

Some Topics in Mathematical Optimization

Nguyễn Quân Bá Hồng¹

July 7, 2022

¹Independent Researcher, Ben Tre City, Vietnam
e-mail: nguyenquanbahong@gmail.com

Contents

1	Wikipedia's	5
1.1	Wikipedia/Shape Optimization	5
1.1.1	Definition	5
1.1.2	Examples	5
1.1.3	Techniques	5
1.1.3.1	Keeping track of the shape	5
1.1.3.2	Iterative methods using shape gradients	6
1.1.3.3	Geometry parametrization	6
1.2	Wikipedia/Topology Optimization	7
1.2.1	Problem statement	7
1.2.2	Implementation methodologies	7
1.2.2.1	Discrete	7
1.2.2.2	Solving the problem with continuous variables	8
1.2.2.3	Shape derivatives	8
1.2.2.4	Topological derivatives	8
1.2.2.5	Level set	8
1.2.2.6	Phase field	8
1.2.2.7	Evolutionary structural optimization	8
1.2.2.8	Commercial software	8
1.2.3	Examples	8
1.2.3.1	Structural compliance	8
1.2.3.2	Multiphysics problems	9
1.2.3.3	3F3D Form Follows Force 3D Printing	9
1.2.3.4	Design-dependent loads	9
1.3	Wikipedia/Water Supply Network	9
1.3.1	Water abstraction & raw water transfer	10
1.3.2	Water treatment	10
1.3.3	Water distribution network	10
1.3.3.1	Topologies	11
1.3.4	Water network maintenance	11
1.3.5	Sustainable urban water supply	11
1.3.5.1	Population growth	11
1.3.5.2	Water scarcity	12
1.3.5.3	Governmental issues	12
1.3.6	Optimizing the water supply network	12
1.3.6.1	Single-objective optimization	12
1.3.6.2	Multi-objective optimization	12
1.3.6.3	Weighted sum method	13
1.3.6.4	The constraint method	13
1.3.6.5	Sensitivity analysis	13
1.3.6.6	Operational constraints	13
1.3.7	Sustainable development	13
1.3.8	Future approaches	14

I	Optimal Control	15
----------	------------------------	-----------

2	Lions, 1971. Optimal Control of Systems Governed by PDEs	16
----------	---	-----------

2.1	Minimization of Functions & Unilateral BVPs	19
2.1.1	Minimization of Coercive Forms	19
2.1.1.1	Notation	19
2.1.1.2	The Case when π is Coercive	20
2.1.1.3	Characterization of the Minimizing Element. Variational Inequalities	20
2.1.1.4	Alternative Form of Variational Inequalities	21
2.1.1.5	Function f being the Sum of a Differentiable & Non-Differentiable Function	21
2.1.1.6	The Convexity Hypothesis on \mathcal{U}_{ad}	22
2.1.2	A Direct Solution of Certain Variational Inequalities	23
2.1.3	Examples	23
2.1.4	A Comparison Theorem	23
2.1.5	Non Coercive Forms	23
2.2	Control of Systems Governed by Elliptic PDEs	23
2.2.1	Control of Elliptic Variational Problems	23
2.2.2	1st Applications	23
2.2.3	A Family of Examples with $N = 0$ & \mathcal{U}_{ad} Arbitrary	23
2.2.4	Observation on the Boundary	23
2.2.5	Control & Observation on the Boundary. Case of the Dirichlet Problem	23
2.2.6	Constraints on the State	23
2.2.7	Existence Results for Optimal Controls	23
2.2.8	1st Order Necessary Conditions	23
2.3	Control of Systems Governed by Parabolic PDEs	23
2.4	Control of Systems Governed by Hyperbolic Equations or by Equations which are well Posed in the Petrowsky Sense	23
2.5	Regularization, Approximation & Penalization	23
2.5.1	Regularization	23
2.5.1.1	Parabolic Regularization	23
2.5.1.2	Approximation in Terms of Systems of Cauchy–Kowaleska Type	23
2.5.1.3	Penalization	23
3	Tröltzsch, 2010. Optimal Control for PDEs	24
3.1	Introduction & Examples	25
3.1.1	What is Optimal Control?	25
3.1.2	Examples of Convex Problems	25
3.1.2.1	Optimal boundary heating	25
3.1.2.2	Optimal nonstationary boundary control	26
3.1.2.3	Optimal vibrations	27
3.1.3	Examples of Nonconvex Problems	27
3.1.3.1	Problems involving semilinear elliptic equations	27
3.1.3.2	Problems involving semilinear parabolic equations	28
3.1.4	Basic Concepts for the Finite-Dimensional Case	28
3.2	Linear-Quadratic Elliptic Control Problems	28
3.3	Linear-Quadratic Parabolic Control Problems	28
3.4	Optimal Control of Semilinear Elliptic Equations	28
3.5	Optimal Control of Semilinear Parabolic Equations	28
3.6	Optimization Problems in Banach Spaces	28
3.6.1	The Karush–Kuhn–Tucker Conditions	28
3.6.1.1	Convex problems	28
3.6.1.2	Differentiable problems	31
3.6.1.3	A semilinear elliptic problem	33
3.6.2	Control Problems with State Constraints	34
3.6.2.1	Convex problems	34
3.7	Supplementary Results on PDEs	34
4	Shape Optimization	35
4.1	Introduction	35
4.1.1	Classical & moving shape analysis	36
4.1.2	Fluid-Structure Interaction Problems	36
4.1.2.1	Control of a fluid flow around a fixed body	38

