ds = \frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot v_j \right] ds + \frac{\gamma}{h} \int_{\partial \widetilde{\Omega}} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_j \right] ds$$

Using the results obtained from the computation of matrix $E$, we can easily find all coefficients for the local $L$ matrices:

**34**

**For the left boundary:**

- For the $x$ component we have that:

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial x}\cdot d_x \cdot v_1\right]ds = \frac{\gamma}{12}\frac{h_y}{h_x^2}\cdot\left[(u_3-u_1)(3\cdot d_{x1}+d_{x2})+(u_4-u_2)(d_{x1}+d_{x2})\right]$$

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial x}\cdot d_x \cdot v_2\right]ds = \frac{\gamma}{12}\frac{h_y}{h_x^2}\cdot\left[(u_3-u_1)(d_{x1}+d_{x2})+(u_4-u_2)(d_{x1}+3\cdot d_{x2})\right]$$

- For the $y$ component we have that:

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial y}\cdot d_y \cdot v_1\right]ds = \frac{\gamma}{6\cdot h_x}\cdot(u_2-u_1)\cdot(2\cdot d_{y1}+d_{y2})$$

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial y}\cdot d_y \cdot v_2\right]ds = \frac{\gamma}{6\cdot h_x}\cdot(u_2-u_1)\cdot(d_{y1}+2\cdot d_{y2})$$

$$L_{local}^{left} = \frac{\gamma}{12}\frac{h_y}{h_x^2}\cdot\begin{pmatrix}-(3\cdot d_{x1}+d_{x2}) & -(d_{x1}+d_{x2}) & 3\cdot d_{x1}+d_{x2} & d_{x1}+d_{x2}\\ -(d_{x1}+d_{x2}) & -(d_{x1}+3\cdot d_{x2}) & d_{x1}+d_{x2} & d_{x1}+3\cdot d_{x2}\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0\end{pmatrix}$$

$$\dots + \frac{\gamma}{6\cdot h_x}\cdot\begin{pmatrix}-(2\cdot d_{y1}+d_{y2}) & 2\cdot d_{y1}+d_{y2} & 0 & 0\\ -(d_{y1}+2\cdot d_{y2}) & d_{y1}+2\cdot d_{y2} & 0 & 0\\ 0 & 0 & 0 & 0\\ 0 & 0 & 0 & 0\end{pmatrix}$$

**For the right boundary:**

- For the $x$ component we have that:

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial x}\cdot d_x \cdot v_3\right]ds = \frac{\gamma}{12}\frac{h_y}{h_x^2}\cdot\left[(u_3-u_1)(3\cdot d_{x3}+d_{x4})+(u_4-u_2)(d_{x3}+d_{x4})\right]$$

$$\frac{\gamma}{h}\int_{\partial R}\left[\frac{\partial \widetilde{u}}{\partial x}\cdot d_x \cdot v_4\right]ds = \frac{1}{12}\frac{h_y}{h_x^2}\cdot\left[(u_3-u_1)(d_{x3}+d_{x4})+(u_4-u_2)(d_{x3}+3\cdot d_{x4})\right]$$

- For the $y$ component we have that:

**35**

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_3 \right] ds = \frac{\gamma}{6 \cdot h_x} \cdot (u_4 - u_3) \cdot (2 \cdot d_{y3} + d_{y4})$$

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_4 \right] ds = \frac{\gamma}{6 \cdot h_x} \cdot (u_4 - u_3) \cdot (d_{y3} + 2 \cdot d_{y4})$$

$$L_{local}^{right} = \frac{\gamma}{12} \frac{h_y}{h_x^2} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -(3 \cdot d_{x3} + d_{x4}) & -(d_{x3} + d_{x4}) & 3 \cdot d_{x3} + d_{x4} & d_{x3} + d_{x4} \\ -(d_{x3} + d_{x4}) & -(d_{x3} + 3 \cdot d_{x4}) & d_{x3} + d_{x4} & d_{x3} + 3 \cdot d_{x4} \end{pmatrix}$$

$$\dots + \frac{\gamma}{6 \cdot h_x} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -(2 \cdot d_{y3} + d_{y4}) & 2 \cdot d_{y3} + d_{y4} \\ 0 & 0 & -(d_{y3} + 2 \cdot d_{y4}) & d_{y3} + 2 \cdot d_{y4} \end{pmatrix}$$

**For the upper boundary:**

- For the $x$ component we have that:

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot v_2 \right] ds = \frac{\gamma}{6 \cdot h_y} \cdot (u_4 - u_2) \cdot (2 \cdot d_{x2} + d_{x4})$$

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot v_4 \right] ds = \frac{\gamma}{6 \cdot h_y} \cdot (u_4 - u_2) \cdot (d_{x2} + 2 \cdot d_{x4})$$

- For the $y$ component we have that:

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_2 \right] ds = \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(3 \cdot d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + d_{y4})]$$

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_4 \right] ds = \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot [(u_2 - u_1)(d_{y2} + d_{y4}) + (u_4 - u_3)(d_{y2} + 3 \cdot d_{y4})]$$

$$L_{local}^{upper} = \frac{\gamma}{6 \cdot h_y} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -(2 \cdot d_{x2} + d_{x4}) & 0 & 2 \cdot d_{x2} + d_{x4} \\ 0 & 0 & 0 & 0 \\ 0 & -(d_{x2} + 2 \cdot d_{x4}) & 0 & d_{x2} + 2 \cdot d_{x4} \end{pmatrix}$$

$$\dots + \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot \begin{pmatrix} 0 & 0 & 0 & 0 \\ -(3 \cdot d_{y2} + d_{y4}) & 3 \cdot d_{y2} + d_{y4} & -(d_{y2} + d_{y4}) & d_{y2} + d_{y4} \\ 0 & 0 & 0 & 0 \\ -(d_{y2} + d_{y4}) & d_{y2} + d_{y4} & -(d_{y2} + 3 \cdot d_{y4}) & d_{y2} + 3 \cdot d_{y4} \end{pmatrix}$$

**36**

**For the lower boundary:**

- For the $x$ component we have that:

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot v_1 \right] ds = \frac{\gamma}{6 \cdot h_y} \cdot (u_3 - u_1) \cdot (2 \cdot d_{x1} + d_{x3})$$

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot d_x \cdot v_3 \right] ds = \frac{\gamma}{6 \cdot h_y} \cdot (u_3 - u_1) \cdot (d_{x1} + 2 \cdot d_{x3})$$

- For the $y$ component we have that:

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_1 \right] ds = \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot \left[ (u_2 - u_1)(3 \cdot d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + d_{y3}) \right]$$

$$\frac{\gamma}{h} \int_{\partial R} \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot d_y \cdot v_3 \right] ds = \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot \left[ (u_2 - u_1)(d_{y1} + d_{y3}) + (u_4 - u_3)(d_{y1} + 3 \cdot d_{y3}) \right]$$

$$L_{local}^{lower} = \frac{\gamma}{6 \cdot h_y} \cdot \begin{pmatrix} -(2 \cdot d_{x1} + d_{x3}) & 0 & 2 \cdot d_{x1} + d_{x3} & 0 \\ 0 & 0 & 0 & 0 \\ -(d_{x1} + 2 \cdot d_{x3}) & 0 & d_{x1} + 2 \cdot d_{x3} & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\cdots + \frac{\gamma}{12} \frac{h_x}{h_y^2} \cdot \begin{pmatrix} -(3 \cdot d_{y1} + d_{y3}) & 3 \cdot d_{y1} + d_{y3} & -(d_{y1} + d_{y3}) & d_{y1} + d_{y3} \\ 0 & 0 & 0 & 0 \\ -(d_{y1} + d_{y3}) & d_{y1} + d_{y3} & -(d_{y1} + 3 \cdot d_{y3}) & d_{y1} + 3 \cdot d_{y3} \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

# Examples

**Problem 3:**

- Consider region $\Omega$ as the set of points $(x, y)$ such that $0 \leq y \leq 4 + \sin(\pi \cdot x)$ and $0 \leq x \leq 5$, and that the physical quantity under investigation is given by $u = 2 + sin\left(\frac{2 \cdot \pi}{5} \cdot x\right) \cdot sin\left(\frac{3 \cdot \pi}{5} \cdot y\right)$, namely the same quantity in Problem 1. We wish to apply our FEA to check how do the errors on $\widetilde{u}$ and $\nabla \widetilde{u}$ decrease when the mesh is refined.

As explained previously, now region $\Omega$ can not be completely filled with finite elements, the idea in order to get the best possible approximation of the physical variable $u$, will be to fill region $\Omega$ with as many finite elements as possible without crossing the physical boundary $\partial \Omega$. This will give a mesh of finite elements with a set of nodes, for which we need to calculate vector $d$ for the nodes on the numerical boundary $\partial \widetilde{\Omega}$, and from there

we can perform our FEA to compute the errors in our approximation. By doing so, we get the following results:
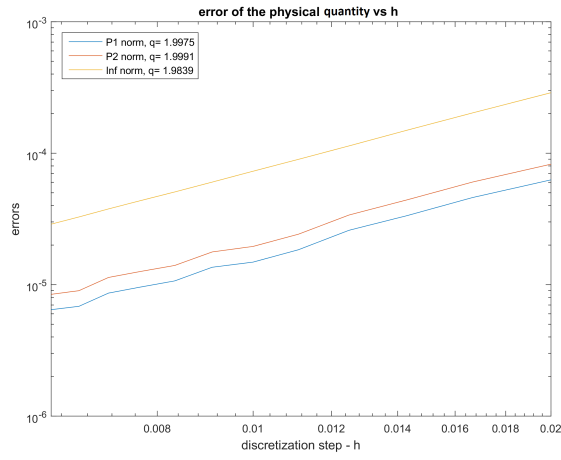


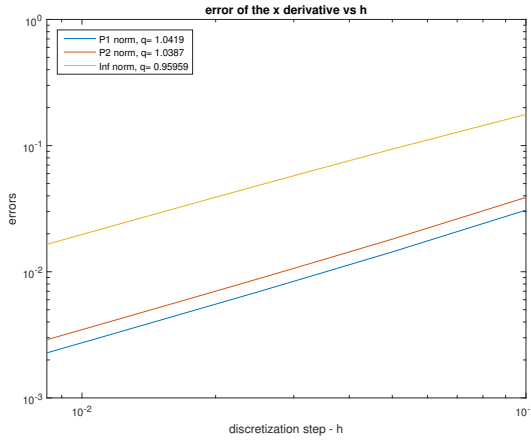Figure 3.5: Errors of $\widetilde{u}$ vs $h$ for problem 3



Figure 3.6: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 3
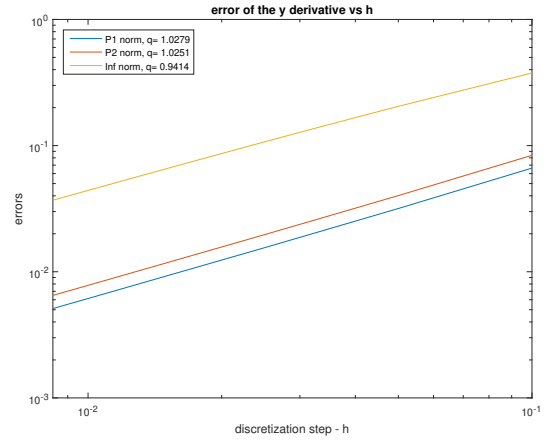


