

UNIVERSITÄT BIELEFELD

FAKULTÄT FÜR MATHEMATIK

**Bachelorarbeit**

im Studiengang Mathematik

zum Thema:

---

Primal-Dual Active Set Strategy for  
Optimal Control Problems with Partial  
Differential Equations and Control  
Constraints

---

vorgelegt von

**Juan Carlos Chavarría Morales**

Matrikel-Nr: 4016452

1. Erstgutachterin: Dr. Tabea Tscherpel
2. Zweitgutachter: Prof. Dr. Lars Diening

Bielefeld, 6. August 2021

## ABSTRACT

This work is devoted to the implementation of a numerical strategy to solve Optimal Control Problems containing partial differential equations as restrictions. In addition, we have considered an admissible set restriction for the control variable.

Optimal Control Problems involving the type of restrictions above mentioned have gained importance in the last two decades and have been extensively studied [15],[24],[14], etc. Applications could be found in fields such diverse as Physics, Optimal Design, Pure and Applied Mathematics and Medicine [24], [14], [16]. Therefore, the scope of this Bachelor-Thesis is wide and oriented to audiences with a solid background in Mathematics, but whose main interests could lie in other disciplines.

The Literature [24], [14] consulted normally treats Optimal Control Problems (OCP) starting with classical Poisson equation and then jumping to more complicated models such as Navier-Stokes system. The first goal of this work is to fill this gap introducing an intermediate model in order to simplify the study of more complicated OCP using accessible mathematical language. Nevertheless, the OCP considered is *per se* interesting since it arises in applied sciences such as Fluid Dynamics [11].

Nowadays there exist several computer tools to solve OCP like the one described in this work [3], [21], [1]. In addition, knowledge of concepts from Functional and Numerical Analysis are necessary to provide an adequate theoretical context [24], [14]. The interface from theory to computer implementation requires Finite Element Method (FEM) also necessary as a prerequisite [24], [14]. The second and most important contribution of this work is precisely to use this mathematical framework to computationally implementate a numerical scheme and then compare with the results from a new software released this year [3]. In the same direction, there are good chances after this work to implementate our approach to a wide type of OCP.

Future research include also the study of another type of OCP, for example with new EDP restrictions, or the implementation of different resolution algorithms. The question that remains is the possibility of extending our approach study its limitations. Whatever the case, the learning achieved in this work allows tackling new challenges in the area of Optimal Control with the tools acquired in the development of this Thesis and the support of the Bachelor courses.

## KURZFASSUNG

Diese Arbeit widmet sich der Anwendung einer numerischen Strategie, um optimale Steuerungsprobleme zu lösen, die partielle Differentialgleichungen als Beschränkungen enthalten. Außerdem haben wir eine zulässige Mengenbeschränkung für die Steuerungsvariable betrachtet.

Optimale Steuerungsprobleme (OSP), die die oben genannten Beschränkungen enthalten, haben in den letzten zwei Jahrzehnten an Bedeutung gewonnen und sind intensiv untersucht worden [15],[24],[14] etc. Anwendungen konnten in diversen Bereichen wie z.B. Physik, optimales Design, reine und angewandte Mathematik und Medizin gefunden werden [24], [14], [16]. Daher ist der Umfang dieser Bachelor-Arbeit weit und für Leserinnen und Leser mit einem soliden mathematischen Hintergrundwissen ausgelegt, deren Hauptforschungsinteresse jedoch in anderen Disziplinen liegt.

Die gängige Literatur zum Thema OSP [24], [14] umfasst normalerweise optimale Steuerungsprobleme, die mit der klassischen Poisson-Gleichung beginnen und dann zu komplizierteren Modellen wie dem Navier-Stokes-System übergehen. Das Hauptziel dieser Arbeit ist es, diese Lücke durch Einführen eines Zwischenmodells zu schließen und so die Untersuchung komplizierterer OSP zu erleichtern. Dabei soll eine zugängliche mathematische Sprache genutzt werden. Gleichwohl ist das hier betrachtete OSP an sich interessant, weil es in angewandten Wissenschaften wie z.B. der Strömungsdynamik vorkommt [11].

Heutzutage existieren verschiedene Computerprogramme zur Lösung von wie dem in dieser Arbeit beschriebenen [3], [21], [1] OSP. Darüber hinaus ist als theoretischer Hintergrund die Kenntnis von Konzepten der funktionellen und numerischen Analyse erforderlich. Die Schnittstelle von der Theorie zur Computerimplementierung erfordert als Voraussetzung die Finite-Elemente-Methode (FEM) [24], [14]. Der zweite und wichtigste Beitrag dieser Arbeit ist es, dieses mathematische Rahmenwerk zu nutzen, um ein numerisches Schema anzuwenden und dann mit den Ergebnissen eines in diesem Jahr erschienenen Software-Programms zu vergleichen [3]. So bestehen auch gute Möglichkeiten, mithilfe des Ansatzes dieser Arbeit weitere Arten von OSP zu bearbeiten.

Zukünftige Forschungen könnten auch die Untersuchung anderer Arten von OSP beinhalten, z.B. mit neuen EDP-Beschränkungen oder der Anwendung verschiedener Lösungsalgorithmen. Es bleibt die Frage, inwieweit sich der erarbeitete Ansatz erweitern lässt und wo er an seine Grenzen stößt. In jedem Fall ermöglichen die Erkenntnisse dieser Arbeit, neue Herausforderungen im Bereich der optimalen Steuerung mithilfe der erarbeiteten Hilfsmittel anzunehmen. Nicht zuletzt können diese ein entlastendes Hilfsmittel in Bachelorkursen darstellen.

## DECLARATION OF ORIGINALITY

I hereby declare, Juan Carlos Chavarría Morales, that I wrote the present work independently and without outside help. I also assure that I have not used any other sources than those given and that I have marked all statements taken literally or analogously from other works as such, and that the submitted work was neither completely nor in essential parts the subject of another examination procedure.

Bielefeld, 6 August 2021

## EIGENSTÄNDIGKEITSERKLÄRUNG

Hiermit erkläre ich, Juan Carlos Chavarría Morales, dass ich die hier vorliegende Arbeit eigenständig und ohne fremde Hilfe verfasst habe. Ebenso versichere ich, dass ich keine anderen als die angegebenen Quellen und alle wörtlich oder sinngemäß aus anderen Werken übernommenen Aussagen als solche gekennzeichnet habe, und dass die eingereichte Arbeit weder vollständig noch in wesentlichen Teilen Gegenstand eines anderen Prüfungsverfahrens gewesen ist.

Bielefeld, 6. August 2021

# CONTENTS

|   |           |
|---|-----------|
| <b>1. Introduction and Motivation</b>                                       | <b>1</b>  |
| 1.1. Motivation . . . . .   | 2         |
| <b>2. Optimal Control Problems</b>  | <b>4</b>  |
| 2.1. Introduction . . . . .   | 4         |
| 2.2. Optimal Control Problem for ADR Equation . . . . .                     | 5         |
| <b>3. Optimality Theory</b>   | <b>8</b>  |
| 3.1. Existence of Optimal Controls . . . . .                                | 8         |
| 3.2. First Order Optimality Conditions . . . . .                            | 9         |
| 3.2.1. Quadratic Optimization in Hilbert Spaces . . . . .                   | 9         |
| 3.3. Optimality for the ADR problem . . . . .                               | 10        |
| 3.3.1. On the pointwise optimality conditions . . . . .                     | 12        |
| <b>4. Finite Element Approximation of the Optimal Control Problem</b>       | <b>14</b> |
| 4.1. A Finite Element Method Scheme . . . . .                               | 14        |
| 4.2. Convergence of the Finite Element Method . . . . .                     | 14        |
| <b>5. Primal-Dual Active Set (PDAS) Strategy to Solve Optimality System</b> | <b>19</b> |
| 5.1. PDAS Strategy in Continuous Problem . . . . .                          | 19        |
| 5.2. Mixed formulation of Optimality System . . . . .                       | 20        |
| <b>6. Numerical Results</b>   | <b>23</b> |
| 6.1. Parameters and exact solutions . . . . .                               | 23        |
| 6.2. Implementation in <b>FENICS</b> . . . . .                              | 23        |
| 6.2.1. Numerical Results . . . . .  | 24        |
| 6.3. Implementation in <b>cashocs</b> . . . . .                             | 25        |
| <b>7. Conclusions and Outlook</b>   | <b>28</b> |
| <b>A. Implementation in <b>fenics</b></b>                                   | <b>B</b>  |
| A.1. Main Program . . . . .   | B         |
| A.2. Optimality system Solver . . . . .                                     | D         |
| A.3. $L^2$ -Projection . . . . .  | E         |
| A.4. Indicator Functions . . . . .  | E         |
| A.4.1. $\chi_a$ . . . . .   | E         |
| A.4.2. $\chi_b$ . . . . .   | E         |
| <b>B. Implementation in <b>cashocs</b></b>                                  | <b>F</b>  |
| B.1. Main Program . . . . .   | F         |

# 1. INTRODUCTION AND MOTIVATION

Most of the human endeavours involve, or at least try, some degree of control in the behavior of some phenomena or system. This generates the question: Since we know how a system will behave, how can we *change it* to obtain a *desired* response? The goal of the Control Theory is to answer this question. In particular, this work deals on a system modeled by a partial differential equation (PDE).

A careful reading shows that the *Control Theory* and the *Optimal Control Theory* answer two different but related, questions. Let's analyze this: A general form to formulate a control problem is the following. Given a system, and a way to control it, which can be formulated as

$$A.y = B.u,$$

where  $A$  and  $B$  are operators, which specify the model of the system and how the control acts on it respectively;  $y$  is the state of the system and  $u$  is the control, taken from an admissible control set  $U_{ad}$ . If in addition a *cost functional*  $J(u)$  is given, which expresses the *cost of using the control*  $u$ , the question that answers the Optimal Control Theory is to find a control  $u$ , in the space of admissible controls  $U_{ad}$ , such that

$$\min_{u \in U_{ad}} J(u).$$

In this case, the relation between the control  $u$  and the state  $y$  is given by a PDE and the boundary conditions. We must proceed carefully in the election of the set of admissible controls. A poor decision will minimise the cost functional, but if this control does not produce the desired results on the system, then the resolution effort is pointless.

Our goal is then to determine an admissible control, called *optimal control* which gives us a satisfactory state and minimizes the functional  $J$ . The next question to answer is the existence (and unicity) of a pair *control, state* and its respective calculation. To obtain the general solution we must use normally numerical methods. The first step to solve numerically requires a discretisation of the control problem. To this end we proceed using finite elements (see, for example, [24] or [2]). An immediate consequence of this approach is the error- and convergence analysis of the method. To perform this analysis we will see that it is necessary to impose certain regularity conditions on the optimal control. We can derive this framework from the first order optimality conditions. Another important tool in the error estimates are the second order conditions, which in our case (linear differential equation and quadratic functional cost) are determined by the problem. Once we have a discrete control problem, we must use some optimisation algorithm to solve. In summary, we can list the difficulties in carrying out the study of our optimal control problem as follows [5]:

1. Existence of a solution.
2. Optimality conditions of first and second order (System of PDE's).
3. Discretization of Optimality System.
4. Numerical solution of the associated discrete problem and convergence analysis.

It is important to mention, that this framework corresponds to (1) Continuous Optimisation Problem (2) Optimise (3) Discretise and (4) Solve scheme [10]. However another schemes are possible by switching steps (2) and (3).

## 1.1. MOTIVATION

This work is devoted to a control problem for the Advection–Diffusion–Reaction (ADR) equation, whose derivation comes from fluid mechanics [25]. Let  $y$  be a scalar function which represents a physical property, such as heat or concentration. If  $y$  is conserved and can only vary through exchange between material particles, reaction of the system itself or through external sources. Then the conservation law for  $y$  is given by

$$\partial_t y - \epsilon \Delta y + \underline{b} \cdot \nabla y + \kappa y = g, \quad \text{in some domain } \Omega, \quad (1.1)$$

where  $\epsilon$  is the diffusion coefficient,  $\kappa$  is the reaction parameter (which is assumed positive, but could be negative) and  $\underline{b}$  is the velocity field associated with the advective process. Although the ADR equation arises naturally in fluid mechanics, is often found in several applications such as Economy (Black-Scholes equation), acoustics, heating transfer, medicine and many others [11], [16]. Also, we normally need to prescribe some boundary conditions in  $\Omega$ . These conditions have a capital importance, because depending on the type of boundary condition (Dirichlet, Neumann or Robin) will determine the kind of optimal control problem to be solved and also the numerical scheme. We will give some directions in the next section. We will consider in particular the stationary ADR equation, which leads to an elliptic partial differential equation. Roughly speaking, the control  $u$  will be inserted into the source term in order to produce a desire response, or state of the system  $y$ . When is possible to do minimising the value of the cost functional  $J$  we say that our solution pair  $(y, u)$  is optimal.

We have extensively used two standard texts from the optimal control literature with PDE's constraints, [15], [24]. In these texts the study of optimal control problems begins with the study of Poisson equation:

$$\begin{aligned} & \min J(y, u) \\ \text{subject to: } & -\Delta y = u, \quad \text{in some domain } \Omega \\ & u \in U_{ad}, \end{aligned}$$

plus some boundary condition for the state variable  $y$ . Here  $U_{ad}$  denotes the admissible set of the control variable  $u$ . However, the path to study other optimal control problems governed by more complex equations of state can often be abrupt for the inexperienced reader. Nonetheless, it is possible to find literature dedicated to the Optimal Control Problem for the ADR equation [17], but normally the prerequisites needed could be too technical for an undergraduate mathematic student or a non mathematic reader. This work tries to fill this gap by providing a basic language to explain and solve an Optimal Control Problem which is halfway between the aforementioned Poisson model and more complex problems.

Pseudo Newton methods like the *Primal-dual Active Set Strategy* [15], [24], are extensively used in Optimal Control Problems because their computational efficiency

performing the iterations to find the optimal solution, where the control problem is formulated as a quadratic minimisation problem in which the restrictions are (1) the state equation, in this case a PDE; (2) the boundary conditions and (3) the possible values that the control can reach, the so-called *Admissible Set*. Generally speaking, following the list above and [15], in order to apply this strategy we must first have existence and uniqueness of a solution of the PDE and an optimality condition for the control. This has determined the outline of this work: First section deals with existence and uniqueness of a solution determined by the restriction given by the PDE and the boundary condition. This is achieved using standard arguments from Functional Analysis,[24] ,[4], [2]. In the next section, we determined optimality conditions, where we replicated the scheme showed in [24]. In the last two sections we implementate a Finite Element Method (FEM) to solve our problem together the Primal-dual Active Set Strategy, perform the respective error analysis and show the numerical resolution . For a detailed discussion on this issue we suggest [6] or [24]. In the implementation stage (Section 5.2)cwe have introduced a Mixed Formulation to simplify the computational effort.