4.1.2.2	Shape design of a fixed solid inside a fluid flow	40
4.1.2.3	Dynamical shape design of a solid inside a fluid flow	43
4.2	Inverse Stefan Problem	45
4.2.1	The inverse problem setting	45
4.2.2	The Eulerian derivative & the transverse field	45
4.3	Dynamical Shape Control of NSEs	46
4.3.1	Problem Statement	47
4.3.2	Elements of Non-cylindrical Shape Calculus	49
4.3.2.1	Non-cylindrical speed method	49
4.3.3	Elements of tangential calculus	52
4.3.3.1	Oriented distance function	52
5	Topology Optimization	53
	Bibliography	54

Foreword

A collection of & some personal notes on Mathematical Optimization, especially the 3 major topics: Optimal Control, Shape Optimization, & Topology Optimization.

Keywords. Optimal control; Shape optimization; Topology optimization.

Chapter 1

Wikipedia's

1.1 Wikipedia/Shape Optimization

“*Shape optimization* is part of the field of **optimal control** theory. The typical problem is to find the **shape** which is optimal in that it minimizes a certain cost **functional** while satisfying given **constraints**. In many cases, the functional being solved depends on the solution of a given **PDE** defined on the variable domain.

Topology optimization is, in addition, concerned with the number of connected components/boundaries belonging to the domain. Such methods are needed since typically shape optimization methods work in a subset of allowable shapes which have fixed topological properties, such as having a fixed number of holes in them. Topological optimization techniques can then help work around the limitations of pure shape optimization.” – [Wikipedia/shape optimization](#)

1.1.1 Definition

“Mathematically, shape optimization can be posed as the problem of finding a **bounded set** Ω , **minimizing a functional** $\mathcal{F}(\Omega)$, possibly subject to a **constraint** of the form $\mathcal{G}(\Omega) = 0$. Usually we are interested in sets Ω which are **Lipschitz** or C^1 **boundary** & consist of finite many **components**, which is a way of saying that we would like to find a rather pleasing shape as a solution, not some jumble of rough bits & pieces. Sometimes additional constraints need to be imposed to that end to ensure well-posedness of the problem & uniqueness of the solution.

Shape optimization is an **infinite-dimensional optimization** problem. Furthermore, the space of allowable shapes over which the optimization is performed does not admit a **vector space** structure, making application of traditional optimization methods more difficult.” – [Wikipedia/shape optimization/definition](#)

1.1.2 Examples

- “among all 3D shapes of given volume, find the one which has minimal surface area. Here: $\mathcal{F}(\Omega) = \text{Area}(\partial\Omega)$, with $\mathcal{G}(\Omega) = \text{Volume}(\Omega) = \text{const.}$ The answer, given by the **isoperimetric inequality**, is a **ball**.
- Find the shape of an airplane wing which minimizes **drag**. Here the constraints could be the wing strength, or the wing dimensions.
- Find the shape of various mechanical structures, which can resist a given **stress** while having a minimal mass/volume.
- Given a known 3D object with a fixed radiation source inside, deduce the shape & size of the source based on measurements done on part of the boundary of the object. A formulation of this **inverse problem** using **least squares** fit leads to a shape optimization problem.” – [Wikipedia/shape optimization/examples](#)

1.1.3 Techniques

“Shape optimization problems are usually solved **numerically**, by using **iterative methods**. I.e., one starts with an initial guess for a shape, & then gradually evolves it, until it morphs into the optimal shape.” – [Wikipedia/shape optimization/techniques](#)

1.1.3.1 Keeping track of the shape

Fig. Example: Shape optimization as applied to building geometry. Example provided courtesy of [Formsolver.com](#).