Figure 3.7: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 3
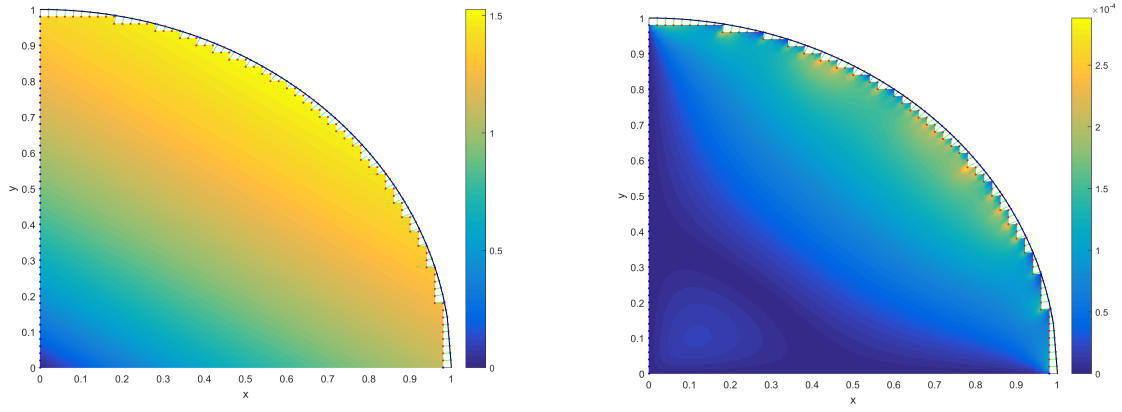
We can appreciate quadratic convergence, or at least super-linear for the Infinity norm, for $\widetilde{u}$, and linear convergence for $\nabla \widetilde{u}$. Graphs of the numerical solution and their errors can be appreciated next:

**38**

Figure 3.8: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.1$ for problem 3



Figure 3.9: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.05$ for problem 3



Figure 3.10: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.025$ for problem 3

**39**

As we can see, the maximum errors occur near the upper boundary, and they decrease rapidly as the discretization increases.

**Problem 4:**

- Consider region $\Omega$ as the set of points $(x, y)$ such that $x^2 + y^2 \leq 1$, $0 \leq x$ and $0 \leq y$, and that the physical quantity under investigation is given by $u = \ln\left(1 + 2 \cdot x + 3 \cdot y\right)$. We wish to apply our FEA to check how do the errors on $\widetilde{u}$ and $\nabla\widetilde{u}$ decrease when the mesh is refined. Applying our FEA we get to the following results:



Figure 3.11: Errors of $\widetilde{u}$ vs $h$ for problem 4



Figure 3.12: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 4



Figure 3.13: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 4

**40**

Figure 3.14: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.02$ for problem 4



Figure 3.15: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.01$ for problem 4



Figure 3.16: Graph of $\widetilde{u}$ and $|u - \widetilde{u}|$ for $h = 0.0063$ for problem 4

Also in this case, we can appreciate quadratic convergence for $\widetilde{u}$, and linear convergence for $\nabla\widetilde{u}$. Graphs of the mesh and numerical solution can be appreciated as well, where the maximum errors concentrate near the boundary and decrease rapidly as the discretization increases.

# Chapter 4

# Extension to a more general problem

## Diffusion: introducing porosity $k$

So far, we have made great advances on a very simple case of Poisson's equation, however a more general case is when we have the following:

$$-\langle \nabla, k \cdot \nabla u \rangle = f \ in \ \Omega \tag{4.1}$$

$$u = g \ on \ \partial\Omega$$

where coefficient $k = k(x, y)$ is a function usually called "porosity" or "conductivity", which satisfies $k \geq k_0 > 0$ in $\Omega$, where $k_0$ is a constant. Equation (4.1) is particularly applied on heat transfer phenomenon and fluid flow situations.

We wish to find a solution for equation (4.1), while still maintaining several of the results obtained for equation (1.1). To do this, we will approximate function $k$, to a certain degree, over each finite element.

To apply the FEA to equation (4.1), first we transform it into its equivalent weak formulation , therefore we get to:

$$\int_{\widetilde{\Omega}} \langle k \cdot \nabla u, \nabla v_j \rangle \, d\Omega - \int_{\partial\widetilde{\Omega}} \left( v_j \cdot k \cdot \frac{du}{dn} \right) ds - \int_{\partial\widetilde{\Omega}} \left( u \cdot k \cdot \frac{dv_j}{dn} \right) ds$$
$$+ \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} (uv_j) \, ds - \int_{\partial\widetilde{\Omega}} \left[ \langle k \cdot \nabla u, d \rangle \cdot \frac{dv_j}{dn} \right] ds + \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} \left[ \langle \nabla u, d \rangle \cdot v_j \right] ds$$
$$= \int_{\widetilde{\Omega}} (f \cdot v_j) \, d\Omega - \int_{\partial\widetilde{\Omega}} \left( g \cdot k \cdot \frac{dv_j}{dn} \right) ds + \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}} (g \cdot v_j) \, ds \quad (4.2)$$

Approximating $u$, $f$, and $g$ as a linear combination of the test functions we immediately notice the problem for this new general case: Having the integral $\int_R I(x, y) \cdot d\Omega$, how can

we compute the integral $\int_R k \cdot I(x, y) \cdot d\Omega$ for a given function $k$? To avoid complications and still use our previous results, we would like to obtain something like $\int_R k \cdot I(x, y) \cdot d\Omega = \bar{k} \cdot \int_R I(x, y) \cdot d\Omega$. This happens only when $k$ is constant, therefore we will approximate $k$ as almost constant along each finite element.

To compute each of the local matrices for this new problem, our solution will be to evaluate porosity $\bar{k}$ at the geometrical center of each finite element and multiply it by the expressions of the local matrices we obtained before. This does not imply any major mistake since the error of our FEA, theoretically, decreases quadratically.

To properly express this, we redefine matrices $A$, $C$, and $E$ as:

$$A = \int_{\widetilde{\Omega}} \langle k \cdot \nabla v_i, \nabla v_j \rangle \, d\Omega$$

$$C = \int_{\partial\widetilde{\Omega}} \left( v_j \cdot k \cdot \frac{dv_i}{dn} \right) ds$$

$$E = \int_{\partial\widetilde{\Omega}} \left[ \langle k \cdot \nabla v_i, d \rangle \cdot \frac{dv_j}{dn} \right] ds$$

And, for example, now the local matrix $A$ is given by:

$$A_{local} = \bar{k} \cdot \frac{\beta}{6} \cdot \begin{pmatrix} 2 & 1 & 1 & -1 \\ 1 & 2 & -1 & 1 \\ 1 & -1 & 2 & 1 \\ -1 & 1 & 1 & 2 \end{pmatrix} - \bar{k} \cdot \frac{1}{2} \cdot \begin{pmatrix} 0 & \frac{h_x}{h_y} & \frac{h_y}{h_x} & 0 \\ \frac{h_x}{h_y} & 0 & 0 & \frac{h_y}{h_x} \\ \frac{h_y}{h_x} & 0 & 0 & \frac{h_x}{h_y} \\ 0 & \frac{h_y}{h_x} & \frac{h_x}{h_y} & 0 \end{pmatrix}$$

The same procedure must be done with matrices $C$ and $E$.

## Advection: introducing velocity field $V$

Equation (4.1) can be further expanded by introducing an advection term. Advection is a phenomenon caused by the transport of a certain quantity along with its mass due to a velocity field $V = V(x, y)$, like energy or momentum. To add this into our formulation we have that:

$$- \langle \nabla, k \cdot \nabla u \rangle + \langle \nabla, V \cdot u \rangle = f \ in \ \Omega \tag{4.3}$$

$$u = g \ on \ \partial\Omega$$

To deal with this new term $\langle \nabla, V \cdot u \rangle$, we do the same procedure as with the others (multiply it by test function $v_j$, integrate over the whole region, and approximate $u$, $f$, and $g$ as a linear combination of the test functions). By doing so and developing further we get that:

$$\int_{\widetilde{\Omega}} \left[ \left\langle \nabla, V \cdot \sum (u_i \cdot v_i) \right\rangle \cdot v_j \right] d\Omega = \sum u_i \int_{\widetilde{\Omega}} \left[ \langle \nabla, V \rangle \cdot v_i \cdot v_j \right] d\Omega + \sum u_i \int_{\widetilde{\Omega}} \left[ \langle V, \nabla v_i \rangle \cdot v_j \right] d\Omega$$

Therefore we need to get the above 2 integrals.

**44**

## Concerning the term $\int_{\widetilde{\Omega}} \left[ \langle \nabla, V \rangle \cdot v_i \cdot v_j \right] d\Omega$

The above term represents a matrix, which will call $P$, whose entries will be $p_{ji} = \int_{\widetilde{\Omega}} \left[ \langle \nabla, V \rangle \cdot v_i \cdot v_j \right] d\Omega$. Due to the shape of the test functions, most of the entries of matrix $P$, will be zero, specifically when nodes $i$ and $j$ do not belong to the same finite element. As we did with the integrals involving porosity $k$, we will approximate the above integral in each finite element as:

$$\int_R \left[ \langle \nabla, V \rangle \cdot v_i \cdot v_j \right] d\Omega \approx \langle \nabla, V \rangle (\epsilon, \eta) \int_R (v_i \cdot v_j) \, d\Omega$$

where $\epsilon$ and $\eta$ represent the coordinates of the geometric center of region $R$. The different coefficients $\int_R (v_i \cdot v_j) \, d\Omega$ are the entries of matrix $K$ introduced in chapter 1, so $P_{local} = \langle \nabla, V \rangle (\epsilon, \eta) \cdot K_{local}$.

## Concerning the term $\int_{\widetilde{\Omega}} \left[ \langle V, \nabla v_i \rangle \cdot v_j \right] d\Omega$

The above term represents a matrix, which we will call $Q$, whose entries will be $q_{ji} = \int_{\widetilde{\Omega}} \left[ \langle V, \nabla v_i \rangle \cdot v_j \right] d\Omega$. Most of the entries of matrix $Q$ will be zero, and this happens when nodes $i$ and $j$ do not belong to the same finite element. For the computation of matrix $Q$, in order to simplify calculations, it is better to calculate the integral $\int_{\widetilde{\Omega}} \left[ \langle V, \nabla \widetilde{u} \rangle \cdot v_j \right] d\Omega$.