## 2. OPTIMAL CONTROL PROBLEMS

### 2.1. INTRODUCTION

Our problem belongs to a large class of minimisation problems. In an abstract form we formulate

$$\min_{u \in U_{ad}} J(y, u), \quad (2.1)$$

where  $U_{ad}$  is a closed subset of admissible control functions [22] and  $y \in Y$  for some function space to specify which solves the problem

$$\mathcal{D}(y, u) = 0. \quad (2.2)$$

Here  $\mathcal{D}$  is some differential operator. We assume that  $J$  is bounded from below and that 2.2 is well posed. Our attention lies in a particular subclass of problems, namely

$$\mathcal{D}(y) = u + g, \quad (2.3)$$

where  $g$  is a given function in the image space of the operator  $\mathcal{D}$ . In other optimization problems the control  $u$  may also enter the differential operator (for example shape optimization, [19]). We are not concerned in this type of problems, but the simplest case, when  $y = 0$  in  $\Gamma$  and 2.3 is obtained solving: find

$$y \in H_0^1(\Omega) = \{y \in H^1(\Omega) : y = 0 \text{ on } \Gamma\} \quad (2.4)$$

such that

$$B(y, v) = (u + g, v), \quad \forall v \in H_0^1(\Omega). \quad (2.5)$$

Here  $\Omega$  is a bounded connected domain in  $\mathbb{R}^d$  with Lipschitz boundary  $\Gamma$  [24] and the functions  $u$  and  $g$  belong to  $L^2(\Omega)$ . We have considered Dirichlet boundary condition in this work, but other boundary conditions could be in an analog way presented modifying the cost functional  $J$  [24].  $B$  is a bilinear form which represents the weak formulation of an elliptic differential operator. For this problems, we consider integral type cost functionals

$$J(y, u) := \frac{1}{2} \|y - y_\Omega\|_{L^2(\Omega)}^2 + \frac{\lambda}{2} \|u\|_{L^2(\Omega)}^2, \quad \lambda > 0,$$

and the set of admissible control functions includes *box* constraints, namely,

$$U_{ad} := \{u \in L^2(\Omega) : u_a \leq u(x) \leq u_b, \text{ a.e. } x \in \Omega\}.$$

It is well known that under the above assumptions these optimal control problems have unique solutions (e.g., see Section 2, Chapter 2 in [18]). Although literature [24], [14], etc. normally enfold the theory through the optimal control problem for the Poisson equation, we focus our efforts in the optimal control problem for the ADR equation. This problem has *per se* a physical relevance [11], but also from the pure mathematic point of view is interesting because it is useful to introduce more complicated control problems like Navier Stokes equations [7].

## 2.2. OPTIMAL CONTROL PROBLEM FOR ADR EQUATION

Our work is devoted to the *distributed* Optimal Control Problem for the stationary ADR equation

$$\min J(y, u) \quad \text{where} \quad J(y, u) := \frac{1}{2} \|y - y_\Omega\|_{L^2(\Omega)}^2 + \frac{\lambda}{2} \|u\|_{L^2(\Omega)}^2, \quad (2.6)$$

subject to the restrictions:

1. The *system to state* condition, given by the Advection-Diffusion-Reaction (ADR) equation, with homogeneous Dirichlet boundary condition

$$\begin{cases} -\varepsilon \Delta y + \underline{b} \cdot \nabla y + \kappa y &= u + g & \text{in } \Omega, \\ y &= 0 & \text{on } \Gamma, \end{cases} \quad (2.7)$$

2. The *Admissibility or Box* condition for the control variable  $u$

$$u \in U_{ad} := \{u \in L^2(\Omega) : u_a \leq u(x) \leq u_b, \quad \text{a.e. } x \in \Omega\}. \quad (2.8)$$

We say that the control problem is *distributed* because the control  $u$  acts in the whole domain  $\Omega$ . Depending on the boundary condition the cost functional  $J$  could be adapted to the control requirements. See for example [24], [14], etc. We make the following assumptions:

**SUPPOSITION 2.1.** *For the problem (2.6) – (2.8) we have that:*

- $\lambda > 0, \varepsilon > 0, \kappa \geq 0$  are given scalars.  $y_\Omega \in L^2(\Omega)$ , represents the desire state,  $g \in L^2(\Omega)$  and  $\underline{b} \in W^{1,\infty}(\Omega)$  is a null divergence vector field in  $\Omega$ .
- The bounds  $u_a < u_b$  define the sets of admissible controls acting on  $\Omega$ . The elements of  $U_{ad}$  are called *admissible Controls*.

We define the norm of our solution space  $H_0^1(\Omega)$  as

$$\|v\|_\Omega := \left( \varepsilon \|\nabla v\|_{L^2(\Omega)}^2 + \kappa \|v\|_{L^2(\Omega)}^2 \right)^{1/2},$$

which is equivalent to the norm  $H^1(\Omega)$ . In fact, the next existence result for (2.7) explains this claim.

**THEOREM 2.1 Existence.** *Under the assumptions (2.1), for each  $u \in U_{ad}$  there exists an unique solution  $y \in H_0^1(\Omega)$  of (2.7). Moreover, there exists a constant  $c_P > 0$ , which only depends on  $\Omega$  such that*

$$\|y\|_\Omega \leq c_P \left( \|g\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)} \right). \quad (2.9)$$

*Proof.* Following [18] the idea of the proof is to check Lax–Milgram’s Lemma in order to obtain the well-posedness of the problem and the uniqueness of a solution when the

data  $u + g$  is given.<sup>1</sup> To do this, we use the weak problem associated to (2.7). In other words, we want to find  $y \in H_0^1(\Omega)$  such that  $B(y, v) = F(v)$ , for all  $v \in H_0^1(\Omega)$ , where

$$B(y, v) = \int_{\Omega} (\varepsilon \nabla y \cdot \nabla v + \underline{b} \cdot \nabla y v + \kappa y v) dx \quad (2.10)$$

$$G(v) = \int_{\Omega} (u + g) v dx. \quad (2.11)$$

In order to apply the result of Lax-Milgram we must verify that: (1)  $B$  is continuous, which is  $|B(u, v)| \leq C_{cont} \|y\|_{\Omega} \|v\|_{\Omega}$ ,  $\forall y, v \in H_0^1(\Omega)$ . (2)  $B$  is coercive. That is  $B(v, v) \geq C_{coer} \|v\|_{\Omega}^2$ ,  $\forall v \in H_0^1(\Omega)$  and (3)  $G$  is continuous in the sense:  $G(v) \leq c \|v\|_{\Omega}$ ,  $\forall v \in H_0^1(\Omega)$ . First, we have

$$\begin{aligned} |B(y, v)| &= \left| \int_{\Omega} (\varepsilon \nabla y \cdot \nabla v + \underline{b} \cdot \nabla y v + \kappa y v) dx \right| \\ &\leq \varepsilon \|\nabla y\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} + \|\underline{b}\|_{L^\infty(\Omega)} \|\nabla y\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} + \kappa \|y\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \\ &= \sqrt{\varepsilon} \|\nabla y\|_{L^2(\Omega)} \sqrt{\varepsilon} \|\nabla v\|_{L^2(\Omega)} + \frac{\|\underline{b}\|_{L^\infty(\Omega)}}{\sqrt{\varepsilon \kappa}} \sqrt{\varepsilon} \|\nabla y\|_{L^2(\Omega)} \sqrt{\kappa} \|v\|_{L^2(\Omega)} \\ &\quad + \kappa \|y\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \\ &\leq \max \left\{ 1, \frac{\|\underline{b}\|_{L^\infty(\Omega)}}{\sqrt{\varepsilon \kappa}} \right\} (\sqrt{\varepsilon} \|\nabla y\|_{L^2(\Omega)} + \sqrt{\kappa} \|y\|_{L^2(\Omega)}) (\sqrt{\varepsilon} \|\nabla v\|_{L^2(\Omega)} + \sqrt{\kappa} \|v\|_{L^2(\Omega)}) \\ &\leq 2 \max \left\{ 1, \frac{\|\underline{b}\|_{L^\infty(\Omega)}}{\sqrt{\varepsilon \kappa}} \right\} \|y\|_{\Omega} \|v\|_{\Omega}. \end{aligned}$$

We have used Cauchy-Schwarz and  $(a + b)^2 \leq 2(a^2 + b^2)$  inequalities and *built* the norm of the state space. Defining  $C_{cont} = 2 \max \left\{ 1, \frac{\|\underline{b}\|_{L^\infty(\Omega)}}{\sqrt{\varepsilon \kappa}} \right\}$  we are done.

To proof the second statement, we note that  $\forall v \in H_0^1(\Omega)$ :

$$B(v, v) = \int_{\Omega} (\varepsilon \nabla v \cdot \nabla v + \kappa v^2) dx + \int_{\Omega} \underline{b} \cdot \nabla v v dx = \|v\|_{\Omega}^2 + \int_{\Omega} \underline{b} \cdot \nabla v v dx.$$

We claim that the second term in the right hand side is null. To see this, we use the identity  $\nabla \cdot (v^2 \underline{b}) = 2v \nabla v \cdot \underline{b} + v^2 \nabla \cdot \underline{b}$ , which gives

$$\int_{\Omega} \underline{b} \cdot \nabla v v dx = \frac{1}{2} \int_{\Omega} \nabla \cdot (v^2 \underline{b}) dx - \frac{1}{2} \int_{\Omega} v^2 \nabla \cdot \underline{b} dx.$$

The second integral in the right hand side is zero since  $\underline{b}$  is a null divergence vector field. We use Gauß theorem for the remaining integral to obtain

$$\int_{\Omega} \underline{b} \cdot \nabla v v dx = \frac{1}{2} \int_{\Gamma} (v^2 \underline{b}) \cdot d\sigma = 0, \quad \text{since } v \in H_0^1(\Omega).$$

---

<sup>1</sup>The result can be found in numerous texts of Functional Analysis. For the sake of simplicity we use the standard reference of this report, [24].

Then  $B(v, v) = \|v\|_\Omega^2 = C_{coerv} \|v\|_\Omega^2$ , with  $C_{coerv} = 1$  and we have our result. Finally, using Cauchy–Schwarz inequality and the definition of norm of the state space we obtain

$$G(v) \leq \|g + u\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} = \frac{\|g + u\|_{L^2(\Omega)}}{\sqrt{\kappa}} \sqrt{\kappa} \|v\|_{L^2(\Omega)} \leq \frac{\|g\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)}}{\sqrt{\kappa}} \|v\|_\Omega.$$

Choosing  $c = \frac{\|g\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)}}{\sqrt{\kappa}}$  we have that  $G$  is continuous.

For the last part, if we set  $v = y$  and use the previous results in the weak formulation we obtain

$$\begin{aligned} C_{coerv} \|y\|_\Omega^2 \leq B(y, y) &= G(y) \leq \frac{\|g\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)}}{\sqrt{\kappa}} \|y\|_\Omega \\ \Rightarrow \|y\|_\Omega &\leq \frac{1}{\sqrt{\kappa}} (\|g\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)}), \end{aligned}$$

which gives the continuous dependence on data. The proof is then complete.  $\square$

**REMARK 2.1.** *The theorem (2.1) gives a relation between the state  $y$  and the control  $u$ , say  $y(u)$ . This allows us to write 2.6 in terms of  $u$ .*

Now we focus our attention on obtain optimality conditions for the optimal control problem.

**DEFINITION 2.1.** *We call a control  $\bar{u} \in U_{ad}$  optimal and its associated state optimal state  $\bar{y} = y(\bar{u})$  if*

$$J(\bar{y}, \bar{u}) \leq J(y(u), u), \quad \forall u \in U_{ad}.$$

Theorem 2.1 guarantees the existence of the following continuous linear operator in view of 2.9:

**DEFINITION 2.2.** *The Control to State operator is defined as  $\Xi : L^2(\Omega) \rightarrow H_0^1(\Omega)$ , which assigns a state  $y(u)$  to each control  $u$ .*

Since  $\|y\|_{L^2} \leq \|y\|_{H^1}$  then the state space  $H_0^1(\Omega) \subset H^1(\Omega)$  is linear embedded in  $L^2(\Omega)$  (see [24]). Therefore, we can *extend* the range of  $\Xi$  to  $L^2(\Omega)$  through the following operator  $E_Y : H^1(\Omega) \rightarrow L^2(\Omega)$  which to each function of  $H^1$  assigns the same element in  $L^2$ , which is also linear and continuous. In this way, through composition we can define the Solution Operator

$$S = E_Y \Xi : L^2(\Omega) \rightarrow L^2(\Omega),$$

where  $S(u) = y(u)$ . This operator allows us to work with admissible controls and reduce (2.9) to the following minimisation problem

$$\min_{u \in U_{ad}} f(u) := \frac{1}{2} \|Su - y_\Omega\|_{L^2(\Omega)}^2 + \frac{\lambda}{2} \|u\|_{L^2(\Omega)}^2. \quad (2.12)$$

We say that  $f$  is the *Reduce Functional*.

### 3. OPTIMALITY THEORY

Taking 2.12 as starting point, we develop the optimality conditions for our Optimal Control Problem. Like in Calculus in one or several variables, when one optimises a function some conditions for existence of critical points are needed. Further, one seeks for a classification of these critical points using to that end Hessian matrices and eventually Karush-Kuhn-Tucker conditions. A relief in this context is obtain when the objective function is convex, which ensures that the founded critical points are minimums. In this new context, we try to proof existence and unicity of functions (control functions, and through  $y = Su$  state functions) which minimise the reduced functional. Therefore, we need to extend the definitions from classical Calculus to function spaces and to profit from the convexity of the reduced functional 2.12. The first two parts of this section deal with these matters, while in the last one we apply the results of the previous subsections to the ADR problem, which leads to a first version of the optimality system which will allow us to solve our problem.