“To solve a shape optimization problem, one needs to find a way to represent a shape in the **computer memory**, & follow its evolution. Several approaches are usually used.

1 approach is to follow the boundary of the shape. For that, one can sample the shape boundary in a relatively dense & uniform manner, i.e., to consider enough points to get a sufficiently accurate outline of the shape. Then, one can evolve the shape by gradually moving the boundary points. This is called the *Lagrangian approach*.

Another approach is to consider a **function** defined on a rectangular box around the shape, which is positive inside of the shape, zero on the boundary of the shape, & negative outside of the shape. One can then evolve this function instead of the shape itself. One can consider a rectangular grid on the box & sample the function at the grid points. As the shape evolves, the grid points do not change; only the function values at the grid points change. This approach, of using a fixed grid, is called the *Eulerian approach*. The idea of using a function to represent the shape is at the basis of the **level set method**.

Fig. Example: Optimization shape families resulting from differing goal parameters. Example provided courtesy of **Formsolver.com**

A 3rd approach is to think of the shape evolution as of a flow problem. I.e., one can imagine that the shape is made of a plastic material gradually deforming s.t. any point inside or on the boundary of the shape can be always traced back to a point of the original shape in a 1-1 fashion. Mathematically, if Ω_0 is the initial shape, & Ω_t is the shape at time t , one considers the **diffeomorphisms** $f_t : \Omega_0 \rightarrow \Omega_t$, for $t \leq t \leq t_0$. The idea is again that shapes are difficult entities to be dealt with directly, so manipulate them by means of a function.” – [Wikipedia/shape optimization/techniques/keeping track of the shape](#)

1.1.3.2 Iterative methods using shape gradients

“Consider a smooth velocity field V & the family of transformations T_s of the initial domain Ω_0 under the velocity field V : $x(0) = x_0 \in \Omega_0$, $x'(s) = V(x(s))$, $T_s(x_0) = x(s)$, $s \geq 0$, & denote $\Omega_0 \mapsto T_s(\Omega_0) = \Omega_s$. Then the Gâteaux or shape derivative of $\mathcal{F}(\Omega)$ at Ω_0 w.r.t. the shape is the limit of

$$d\mathcal{F}(\Omega_0; V) = \lim_{s \rightarrow 0} \frac{\mathcal{F}(\Omega_s) - \mathcal{F}(\Omega_0)}{s}$$

if this limit exists. If in addition the derivative is linear w.r.t. V , there is a unique element of $\nabla \mathcal{F} \in L^2(\partial\Omega)$ & $d\mathcal{F}(\Omega_0; V) = \langle \nabla \mathcal{F}, V \rangle_{\partial\Omega_0}$ where $\nabla \mathcal{F}$ is called the *shape gradient*. This gives a natural idea of **gradient descent**, where the boundary $\partial\Omega$ is evolved in the direction of negative shape gradient in order to reduce the value of the cost functional. Higher order derivatives can be similarly defined, leading to Newtonlike methods.

Typically, gradient descent is preferred, even if requires a large number of iterations, because, it can be hard to compute the 2nd-order derivative (i.e., the **Hessian**) of the objective functional \mathcal{F} .