To deal with the above term, we separate it by performing the dot product. Therefore assuming that $V$ has components $V_x$ and $V_y$ in the $x$ and $y$ direction respectively, we get:

$$\int_R \left[ \langle V, \nabla \widetilde{u} \rangle \cdot v_j \right] d\Omega \approx V_x(\epsilon, \eta) \cdot \int_R \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot v_j \right] d\Omega + V_y(\epsilon, \eta) \cdot \int_R \left[ \frac{\partial \widetilde{u}}{\partial y} \cdot v_j \right] d\Omega$$

where $\epsilon$ and $\eta$ are again the coordinates of the geometric center of the finite element. Calculating the above integrals for a local numbering of the nodes, where the derivatives of $\widetilde{u}$ are also expressed under the local numbering like:

$$\frac{\partial \widetilde{u}}{\partial x} = \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y}$$

$$\frac{\partial \widetilde{u}}{\partial y} = \frac{(u_2 - u_1)}{h_y} \cdot \left( 1 - \frac{x}{h_x} \right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}$$

we get:

$$\int_R \left[ \frac{\partial \widetilde{u}}{\partial x} \cdot v_1 \right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[ \frac{(u_3 - u_1)}{h_x} \cdot \left( 1 - \frac{y}{h_y} \right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y} \right] \left[ \left( 1 - \frac{x}{h_x} \right) \left( 1 - \frac{y}{h_y} \right) \right] \cdot dx dy$$

$$= \frac{h_y}{12} \left[ 2 \cdot (u_3 - u_1) + (u_4 - u_2) \right]$$

**45**

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial y} \cdot v_1\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_2 - u_1)}{h_y} \cdot \left(1 - \frac{x}{h_x}\right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}\right] \left[\left(1 - \frac{x}{h_x}\right)\left(1 - \frac{y}{h_y}\right)\right] \cdot dxdy$$

$$= \frac{h_x}{12} \left[2 \cdot (u_2 - u_1) + (u_4 - u_3)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial x} \cdot v_2\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_3 - u_1)}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y}\right] \left[\left(1 - \frac{x}{h_x}\right) \cdot \frac{y}{h_y}\right] \cdot dxdy$$

$$= \frac{h_y}{12} \left[(u_3 - u_1) + 2 \cdot (u_4 - u_2)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial y} \cdot v_2\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_2 - u_1)}{h_y} \cdot \left(1 - \frac{x}{h_x}\right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}\right] \left[\left(1 - \frac{x}{h_x}\right) \cdot \frac{y}{h_y}\right] \cdot dxdy$$

$$= \frac{h_x}{12} \left[2 \cdot (u_2 - u_1) + (u_4 - u_3)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial x} \cdot v_3\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_3 - u_1)}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y}\right] \left[\frac{x}{h_x} \cdot \left(1 - \frac{y}{h_y}\right)\right] \cdot dxdy$$

$$= \frac{h_y}{12} \left[2 \cdot (u_3 - u_1) + (u_4 - u_2)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial y} \cdot v_3\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_2 - u_1)}{h_y} \cdot \left(1 - \frac{x}{h_x}\right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}\right] \left[\frac{x}{h_x} \cdot \left(1 - \frac{y}{h_y}\right)\right] \cdot dxdy$$

$$= \frac{h_x}{12} \left[(u_2 - u_1) + 2 \cdot (u_4 - u_3)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial x} \cdot v_4\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_3 - u_1)}{h_x} \cdot \left(1 - \frac{y}{h_y}\right) + \frac{(u_4 - u_2)}{h_x} \cdot \frac{y}{h_y}\right] \left[\frac{x}{h_x} \cdot \frac{y}{h_y}\right] \cdot dxdy$$

$$= \frac{h_y}{12} \left[(u_3 - u_1) + 2 \cdot (u_4 - u_2)\right]$$

$$\int_R \left[\frac{\partial \widetilde{u}}{\partial y} \cdot v_4\right] d\Omega = \int_0^{h_y} \int_0^{h_x} \left[\frac{(u_2 - u_1)}{h_y} \cdot \left(1 - \frac{x}{h_x}\right) + \frac{(u_4 - u_3)}{h_y} \cdot \frac{x}{h_x}\right] \left[\frac{x}{h_x} \cdot \frac{y}{h_y}\right] \cdot dxdy$$

$$= \frac{h_x}{12} \left[(u_2 - u_1) + 2 \cdot (u_4 - u_3)\right]$$

Using all of the above information, we can build the local matrix $Q$ as:

$$Q_{local} = V_x(\epsilon, \eta) \frac{h_y}{12} \cdot \begin{pmatrix} -2 & -1 & 2 & 1 \\ -1 & -2 & 1 & 2 \\ -2 & -1 & 2 & 1 \\ -1 & -2 & 1 & 2 \end{pmatrix} + V_y(\epsilon, \eta) \frac{h_x}{12} \cdot \begin{pmatrix} -2 & 2 & -1 & 1 \\ -2 & 2 & -1 & 1 \\ -1 & 1 & -2 & 2 \\ -1 & 1 & -2 & 2 \end{pmatrix}$$

# Examples

**Problem 5:**

- Consider region $\Omega$ as the set of points $(x, y)$ such that $0 \le y \le 3 + sin(2 \cdot \pi \cdot x) \cdot tan^{-1}(x - 4)$ and that $0 \le x \le 5$. Assume that function $u$ is some polynomial of the shape $u = c_1 \cdot (x + 1)^2 + c_2 \cdot (y + 1)^3$, that porosity $k = 4$, and velocity field as $V = (2,1)$. In the following images we can see how the errors behave when the mesh is refined:
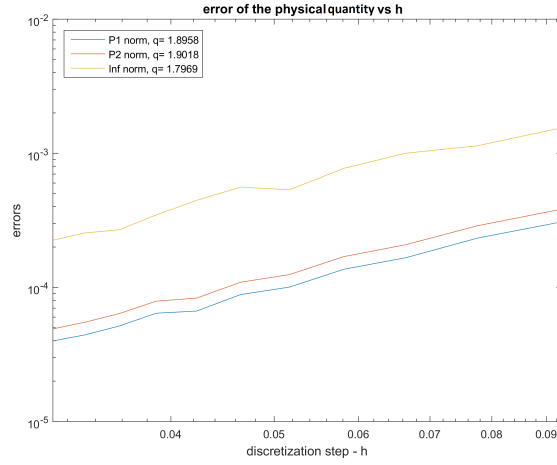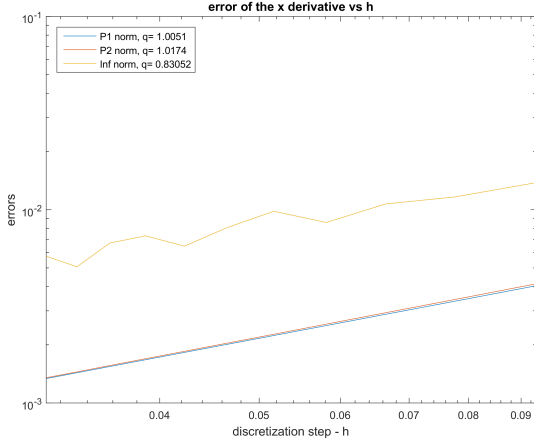


Figure 4.1: Errors of $\widetilde{u}$ vs $h$ for problem 5



Figure 4.2: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 5



Figure 4.3: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 5

We can appreciate near quadratic convergence for $\widetilde{u}$ and near linear convergence for $\nabla \widetilde{u}$. Also we can show different discretizations for our FEA:

**47**

Figure 4.4: Graphs of $\widetilde{u}$ for several discretizations for problem 5

**Problem 6**

- As a more general case, we take as region $\Omega$ the set of points $(x, y)$ such that $0 \leq y \leq 5 - e^{\frac{x}{2}}$ and $0 \leq x \leq 2 \cdot \ln(5)$, porosity $k = (x+1)^2 + (y+1)^2$, velocity field $V = (x + y + 2, sin(x \cdot y))$, and solution $u$ as $u = 2 + sin\left(\frac{\pi}{\ln(5)} \cdot x\right) \cdot sin\left(\frac{3 \cdot \pi}{4} \cdot y\right)$. Note that for this example $k_0 = 2$, since $k \geq 2 > 0$ In the following images we can see the behavior of the errors with respect to the discretization step:

Figure 4.5: Errors of $\widetilde{u}$ vs $h$ for problem 6



Figure 4.6: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 6



Figure 4.7: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 6

Again, we can appreciate near quadratic convergence for $\widetilde{u}$ and near linear convergence for $\nabla \widetilde{u}$. Graphs of the discretizations are shown next:

**49**

Figure 4.8: Graphs of $\widetilde{u}$ for several discretizations for problem 6

## Problem 7

- So far, we have only taken on problems where we knew the solution, condition which was necessary to compute the errors in the approximation. Now we decide to study another type of problem just to show convergence as the discretization increases, which is:

$$- \langle \nabla, k \cdot \nabla u \rangle + \langle \nabla, V \cdot u \rangle = 1 \; in \; \Omega$$

$$u = \sqrt{x^2 + y^2} \; on \; \partial\Omega$$

where we take region $\Omega$ as the set of points $(x, y)$ such that $0 \leq y \leq \sqrt{1 - x^2}$ and $0 \leq x \leq 1$, porosity $k = \frac{1}{1 + x^2 + y^2}$, and velocity field with components $V_x = ln(1 + x + y)$ and $V_y = 5 + e^{x-y}$. Note that for this example $k_0 = \frac{1}{2}$.

Performing our FEA for several discretization steps, we get the following graphs for $\widetilde{u}$:

**50**

Figure 4.9: Graphs of $\widetilde{u}$ for discretization steps $h = 0.05$, $h = 0.025$, $h = 0.0125$, and $h = 0.00625$ (left to right) for problem 7

Notice how as the the discretization is enriched, the resulting graph converges more and more. To show this, we show the value of $\widetilde{u}$ at point $(x, y) = (0.8, 0.2)$ for the different values of $h$:

| $h$ | $\widetilde{u}$ |
|---------|------------------------|
| 0.05 | 0.824376171840539 |
| 0.025 | 0.824577483938598 |
| 0.0125 | 0.824658354934295 |
| 0.00625 | 0.824674019954887 |

Table 4.1: Values of $\widetilde{u}$ at point $(x, y) = (0.8, 0.2)$ for several values of $h$

Using the last value of $\widetilde{u}$ on the above table as a reference for the actual solution, we can get that the overall convergence rate is 2.12, meaning quadratically.

# Chapter 5

# Introducing time dependence

We have made continuous progress in solving Poisson's equation and its extensions through the FEA, however, so far we only handled steady state cases, which translates to all time derivatives being null.

Therefore, in this chapter we deal with the time dependent case of Poisson's equation and its extensions, or how it is widely known, **the heat equation**. For this we introduce the time variable denoted by $t$. This implies that variable $u$ might change over time, as well as other properties like porosity $k$ and velocity field $V$. We will also admit that region $\Omega$ may change over time, i.e., $\Omega = \Omega(t) = \Omega_t$, according to a given law.

To study this, we further extend equation (4.3) by adding a time dependent term in the following way:

$$\frac{\partial}{\partial t}\left(s \cdot u\right) - \langle \nabla, k \cdot \nabla u \rangle + \langle \nabla, V \cdot u \rangle = f \ in \ \Omega_t \tag{5.1}$$

$$u = g \ on \ \partial\Omega_t \qquad u = u_0 \ in \ \Omega_0$$

where $u_0$ are the initial conditions and $s$ is called "volumetric capacity". Volumetric capacity $s$ is a property which measures the "resistance" of a small volume of mass to keep its current value of quantity $u$ even when subjected to different conditions over time, a concept similar to inertia. It is to be noted that $s = s(x, y, t)$

To deal with this new problem, we apply our FEA as previously: we multiply equation (5.1) by the test functions $v_j$, , integrate over the whole domain $\widetilde{\Omega_t}$, and assume that $u$, $f$, and $g$ are approximated by a linear combination of the test functions. Doing so will give similar expressions to the matrices we computed previously, however now time dependant. Because of this, we redefine all previously computed matrices as:

$$A^t = \int_{\widetilde{\Omega}_t} \langle k \cdot \nabla v_i, \nabla v_j \rangle \, d\Omega \qquad\qquad C^t = \int_{\partial\widetilde{\Omega}_t} \left( v_j \cdot k \cdot \frac{dv_i}{dn} \right) ds$$

$$D^t = \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}_t} (v_i v_j) \, ds \qquad\qquad E^t = \int_{\partial\widetilde{\Omega}_t} \left[ \langle k \cdot \nabla v_i, d \rangle \cdot \frac{dv_j}{dn} \right] ds$$

$$L^t = \frac{\gamma}{h} \int_{\partial\widetilde{\Omega}_t} [\langle \nabla v_i, d \rangle \cdot v_j] \, ds \qquad\qquad P^t = \int_{\widetilde{\Omega}_t} [\langle \nabla, V \rangle \cdot v_i \cdot v_j] \, d\Omega$$

$$Q^t = \int_{\widetilde{\Omega}_t} [\langle V, \nabla v_i \rangle \cdot v_j] \, d\Omega \qquad\qquad K^t = \int_{\widetilde{\Omega}_t} (v_i \cdot v_j) \, d\Omega$$

where the superscript $^t$ means that all matrices were computed at time $t$.