#### 3.1. EXISTENCE OF OPTIMAL CONTROLS

Our first result deals with the existence of an Optimal Control for the minimisation problem of the reduced functional (2.12). This argument will be replicated in other contexts and therefore we present the most general version. We refer to [24] or [14] for a detailed proof

**THEOREM 3.1.** *Let  $\{U, \|\cdot\|_U\}$  and  $\{H, \|\cdot\|_H\}$  two real Hilbert spaces and  $U_{ad} \subset U$  a bounded closed convex non-empty set, as well as  $y_d \in H$  and a constant  $\lambda \geq 0$  given. Moreover, let  $S : U \rightarrow H$  be a linear continuous operator. Then the quadratic optimization problem on Hilbert spaces*

$$\min_{u \in U_{ad}} f(u) := \frac{1}{2} \|Su - y_d\|_H^2 + \frac{\lambda}{2} \|u\|_U^2 \quad (3.1)$$

*admits an optimal solution  $\bar{u}$ . If besides  $\lambda > 0$  or  $S$  is injective, then the solution is uniquely determined.*

As a direct consequence of the preceding result, we obtain the existence and uniqueness for the elliptic optimal control problem (2.6)–(2.8) [24]:

**THEOREM 3.2.** *Suppose that the conditions of (2.1) are satisfied. Then, the problem (2.6)–(2.8) with  $g = 0$  possesses at least one optimal control  $\bar{u}$ . If, in addition  $\lambda > 0$ , the solution is unique.*

*Proof.* We apply the previous theorem for  $U = H = L^2(\Omega)$ ,  $y_d = y_\Omega$  and  $S = E_Y G$ . The set  $U_{ad} = \{u \in L^2(\Omega) : u_a \leq u \leq u_b \text{ a.e. in } \Omega\}$  is closed, bounded and convex. Then, from the theorem (3.1) it follows that the minimisation problem (2.12) admits at least one solution  $\bar{u}$  which is unique if  $\lambda > 0$ . If  $Su = 0$ , we have that the operator  $S$  for this problem is injective. In fact, if  $Su = 0$ , then  $y = 0$  for (2.9) and inserting this into the differential equation follows that  $u = 0$  almost everywhere in  $\Omega$ . This gives the injectivity of  $S$  and we have uniqueness for this case.  $\square$

**REMARK 3.1.** *In the proof of last theorem,  $\bar{u}$  we obtain the limit of a weakly convergent sequence  $\{u_{n_k}\}$ . Since the control-to-state operator  $\Xi : L^2(\Omega) \rightarrow H_0^1(\Omega)$  is a continuous linear operator, it is also weakly continuous. This implies that the series of states  $\{y_{n_k}\}$  converges weakly in  $H_0^1(\Omega)$  to  $\bar{y} = \Xi\bar{u}$ . Also, the case  $g \neq 0$  could be in an analog way treated simply defining  $w := u + g \in L^2(\Omega)$ , rescaling  $U_{ad}$  and repeating the previous proof.*

It is possible [24] extend the set  $U_{ad}$  allowing  $u_a = -\infty$  and/or  $u_b = +\infty$ . In this case,  $U_{ad}$  is not bounded and therefore is not sequentially weak compact. However, we still have existence for the case  $\lambda > 0$ .

### 3.2. FIRST ORDER OPTIMALITY CONDITIONS

The aim of this section is to provide similar arguments to the Calculus in several variables for optimality conditions in the context of our cost functional. The necessary conditions that we will find allow us to go deeper into the structure of an optimal control and the numerical verification of its optimality.

#### 3.2.1. QUADRATIC OPTIMIZATION IN HILBERT SPACES

In order to prove the existence of optimal controls, we have transformed the control problems under analysis into quadratic optimization problems reduced in terms of  $u$  to

$$\min_{u \in U_{ad}} f(u) := \frac{1}{2} \|Su - y_d\|_H^2 + \frac{\lambda}{2} \|u\|_U^2. \quad (3.2)$$

For this minimisation problem, the following fundamental result can be applied. It is the key to deriving first-order necessary conditions in the presence of control constraints and is in complete analogy to the well known theory for optimality conditions for functions in  $\mathbb{R}^d$ . The proof could be found in [24] or [14].

**LEMMA 3.1.** *If  $C$  denotes a non-empty convex subset of a real Banach space  $U$ , and  $f$  is a real Gâteaux-differentiable application on an open subset of  $U$  containing  $C$ . If  $\bar{u}$  is the solution of the problem*

$$\min_{u \in C} f(u),$$

*then it solves the inequality*

$$f'(\bar{u})(u - \bar{u}) \geq 0 \quad u \in C. \quad (3.3)$$

*Reciprocally, if  $\bar{u} \in C$  solves the above inequality and  $f$  is convex, then  $\bar{u}$  is a solution of the minimisation problem  $\min_{u \in C} f(u)$ .*

The above result gives a *necessary* condition, which is *sufficient* in the case of convexity. It is also called *first-order optimality condition*. Let us apply this result to the problem (3.2).

**THEOREM 3.3.** *Suppose that we have  $U, H$  Hilbert spaces, a non-empty convex subset  $U_{ad} \subset U$ , some  $y_d \in H$  and  $\lambda \geq 0$  are given. Moreover, let  $S : U \rightarrow H$  be a linear and*

continuous operator. Then  $\bar{u} \in U_{ad}$  is a solution of the minimisation problem (3.2) if and only if it satisfies the variational inequality

$$(S^*(S\bar{u} - y_d) + \lambda\bar{u}, u - \bar{u})_U \geq 0, \quad \forall u \in U_{ad}. \quad (3.4)$$

*Proof.* Usig the chain rule for Gâteaux and Fréchet derivatives (see, for example Theorem 2.20 in [24]), the gradient of the functional defined by (3.2) has the form

$$f'(\bar{u}) = S^*(S\bar{u} - y_d) + \lambda\bar{u}, \quad (3.5)$$

where  $S^*$  denotes the *adjoint* operator of  $S$ . Then our result follows as consequence of the previous Lemma.  $\square$

Next, we focus our attention in this operator  $S^*$ .

### 3.3. OPTIMALITY FOR THE ADR PROBLEM

The ADR problem defined in the first section states that under the assumptions (2.1) we must find a pair  $(y, u)$  such that minimises  $J(y, u)$  defined by 2.6 subject to 2.7 and 2.8.

We will denote by  $S$  to the solution operator of the boundary value problem, which according to the previous section muss satisfies the variational inequality (3.4) adapted to our case, that is

$$(S^*(S\bar{u} - y_\Omega) + \lambda\bar{u}, u - \bar{u})_U \geq 0, \quad \forall u \in U_{ad}, \quad (3.6)$$

where the adjoint operator  $S^*$  is yet to be determined. We replicate the analysis of [24] for our case in order to determine  $S^*$ . To that end, we refer to the next pair of problems: Find  $y, p \in H_0^1(\Omega)$ , such that for  $u, z \in L^2(\Omega)$  given, solve

$$\begin{cases} -\varepsilon\Delta y + \underline{b} \cdot \nabla y + \kappa y &= u + g & \text{in } \Omega, \\ y &= 0 & \text{on } \Gamma, \end{cases} \quad \begin{cases} -\varepsilon\Delta p - \underline{b} \cdot \nabla p + \kappa p &= z & \text{in } \Omega, \\ p &= 0 & \text{on } \Gamma, \end{cases} \quad (3.7)$$

We claim that the problem on the right is the dual problem associated to our state equation. To see this, we consider  $p \in H_0^1(\Omega)$  as test function for the problem of the left side to write the weak formulation of this problem, which yields

$$\begin{aligned} \int_{\Omega} \{ \varepsilon \nabla y \cdot \nabla p + (\underline{b} \cdot \nabla y + \kappa y) p \} dx &= \\ \int_{\Omega} \{ \varepsilon \nabla p \cdot \nabla y - \underline{b} \cdot \nabla p y + \kappa p y \} dx &= \int_{\Omega} (u + g) p dx. \end{aligned}$$

The later is the weak formulation of the problem at the right hand side in 3.7, where  $y \in H_0^1(\Omega)$  is used this time as test function. This allow us to conclude

$$\int_{\Omega} zy dx = \int_{\Omega} p(u + g) dx. \quad (3.8)$$

In view of our previous discussion  $S : L^2(\Omega) \rightarrow L^2(\Omega)$ , we have that  $S^* : L^2(\Omega) \rightarrow L^2(\Omega)$  and this is given by

$$(z, Su)_{L^2(\Omega)} = (S^*z, u)_{L^2(\Omega)}, \quad \forall z, u \in L^2(\Omega).$$

Now, taking  $u = u + g$ , we have that  $y = Su$ , and we find that

$$(z, Su)_{L^2(\Omega)} = (z, y)_{L^2(\Omega)} = (p, u + g)_{L^2(\Omega)},$$

according to (3.8). The mapping  $z \mapsto p$  is linear and continuous from  $L^2(\Omega) \rightarrow L^2(\Omega)$  and given that  $z, u$  could be arbitrary choose and  $S^*$  is uniquely determined. Summarizing:

**LEMMA 3.2.** *For the boundary value problem (2.7), the adjoint operator  $S^* : L^2(\Omega) \rightarrow L^2(\Omega)$  is given by*

$$S^*z = p,$$

where  $p \in H_0^1(\Omega)$  is solution of the boundary value problem

$$\begin{aligned} -\varepsilon \Delta p - \underline{b} \cdot \nabla p + \kappa p &= z, & \text{in } \Omega \\ p &= 0, & \text{on } \Gamma \end{aligned} \tag{3.9}$$

**DEFINITION 3.1.** *The weak solution  $p \in H_0^1(\Omega)$  of the adjoint problem 3.9 is called adjoint state associated to  $\bar{y}$ .*

The right hand side of the adjoint equation belongs to  $L^2(\Omega)$  since  $y_\Omega \in L^2(\Omega)$  and  $\bar{y} \in Y = H_0^1(\Omega) \hookrightarrow L^2(\Omega)$ . By (2.1), we have that (3.9) admits a unique solution  $p \in H_0^1(\Omega)$ .

Making  $z = y - y_\Omega$ , we conclude from the previous Lemma that

$$S^*(S\bar{u} - y_\Omega) = S^*(\bar{y} - y_\Omega) = p, \tag{3.10}$$

which, applied to (3.6) gives

$$(p + \lambda \bar{u}, u - \bar{u})_{L^2(\Omega)} \geq 0, \quad \forall u \in U_{ad}.$$

Thus, the next result follows directly from the variational inequality

**THEOREM 3.4.** *Suppose that  $\bar{u}$  is an optimal control for the ADR problem and let  $\bar{y}$  their associated state. Then the adjoint equation (3.9) has an unique weak solution which satisfies the variational inequality*

$$\int_{\Omega} (p(x) + \lambda \bar{u}(x)) (u(x) - \bar{u}(x)) dx \geq 0, \quad \forall u \in U_{ad}. \tag{3.11}$$

Reciprocally, any control  $\bar{u} \in U_{ad}$  which together with the associated state  $\bar{y} = y(\bar{u})$  and the solution  $p$  of (3.9), satisfies the variational inequality (3.11) is optimal.

The sufficiency of the previous theorem follows from the convexity of  $f$ . Therefore, a control  $u$ , together with the optimal state  $y$  and the adjoint state  $p$ , for the ADR problem is optimal if the triple  $(y, p, u) \in H_0^1(\Omega) \times H_0^1(\Omega) \times U_{ad}$  satisfies the following optimal system

$$\begin{aligned} -\varepsilon \Delta y + \underline{b} \cdot \nabla y + \kappa y &= u + g \\ -\varepsilon \Delta p - \underline{b} \cdot \nabla p + \kappa p &= y - y_\Omega \\ (p + \lambda \bar{u}, u - \bar{u})_{L^2(\Omega)} &\geq 0, \quad \forall u \in U_{ad}. \end{aligned}$$

(3.12)



### 3.3.1. ON THE POINTWISE OPTIMALITY CONDITIONS

The last inequality in our optimality system 3.12 is still not enough to perform any numerical intent to solve our OCP. However, we will see that using some projection operator we can establish an optimality system consisting of 3 equations. To that end, we formulate the variational inequality (3.11) from the *pointwise* point of view. The next result marks the first step. We follow the outline proposed in [24].

**LEMMA 3.3.** *A necessary and sufficient condition to satisfy the variational inequality (3.11) is that for almost every  $x \in \Omega$ ,*

$$\bar{u}(x) = \begin{cases} u_a & \text{if } p(x) + \lambda \bar{u}(x) > 0 \\ \in [u_a, u_b] & \text{if } p(x) + \lambda \bar{u}(x) = 0 \\ u_b & \text{if } p(x) + \lambda \bar{u}(x) < 0. \end{cases} \quad (3.13)$$

An equivalent condition is given by the variational inequality in  $\mathbb{R}$ ,

$$(p(x) + \lambda \bar{u}(x)) (v - \bar{u}(x)) \geq 0 \quad \forall v \in [u_a, u_b], \text{ a.e. } x \in \Omega. \quad (3.14)$$

Next, we note that a simple rearrangement of terms in (3.14) gives

$$(p(x) + \lambda \bar{u}(x)) \bar{u}(x) \leq (p(x) + \lambda \bar{u}(x)) v \quad \forall v \in [u_a, u_b], \quad (3.15)$$

for almost  $x \in \Omega$ . Here like in (3.14),  $v$  is understood as a constant.

**THEOREM 3.5.** *A control  $\bar{u} \in U_{ad}$  is optimal for (2.6)–(2.8) if and only if satisfies together with the adjoint state  $p$  of (3.9) one of the following minimisation conditions for almost all  $x \in \Omega$ :*

- Weak minimum principle

$$\min_{v \in [u_a, u_b]} \{(p(x) + \lambda \bar{u}(x)) v\} = (p(x) + \lambda \bar{u}(x)) \bar{u}(x).$$

- Or The minimum principle

$$\min_{v \in [u_a, u_b]} \left\{ p(x)v + \frac{\lambda}{2} v^2 \right\} = p(x)\bar{u}(x) + \frac{\lambda}{2} \bar{u}(x)^2.$$

The point conditions obtained can be used to obtain more information if we know the sign of  $\lambda$ . The consequences are

**Case 1:  $\lambda = 0$ .** In this case, using (3.13) we have that a.e.,

$$\bar{u}(x) = \begin{cases} u_a & \text{if } p(x) > 0 \\ u_b & \text{if } p(x) < 0. \end{cases} \quad (3.16)$$

In the points  $x \in \Omega$  where  $p(x) = 0$ , there is no information of  $\bar{u}(x)$  that we can obtain. Otherwise, we say that  $\bar{u}$  is a *bang-bang control*, which means that the values of  $\bar{u}$  coincide a.e. with some of the bounds  $u_a$  or  $u_b$ .

**Case 2:  $\lambda > 0$ .** We could interpretate that the second of the relations in (3.13) noting that  $\bar{u}$  is undetermined if  $\lambda\bar{u} + p = 0$ . However, this is incorrect, since these relation immediatly gives that

$$\bar{u}(x) = -\frac{1}{\lambda}p(x)$$

and therefore is an indicator to use the minimum condition. This is specified with the following result.

**THEOREM 3.6.** *If  $\lambda > 0$ , then  $\bar{u}$  is optimal for (2.6)–(2.8) if and only if satisfies, together with the adjoint state  $p$ , the projection formula*

$$\bar{u}(x) = \Pi_{[u_a, u_b]}^\perp \left( -\frac{1}{\lambda}p(x) \right) \quad a.e. \ x \in \Omega, \quad (3.17)$$

where, for real numbers  $a \leq b$ ,  $\Pi_{[a, b]}^\perp$  denotes the projection of  $\mathbb{R}$  on  $[a, b]$ ,

$$\Pi_{[a, b]}^\perp(u) := \min\{b, \max\{a, u\}\}.$$

*Proof.* This is a direct consequence of the previous theorem. In fact, the solution of the quadratic optimisation problem in  $\mathbb{R}$  formulated in terms of the minimum principle

$$\min_{v \in [u_a, u_b]} \left\{ p(x)v + \frac{\lambda}{2}v^2 \right\} = p(x)\bar{u}(x) + \frac{\lambda}{2}\bar{u}(x)^2$$

is given by the projection formula (3.17). □

Cases  $\lambda > 0$  and  $U_{ad} = L^2(\Omega)$  are not considered in this work. We suggest [24] for a study of this situations. We will implementate the finite element method in the next sections on the case  $\lambda > 0$  and  $U_{ad}$  given by 2.8.

## 4. FINITE ELEMENT APPROXIMATION OF THE OPTIMAL CONTROL PROBLEM

Following our outline, after establishing the foundations of existence and uniqueness of the triple  $(y, p, u)$ , we focus our attention in the discretisation of the optimal system 3.12 and its convergence. This is the first stage of our resolution algorithm. Later we will solve the discrete optimization problem with the help of the strategy shown in the next section.

### 4.1. A FINITE ELEMENT METHOD SCHEME

From previous sections, we know that an equivalent formulation of the control problem is given by the non-linear system of equations:

*Find  $(y, p, u) \in H_0^1(\Omega) \times H_0^1(\Omega) \times U_{ad}$  such that*

$$\begin{cases} B(y, v) - (u, v)_{L^2(\Omega)} = (g, v)_{L^2(\Omega)} & \forall v \in H_0^1(\Omega), \\ -(y, v)_{L^2(\Omega)} + B^*(p, v) = -(y_\Omega, v)_{L^2(\Omega)} & \forall v \in H_0^1(\Omega), \\ (p + \lambda u, \bar{u} - u)_{L^2(\Omega)} \geq 0 & \forall \bar{u} \in U_{ad}, \end{cases} \quad (4.1)$$

where the bilinear forms are given by

$$\begin{aligned} B(y, v) &= \int_{\Omega} (\varepsilon \nabla y \cdot \nabla v + \underline{b} \cdot \nabla y v + \kappa y v) \, dx \quad \text{and} \\ B^*(p, v) &= \int_{\Omega} (\varepsilon \nabla p \cdot \nabla v - \underline{b} \cdot \nabla p v + \kappa p v) \, dx. \end{aligned}$$

In order to introduce a FEM scheme we consider the following Lagrange finite element space:

$$V_h := \{v \in \mathcal{C}^0(\bar{\Omega}) : v|_K \in \mathbb{P}_1(K), \forall K \in \mathcal{P}\} \cap H_0^1(\Omega), \quad (4.2)$$

where  $\mathcal{P}$  refers to a triangulation of  $\Omega$ , formed by triangular elements  $K$ . Also  $\mathbb{P}_l(K)$  is the set of all polynomials of degree less or equal to  $l \in \mathbb{N}_0$  and  $v|_K$  represents the restriction of the function  $v$  to the element  $K$ . To approximate the control variable we define in the same way the set,

$$U_{ad}^h = \{t \in L^\infty(\Omega) : t|_K \in \mathbb{P}_0(K), \forall K \in \mathcal{P}\} \cap U_{ad}. \quad (4.3)$$

Now, the finite element discretisation for our control problem is the following:

*Find  $(y_h, p_h, u_h) \in V_h \times V_h \times U_{ad}^h$  such that*

$$\begin{cases} B(y_h, v_h) - (u_h, v_h)_{L^2(\Omega)} = (g_h, v_h)_{L^2(\Omega)} & \forall v_h \in V_h, \\ -(y_h, w_h)_{L^2(\Omega)} + B^*(p_h, w_h) = -(y_\Omega, w_h)_{L^2(\Omega)} & \forall w_h \in V_h, \\ (p_h + \lambda u_h, \bar{u}_h - u_h)_{L^2(\Omega)} \geq 0 & \forall \bar{u}_h \in U_{ad}^h. \end{cases} \quad (4.4)$$

### 4.2. CONVERGENCE OF THE FINITE ELEMENT METHOD

In what follows we will deal with the study of the convergence of the finite element scheme. For this purpose, we will rewrite the previous problems in a convenient way:

$$\begin{cases} B(y, v) = (g + u, v)_{L^2(\Omega)}, \\ B^*(p, v) = (y - y_\Omega, v)_{L^2(\Omega)}, \\ (p + \lambda u, \bar{u} - u)_{L^2(\Omega)} \geq 0, \end{cases} \quad \begin{cases} B(y_h, v_h) = (g + u_h, v_h)_{L^2(\Omega)}, \\ B^*(p_h, v_h) = (y_h - y_\Omega, v_h)_{L^2(\Omega)}, \\ (p_h + \lambda u_h, \bar{u}_h - u_h)_{L^2(\Omega)} \geq 0. \end{cases} \quad (4.5)$$

It is useful to recall that in the previous systems that each state, weak or discrete depends on its respective control, that is,  $y = y(u)$  or  $y_h = y_h(u_h)$ . In addition we define the following auxiliary problems: Given the optimal control  $u \in U_{ad}$  and its discretisation  $u_h \in U_{ad}^h$ , we define the pair  $(y_h(u), r_h) \in V_h \times V_h$ , as the solution of the system

$$\begin{cases} B(y_h(u), v_h) &= (g + u, v_h), \\ B^*(r_h, v_h) &= (y_h(u) - y_\Omega, v_h). \end{cases} \quad (4.6)$$

Besides we define  $q_h \in V_h$  as the solution of

$$B^*(q_h, v_h) = (y(u) - y_\Omega, v_h). \quad (4.7)$$

The previous equalities are valid for any  $v_h \in V_h$ . Finally we define the following orthogonal projection,

$$\Pi^\perp : L^2(\Omega) \rightarrow \mathbb{P}_0(\mathcal{P}),$$

which for each element  $K \in \mathcal{P}$ , fulfills

$$\int_K (j - \Pi^\perp(j)) = 0.$$

We note that  $\Pi^\perp(u) \in U_{ad}$  if  $u \in U_{ad}$ . In fact, since  $u_a \leq u \leq u_b$ , then

$$u_a \leq \frac{1}{|K|} \int_K u \leq u_b,$$

and we conclude that  $\Pi^\perp(u)|_K = \frac{1}{|K|} \int_K u$ , for all  $K \in \mathcal{P}$ . Besides if we recall the following optimal Poincaré inequality (see [2]):

$$\|j - \Pi^\perp(j)\|_{L^2(K)} \leq \frac{h_K}{\pi} \|\nabla(j)\|_{L^2(K)} \quad \forall j \in H^1(K). \quad (4.8)$$

If we take  $\bar{u} = u_h$  and  $\bar{u}_h = \Pi^\perp(u)$  in the variational inequalities (4.5), we have that

$$(p + \lambda u, u_h - u)_{L^2(\Omega)} \geq 0, \quad (4.9)$$

$$(p_h + \lambda u_h, \Pi^\perp(u) - u + u - u_h)_{L^2(\Omega)} \geq 0. \quad (4.10)$$

Adding this optimal inequalities we note that

$$\lambda \|u - u_h\|_{L^2(\Omega)}^2 \leq (p - p_h, u_h - u)_{L^2(\Omega)} + (p_h + \lambda u_h, \Pi^\perp(u) - u)_{L^2(\Omega)}. \quad (4.11)$$

Now we will bound the second term of the right hand side of the inequality above. Considering the definition of the orthogonal projector and using that  $u_h|_K \in \mathbb{P}_0(K)$

and  $p_h \in L^2(\Omega)$ , for all  $K \in \mathcal{P}$ , it follows that

$$\begin{aligned}
(p_h + \lambda u_h, \Pi^\perp(u) - u)_{L^2(\Omega)} &= \sum_{K \in \mathcal{P}} \left( (p_h, \Pi^\perp(u) - u)_{L^2(K)} + \lambda \underbrace{(u_h, \Pi^\perp(u) - u)_{L^2(K)}}_{=0} \right) \\
&= \sum_{K \in \mathcal{P}} (p_h, \Pi^\perp(u) - u)_{L^2(K)} \\
&= \sum_{K \in \mathcal{P}} (p_h - \Pi^\perp(p_h), \Pi^\perp(u) - u)_{L^2(K)} \\
&\leq \sum_{K \in \mathcal{P}} \left( \frac{h_K}{\pi} \right)^2 \|\nabla(p_h)\|_{L^2(K)} \|\nabla(u)\|_{L^2(K)}, \\
&\leq \max_{K \in \mathcal{P}} \{h_K^2 \pi^{-2}\} \left( \sum_{K \in \mathcal{P}} \|\nabla(p_h)\|_{L^2(K)}^2 \right)^{1/2} \left( \sum_{K \in \mathcal{P}} \|\nabla(u)\|_{L^2(K)}^2 \right)^{1/2} \\
&\leq C h^2 \|\nabla(p_h)\|_{L^2(\Omega)} \|\nabla(u)\|_{L^2(\Omega)},
\end{aligned}$$

where we have used the fact that the optimal control  $u$  has  $H^1(\Omega)$  regularity (see Lemma 1.1 in [6]). Taking now  $w_h = p_h$  in (4.5), and using the coercivity of  $B^*$ , we have that

$$\begin{aligned}
\varepsilon \|\nabla(p_h)\|_{L^2(\Omega)}^2 &\leq B^*(p_h, p_h) \\
&= (y_h - y_\Omega, p_h)_{L^2(\Omega)} \\
&\leq (\|y_h\|_{L^2(\Omega)} + \|y_\Omega\|_{L^2(\Omega)}) \|p_h\|_{L^2(\Omega)} \\
&\leq \frac{c_P}{\sqrt{\varepsilon}} \left( \frac{1}{\sqrt{\kappa}} \sqrt{\kappa} \|y_h\|_{L^2(\Omega)} + \|y_\Omega\|_{L^2(\Omega)} \right) \sqrt{\varepsilon} \|\nabla(p_h)\|_{L^2(\Omega)},
\end{aligned}$$

where we used the Cauchy–Schwarz inequality. Hence

$$\sqrt{\varepsilon} \|\nabla(p_h)\|_{L^2(\Omega)} \leq \frac{c_P}{\sqrt{\varepsilon}} \left( \frac{1}{\sqrt{\kappa}} \sqrt{\kappa} \|y_h\|_{L^2(\Omega)} + \|y_\Omega\|_{L^2(\Omega)} \right).$$

In the same way, we have

$$\kappa \|y_h\|_{L^2(\Omega)}^2 \leq B(y_h, y_h) = (g + u_h, y_h)_{L^2(\Omega)} \leq \frac{1}{\sqrt{\kappa}} (\|g\|_{L^2(\Omega)} + \|u_h\|_{L^2(\Omega)}) \sqrt{\kappa} \|y_h\|_{L^2(\Omega)}.$$

Then,

$$\sqrt{\kappa} \|y_h\|_{L^2(\Omega)} \leq \frac{1}{\sqrt{\kappa}} (\|g\|_{L^2(\Omega)} + \|u_h\|_{L^2(\Omega)}),$$

which allow us to obtain

$$\sqrt{\varepsilon} \|\nabla(p_h)\|_{L^2(\Omega)} \leq \frac{c_P}{\sqrt{\varepsilon}} \left( \frac{1}{\kappa} (\|g\|_{L^2(\Omega)} + \|u_h\|_{L^2(\Omega)}) + \|y_\Omega\|_{L^2(\Omega)} \right),$$

and finally, as  $u_h \in U_{ad}^h$ , we have that  $u_a \leq u_h \leq u_b$ , which leads us to conclude

$$\sqrt{\varepsilon} \|\nabla(p_h)\|_{L^2(\Omega)} \leq C (\|g\|_{L^2(\Omega)} + \|y_\Omega\|_{L^2(\Omega)} + \max\{|u_a|, |u_b|\}).$$

This allow us to bound (4.11), as

$$\lambda \|u - u_h\|_{L^2(\Omega)}^2 \leq (p - p_h, u_h - u)_{L^2(\Omega)} + Ch^2. \quad (4.12)$$

Now, we add the function  $q_h$  to the first term in the previous right hand side, defined in (4.7), to obtain

$$(p - p_h, u_h - u) = (p - q_h, u_h - u) + (q_h - p_h, u_h - u).$$

Recalling that  $p \in H_0^1(\Omega)$  and  $q_h \in V_h$ , are solution of the following problems

$$\begin{aligned} B^*(p, v) &= (y(u) - y_\Omega, v) \quad \forall v \in H_0^1(\Omega), \\ B^*(q_h, v_h) &= (y(u) - y_\Omega, v_h) \quad \forall v_h \in V_h, \end{aligned}$$

respectively. This tell us that  $q_h$  corresponds to a discretisation of finite elements for  $p$ , and from the classic theory of FEM (see for example Theorem 3.16 in [9]), we have that

$$\|p - q_h\|_\Omega \leq Ch.$$

For the second term we repeat slightly the procedure in the next fashion, with  $r_h$  defined in (4.6),

$$(q_h - p_h, u_h - u) = (q_h - r_h, u_h - u) + (r_h - p_h, u_h - u).$$

This time, using the discrete and auxiliar problems, making explicit  $y_h = y_h(u_h)$ , we have that

$$\begin{aligned} B(y_h(u_h), v_h) &= (g + u_h, v_h) \Rightarrow B(y_h(u_h) - y_h(u), v_h) = (u_h - u, v_h) \\ B(y_h(u), v_h) &= (g + u, v_h) \end{aligned}$$

If we choose  $v_h = r_h - p_h \in V_h$ , then  $B(y_h(u_h) - y_h(u), r_h - p_h) = (u_h - u, r_h - p_h) = (r_h - p_h, u_h - u)$  because of the symmetry of the inner product in  $L^2(\Omega)$ . Now, using the dual problem we find that

$$\begin{aligned} (r_h - p_h, u_h - u) &= B^*(r_h - p_h, y_h(u_h) - y_h(u)) \\ &= (y_h(u) - y_h(u_h), y_h(u_h) - y_h(u)) \\ &= -\|y_h(u) - y_h(u_h)\|_{L^2(\Omega)}^2 \leq 0. \end{aligned} \quad (4.13)$$

On the other hand, by Cauchy-Schwarz and Poincaré inequalities we have that

$$(q_h - r_h, u_h - u) \leq \|q_h - r_h\|_{L^2(\Omega)} \|u_h - u\|_{L^2(\Omega)} \leq \frac{1}{\sqrt{\kappa}} \|q_h - r_h\|_\Omega \|u_h - u\|_{L^2(\Omega)}.$$