Additional to the above matrices, our FEA on problem (5.1) will give a new contribution which is expressed by $\int_{\widetilde{\Omega}_t} \frac{\partial\left(s \cdot \sum u_i^t \cdot v_i\right)}{\partial t} \cdot v_j \cdot d\Omega$, where $u_i^t$ means quantity $\widetilde{u}$ at node $i$ and at time $t$.

To deal with the time dependence, we will assume that the different values of $u_i^t$ at the different nodes are time dependent. Focusing on the above single term, we can rewrite it as:

$$\int_{\widetilde{\Omega}_t} \frac{\partial\left(s \cdot \sum u_i^t \cdot v_i\right)}{\partial t} \cdot v_j \cdot d\Omega = \sum u_i^t \cdot \int_{\widetilde{\Omega}_t} \left( \frac{\partial s}{\partial t} \cdot v_i \cdot v_j \right) \cdot d\Omega + \sum \frac{\partial u_i^t}{\partial t} \cdot \int_{\widetilde{\Omega}_t} (s \cdot v_i \cdot v_j) \cdot d\Omega$$

Now we define two new matrices as:

$$T^t = \int_{\widetilde{\Omega}_t} \left( \frac{\partial s}{\partial t} \cdot v_i \cdot v_j \right) \cdot d\Omega \qquad\qquad W^t = \int_{\widetilde{\Omega}_t} (s \cdot v_i \cdot v_j) \cdot d\Omega$$

whose local matrices will be approximated as:

$$T^t_{local} \approx \overline{\frac{\partial s}{\partial t}} \cdot K_{local} \qquad\qquad W^t_{local} \approx \bar{s} \cdot K_{local}$$

where matrix $K_{local}$ was presented in chapter 1, and $\bar{s}$ and $\overline{\frac{\partial s}{\partial t}}$ are the values of $s$ and $\frac{\partial s}{\partial t}$, respectively, at the geometrical center of finite element $R$. Notice how matrices $T^t$ and $W^t$ are similar in structure to matrix $P^t$.

Therefore, our FEA results in the following system of equations:

$$W^t \cdot \frac{\partial \vec{u^t}}{\partial t} + \left( A^t - C^t - C^{t^T} + D^t - E^t + L^t + P^t + Q^t + T^t \right) \cdot \vec{u^t} = K^t \cdot \vec{f^t} + \left( D^t - C^{t^T} \right) \cdot \vec{g^t}$$

In order to simplify the above expression, it is useful to group some of the above terms like:

**54**

$$Z^t = A^t - C^t - C^{t^T} + D^t - E^t + L^t + P^t + Q^t + T^t$$
$$\vec{F^t} = K^t \cdot \vec{f^t} + \left(D^t - C^{t^T}\right) \cdot \vec{g^t}$$

so:

$$W^t \cdot \frac{\partial \vec{u^t}}{\partial t} + Z^t \cdot \vec{u^t} = \vec{F^t} \tag{5.2}$$

It's important to remember that our goal with this procedure is to approximate function $u$ across domain $\Omega_t$ while we are able to advance in time. For this we can apply any time time advancing scheme to equation (5.2).

## Change of region $\Omega_t$

From problem (5.1) we managed to apply our FEA to reach equation (5.2), and choosing a time advancing scheme for equation (5.2) allows to approximate a solution for problem (5.1) at any time instant. However there are still some details to discuss for the case in which region $\Omega_t$ changes in time.

To understand this idea, let's concentrate first in the case in which region $\Omega_t$ is fixed over time. In this case, for fixed spatial discretization steps, there is a constant number of nodes on $\Omega_t$ at any time $t$. Besides that, the initial condition of each node is contained within the information given by $u_0$.

Now let's concentrate in the case in which region $\Omega_t$ does change over time, specifically when $\Omega_t$ grows larger. In this case, for fixed spatial discretization steps, the number of nodes in $\Omega_t$ increases over time, and $u_0$ only gives information of initial conditions on the nodes appearing at time $t = 0$. Now the question is, how to get "initial conditions" for newly created nodes?.

As a reference, look at the following image:



Figure 5.1: Description of the creation of new nodes during the time advancing scheme

therein, we have a set of nodes (black and green), and two curves (blue and red). The blue curve is the physical boundary at a time $t$, meaning $\partial\Omega_t$, and the red curve is the physical boundary at a later time $t + \Delta t$, meaning $\partial\Omega_{t+\Delta t}$. The black nodes are nodes which participate on the FEA at time $t$ and $t + \Delta t$, but the green node (identified as point $C$) only participates on the FEA at time $t + \Delta t$.

**55**

For the FEA to be able to advance in time in a situation similar to the above, we need to find two things:

- The moment at which point $C$ was "created", which we will name time $t + \delta t$.

- The value of quantity $u$ at point $C$ at the moment $t + \delta t$, which following the notation will be named as $u_C^{t+\delta t}$.

For that, we find two points, $A$ and $B$, such that points $A$, $B$ and $C$ are co-linear. Point $A$ belongs to $\partial \Omega_{t+\Delta t}$, and point $B$ belongs to $\partial \Omega_t$. Now we write first order Taylor expansions centered at point $C$, neglecting higher order terms, which relate to points $A$ and $B$, so:

$$u_A^{t+\Delta t} = u_C^{t+\delta t} + \langle \nabla u_C^{t+\delta t}, \vec{e} \rangle \cdot d_2 + \frac{\partial u_C^{t+\delta t}}{\partial t} \cdot (\Delta t - \delta t)$$

$$u_B^t = u_C^{t+\delta t} - \langle \nabla u_C^{t+\delta t}, \vec{e} \rangle \cdot d_1 - \frac{\partial u_C^{t+\delta t}}{\partial t} \cdot \delta t$$

where:

- $\vec{e}$ is a unit vector which connects point $C$ to point $A$.

- $d_1$ is the distance between points $B$ and $C$.

- $d_2$ is the distance between points $C$ and $A$.

In the above 2 equations, our unknowns are $u_C^{t+\delta t}$ and $\delta t$. To find these two unknowns, we multiply the first Taylor expansion by $d_1$, the second Taylor expansion by $d_2$, and add them:

$$d_1 \cdot u_A^{t+\Delta t} + d_2 \cdot u_B^t = (d_1 + d_2) \cdot u_C^{t+\delta t} + \frac{\partial u_C^{t+\delta t}}{\partial t} \cdot [d_1 \cdot \Delta t - (d_1 + d_2) \cdot \delta t]$$

Since $\frac{\partial u_C^{t+\delta t}}{\partial t}$ is not known, we define $\delta t$ as:

$$\delta t = \frac{d_1 \cdot \Delta t}{d_1 + d_2} \tag{5.3}$$

and we get that:

$$u_C^{t+\delta t} = \frac{d_1 \cdot u_A^{t+\Delta t} + d_2 \cdot u_B^t}{d_1 + d_2} \tag{5.4}$$

With this information, we can build the time advancing scheme for any situation. Notice how for the case in which $\Omega_t$ grows smaller over time, this question about "initial conditions" does not happen, since no nodes are created. Before proceeding, it is important to comment some details for this methodology:

**56**

- The geometrical interpretation of this methodology is that there is a constant velocity of deformation for boundary $\partial\Omega_t$ between times $t$ and $t + \Delta t$, so the time increment $\delta t$ at which point $C$ is reached is proportional to the space increment $d_1$ between points $B$ and $C$. Furthermore, it is implicit that we assume that $u$ varies linearly between $u_A^{t+\Delta t}$ and $u_B^t$.

- For a given time step and assuming that $\partial\Omega_t$ grows over time, if the difference between $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$ is large and new nodes are created, then the actual behaviour of function $u$ on the newly created nodes might not be correctly modelled by our linear approach. This is not convenient, therefore we need to introduce a CFL condition to avoid this. This condition can be written as:

$$\frac{V_{\partial\Omega_t} \cdot \Delta t}{h} < C$$

where $V_{\partial\Omega_t}$ is the velocity of deformation of boundary $\partial\Omega_t$, and $C$ is a certain constant. Since $V_{\partial\Omega_t}$ is considered constant between times $t$ and $t + \Delta t$, this relation can be rewritten for every time step as:

$$\frac{(d_1 + d_2)_{max}}{h} < C \tag{5.5}$$

where $(d_1 + d_2)_{max}$ is the maximum distance between points $A$ and $B$ measured along one of the newly created nodes. From now on, we will refer to coefficient $\frac{(d_1+d_2)_{max}}{h}$ as the CFL number.

- Obviously, the correct choice of points $A$ and $B$ for every newly created node affects the performance of the method. There are infinite combinations for points $A$ and $B$ which will work, but it is better to choose the **best** combination of them.

## Computation of points $A$ and $B$

When boundary $\partial\Omega_t$ grows in time, new nodes are created and we must find points $A$ and $B$ to compute the "initial conditions" for such nodes. The best choice for them is so that the distance between points $A$ and $B$ is minimum.

This is a purely geometrical problem, since having a newly created node, which we have called point $C$, and boundaries $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$, the goal is to find points $A \in \partial\Omega_{t+\Delta t}$ and $B \in \partial\Omega_t$, such that points $A$, $B$, and $C$ are colinear, and the distance between points $A$ and $B$ is minimal. Due to the physical interpretation of the circumstances, another condition is that point $C$ must be between points $A$ and $B$.

Figure 5.2: Representation of how boundaries $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$ are linearized over a small interval

Since finding a solution for a general problem like this is highly difficult (and in some cases unfeasible), we will study this problem over small segments of boundaries $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$ which are approximated as straight lines, as it can be seen in figure (5.2). However this is not an obstacle, since boundaries $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$ can be thought as to be made of consecutive segments of straight lines.

For this, assume we have 2 line segments which represent boundaries $\partial\Omega_t$ and $\partial\Omega_{t+\Delta t}$:

- $L_1$ which represents $\partial\Omega_{t+\Delta t}$, and has extreme points $P_1 = (x_1, y_1)$ and $P_2 = (x_2, y_2)$. Remember that $A \in L_1$.

- $L_2$ which represents $\partial\Omega_t$, and has extreme points $P_3 = (x_3, y_3)$ and $P_4 = (x_4, y_4)$. Remember that $B \in L_2$

From the above ideas, we can say that:

$$A = P_1 + (P_2 - P_1) \cdot a$$
$$B = P_3 + (P_4 - P_3) \cdot b$$

where $a$ and $b$ are scalars such that $0 \le a \le 1$ and $0 \le b \le 1$. The key now is finding $a$ and $b$, and from there finding points $A$ and $B$.