Now, if we use the coercivity of  $B^*$ , it follows that

$$\begin{aligned} \|q_h - r_h\|_\Omega^2 &\leq B^*(q_h - r_h, q_h - r_h) = (y(u) - y_h(u), q_h - r_h)_{L^2(\Omega)} \\ &\leq \|y(u) - y_h(u)\|_{L^2(\Omega)} \|q_h - r_h\|_{L^2(\Omega)} \leq \frac{1}{\kappa} \|y(u) - y_h(u)\|_\Omega \|q_h - r_h\|_\Omega, \end{aligned}$$

and noting that  $y_h(u)$  is a finite element discretisation for  $y(u)$  which contains the same right hand side allows us to use similar arguments as before to conclude that  $\|y(u) - y_h(u)\|_\Omega \leq Ch$ . In resume, we have that

$$\|q_h - r_h\|_\Omega \leq Ch.$$

Using again Cauchy–Schwarz inequality,  $ab \leq \frac{1}{2\lambda}a^2 + \frac{\lambda}{2}b^2$  and previous results we obtain

$$\begin{aligned} \lambda \|u - u_h\|_{L^2(\Omega)}^2 &\leq C (\|p - q_h\|_\Omega + \|q_h - r_h\|_\Omega) \|u - u_h\|_{L^2(\Omega)} + \underbrace{(r_h - p_h, u_h - u)_{L^2(\Omega)}}_{\leq 0} \\ &\leq C (\|p - q_h\|_\Omega^2 + \|q_h - r_h\|_\Omega^2) + \frac{\lambda}{2} \|u - u_h\|_{L^2(\Omega)}^2 \\ &\leq Ch^2 + \frac{\lambda}{2} \|u - u_h\|_{L^2(\Omega)}^2. \end{aligned}$$

With this, we have obtained the following:

**THEOREM 4.1 Control Convergence.** *Let  $(y, p, u)$  be the solution of the weak problem (4.1) and  $(y_h, p_h, u_h)$  the discrete FEM solution of (4.4). Then, under the assumptions of Lemma 1.1 in [6] and Theorem 3.16 in [9], there exist a constant  $C$ , which only depends on  $\Omega$  and the physical parameters such that*

$$\lambda \|u - u_h\|_{L^2(\Omega)} \leq Ch. \quad (4.14)$$

Now, to obtain the convergence for the state and its adjoint, we notice that using the coercivity of  $B$  and substracting the continuous and discrete state equations in (4.5) and then using Cauchy–Schwarz that

$$\begin{aligned} \|y - y_h\|_\Omega^2 &\leq B(y - y_h, y - y_h) = (u - u_h, y - y_h)_{L^2(\Omega)} \\ &\leq \frac{1}{\sqrt{\kappa}} \|u - u_h\|_{L^2(\Omega)} \sqrt{\kappa} \|y - y_h\|_{L^2(\Omega)} \\ &\leq \frac{1}{\sqrt{\kappa}} \|u - u_h\|_{L^2(\Omega)} \|y - y_h\|_\Omega, \\ \Rightarrow \|y - y_h\|_\Omega &\leq \frac{1}{\sqrt{\kappa}} \|u - u_h\|_{L^2(\Omega)}. \end{aligned}$$

In the same way, we conclude that

$$\begin{aligned} \|p - p_h\|_\Omega^2 &\leq B^*(p - p_h, p - p_h) = (y - y_h, p - p_h)_{L^2(\Omega)} \\ &\leq \frac{1}{\kappa} \sqrt{\kappa} \|y - y_h\|_{L^2(\Omega)} \sqrt{\kappa} \|p - p_h\|_{L^2(\Omega)} \leq \frac{1}{\kappa} \|y - y_h\|_\Omega \|p - p_h\|_\Omega \\ &\leq \frac{1}{\kappa^{3/2}} \|u - u_h\|_{L^2(\Omega)} \|p - p_h\|_\Omega, \\ \Rightarrow \|p - p_h\|_\Omega &\leq \frac{1}{\kappa^{3/2}} \|u - u_h\|_{L^2(\Omega)}. \end{aligned}$$

Summarizing we have,

**THEOREM 4.2 State and Adjoint Convergence.** *Let  $(y, p, u)$  the solution of the weak problem (4.1) and  $(y_h, p_h, u_h)$  the FEM solution of (4.4). Then, under the assumptions of Lemma 1.1 in [6] and Theorem 3.16 in [9], there exists constants  $C$  and  $\tilde{C}$ , which only depend on  $\Omega$  and the physical parameters such that*

$$\|p - p_h\|_\Omega \leq Ch \quad \text{and} \quad \|y - y_h\|_\Omega \leq \tilde{C} h. \quad (4.15)$$

## 5. PRIMAL-DUAL ACTIVE SET (PDAS) STRATEGY TO SOLVE OPTIMALITY SYSTEM

There exist several strategies to obtain the solution of our Control Problem (see for example [14], Chapter 2). In this section we present a pseudo Newton method ([14]) which allows the application to a wide number of quadratic minimisation problems, as the one we discuss. These methods begin with a feasible solution and obtain a descent sequence where each point in the sequence is feasible.

This type of methods are an invitation besides to generalize into the solution of control problems through Lagrange multipliers. This treatment is not discuss in the present work, but it can be found for example in [15]. We develop first the method for the continuous case and then we implement it in the discrete case. We have already guaranteed the convergency of the finite element method presented in the previous section.

### 5.1. PDAS STRATEGY IN CONTINUOUS PROBLEM

The Optimal Control problem under consideration is minimize the Reduce functional 2.12, subject to the Box condition 2.8 and the ADR boundary value problem 2.7. According to the projection relation (3.17), the optimal control  $u$  satisfies

$$\bar{u}(x) = \Pi_{[u_a, u_b]}^\perp(-\lambda^{-1}p(x)) = \min\{u_b, \max\{u_a, -\lambda^{-1}p(x)\}\},$$

where the adjoint state  $p$  is solution of (3.9). Therefore, defining

$$\mu = -(-\lambda^{-1}p + \bar{u}) = -\lambda^{-1}f'(\bar{u}),$$

by Lemma 3.2, (3.5) and (3.10). We have then that (3.13) takes the form

$$\bar{u}(x) = \begin{cases} u_a & \text{if } -\lambda^{-1}p(x) < u_a \quad (\Leftrightarrow \mu(x) < 0) \\ -\lambda^{-1}p(x) & \text{if } -\lambda^{-1}p(x) \in [u_a, u_b] \quad (\Leftrightarrow \mu(x) = 0) \\ u_b & \text{if } -\lambda^{-1}p(x) < u_b \quad (\Leftrightarrow \mu(x) > 0). \end{cases} \quad (5.1)$$

In the first case,  $\mu(x) < 0$  and thus  $\bar{u}(x) + \mu(x) < u_a$ , because the definition of  $\mu$  and the fact that  $\bar{u} = u_a$ . The same argument could be use in the third case. In the second case,  $\mu(x) = 0$ , hence  $\bar{u}(x) + \mu(x) = -\lambda^{-1}p(x) \in [u_a, u_b]$ . In summary, if we make  $u = \bar{u}$  the following relations are fulfilled

$$u(x) = \begin{cases} u_a & \text{if } u(x) + \mu(x) < u_a \\ -\lambda^{-1}p(x) & \text{if } u(x) + \mu(x) \in [u_a, u_b] \\ u_b & \text{if } u(x) + \mu(x) < u_b. \end{cases} \quad (5.2)$$

Reciprocally, if  $u \in U_{ad}$  fulfills the previous relation, then it satisfies the projection condition and then it is optimal. For example, let's consider  $u = u_a$ , from where  $\mu(x) < 0$  and then

$$0 > -\lambda^{-1}p - u = -\lambda^{-1}p - u_a \Rightarrow -\lambda^{-1}p < u_a,$$



and thus  $u(x) = \Pi_{[u_a, u_b]}^\perp(-\lambda^{-1}p(x))$ . We can repeat the procedure for the other two cases.

We can conclude that  $u + \mu$  indicates when the inequality restrictions are active or not. This motivates the *primal-dual active set strategy*, which is described now.

The optimality system is: find  $(y, p, u) \in H_0^1(\Omega) \times H_0^1(\Omega) \times U_{ad}$  such that  $y, p$  solve 2.7, 3.9 respectively and

$$u = \Pi_{[u_a, u_b]}^\perp(-\lambda^{-1}p(x)) = \begin{cases} u_a & \text{if } u(x) + \mu(x) < u_a \\ -\lambda^{-1}p(x) & \text{if } u(x) + \mu(x) \in [u_a, u_b] \\ u_b & \text{if } u(x) + \mu(x) > u_b, \end{cases} \quad (5.3)$$

with  $\mu(x) = -(-\lambda^{-1}p(x) + \bar{u}(x))$ . Now, we rewrite the last projection, using the following characteristic functions:

$$\chi^a(x) := \begin{cases} 1 & \text{if } x : u(x) + \mu(x) < u_a, \\ 0 & \text{otherwise,} \end{cases} \quad \chi^b(x) := \begin{cases} 1 & \text{if } x : u(x) + \mu(x) > u_b, \\ 0 & \text{otherwise,} \end{cases} \quad (5.4)$$

which allows to write the projection as

$$(1 - \chi^a(x) - \chi^b(x))\lambda^{-1}p(x) + u(x) = \chi^a(x)u_a + \chi^b(x)u_b. \quad (5.5)$$

This relation should clearly understand in the weak form. Then, in what follows we will be concerned with the following, and final, optimality system:

$$\boxed{\begin{aligned} -\varepsilon \Delta y + \underline{b} \cdot \nabla y + \kappa y &= u + g, & \text{in } \Omega \\ -\varepsilon \Delta p - \underline{b} \cdot \nabla p + \kappa p &= y - y_\Omega, & \text{in } \Omega, \\ (1 - \chi^a(x) - \chi^b(x))\lambda^{-1}p(x) + u(x) &= \chi^a(x)u_a + \chi^b(x)u_b. \end{aligned}} \quad (5.6)$$

## 5.2. MIXED FORMULATION OF OPTIMALITY SYSTEM

In order to simplify the computational implementation of our approach we will consider a mixed formulation of the latter system 5.6. To this end, we will use test functions  $(v, q, w) \in H_0^1(\Omega) \times H_0^1(\Omega) \times L^2(\Omega)$ . As usual, we multiply each equation by  $v, q$  and  $w$  respectively and then integrate on  $\Omega$  to obtain

$$\begin{aligned} -\varepsilon(\nabla y, \nabla v) + (\underline{b} \cdot \nabla y, v) + \kappa(y, v) - (u, v) &= (g, v), \\ -\varepsilon(\nabla p, \nabla q) - (\underline{b} \cdot \nabla p, q) + \kappa(p, q) - (y, q) &= -(y_\Omega, q), \\ (1 - \chi^a(x) - \chi^b(x))\lambda^{-1}(p, w) + (u, w) &= (\chi^a(x)u_a + \chi^b(x)u_b, w). \end{aligned} \quad (5.7)$$

This can be simplified using the bilinear forms  $B$  and  $B^*$  defined in section 4.1 and defining the bilinear form  $e(u, v) = -(u, v)$ , simply the  $L^2$  product of functions  $u, v$ . Also we define

$$\begin{aligned} \beta &= (1 - \chi^a(x) - \chi^b(x)), \\ \eta &= \chi^a(x)u_a + \chi^b(x)u_b, \end{aligned}$$

Then

$$\begin{aligned} B(y, v) + e(u, v) &= (g, v), \\ e(y, q) + B^*(p, q) &= -(y_\Omega, q), \\ -e(\lambda^{-1}\beta p + u, w) - e(u, w) &= (\eta, w). \end{aligned} \tag{5.8}$$

To state the mixed formulation we will add the three left hand sides of the previous system to define

$$\begin{aligned} a((y, p, u); (v, q, w)) &= (B(y, v) + e(u, v)) + (e(y, q) + B^*(p, q)) \\ &\quad - (e(\lambda^{-1}\beta p + u, w) + e(u, w)). \end{aligned} \tag{5.9}$$

In the same way we define the sum of each right hand side in 5.6:

$$R((g, y_\Omega, \eta); (v, q, w)) = (g, v) + (y_\Omega, q) + (\eta, w). \tag{5.10}$$

Then, our problem is formulated as follows: Find a triple  $(y, p, u) \in H_0^1(\Omega) \times H_0^1(\Omega) \times L^2(\Omega)$  such that, for any test functions  $(v, q, w) \in H_0^1(\Omega) \times H_0^1(\Omega) \times L^2(\Omega)$ :

$$a((y, p, u); (v, q, w)) = R((g, y_\Omega, \eta); (v, q, w)) \tag{5.11}$$

We need to show that this formulation preserves the convergence orders obtained in theorems 4.1 and 4.2. To see this, we note that coercivity and continuity of  $a$  are an immediate consequence of theorem 2.1 applied to  $B$  and  $B^*$  and the properties of the inner product in  $L^2(\Omega)$ . To recover the convergence orders for each variable, we refer to the article of Leykekhman [17] and the Lecture Notes on Stokes System from the Numerics of EDP 2 course from Lars Diening [8]. In the same direction, we recommend [23] or [12] for an introductory approach.

With all these elements, we can formulate the following algorithm to solve our OCP:

---

**Algorithm 1** PDAS algorithm for ADR control problem, mixed formulation

---

- 1: **Input:** Initial guess  $u_0, \mu_0$ , max. iteration number  $K$
- 2: **Set**  $k = 0$
- 3: Determine the indicator functions  $\chi_0^a, \chi_0^b$  of the active sets

$$\mathcal{A}_0^a = \{x \in \Omega : \mu_0(x) + u_0(x) < u_a\}$$

$$\mathcal{A}_0^b = \{x \in \Omega : \mu_0(x) + u_0(x) > u_b\}.$$

- 4: **while**  $1 \leq k \leq K$  **do**
- 5:   **Define**

$$\beta_{k-1} = (1 - \chi_{k-1}^a - \chi_{k-1}^b)$$

$$\eta_{k-1} = u_a \chi_{k-1}^a + u_b \chi_{k-1}^b$$

- 6:   **Solve** Find  $y_k, p_k \in H_0^1(\Omega)$  and  $u_k \in L^2(\Omega)$  such that

$$a((y_k, p_k, u_k); (v, q, w)) = R((g, y_\Omega, \eta_{k-1}); (v, q, w)), \quad (5.12)$$

for test functions  $v, q \in H_0^1(\Omega)$  and  $w \in L^2(\Omega)$ .

- 7:   **Actualize**  $\mu_k = p_k - \lambda u_k, \chi_k^a, \chi_k^b$ .
  - 8:   **if**  $\|\chi_k^a - \chi_{k-1}^a\| + \|\chi_k^b - \chi_{k-1}^b\| = 0$ : **then**                       $\triangleright$  This guaranties that  $u_k \in U_{ad}$ .
  - 9:     **Break**
  - 10:   **else**
  - 11:     **Set**  $k = k + 1$
  - 12:   **end if**
  - 13: **end while**
  - 14: **end**
-

## 6. NUMERICAL RESULTS

In this section we develop numerical experiments for the ADR Optimal Control Problem 5.6. Our goal is to check the convergence orders for control, state and adjoint state given by theorems 4.1 and 4.2, as well as to establish a comparison between our method and the novel optimal control software **cashocs** [3]. First we show the parameters and exact solutions, then we provide the convergence analysis for our approach to conclude with a similar analysis using **cashocs** to determine a comparison between the two techniques. We will see that convergence orders using the mixed strategy are in concordance with the convergence results presented in section 4.2, but are strangely improved with the use of **cashocs**. We give a conjecture for this issue in section 5.3.