Before starting, it is useful to define 3 important notations:

- We define vector $V_{ij}$, such that $V_{ij} = P_j - P_i$, which is a vector that goes from point $P_i$ to point $P_j$.

- We define the norm of the cross product between 2 vectors $\vec{a}$ and $\vec{b}$ as $|\vec{a} \times \vec{b}|$.

- The dot product between 2 vectors $\vec{a}$ and $\vec{b}$ is still given by $\langle \vec{a}, \vec{b} \rangle$

To proceed, let's focus on point $A$, and connect it to point $B$ through point $C$. For this, assume that point $C$ has coordinates $(x_C, y_C)$, and let's define line $L_3$ which goes from $A$ to $C$, so:

$$L_3 = A + (C - A) \cdot c$$

**58**

where $c$ is another scalar. Point $B$ follows from colinearity, so we need to find $b$ based on this condition. For this let's intersect lines $L_2$ and $L_3$, and we get the following system of equations to solve for parameters $b$ and $c$ (even though $c$ is not important):

$$V_{34} \cdot b + V_{CA} \cdot c = V_{3A}$$

Using Kramer's rule, we find that:

$$b = \frac{|V_{3A} \times V_{CA}|}{|V_{34} \times V_{CA}|}$$

Let's try and simplify the above into something more understandable. Since $V_{3A} = V_{31} + V_{12} \cdot a$ and $V_{CA} = V_{C1} + V_{12} \cdot a$, we can write:

$$b = \frac{|(V_{31} + V_{12} \cdot a) \times (V_{C1} + V_{12} \cdot a)|}{|V_{34} \times (V_{C1} + V_{12} \cdot a)|} = \frac{|V_{31} \times V_{C1}| + |V_{3C} \times V_{12}| \cdot a}{|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a}$$

The above value of $b$ represents point $B$, such that points $A$, $C$ and $B$ are colinear. As a reference, look at image (5.3) where we can see the colinearity of points $A$, $C$, and $B$, and the different vectors used in our analysis:



Figure 5.3: Graph of the different vectors used in the computations of points $A$ and $B$

Having this, we define the distance $d_{AB}$ between points $A$ and $B$ from the following:

$$d_{AB}^2 = \langle V_{BA}, V_{BA} \rangle$$

Now, points $A$ and $B$ both depend on parameter $a$, so in order to minimize the distance, we derive with respect to $a$ and nullify the derivative. To make things easier, let's find $db/da$ first.

$$\frac{db}{da} = \frac{|V_{3C} \times V_{12}| \cdot (|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a) - (|V_{31} \times V_{C1}| + |V_{3C} \times V_{12}| \cdot a) \cdot |V_{34} \times V_{12}|}{(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a)^2}$$

$$= \frac{|V_{3C} \times V_{12}| \cdot |V_{34} \times V_{C1}| - |V_{31} \times V_{C1}| \cdot |V_{34} \times V_{12}|}{(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a)^2}$$

The above can be further simplified like:

$$\frac{db}{da} = -\frac{|V_{34} \times V_{C3}| \cdot |V_{C1} \times V_{12}|}{(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a)^2}$$

Now we can get the minimum distance by deriving $d_{AB}$ with respect to $a$ and nulling the derivative. This is given by solving for $a$ in the following equation:

$$2 \cdot \left\langle V_{BA}, \frac{dV_{BA}}{da} \right\rangle = 0 \tag{5.6}$$

Computing what is needed:

$$\begin{aligned}
V_{BA} &= A - B \\
&= P_1 + (P_2 - P_1) \cdot a - (P_3 + (P_4 - P_3) \cdot b) \\
&= V_{31} + V_{12} \cdot a - V_{34} \cdot \frac{|V_{31} \times V_{C1}| + |V_{3C} \times V_{12}| \cdot a}{|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a} \\
&= \frac{(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a) \cdot (V_{31} + V_{12} \cdot a) - (|V_{31} \times V_{C1}| + |V_{3C} \times V_{12}| \cdot a) \cdot V_{34}}{|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a}
\end{aligned}$$

Which simplifies to:

$$V_{BA} = \frac{|V_{34} \times V_{31}| + |V_{34} \times V_{12}| \cdot a}{|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a} \cdot (V_{C1} + V_{12} \cdot a)$$

Introducing the above into equation (5.6), we get:

$$\frac{|V_{34} \times V_{31}| + |V_{34} \times V_{12}| \cdot a}{|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a} \cdot \left\langle (V_{C1} + V_{12} \cdot a), \left( V_{12} - V_{34} \cdot \frac{db}{da} \right) \right\rangle = 0$$

Which, already gives the first solution, which is when:

$$|V_{34} \times V_{31}| + |V_{34} \times V_{12}| \cdot a = 0$$

So:

$$a = -\frac{V_{34} \times V_{31}}{V_{34} \times V_{12}}$$

**60**

The above value of $a$ represents the point where lines $L_1$ and $L_2$ would theoretically intersect (meaning $A$ and $B$ would be the same point), giving $d = 0$, a solution which is independent of point $C$ but it is not what we're looking for since $C$ must be between points $A$ and $B$, so we discard it. The remaining solutions are when:

$$\left\langle (V_{C1} + V_{12} \cdot a), \left( V_{12} - V_{34} \cdot \frac{db}{da} \right) \right\rangle = 0$$

Substituting for $db/da$:

$$\left\langle (V_{C1} + V_{12} \cdot a), \left( V_{12} + \frac{|V_{34} \times V_{C3}| \cdot |V_{C1} \times V_{12}|}{(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a)^2} \cdot V_{34} \right) \right\rangle = 0$$

Multiplying by $(|V_{34} \times V_{C1}| + |V_{34} \times V_{12}| \cdot a)^2$, performing the algebra and the dot product, reduces to:

$$C_0 + C_1 \cdot a + C_2 \cdot a^2 + C_3 \cdot a^3 = 0 \tag{5.7}$$

where:

$$C_0 = |V_{34} \times V_{C1}|^2 \langle V_{C1}, V_{12} \rangle + |V_{34} \times V_{C3}| \cdot |V_{C1} \times V_{12}| \cdot \langle V_{C1}, V_{34} \rangle$$

$$C_1 = 2 \cdot |V_{34} \times V_{C1}| \cdot |V_{34} \times V_{12}| \cdot \langle V_{C1}, V_{12} \rangle + |V_{34} \times V_{C1}|^2 \cdot \langle V_{12}, V_{12} \rangle + \dots$$
$$\dots + |V_{34} \times V_{C3}| \cdot |V_{C1} \times V_{12}| \cdot \langle V_{12}, V_{34} \rangle$$

$$C_2 = |V_{34} \times V_{12}|^2 \cdot \langle V_{C1}, V_{12} \rangle + 2|V_{34} \times V_{C1}| \cdot |V_{34} \times V_{12}| \cdot \langle V_{12}, V_{12} \rangle$$

$$C_3 = |V_{34} \times V_{12}|^2 \cdot \langle V_{12}, V_{12} \rangle$$

Solving for $a$ in equation (5.7) gives information about points $A$ and $B$. It's valid to mention that equation (5.7) is a cubic polynomial, so it has 3 solutions. Since coefficients $C_0$, $C_1$, $C_2$, and $C_3$ are real scalars, then (5.7) has either 1 real and 2 complex solutions, or 3 real solutions. Both cases are possible but complex solutions are discarded.

Having a solution for equation (5.7), it has to be tested it with our restrictions, meaning that $0 \le a \le 1$, $0 \le b \le 1$, and point $C$ must be between $A$ and $B$. This last conditions translates to:

$$0 \le \frac{\langle V_{AC}, V_{AB} \rangle}{\langle V_{AB}, V_{AB} \rangle} \le 1$$

In case **any** of these restrictions is not met, then other segments of boundaries $\partial \Omega_t$ and $\partial \Omega_{t+\Delta t}$ have to be considered until all conditions are fulfilled.

It is important to mention that all possible combination of segments between boundaries $\partial \Omega_t$ and $\partial \Omega_{t+\Delta t}$ have to be considered in order to find the global minimum for $d_{AB}$.

## Assembling all pieces together

Suppose we wish to solve the general problem (5.1). For that we apply our FEA and reach to equation (5.2). Suppose that our choice for the time advancing scheme is Backward Euler, therefore:

$$\frac{\partial u_i^t}{\partial t} \approx \frac{u_i^{t+\Delta t} - u_i^{t+\delta t_i}}{\Delta t - \delta t_i}$$

where $\delta t_i$ is the time interval of creation of node $i$ after time $t$, and is given by equation (5.3). In case node $i$ is not a newly created node, then $\delta t_i = 0$. It is to be noted that during a time step, several nodes can be created, as well as others can be "destroyed".

The above derivative can be expressed for all nodes on $\Omega_t$ as a matrix vector product like:

$$\frac{\partial \vec{u^t}}{\partial t} \approx \frac{1}{\Delta T} \cdot \left( \overrightarrow{u^{t+\Delta t}} - \overrightarrow{u^{t+\delta t_i}} \right)$$

where matrix $\frac{1}{\Delta T}$ is a diagonal matrix whose entries are the different $\frac{1}{\Delta t - \delta t_i}$ coefficients.

From here equation (5.2) becomes:

$$W^{t+\Delta t} \cdot \frac{1}{\Delta T} \cdot \left( \overrightarrow{u^{t+\Delta t}} - \overrightarrow{u^{t+\delta t_i}} \right) + Z^{t+\Delta t} \cdot \overrightarrow{u^{t+\Delta t}} = \overrightarrow{F^{t+\Delta t}}$$

Solving for $\overrightarrow{u^{t+\Delta t}}$ allows advancing in time, having a set of initial conditions.

It is worthwhile mentioning that in case $\partial \Omega_t$ is constant over time, it is easier to use higher order methods than Backward Euler on equation (5.2) since the issue of "initial conditions" on newly created nodes does not happen since there are not any, and $\delta t_i = 0$ for every node. For example, we can use the Trapezoidal rule (which is a second order method), which starting from equation (5.2) can be written as:

$$\frac{W^{t+\Delta t} + W^t}{2} \cdot \frac{\left( \overrightarrow{u^{t+\Delta t}} - \overrightarrow{u^t} \right)}{\Delta t} + \frac{Z^{t+\Delta t} + Z^t}{2} \cdot \frac{\left( \overrightarrow{u^{t+\Delta t}} - \overrightarrow{u^t} \right)}{2} = \frac{\overrightarrow{F^{t+\Delta t}} + \overrightarrow{F^t}}{2}$$

where, as for Backward Euler, solving for $\overrightarrow{u^{t+\Delta t}}$ allows advancing in time.

## Practical implementation of the method

In previous chapters, our scheme for the FEA was straight forward:

- We started from a known function $u$, together with porosity $k$ and velocity field $V$, which allowed to compute functions $f$ and $g$.

- Then we built a mesh of finite elements, which allowed to compute function $d$ across $\partial \widetilde{\Omega}$.

**62**

- With a proper choice of parameter $\gamma$, we were able to compute all necessary matrices for the simulation.

- Solving the system of equations allowed to get $\widetilde{u}$, and together with function $u$, we measured the error in our approximation for a given norm.

- Then post-process the results to check certain conditions, as well as the convergence of the method, which we have numerically proofed to be quadratic with respect to the spatial step $h$ for a wide variety of cases.

Since we are now dealing with time dependence, for our FEA we will have to do all of the above steps for every time step of our simulations, in addition to the following:

- Include volumetric capacity $s$ to compute its contribution to function $f$, and to compute the matrices which depend on it, meaning matrices $W^t$ and $T^t$.

- In case $\Omega_t$ grows at a time step, compute the "initial conditions" for the newly created nodes.

Since we already demonstrated quadratic convergence with respect to the spatial discretization, for this new case we will be interested in studying how the error behaves with respect to the time step. For this, now we consider that the error will be approximated by the following relationship:

$$error \approx C_{space} \cdot \bar{h}^2 + C_{time} \cdot \Delta t^q \tag{5.8}$$

where $C_{time}$ is a certain scalar constant, and now $q$ is the rate of decay of the error with respect to the temporal discretization step.

The current interpretation of the error is that it is composed by a sum of an error due to the discretization in space and another error due to the discretization in time. This is valid since the differential equation we are trying to solve, meaning the heat equation, is linear.

By performing some simulations we expect that $q$ will be similar to the of order of the method used to advance in time. For example, we expect $q \approx 1$ for the Backward Euler scheme, and $q \approx 2$ for the Trapezoidal rule.

We will also modify the expressions we used to measure the errors to include the time variation. For this, we redefine the error as:

$$e_p = \left( \frac{1}{t_f} \int_0^{t_f} \frac{\int_{\Omega_t} \left| u - \widetilde{u^t} \right|^p d\Omega}{\|\Omega_t\|} dt \right)^{1/p}$$

where $t_f$ is the final time instance of the simulation. To approximate the time integral, we focus on a single time step between times $t$ and $t + \Delta t$ and use the Trapezoidal rule, so:

$$\int_t^{t+\Delta t} \frac{\int_{\Omega_t} \left| u - \widetilde{u^t} \right|^p d\Omega}{\|\Omega_t\|} dt \approx \frac{\frac{\int_{\Omega_t} \left| u - \widetilde{u^t} \right|^p d\Omega}{\|\Omega_t\|} + \frac{\int_{\Omega_{t+\Delta t}} \left| u - \widetilde{u^{t+\Delta t}} \right|^p d\Omega}{\|\Omega_{t+\Delta t}\|}}{2} \cdot \Delta t$$

**63**

A special case to consider, as was also described in chapter 1, is when $p \to \infty$, since:

$$e_\infty = lim_{p \to \infty} \left( \frac{1}{t_f} \int_0^{t_f} \frac{\int_{\Omega_t} \left| u - \widetilde{u^t} \right|^p d\Omega}{\|\Omega_t\|} dt \right)^{1/p} = max_{\Omega_t \, [0,t_f]} |u - \widetilde{u^t}|$$

Meaning that $e_\infty$ is the maximum error occurred during the whole simulation.
Now we have everything necessary to test our FEA. The goals are 2:

1. Test if the the rate of decay of the error with respect to the temporal discretization step, meaning $q$, does match that of the order of the time advancing scheme.

   For this we will compare the errors for different combinations of $h$ and $\Delta t$ in a logical manner. The idea is that, having a simulation with parameters $h_0$ and $\Delta t_0$, we expect to see quadratic convergence if, for example, another simulation has parameters $h$ and $\Delta t$ such that:
   $$\left( \frac{h_0}{h} \right)^2 = \left( \frac{\Delta t_0}{\Delta t} \right)^q$$

2. Study how the errors behave with respect to the CFL number, given in equation (5.5). For this we will perform several simulations for a fixed discretization step $h$ while changing $\Delta t$, and compute the errors and the CFL number.

## Examples

**Problem 8:**

- Consider the following problem:

  - Region $\Omega_t$ as the set of points $(x, y)$ such that $0 \le y \le 2 + \sin(\pi \cdot x) \cdot \sin\left(t + \frac{\pi}{2}\right)$, where $0 \le x \le 3$ and $0 \le t \le 2 \cdot \pi$.
  - Assume that function $u$ is of the shape $u = 2 + \sin\left(\frac{2\pi \cdot x}{3}\right) \cdot \sin(\pi \cdot y) \cdot \sin(t)$
  - Porosity $k$ is given by $k = 1 + \left(\frac{t}{6 \cdot \pi}\right)^2 \cdot (x^2 + y^2)$
  - Velocity field $V$ has components $V_x = x + y - t$ and $V_y = 1 - \exp\left(-\frac{x+y+t}{3}\right)$
  - Volumetric capacity $s$ is given by $s = 2 + sin(x \cdot t - y \cdot t)$

We wish to study how the errors computed by our FEA behave for the above problem. This can be done by performing several simulations of our FEA, measuring the errors, and checking how they behave with respect to the discretization steps $\Delta t$ and $h$. Since boundary $\partial\Omega_t$ changes in time, we prefer to use Backward Euler as a time advancing scheme.

**_Changing_ Δ*t* _while keeping_ *h* _fixed_**

For a fixed spatial discretization step $h \approx 0.05$ and varying $\Delta t$, we get the following results for the CFL number and the errors:



Figure 5.4: Dependence of the CFL number with respect to $\Delta t$ for problem 8



Figure 5.5: Errors of $\widetilde{u}$ vs $\Delta t$ for problem 8



Figure 5.6: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs $\Delta t$ for problem 8
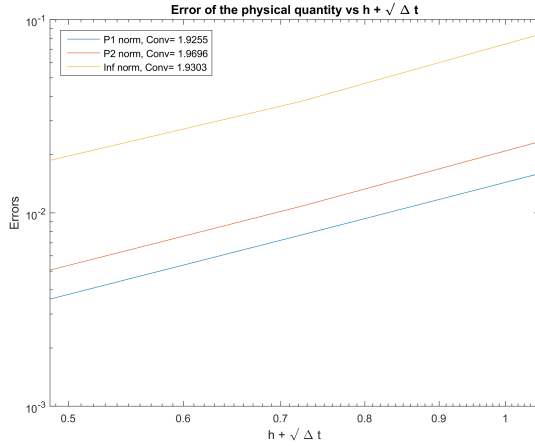
Figure 5.7: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $\Delta t$ for problem 8

Figure 5.8: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $\Delta t$ for problem 8

From the above figures, 4 things are to be noted:

- The rapid decrease of the CFL number as $\Delta t$ becomes smaller. Simulations were made with CFL numbers ranging from as high as 20 to as low as 2, and the method is stable. Obviously as $\Delta t \to 0$ then $CFL \to 0$, though this highly increases the number of time steps for the whole simulations. From the above we can see that the CFL number does not pose a real restriction on the choice of $\Delta t$.

- The errors for $\widetilde{u}$ decrease with a decreasing $\Delta t$ in a trend which is almost linear, as it should be since Backward Euler is a first order method. The reason why the trend is not exactly linear is because of the errors due to the spatial discretization.

- The errors for $\frac{\partial \widetilde{u}}{\partial t}$ also decrease almost linearly with respect to $\Delta t$ for the 1-norm and 2-norm, but not for the infinity norm (which seems to oscillate around a constant value).

  This is the case since a first order time advancing scheme gives a zeroth order approximation for the first time derivative, which means that $\frac{\partial \widetilde{u}}{\partial t}$ might not converge, as shown.

- There is almost no change in the errors for the spatial derivatives, $\frac{\partial \widetilde{u}}{\partial x}$ and $\frac{\partial \widetilde{u}}{\partial y}$, with respect to $\Delta t$.

### *Changing $h$ while keeping $\Delta t$ fixed*

We will also be interested in studying how the errors behave with respect to $h$ while keeping $\Delta t$ fixed. For a fixed time discretization step $\Delta t \approx 0.42$ and varying $h$, we get the following results for the errors:
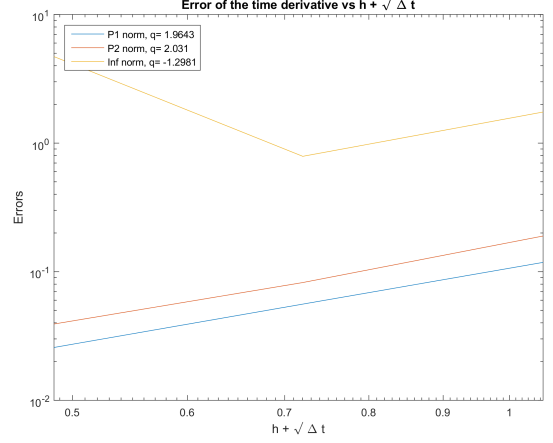
Figure 5.9: Errors of $\widetilde{u}$ vs $h$ for problem 8



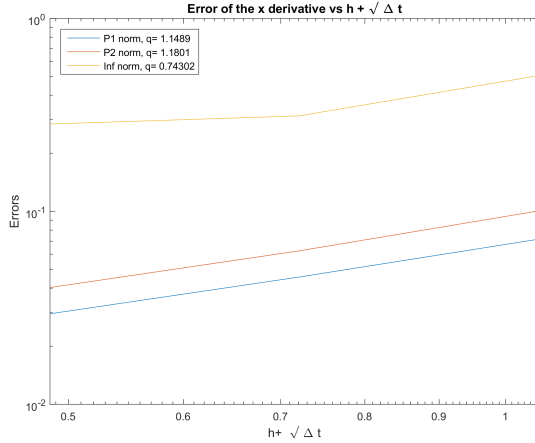Figure 5.10: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs $h$ for problem 8



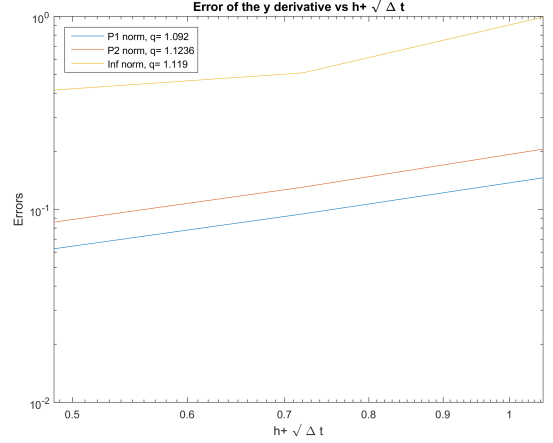Figure 5.11: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 8



Figure 5.12: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 8

From the above figures, 3 things are to be noted:

- The errors for $\widetilde{u}$ decrease quadratically with a decreasing $h$ in the first part of the plot, but then they flatten. This is because of the predominance of the errors due to the time discretization.

- We still have linear convergence for $\nabla u$ for the 1-norm and 2-norm, however not for the infinity norm. This is due to the errors involved in the time discretization.

- There is no change in the errors for the time derivative, $\frac{\partial \widetilde{u}}{\partial t}$, with respect to $h$.

### *Changing $h$ and $\Delta t$ together*

Now we would like to check if the rate of decay of the error with respect to the temporal discretization step, $q$, does match that of the order of the time advancing scheme. In shorter words, we would like to see if $q \approx 1$ since we are using Backward Euler.

For this we start with a simulation with parameters $h_0$ and $\Delta t_0$, perform other simulations by dividing $h_0$ by 2 and $\Delta t_0$ by 4, and compare the errors to check if we have quadratic convergence. If yes, then indeed $q \approx 1$.