### 6.1. PARAMETERS AND EXACT SOLUTIONS

For the parameters and data showed in Table 6.1 we set the desire state  $y_\Omega$  and the external function  $g$  as

| Parameter       | Value  |
|-----------------|--------|
| $\varepsilon$   | 1.0    |
| $\underline{b}$ | (3, 4) |
| $\kappa$        | 1      |
| $\lambda$       | 1      |
| $u_a$           | -2     |
| $u_b$           | 2      |

Table 6.1: Parameters for example problem 5.6 in domain  $\Omega = [0, 1]^2$ .

$$\begin{aligned}
 y_\Omega &= -x_1^2 x_2^2 + 9x_1^2 x_2 - 2x_1^2 + 7x_1 x_2^2 - 15x_1 x_2 + 2x_1 - x_2^2 + x_2 \\
 g &= -6x + 6x_1^2 - 5x_2 + 12x_1 x_2 - 6x_1^2 x_2 + 5x_2^2 - 4x_1 x_2^2 - 2x_1^2 x_2^2 \\
 &\quad + \min(-2, \max(-2, -(1-x_1)x_1(1-x_2)x_2))
 \end{aligned}$$

We consider the exact solution triple, with the control given by 3.17

$$\begin{aligned}
 y &= xy(x-1)(y-1) \\
 p &= xy(x-1)(y-1), \\
 u &= \min(u_b, \max(u_a, -p(x)/\lambda)).
 \end{aligned}$$

For both implementations we have chosen a tolerance  $\text{tol} = 10^{-9}$ .

### 6.2. IMPLEMENTATION IN FENICS

We implementate the algorithm 1 in the **Python** based dedicated softwares **FENICS** [1], [20] and **Dolfin** [21]. The main advantage of these softwares relies on the hybrid resolution scheme they applied. We just need to input the domain, mesh, finite element spaces, parameters, weak formulation and boundary conditions and the software will internally resolve and return the numerical solutions. The programming was divided in four parts: a main program containing the different meshes, parameters and the PDAS

loop and four subroutines containing a solver for the mixed weak formulation 5.12, the  $L^2(\Omega)$  projection implementation needed by equation 5.5 and the indicator functions  $\chi_k^a, \chi_k^b$  respectively. At the end of the main program, we develop the postprocessing analysis in order to provide figures of the numerical solutions as well as convergence rates of our algorithm. Details can be found in the Appendix.

### 6.2.1. NUMERICAL RESULTS

We have consider 5 different meshes with 16, 32, 64, 128 and 256 elements. Like in section 4, we choose piecewise Lagrange discontinuos triangular elements for the control function  $u$  and linear continuous Lagrange triangular elements for state and adjoint state  $y, p$  respectively. A two dimensional representation of the numerical solutions is showed in Figure 6.1

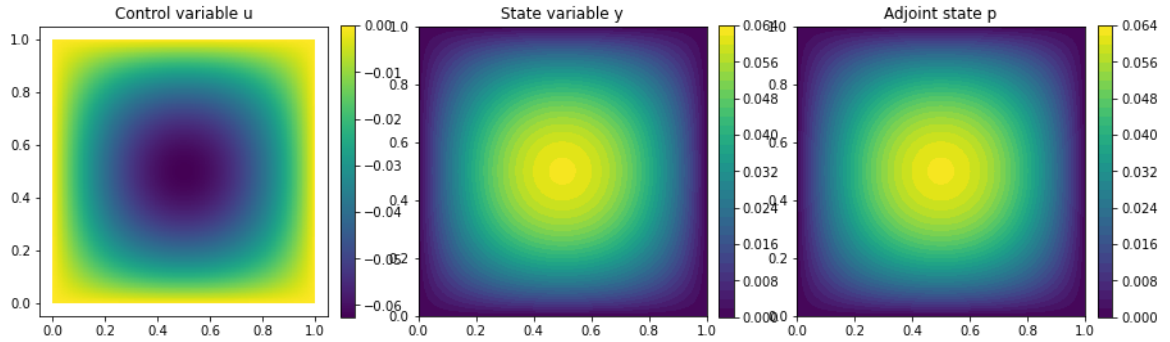


Figure 6.1: Numerical solutions for control (left), state (middle) and adjoint state for 5.12 problem and parameters given by Table 6.1. We have used 256 elements for this representation.

Next, we consider the convergence rates of our approximation technique. The numerical results essentially confirm the theoretical estimate of Theorems 4.1 and 4.2, with a convergence rate of  $h^1$  observed for the  $L^2$ -norm for control and  $H^1$ -norm for both state and adjoint state. Another important point corresponds to the mesh-independance of the iterations. This is concordance with the literature [13]. This means that for each mesh, our method requires the same number of iterations to find the optimal solution. The low number of iterations, just 2 for each mesh, is probably because of the relative simplicity of our problem test selection and the type of our state and adjoint state equations, since both are linear PDE's. The election of  $\underline{b}$  could be also relevant (see, for example [17]). The results are showed in Table 6.2 and Figure 6.2. The adjusted convergence rates were made with relative coarse meshes (10, 20, 30, 40, 50, 60 elements), but they succeed to describe correctly the behaviour of each error.

| NDOF | $h_{\max}$               | Iterations | $\ u - u_h\ _{L^2(\Omega)}$ | $\ y - y_h\ _{H^1(\Omega)}$ | $\ p - p_h\ _{H^1(\Omega)}$ |
|------|--------------------------|------------|-----------------------------|-----------------------------|-----------------------------|
| 16   | $3.125 \times 10^{-2}$   | 2          | $2.23198 \times 10^{-3}$    | $1.51916 \times 10^{-2}$    | $1.51946 \times 10^{-2}$    |
| 32   | $1.5625 \times 10^{-2}$  | 2          | $1.10257 \times 10^{-3}$    | $7.6044 \times 10^{-3}$     | $7.60478 \times 10^{-3}$    |
| 64   | $7.8125 \times 10^{-3}$  | 2          | $5.49579 \times 10^{-4}$    | $3.80327 \times 10^{-3}$    | $3.80332 \times 10^{-3}$    |
| 128  | $3.90625 \times 10^{-3}$ | 2          | $2.74575 \times 10^{-4}$    | $1.90177 \times 10^{-3}$    | $1.90178 \times 10^{-3}$    |
| 256  | $1.95312 \times 10^{-3}$ | 2          | $1.37261 \times 10^{-4}$    | $9.50902 \times 10^{-4}$    | $9.50902 \times 10^{-4}$    |

Table 6.2: Error approximations for test problem 5.6 in domain  $\Omega = [0, 1]^2$ . Although we showed that norms  $\|\cdot\|_{\Omega}$  and  $\|\cdot\|_{H^1(\Omega)}$  are equivalent in Section 1, the election of the parameters sets this equivalence as an equality.

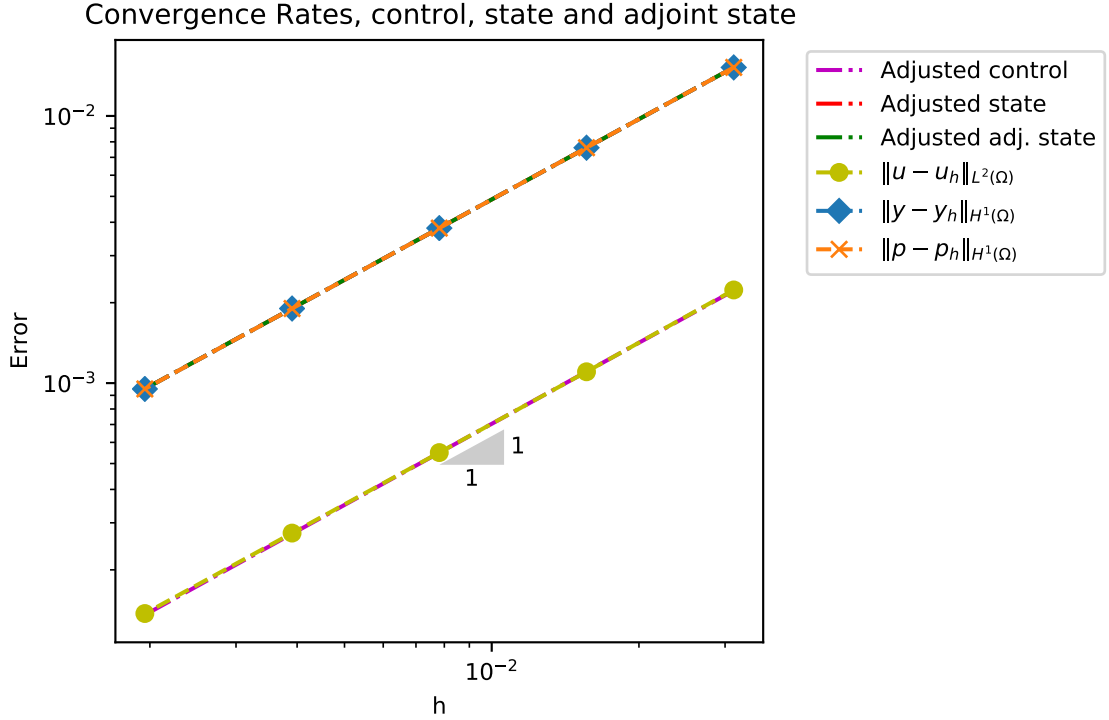


Figure 6.2: Convergence Rates for control  $u$  in  $L^2(\Omega)$  error, state  $y$  and adjoint state  $p$  in  $H^1(\Omega)$  error. The adjusted convergence rates were with least squares method calculated and are 1.00838 for control variable and  $9.9923 \times 10^{-1}$ ,  $9.99315 \times 10^{-1}$  for state and adjoint state variables respectively.

### 6.3. IMPLEMENTATION IN CASHOCS

We repeat our analysis for the problem under consideration using this time the dedicated optimal control software **cashocs** [3]. This software uses the framework given by **FENICS** and **Dolfin** to implementate its internal solvers. We have chosen the PDAS Strategy as solver in order to compare with our results. All the settings are the same as in previous section with a single difference: by software restriction, we must choose the same finite element space for the three variables. That is, linear continuous Lagrange triangular elements for control, state and adjoint state. This extra regularity in the election of the finite element space for the control variable has implied an extra regularity for the numerical control solution as we can see in Figure 6.4. In this case

we can conjecture that the optimal control lies in a subset of  $U_{ad} \subset L^2(\Omega)$ . On the other hand, convergence rates for both state and adjoint state were preserved and were not affected by the election of the finite element space for the control variable. We can see this behaviour in Figure 6.4.

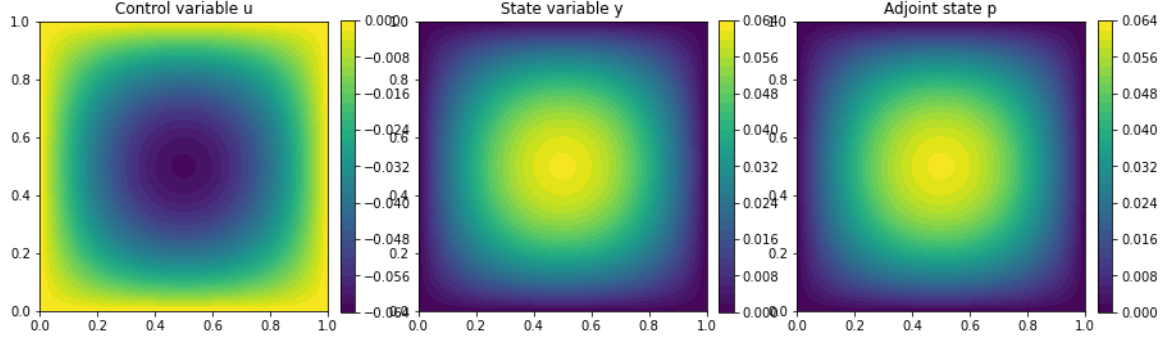


Figure 6.3: Numerical solutions for control (left), state (middle) and adjoint state for 5.12 problem and parameters given by Table 6.1 using `cashocs`. It is important to remark that the scheme to solve the adjoint equation used by this software follows the approach of [14] and therefore will solve the optimality system for  $-p$ . We have used 256 elements for this representation.

| NDOF | $h_{\max}$               | Iterations | $\ u - u_h\ _{L^2(\Omega)}$ | $\ y - y_h\ _{H^1(\Omega)}$ | $\ p - p_h\ _{H^1(\Omega)}$ |
|------|--------------------------|------------|-----------------------------|-----------------------------|-----------------------------|
| 16   | $3.125 \times 10^{-2}$   | 2          | $4.38845 \times 10^{-4}$    | $1.51937 \times 10^{-2}$    | $1.51959 \times 10^{-2}$    |
| 32   | $1.5625 \times 10^{-2}$  | 2          | $1.10054 \times 10^{-4}$    | $7.60466 \times 10^{-3}$    | $7.60494 \times 10^{-3}$    |
| 64   | $7.8125 \times 10^{-3}$  | 2          | $2.7535 \times 10^{-5}$     | $3.8033 \times 10^{-3}$     | $3.80334 \times 10^{-3}$    |
| 128  | $3.90625 \times 10^{-3}$ | 2          | $6.88508 \times 10^{-6}$    | $1.90177 \times 10^{-3}$    | $1.90178 \times 10^{-3}$    |
| 256  | $1.95312 \times 10^{-3}$ | 2          | $1.72135 \times 10^{-6}$    | $9.50902 \times 10^{-4}$    | $9.50903 \times 10^{-4}$    |

Table 6.3: Error approximations for test problem 5.6 in domain  $\Omega = [0, 1]^2$  using `cashocs` PDAS solver. As we can see, the error  $L^2(\Omega)$  is less than the approximation obtained in our algorithm, but is almost equal for the  $H^1(\Omega)$  error for state and adjoint variables.