To include the contribution of both discretization steps, we graph the errors vs parameter $h + \sqrt{\Delta t}$. By doing so, we get the following graph of the errors:



Figure 5.13: Errors of $\widetilde{u}$ vs parameter $h + \sqrt{\Delta t}$ for problem 8



Figure 5.14: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs parameter $h + \sqrt{\Delta t}$ for problem 8



Figure 5.15: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs parameter $h + \sqrt{\Delta t}$ for problem 8



Figure 5.16: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs parameter $h + \sqrt{\Delta t}$ for problem 8

**68**

From the above figures, 3 things are to be noted:

- As we can see at the top left corner of figure (5.13), we do get quadratic convergence for $\widetilde{u}$ for all norms, so indeed $q \approx 1$. This also means that equation (5.8) does give a very good approximation of how errors of $\widetilde{u}$ behave with respect to discretization steps $h$ and $\Delta t$.

- We have complete linear convergence for $\nabla u$ for the 1-norm, 2-norm, and the infinity norm. Even though the convergence for the infinity norm on $\frac{\partial \widetilde{u}}{\partial x}$ is a bit low (0.743), we can see that for $\frac{\partial \widetilde{u}}{\partial y}$ it is actually pretty high (1.119).

- The errors of $\frac{\partial \widetilde{u}}{\partial t}$ for the 1-norm and 2-norm become quadratic, but since we're using a first order time advancing scheme, the infinity norm still diverges.

### *Graphs of interest*

Together with this file, there are attached 3 videos titled **"Problem 8 - u"**, **"Problem 8 - f"**, and **"Problem 8 - e"**, which represent a live graphical representation of the approximation of $\widetilde{u}$, $\widetilde{f}$, and $|\widetilde{u} - u|$ respectively. Also, we show several graphs of $\widetilde{u}$ for different time instances:



Figure 5.17: Graph of $\widetilde{u}$ at time $t = \frac{10}{33}\pi$ for problem 8

Figure 5.18: Graph of $\widetilde{u}$ at time $t = \frac{28}{33}\pi$ for problem 8

Figure 5.19: Graph of $\widetilde{u}$ at time $t = \frac{46}{33}\pi$ for problem 8

Figure 5.20: Graph of $\widetilde{u}$ at time $t = \frac{64}{33}\pi$ for problem 8

**Problem 9:**

- Consider the following problem:

  - Region $\Omega_t$ as the set of points $(x, y)$ such that $0 \leq y \leq \frac{6}{x}$, where $1 \leq x \leq 6$ and $0 \leq t \leq 6$.

  - Assume that function $u$ is of the shape $u = (1 + 2 \cdot x + 3 \cdot y) \cdot (t - 3)$

  - Porosity $k$ is given by $k = 1 + \left(\frac{t}{6 \cdot \pi}\right)^2 \cdot \left(x^2 + y^2\right)$

  - Velocity field $V$ has components $V_x = x + y - t$ and $V_y = 1 - \exp\left(-\frac{x+y+t}{3}\right)$

  - Volumetric capacity $s$ is given by $s = 2 + sin\left(x \cdot t - y \cdot t\right)$

We wish to study how the errors computed by our FEA behave for the above problem. This can be done by performing several simulations of our FEA, measuring the errors, and checking how they behave with respect to the discretization steps $\Delta t$ and $h$. Since region $\Omega_t$ is constant in time, we decide to use the Trapezoidal rule as a time advancing scheme.

In this problem, our main concern is to validate if $q$ in equation (5.8) is approximately equal to the order of the time advancing scheme, meaning that we wish to check if $q \approx 2$. Since region $\Omega_t$ is constant in time, there is not any CFL condition associated to this problem.

We will also be interested in studying how to the errors behave while changing parameters $h$ and $\Delta t$ separately.

### *Changing $\Delta t$ while keeping $h$ fixed*

For a fixed spatial discretization step $h \approx 0.2$ and varying $\Delta t$, we get the following results for the errors:

Figure 5.21: Errors of $\widetilde{u}$ vs $\Delta t$ for problem 9



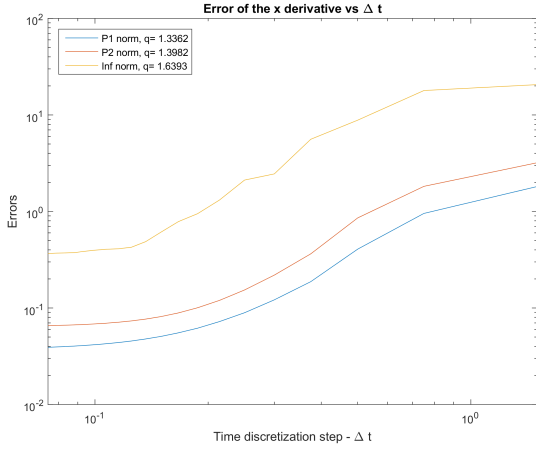Figure 5.22: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs $\Delta t$ for problem 9



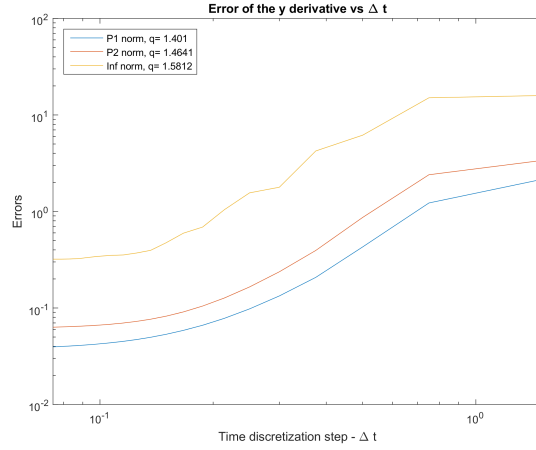Figure 5.23: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $\Delta t$ for problem 9
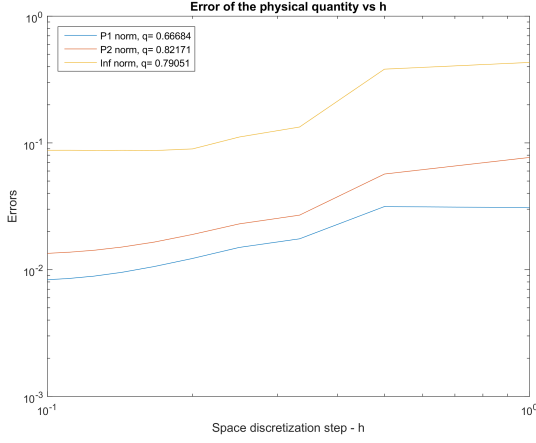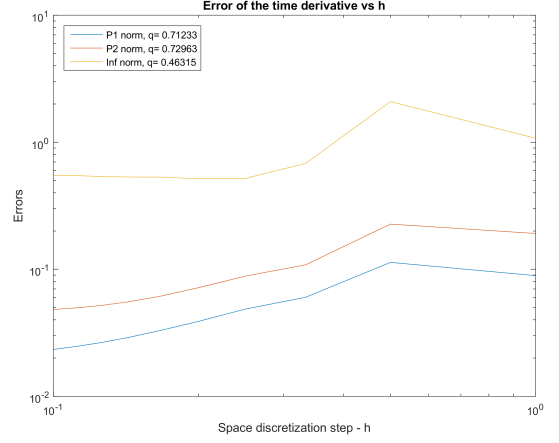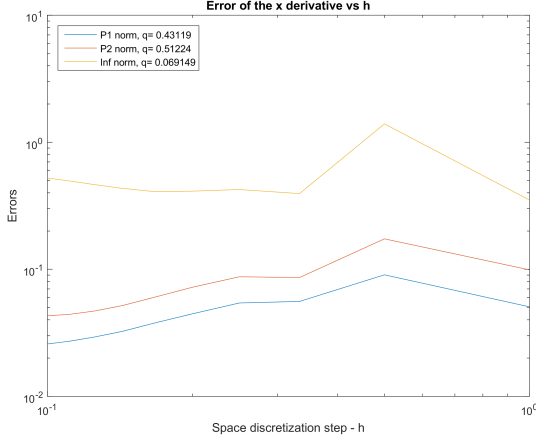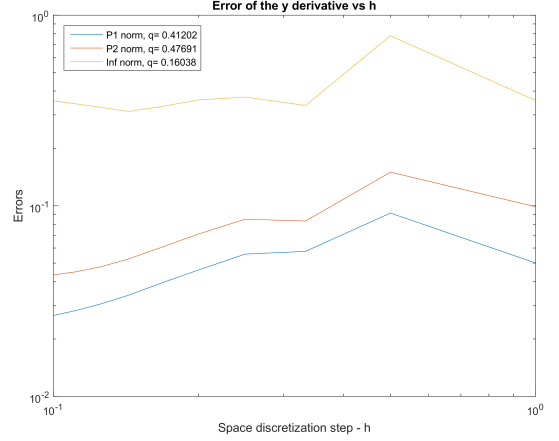


Figure 5.24: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $\Delta t$ for problem 9

From the above figures, we can see that all of the errors decrease with a decrease in $\Delta t$, but up to a certain extent since, when $\Delta t$ becomes very small, the curves become flat. This is because the errors of the spatial discretization become more prominent than those of the temporal discretization.

### *Changing $h$ while keeping $\Delta t$ fixed*

For a fixed temporal discretization step $\Delta t \approx 0.3$ and varying $h$, we get the following results for the errors:

**71**

Figure 5.25: Errors of $\widetilde{u}$ vs $h$ for problem 9



Figure 5.26: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs $h$ for problem 9



Figure 5.27: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs $h$ for problem 9



Figure 5.28: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs $h$ for problem 9

From the above figures, we can see that there is not much change for the errors with respect to $h$, either for $\widetilde{u}$ or for its derivatives.

### *Changing $h$ and $\Delta t$ together*

Now we decide to change parameters $h$ and $\Delta t$ together. We wish to study how the errors behave for this case, and check if $q \approx 2$. For this, we start with a simulation with discretization parameters $h_0$ and $\Delta t_0$, perform other simulations by dividing $h_0$ by 2 and $\Delta t_0$ by 2, and compare the errors to check if we have quadratic convergence. If yes, then indeed $q \approx 2$.