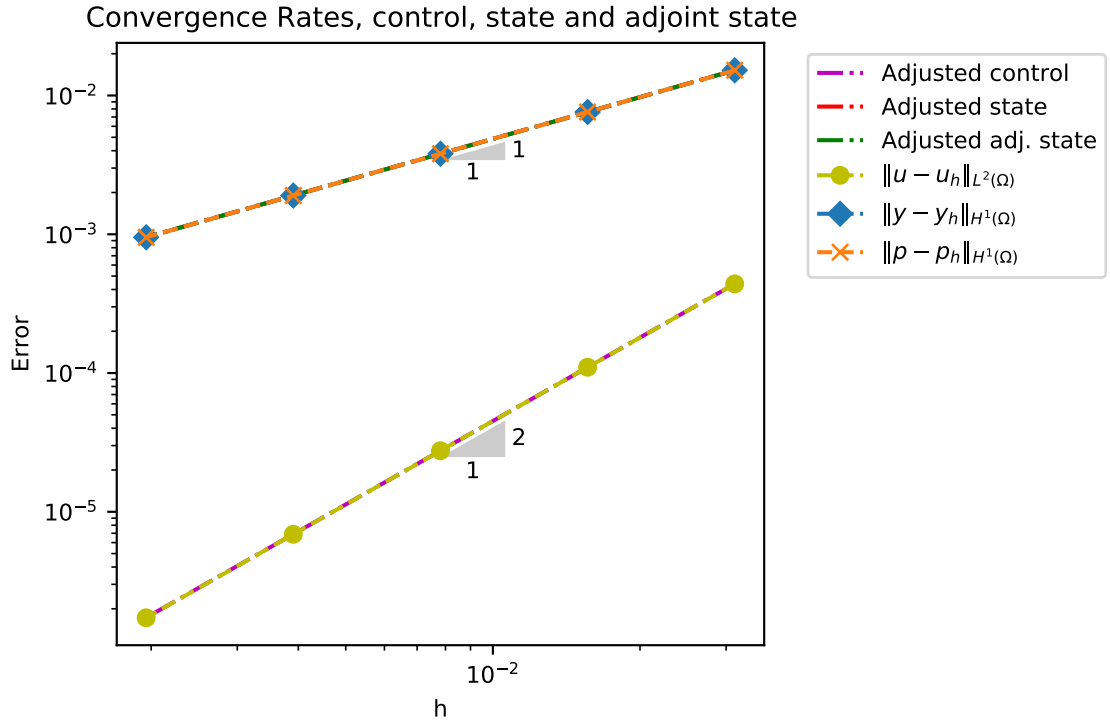


Figure 6.4: Convergence Rates for control  $u$  in  $L^2(\Omega)$  error, state  $y$  and adjoint state  $p$  in  $H^1(\Omega)$  error. The adjusted convergence rates were with least squares method calculated and are  $1.99787$  for control variable and  $9.993 \times 10^{-1}$ ,  $9.99375 \times 10^{-1}$  for state and adjoint state variables respectively.



## 7. CONCLUSIONS AND OUTLOOK

### Summary of the work

Following the guideline given by the literature [24], [18], we have managed to prove the existence and uniqueness of the optimal state-control pair. In addition, the active primal-dual set strategy has been successfully implemented using a mixed formulation approach of the optimal system. Numerical results of our approach coincide with the theory proposed and demonstrated in this work.

The fact of using  $L^2$ -Projection in the implementation to coincide with the theory shows good results, since the convergence orders stipulated by the theoretical framework are respected and the approximation errors are kept under control. On the other hand, the results provided by `cashocs`, despite having a better order of convergence for the control, do not respect the theory and must be rigorously examined. Interesting is the fact that the convergence orders for state and adjoint state remain as suggested by our theoretical framework.

### Important Results

Without a doubt, after providing the theoretical framework of the work in Sections 2 and 3, results for the approximation have been obtained through the finite element method that coincide with the related literature [17]. When we have our results proved, we did not have such reference, so the fourth section constitutes the soul of our work. The numerical results only confirm the correctness of our analysis.

The mixed formulation of the optimal system has the advantage of a notorious simplification of the programming task and it allows us to parallelize the routines and minimise the computational calculations.

### State of Research and Outlook

Notwithstanding the foregoing, the equivalence between the optimal system and the mixed formulation remains to be shown more clearly. Especially what concerns the projection relationship between adjoint state and control 5.5. This is one of the points where to persevere in our investigation.

The programs developed allow to consider different solvers for other control problems. Remains to study the limitations of this approach and obviously to look the perspective for further applications. For now, we have already found a second test problem in section 5.1 of [17] which seems to offer more difficulties than the one studied in our work. It will be interesting to compare both schemes in the near future.

## REFERENCES

- [1] Martin S. Alnæs, Jan Blechta, Johan Hake, August Johansson, Benjamin Kehlet, Anders Logg, Chris Richardson, Johannes Ring, Marie E. Rognes, and Garth N. Wells. The fenics project version 1.5. *Archive of Numerical Software*, 3(100), 2015.
- [2] M. Bebendorf. A note on the Poincaré inequality for convex domains. *Z. Anal. Anwendungen*, 22(4):751–756, 2003.
- [3] Sebastian Blauth. cashocs: A computational, adjoint-based shape optimization and optimal control software. *SoftwareX*, 13:100646, Jan 2021.
- [4] Susanne C. Brenner and L. Ridgway Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer, New York, third edition, 2008.
- [5] Eduardo Casas. Optimal control of pde theory and numerical analysis. page 136, Aug 2006. Lecture.
- [6] Eduardo Casas and Fredi Tröltzsch. Error estimates for linear-quadratic elliptic control problems. In *Analysis and optimization of differential systems (Constanta, 2002)*, pages 89–100. Kluwer Acad. Publ., Boston, MA, 2003.
- [7] Juan Carlos de los Reyes and Fredi Tröltzsch. Optimal control of the stationary navier–stokes equations with mixed control-state constraints. *SIAM J. Control. Optim.*, 46(2):604–629, 2007.
- [8] Lars Diening. *Numerics of PDE II*. Universität Bielefeld, 2020.
- [9] Alexandre Ern and Jean-Luc Guermond. *Theory and practice of finite elements*, volume 159 of *Applied Mathematical Sciences*. Springer-Verlag, New York, 2004.
- [10] Funke, S.W. and Nordass, M. PDE constrained optimisation in Hilbert Spaces. <https://fenicsproject.org/pub/presentations/fenics14-paris/SimonFunke.pdf>, 2014.
- [11] Carlos Galeano Urueña. Técnicas de solución numérica de la ecuación difusión-advección-reacción para el estudio de dispersión de contaminantes. Master’s thesis, 2009.
- [12] Gabriel Gatica. *A Simple Introduction to the Mixed Finite Element Method. Theory and Applications*. 01 2014.
- [13] Ulbrich M. Hintermüller, M. A mesh-independence result for semismooth newton methods. *Mathematical Programming*, 101:151–184, 08 2004.
- [14] M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE constraints*, volume 23 of *Mathematical Modelling: Theory and Applications*. Springer, New York, 2009.

- [15] Kazufumi Ito and Karl Kunisch. *Lagrange multiplier approach to variational problems and applications*, volume 15 of *Advances in Design and Control*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2008.
- [16] Günter Leugering, Peter Benner, Sebastian Engell, Andreas Griewank, Helmut Harbrecht, Michael Hinze, Rolf Rannacher, and Stefan Ulbrich. *Trends in PDE Constrained Optimization*, volume 165. 01 2014.
- [17] Dmitriy Leykekhman. Investigation of commutative properties of discontinuous galerkin methods in pde constrained optimal control problems. *J. Sci. Comput.*, 53(3):483–511, 12 2012.
- [18] J.-L. Lions. *Optimal control of systems governed by partial differential equations*. Die Grundlehren der mathematischen Wissenschaften, Band 170. Springer-Verlag, New York-Berlin, 1971. Translated from the French by S. K. Mitter.
- [19] W. Liu, Pekka Neittaanmäki, and Dan Tiba. Existence for shape optimization problems in arbitrary dimension. *SIAM J. Control and Optimization*, 41:1440–1454, 01 2002.
- [20] Anders Logg, Kent-Andre Mardal, Garth N. Wells, et al. *Automated Solution of Differential Equations by the Finite Element Method*. Springer, 2012.
- [21] Anders Logg and Garth N. Wells. Dolfin: Automated finite element computing. *ACM Trans. Math. Softw.*, 37(2), April 2010.
- [22] Pekka Neittaanmaki and Sergey Repin. Two-sided guaranteed estimates of the cost functional for optimal control problems with elliptic state equations. *Lecture Notes in Computational Science and Engineering*, 101, 2014.
- [23] Endre Süli. A brief excursion into the mathematical theory of mixed finite element methods.
- [24] Fredi Tröltzsch. *Optimal control of partial differential equations*, volume 112 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2010. Theory, methods and applications, Translated from the 2005 German original by Jürgen Sprekels.
- [25] Pieter Wesseling. *The stationary convection-diffusion equation*, pages 111–161. Springer Berlin Heidelberg, Berlin, Heidelberg, 2001.

## APPENDIX A IMPLEMENTATION IN FENICS

Here the interested reader can find the routines used to implementate the PDAS to the ADR optimal Control Problem in FENICS

### A.1 MAIN PROGRAM

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Thu Jul 15 19:47:55 2021
5
6  \@author: Juan Carlos Chavarr\`ia
7
8  This script will return the solution of an Optimal Control Problem
9  with restrictions given by an PDE and a box condition using the PDAS.
10 The script is general. The user provide the PDE, used as an input (solver),
11 and run this routine to obtain the solution as well as the desired state used
12 to compute de values of the objective functional J.
13
14 --- In this case we consider the ADR equation studied in the Bachellorarbeit
15 """
16 from fenics import *
17 from math import sqrt, inf
18 from dolfin import *
19 from SolverADR import *
20 from project_up import *
21 from project_low import *
22 from L2projection import *
23 import numpy as np
24 import matplotlib.pyplot as plt
25 from plotslopes import *
26 #from mpltools import annotation
27 # Set some algorithmic parameters
28 atol = 1e-9
29 maxiter = 100
30 # Set the lambda parameter
31 lam = Constant(1.0)
32 # Define polynomial degrees for control and state discretization
33 degree_con = 0
34 degree_state = 1
35 # Exact Solutions
36 ye=Expression('x[0]*x[1]*(1-x[0])*(1-x[1])', degree=4)
37 pe=Expression('x[0]*x[1]*(1-x[0])*(1-x[1])', degree=4)
38 ue=Expression('min(2.0, max(-2.0, -((1-x[0])*x[0]*(1-x[1])*x[1])/lam))', lam=lam, degree=3)
39 #Define Mesh
40 # number of uniform refinements
41 Nref = 5
42 #### Store information ####
43 # a counter
44 array=np.arange(Nref)
45 # NDOF of each mesh
46 l=np.empty(Nref, dtype=int)
47 # Maximum diameter
48 h=np.empty(Nref)
49 # Errors
50 L2error = np.empty(Nref)
51 State_error_y = np.empty(Nref)
52 Adjointerror_p = np.empty(Nref)
53 #### Loops ####
54 # exterior loop: set each mesh and calculate the approx. errors
55 for j in array:
56     ## NDOF ##
57     l[j]=2**(j+4)
58     ## mesh size
59     h[j]=np.divide(1,(np.multiply(2,l[j])))
60     ## mesh
61     mesh = UnitSquareMesh(1[j],1[j])
62     ##### Spaces #####
63     ## Finite Element Spaces
64     PU = FiniteElement('DG',mesh,ufl_cell(),degree_con)
65     PY = FiniteElement('Lagrange',mesh,ufl_cell(),degree_state)
66     ## Function Spaces
67     U = FunctionSpace(mesh,PU)
68     Y = FunctionSpace(mesh,PY)
69     # Desired State
70     yom=Expression('(2*x[0] - 2*pow(x[0],2) + x[1] -14*x[0]*x[1] + 8*x[1]*pow(x[0],2) \
71     -pow(x[1],2) + 6*x[0]*pow(x[1],2))', degree=degree_state)
72     yd=interpolate(yom,Y)
73     # RHS state equation
74     exf=Expression('(6*x[0]-6*pow(x[0],2)+5*x[1]-13*x[0]*x[1]\
75     +7*pow(x[0],2)*x[1]-5*pow(x[1],2)+5*x[0]*pow(x[1],2)\
76     +pow(x[0],2)*pow(x[1],2)\
77     -min(2.0, max(-2.0, -(1-x[0])*x[0]*(1-x[1])*x[1]))))'\
78     ,degree=degree_state)
79     g=interpolate(exf,Y)
80     # Set the bounds of box condition
```

```

81     ua = interpolate(Expression('-2.0', degree=degree_con), U)
82     ub = interpolate(Expression('2.0', degree=degree_con), U)
83     # Prepare for the PDAS loop
84     # Set the initial guesses
85     u0 = interpolate(Expression('-100.0', degree=degree_con), U)
86     mu0 = interpolate(Expression('-100.0', degree=degree_con), U)
87     # Create functions holding the iterations
88     u = Function(U)
89     mu = Function(U)
90     toproj=Function(U)
91     # Build the constant 1 to use in beta function
92     uno = interpolate(Expression('1.0', degree=degree_con), U)
93     # Set the initial iterate
94     u.assign(u0)
95     mu.assign(mu0)
96     toproj.assign(u0+mu0)
97     ## Build the initial indicator functions
98     Xa_old = project_up(toproj, U, ua)
99     Xb_old = project_low(toproj, U, ub)
100    # counter & condition
101    iter = 0
102    done = False
103    # inner loop: PDAS
104    while (iter < maxiter) and (done == False):
105        beta = (uno-Xa_old-Xb_old)/lam
106        phi = ua*Xa_old+ub*Xb_old
107        u, y, p = SolverADR(mesh, PU, PY, beta, phi, yd, g, lam)
108        Ph=L2Projection(mesh, p, degree_con) # We need to project, otherwise mu mix different
            function spaces
109        tomu = np.add(np.multiply(1/lam,Ph.vector()[:]), np.multiply(-1, u.vector()[:]))
110        mu.vector()[:] = tomu
111        coefftoproj = np.add(u.vector()[:], mu.vector()[:])
112        toproj.vector()[:] = coefftoproj
113        Xa_new = project_up(toproj, U, ua)
114        Xb_new = project_low(toproj, U, ub)
115        error_a = sqrt(assemble((Xa_old-Xa_new)*(Xa_old-Xa_new)*dx))
116        error_b = sqrt(assemble((Xb_old-Xb_new)*(Xb_old-Xb_new)*dx))
117        print("Error is= ", error_a, error_b)
118        if error_a+error_b < atol:
119            done = True
120            Xa_old = Xa_new
121            Xb_old = Xb_new
122            iter = iter+1
123    # end of pdas loop
124    print('number of iterations',iter)
125    ## Errors: L2 for u, H1 for y and p
126    L2error[j] = errornorm(ue,u,norm_type='L2')
127    State_error_y[j]=errornorm(ye, y, norm_type='H1')
128    Adjointerror_p[j]=errornorm(pe, p, norm_type='H1')
129    # end of loops
130    # print some information
131    print("maximal diameter of mesh is=", h )
132    print("L2 error",L2error)
133    print("H1 error-state", State_error_y )
134    print("H1 error-adjoint", Adjointerror_p)
135    ### Post Processing
136    # Plot the solutions
137    # and export to paraview
138    plt.figure(figsize=(15,5))
139    # Control
140    plt.subplot(1, 3, 1)
141    fig = plot(u)
142    plt.colorbar(fig, fraction=0.046, pad=0.04)
143    plt.title('Control variable u')
144    vtkfile = File('control.pvd')
145    vtkfile << u
146    # State
147    plt.subplot(1,3,2)
148    fig = plot(y)
149    plt.colorbar(fig, fraction=0.046, pad=0.04)
150    plt.title('State variable y')
151    vtkfile = File('state.pvd')
152    vtkfile << y
153    # Adjoint State
154    plt.subplot(1,3,3)
155    fig = plot(p)
156    plt.colorbar(fig, fraction=0.046, pad=0.04)
157    plt.title('Adjoint state p')
158    vtkfile = File('adjoint.pvd')
159    vtkfile << p
160    # Saving
161    plt.savefig('solution.pdf')
162
163    # error plots
164    plt.figure(2)
165    plt.rcParams.update({'font.size': 8.5})
166    # convergence rate (adapted by hand)
167    # adjusted control
168    const_shift = -1.13571 # shift of line in y-direction
169    const_rate = 1.00838 # slope for line in loglog Plot
170    rate = pow(10,const_shift)*np.exp((const_rate*np.log(h))) # function to plot
171    pl=plt.loglog(h,rate,'m-', label='Adjusted control')
172
173    # adjusted state

```

```

174 const_shift = -0.314198 # shift of line in y-direction
175 const_rate = 0.99923 # slope for line in loglog Plot
176 ratey = pow(10,const_shift)*np.exp((const_rate*np.log(h))) # function to plot
177 p2=plt.loglog(h,ratey,'r-.', label='Adjusted state')
178
179 # adjusted adjoint state
180 const_shift = -0.314022 # shift of line in y-direction
181 const_rate = 0.999315 # slope for line in loglog Plot
182 ratep = pow(10,const_shift)*np.exp((const_rate*np.log(h))) # function to plot
183 p3=plt.loglog(h,ratep,'g-.', label='Adjusted adj. state')
184 p4=plt.loglog(h, L2error, 'y-o', label=r'$\Vert u-u_h\Vert_{L^2(\Omega)}$')
185 p5=plt.loglog(h, State_error_y, '-D', label=r'$\Vert y-y_h\Vert_{H^1(\Omega)}$')
186 p6=plt.loglog(h, Adjointerror_p, '-x', label=r'$\Vert p-p_h\Vert_{H^1(\Omega)}$')
187 slope_marker((h[2],L2error[2]-0.1*L2error[2]), (1,1))
188 #slope_marker((h[1],State_error_y[1]-0.1*State_error_y[1]), (1,1))
189 plt.xlabel('h')
190 plt.ylabel('Error')
191 plt.title("Convergence Rates, control, state and adjoint state")
192
193 plt.legend(bbox_to_anchor=(1.05, 1.0), loc='upper left')
194 plt.tight_layout()
195 plt.savefig('convergence_ich.pdf')
196 plt.show()

```

## A.2 OPTIMALITY SYSTEM SOLVER

```

1 from dolfin import *
2
3 import matplotlib.pyplot as plt
4 import sys
5 import numpy as np
6
7 def SolverADR(mesh,PU,PY,beta,phi,yd,g,lam):
8     ### Parameters
9     epsi =Constant(1.0)
10    kappa=Constant(1.0)
11    be=Constant((3.0, 4.0))
12    ### BC parameter for mixed formulation
13    Dbc= Constant(0.0)
14    ##### function spaces: mixed formulation #####
15    U = FunctionSpace(mesh,PU)
16    Y = FunctionSpace(mesh,PY)
17    # ## Interpolation of data
18    # beta=interpolate(beta,U)
19    # phi=interpolate(phi,U)
20    # yd=interpolate(yd,Y)
21    # g=interpolate(g,U)
22    # (u,y,p) \in (U,Y,Y) but we need the element spaces
23    X = FunctionSpace(mesh,MixedElement([PU,PY,PY]))
24    ## Boundary Conditions for State and Dual eqs.
25    bcy = DirichletBC(X.sub(1),Dbc,'on_boundary')
26    bcp = DirichletBC(X.sub(2),Dbc,'on_boundary')
27    BC=[bcy, bcp]
28    ##### solve the ADR equations #####
29    (u,y,p) = TrialFunctions(X)
30    (w,v,q) = TestFunctions(X)
31    ## Bilinear Forms
32    Co= (beta*p*w+u*w)*dx
33    Ay= (inner(epsi*grad(y),grad(v))+inner(be,grad(y)))*v+kappa*y*v-u*v)*dx
34    Ap= (inner(epsi*grad(p),grad(q))-inner(be,grad(p)))*q+kappa*p*q-y*q)*dx
35    a = Co+Ay+Ap
36    L = (phi*w+g*v-yd*q)*dx
37    ## Assign the solution
38    uh_yh_ph = Function(X)
39    ## Solve
40    solve(a == L,uh_yh_ph,BC)
41    ##### Split components of the solution #####
42    uh = Function(U)
43    yh = Function(Y)
44    ph = Function(Y)
45    assign(uh,uh_yh_ph.sub(0))
46    assign(yh,uh_yh_ph.sub(1))
47    assign(ph,uh_yh_ph.sub(2))
48
49    return (uh,yh,ph)

```

## A.3 $L^2$ -PROJECTION

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Sat Jul 24 20:02:38 2021
5
6  @author: juancarlos
7  """
8  from dolfin import *
9  import numpy as np
10
11  def L2Projection(mesh, f, polydegree):
12      V = FunctionSpace(mesh, 'DG', polydegree)
13      u = TrialFunction(V)
14      v = TestFunction(V)
15      r = u*v*dx
16      rhs = f*v*dx
17      Pu = Function(V)
18      solve(r == rhs, Pu)
19
20      return (Pu)
```

## A.4 INDICATOR FUNCTIONS

### A.4.1 $\chi_a$

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Wed Jul 14 21:43:57 2021
5  This script creates the indicatrix function of the Active set A^a.
6  It needs three inputs: the functions to compare, f and ua as well as the
7  function space U.
8  @author: juancarlos
9  """
10
11  from dolfin import *
12  import numpy as np
13  import matplotlib.pyplot as plt
14
15  def project_up(f,U,ua):
16      # Build the vectors to compare
17      fv = f.vector()[:]
18      uperv = ua.vector()[:]
19      # Indicatrix function
20      cond = np.where(fv<uperv,1,0)
21      # Build the vector with the binary values
22      proj_up = Function(U)
23      proj_up.vector()[:] = cond
24
25      return proj_up
```

### A.4.2 $\chi_b$

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Wed Jul 14 21:43:57 2021
5  This script creates the indicatrix function of the Active set A^b.
6  It needs three inputs: the functions to compare, f and ua as well as the
7  function space U.
8  @author: juancarlos
9  """
10
11  from dolfin import *
12  import numpy as np
13  import matplotlib.pyplot as plt
14
15  def project_low(f,U,ub):
16      # Build the vectors to compare
17      fv = f.vector()[:]
18      uperv = ub.vector()[:]
19      # Indicatrix function
20      cond = np.where(fv>uperv,1,0)
21      # Build the vector with the binary values
22      proj_low = Function(U)
23      proj_low.vector()[:] = cond
24
25      return proj_low
```

## APPENDIX B IMPLEMENTATION IN CASHOCS

Here the interested reader can find the code used to implementate the PDAS to the ADR optimal Control Problem in `cashocs`

### B.1 MAIN PROGRAM

```
1  # Copyright (C) 2020–21 Sebastian Blauth
2  #
3  # This file is part of CASHOCS.
4  #
5  # CASHOCS is free software: you can redistribute it and/or modify
6  # it under the terms of the GNU General Public License as published by
7  # the Free Software Foundation, either version 3 of the License, or
8  # (at your option) any later version.
9  #
10 # CASHOCS is distributed in the hope that it will be useful,
11 # but WITHOUT ANY WARRANTY; without even the implied warranty of
12 # MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
13 # GNU General Public License for more details.
14 #
15 # You should have received a copy of the GNU General Public License
16 # along with CASHOCS. If not, see <https://www.gnu.org/licenses/>.
17
18 ## "For the documentation of this demo see
19 ## https://cashocs.readthedocs.io/en/latest/demos/optimal_control/doc_box_constraints.html.
20 ## "
21
22 from fenics import *
23 import numpy as np
24 import cashocs
25 import matplotlib.pyplot as plt
26 import matplotlib.pyplot as plt
27 from plotslopes import *
28 # config file: see https://cashocs.readthedocs.io/en/latest/demos/optimal_control/
29 # doc_box_constraints.html
30 config = cashocs.load_config('config_box.ini')
31 # Maximum number of uniform refinements, we start with 16 elements
32 Nref = 5
33 ### Parameters
34 epsi = Constant(1.0)
35 kappa = Constant(1.0)
36 be = Constant((3.0, 4.0))
37 g = Expression('6*x[0]-6*pow(x[0],2)+5*x[1]-13*x[0]*x[1]\
38 +7*pow(x[0],2)*x[1]-5*pow(x[1],2)+5*x[0]*pow(x[1],2)\
39 +pow(x[0],2)*pow(x[1],2)\
40 -min(2.0, max(-2.0, -(1-x[0])*x[0]*(1-x[1])*x[1]))'\
41 , degree=1)
42 yde = Expression('2*x[0] - 2*pow(x[0],2) + x[1] -14*x[0]*x[1] + 8*x[1]*pow(x[0],2) \
43 -pow(x[1],2) + 6*x[0]*pow(x[1],2) ', degree=1)
44 ### exact solutions
45 ye = Expression('x[0]*x[1]*(1-x[0])*(1-x[1]) ', degree=4)
46 pe = Expression('-x[0]*x[1]*(1-x[0])*(1-x[1]) ', degree=4)
47 ue = Expression('min(2.0, max(-2.0, -(1-x[0])*x[0]*(1-x[1])*x[1])/lam) ', lam=1, degree=3)
48 ### Store information ###
49 # a counter
50 array = np.arange(Nref)
51 # NDOF of each mesh
52 l = np.empty(Nref, dtype=int)
53 # Maximum diameter
54 h = np.empty(Nref)
55 # Errors
56 L2error = np.empty(Nref)
57 State_error_y = np.empty(Nref)
58 Adjointerror_p = np.empty(Nref)
59 ### Loops ###
60 # exterior loop: set each mesh and calculate the approx. errors
61 for j in array:
62     ## NDOF ##
63     l[j] = 2**(j+4)
64     ## mesh size
65     h[j] = np.divide(1, (np.multiply(2, l[j])))
66     ## mesh
67     mesh, subdomains, boundaries, dx, ds, dS = cashocs.regular_mesh(l[j])
68     V = FunctionSpace(mesh, 'CG', 1)
69     #W = FunctionSpace(mesh, 'DG', 0)
70     y = Function(V)
71     p = Function(V)
72     u = Function(V)
73     ## Box condition limits
74     u_a = interpolate(Expression('-2.0', degree=1), V)
75     u_b = interpolate(Expression('2.0', degree=1), V)
76
77     ### Primal (forward) weak equation ###
78
79     e = epsi*inner(grad(y), grad(p))*dx + inner(be, grad(y))*p*dx + kappa*y*p*dx - u*p*dx - g*p*dx
```



```

80
81 bcs = cashocs.create_bcs_list(V, Constant(0), boundaries, [1, 2, 3, 4])
82
83 alpha = Constant(1.0)
84 J = Constant(0.5)*(y - yde)*(y - yde)*dx + Constant(0.5*alpha)*u*u*dx
85
86 ## Box condition
87 cc = [u_a, u_b]
88 ## Solver
89 ocp = cashocs.OptimalControlProblem(e, bcs, J, y, u, p, config, control_constraints=cc)
90 ocp.solve('pdas', 1e-9, 0.0, 100)
91 ## primal dual Active set strategy, rel. tol=1e-9, abs tol.=0, max it.=100 --^
92 ## Check point
93 assert np.alltrue(u_a.vector()[:] <= u.vector()[:]) and np.alltrue(u.vector()[:] <= u_b.vector
94 ([:]))
95 L2error[j] = errornorm(ue, u, norm_type='L2', degree_rise=2)
96 State_error_y[j] = errornorm(ye, y, norm_type='H1', degree_rise=2)
97 Adjointerror_p[j] = errornorm(pe, p, norm_type='H1', degree_rise=2)
98 h[j]=1/(2*1[j])
99 # L2errorList.append(L2error)
100 # State_errorList_y.append(H1errorstate)
101 # AdjointerrorList_p.append(H1erroradjoint)}
102
103 print("maximal diameter of the mesh is =", h)
104 print("L2 error", L2error)
105 print("H1 error-state", State_error_y)
106 print("H1 error-adjoint", Adjointerror_p)
107
108 ## Post processing
109 # Plot the solutions
110 plt.figure(figsize=(15,5))
111 # control
112 plt.subplot(1, 3, 1)
113 fig = plot(u)
114 plt.colorbar(fig, fraction=0.046, pad=0.04)
115 plt.title('Control variable u')
116 vtkfile = File('control.pvd')
117 vtkfile << u
118 # state
119 plt.subplot(1,3,2)
120 fig = plot(y)
121 plt.colorbar(fig, fraction=0.046, pad=0.04)
122 plt.title('State variable y')
123 # adjoint state
124 plt.subplot(1,3,3)
125 fig = plot(-p) # We plot -p because cashocs use the scheme of Hinze to solve
126 plt.colorbar(fig, fraction=0.046, pad=0.04)
127 plt.title('Adjoint state p')
128 plt.show()
129 plt.savefig('solution_cashocs.png')
130 plt.tight_layout()
131
132 # error plots
133 plt.figure(2)
134 plt.rcParams.update({'font.size': 8.5})
135 # convergence rate (adapted by hand)
136 # adjusted control
137 const_shift = -0.350016 # shift of line in y-direction
138 const_rate = 1.99787 # slope for line in loglog Plot
139 rate = pow(10, const_shift)*np.exp((const_rate*np.log(h))) # function to plot
140 p1=plt.loglog(h,rate, 'm-', label='Adjusted control')
141 # adjusted state
142 const_shift = -0.314052 # shift of line in y-direction
143 const_rate = 0.9993 # slope for line in loglog Plot
144 ratey = pow(10, const_shift)*np.exp((const_rate*np.log(h))) # function to plot
145 p2=plt.loglog(h,ratey, 'r-', label='Adjusted state')
146 # adjusted adjoint state
147 const_shift = -0.313897 # shift of line in y-direction
148 const_rate = 0.999375 # slope for line in loglog Plot
149 ratep = pow(10, const_shift)*np.exp((const_rate*np.log(h))) # function to plot
150 p3=plt.loglog(h,ratep, 'g-', label='Adjusted adj. state')
151 p4=plt.loglog(h, L2error, 'y--o', label=r'$\Vert u-u_h\Vert_{L^2(\Omega)}$')
152 p5=plt.loglog(h, State_error_y, 'b--D', label=r'$\Vert y-y_h\Vert_{H^1(\Omega)}$')
153 p6=plt.loglog(h, Adjointerror_p, 'r--x', label=r'$\Vert p-p_h\Vert_{H^1(\Omega)}$')
154 slope_marker((h[2], L2error[2]-0.1*L2error[2]), (2,1))
155 slope_marker((h[2], State_error_y[2]-0.1*State_error_y[2]), (1,1))
156 plt.xlabel('h')
157 plt.ylabel('Error')
158 plt.title("Convergence History Plot, un. refinement")
159 plt.legend(bbox_to_anchor=(1.05, 1.0), loc='upper left')
160 plt.tight_layout()
161 plt.savefig('convergence_cashocs.pdf')
162 plt.show()

```