To include the contribution of both discretization steps, we graph the errors vs parameter $h + \Delta t$. By doing so, we get the following graph of the errors:
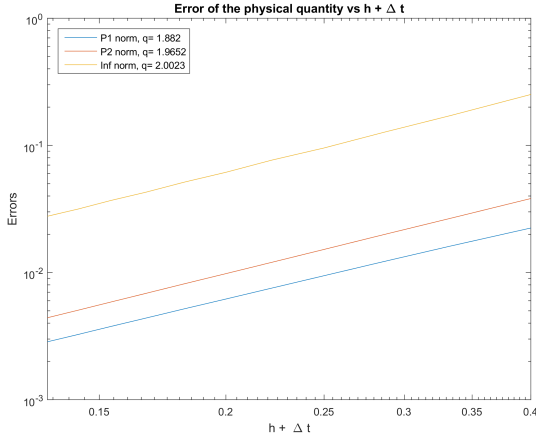
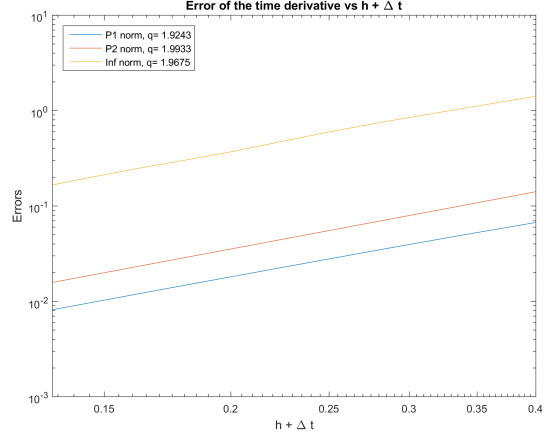Figure 5.29: Errors of $\widetilde{u}$ vs parameter $h+\Delta t$ for problem 9



Figure 5.30: Errors of $\frac{\partial \widetilde{u}}{\partial t}$ vs parameter $h+\Delta t$ for problem 9
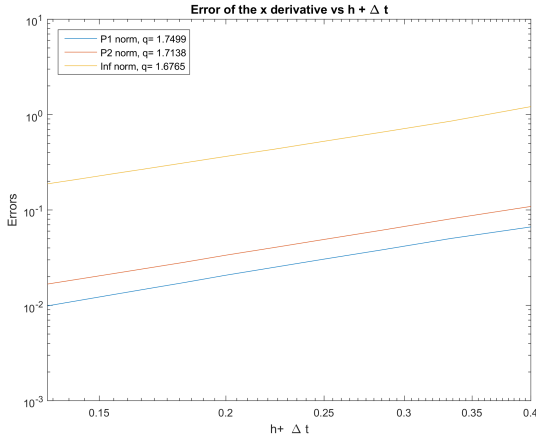


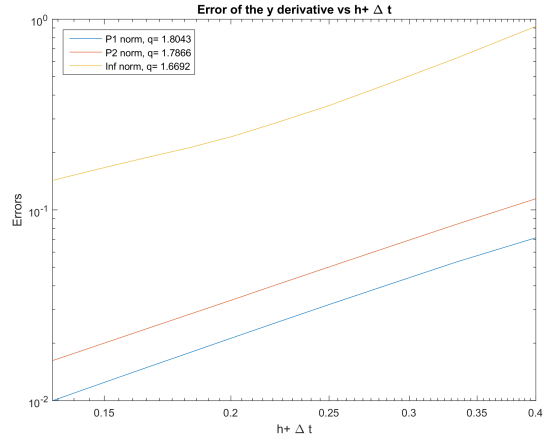Figure 5.31: Errors of $\frac{\partial \widetilde{u}}{\partial x}$ vs parameter $h+\Delta t$ for problem 9



Figure 5.32: Errors of $\frac{\partial \widetilde{u}}{\partial y}$ vs parameter $h+\Delta t$ for problem 9

From the above figures, 3 things are to be noted:

- We do get quadratic convergence for $\widetilde{u}$ for all norms, so indeed $q \approx 2$.

- We have super linear convergence for $\nabla u$ for all norms.

- The errors of $\frac{\partial \widetilde{u}}{\partial t}$ also quadratically converge.

### *Graphs of interest*

Also, we show several graphs of $\widetilde{u}$, $\widetilde{f}$, and the absolute error $|\widetilde{u} - u|$ at several time instances:
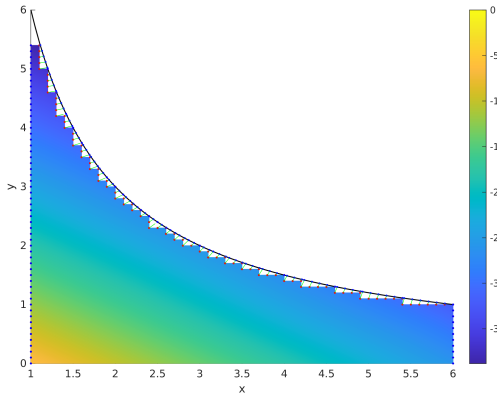
**73**

Figure 5.33: Graph of $\widetilde{u}$ at time $t = 1$ for problem 9
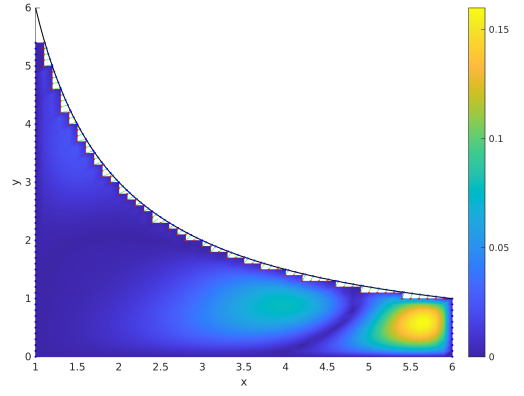


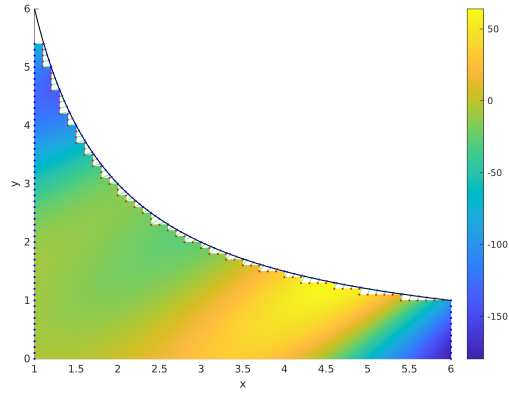Figure 5.34: Graph of $|\widetilde{u} - u|$ at time $t = 1$ for problem 9



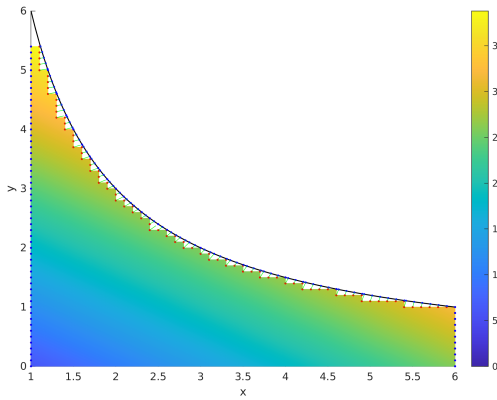Figure 5.35: Graph of $\widetilde{f}$ at time $t = 1$ for problem 9



Figure 5.36: Graph of $\widetilde{u}$ at time $t = 5$ for problem 9
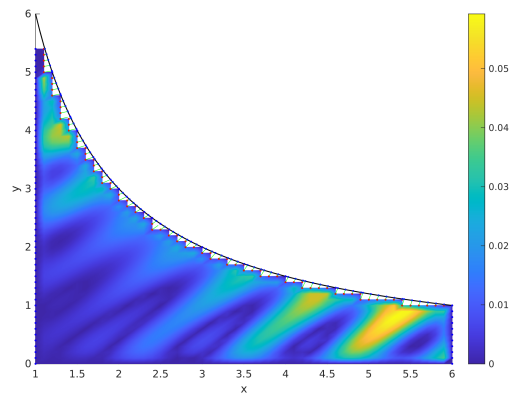


Figure 5.37: Graph of $|\widetilde{u} - u|$ at time $t = 5$ for problem 9
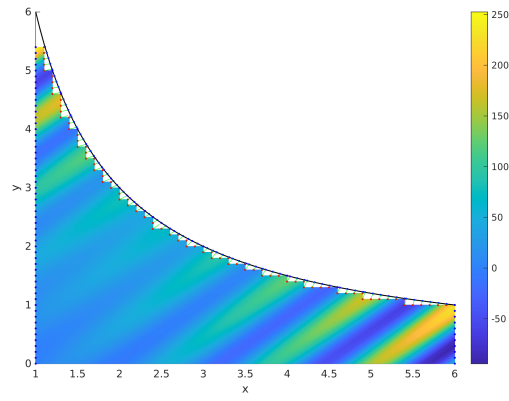
Figure 5.38: Graph of $\widetilde{f}$ at time $t = 5$ for problem 9

# Chapter 6

# Conclusions

The aim of this thesis was to instruct and show the reader the benefits and extents of the Shifted Boundary Method. For this, we made continuous progress through the different chapters, with detailed explanations and testing.

On chapter 1, we started with a simple problem which is common in the literature. On chapter 2, we focused on how to apply boundary conditions through Nitsche's method, necessary for the procedure on chapter 3, in which we showed how to shift boundary conditions. Chapters 4 and 5 were dedicated on the extensions of the main problem from chapter 1.

From our results and studies, we can conclude that:

- Shifting boundary conditions is a useful technique when the numerical boundary does not match with the physical boundary.

- Shifting boundary conditions increases the convergence of the FEA up to the minimum order between the approximation done by the test functions, and function $d$. Because of this, it is useful if both, the test functions and function $d$, have the same order of approximation.

- For all of our case of studies done up to chapter 4, we could appreciate quadratic convergence for $\widetilde{u}$ and linear convergence for $\nabla\widetilde{u}$ with respect to the spatial discretization step.

- For the time dependent cases studied in chapter 5, we can appreciate that the rate of decay of the error with respect to the time discretization matches the order of the time advancing scheme.

- For time dependent cases, since we have spatial and time discretizations, we obtained the best results for convergence when both of the discretization steps where changed accordingly to their respective order of convergence.

- Also we saw that the CFL condition does not restrict the choice of spatial and temporal discretization steps, since there was no correlation between the magnitude of the errors and the CFL number.

# Bibliography

[1] N. M. Atallah, C. Canuto, and G. Scovazzi, *Analysis of the Shifted Boundary Method for the Poisson Problem in General Domains*, arXiv:2006.00872v1 [math.NA] 1 Jun 2020.

[2] N. M. Atallah, C. Canuto, and G. Scovazzi, *Analysis of the Shifted Boundary Method for the Stokes Problem*, Comput. Methods Appl. Mech. Eng., 358:112609, 2020.

[3] N. M. Atallah, C. Canuto, and G. Scovazzi, *The Shifted Boundary Method for the Darcy flow problem*, J. Comput. Phys., 2020. in preparation.

[4] N. M. Atallah, C. Canuto, and G. Scovazzi, *The Second Generation Shifted Boundary Method and Its Numerical Analysis*, Comput. Methods Appl. Mech. Engrg. 372 (2020), 113341

[5] A. Main, and G. Scovazzi, *The shifted boundary method for embedded domain computations. Part I: Poisson and Stokes problems*, J. Comput. Phys. 372 (2018), 972-995.

[6] A. Main, and G. Scovazzi, *The shifted boundary method for embedded domain computations. Part II: Linear advection–diffusion and incompressible Navier–Stokes equations*, J. Comput. Phys. 372 (2018), 996-1026.

[7] Lecture notes on "Numerical Modelling" by C. Canuto, Politecnico di Torino.

[8] K. J. Bathe, *Finite Element Procedures*, Prentice Hall, 1996.

[9] A. Bonito, I. Kyza, and R. Nochetto, *Time Discrete Higher Order ALE Formulations: Stability*, SIAM J. Numer. Anal. 51 (2012), 577-604.

[10] L. Formaggia and F. Nobile, *A stability analysis for the arbitrary Lagrangian Eulerian formulation with finite elements*, East-West J. Numer. Math., 7(2):105–131, 1999.

[11] R. Aris, *Vectors, Tensors, and the Basic equations of Fluid Mechanics*, Dover Publications Inc, Berlin, 1989.