If the shape optimization problem has constraints, i.e., the functional \mathcal{G} is present, one has to find ways to convert the constrained problem into an unconstrained one. Sometimes ideas based on **Lagrange multipliers**, like the **adjoint state method**, can work.” – [Wikipedia/shape optimization/techniques/iterative methods using shape gradients](#)

1.1.3.3 Geometry parametrization

“Shape optimization can be faced using standard optimization methods if a parametrization of the geometry is defined. Such parametrization is very important in CAE field where goal functions are usually complex functions evaluated using numerical models (CFD, FEA, ...). A convenient approach, suitable for a wide class of problems, consists in the parametrization of the CAD model coupled with a full automation of all the process required for function evaluation (meshing, solving & result processing). *Mesh morphing* is a valid choice for complex problems that resolves typical issues associated with *re-meshing* such as discontinuities in the computed objective & constraint functions. In this case the parametrization is defined after the meshing stage acting directly on the numerical model used for calculation that is changed using mesh updating methods. There are several algorithms available for mesh morphing (*deforming volumes*, *pseudosolids*, **radical basis functions**). The selection of the parametrization approach depends mainly on the size of the problem: the CAD approach is preferred for small-to-medium sized models whilst the mesh morphing approach is the best (& sometimes the only feasible one) for large & very large models. The multi-objective Pareto optimization (NSGA II) could be utilized as a powerful approach for shape optimization. In this regard, the Pareto optimization approach displays useful advantages in design method such as the effect of area constraint that other multi-objective optimization cannot declare it. The approach of using a penalty function is an effective technique which could be used in the 1st stage of optimization. In this method the constrained shape design problem is adapted to an unconstrained problem with utilizing the constraints in the objective function as a penalty factor. Most of the time penalty factor is dependent to the amount of constraint variation rather than constrain number. The GA real-coded technique is applied in the present optimization problem. Therefore, the calculations are based on real value of variables.” – [Wikipedia/shape optimization/techniques/geometry parametrization](#)

1.2 Wikipedia/Topology Optimization

“*Topology optimization (TO)* is a mathematical method that optimizes material layout within a given design space, for a given set of **loads**, **boundary conditions** & **constraints** with the goal of maximizing the performance of the system. Topology optimization is different from **shape optimization** & sizing optimization in the sense that the design can attain any shape within the design space, instead of dealing with predefined configurations.

The conventional topology optimization formulation uses a **FEM** to evaluate the design performance. The design is optimized using either gradient-based **mathematical programming** techniques such as the optimality criteria algorithm & the *method of moving asymptotes* or non gradient-based algorithms such as **genetic algorithms**.

Topology optimization has a wide range of applications in aerospace, mechanical, bio-chemical & civil engineering. Currently, engineers mostly use topology optimization at the concept level of a **design process**. Due to the free forms that naturally occur, the result is often difficult to manufacture. For that reason the result emerging from topology optimization is often fine-tuned for manufacturability. Adding constraints to the formulation in order to **increase the manufacturability** is an active field of research. In some cases results from topology optimization can be directly manufactured using **additive manufacturing**; topology optimization is thus a key part of **design for additive manufacturing**.” – [Wikipedia/topology optimization](#)

1.2.1 Problem statement

“A topology optimization problem can be written in the general form of an **optimization problem** as

$$\min_{\rho} F \text{ where } F = F(\mathbf{u}(\rho), \rho) = \int_{\Omega} f(\mathbf{u}(\rho), \rho) dV \text{ subject to } G_0(\rho) = \int_{\Omega} \rho dV - V_0 \leq 0, \quad G_j(\mathbf{u}(\rho), \rho) \leq 0, \quad j = 1, \dots, m.$$

The problem statement includes the following:

- An **objective function** $F(\mathbf{u}(\rho), \rho)$. This function represents the quantity that is being minimized for best performance. The most common objective function is compliance, where minimizing compliance leads to maximizing the stiffness of a structure.
- The *material distribution* as a problem variable. This is described by the density of the material at each location $\rho(\mathbf{x})$. Material is either *present*, indicated by a 1, or *absent*, indicated by a 0. $\mathbf{u} = \mathbf{u}(\rho)$ is a state field that satisfies a linear or nonlinear state equation depending on ρ .
- The *design space* (Ω). This indicates the allowable volume within which the design can exist. Assembly & packaging requirements, human & tool accessibility are some of the factors that need to be considered in identifying this space. With the definition of the design space, regions or components in the model that cannot be modified during the course of the optimization are considered as non-design regions.
- m **constraints** $G_j(\mathbf{u}(\rho), \rho) \leq 0$ a characteristic that the solution must satisfy. Examples are the maximum amount of material to be distributed (volume constraint) or maximum stress values.

Evaluating $\mathbf{u}(\rho)$ often includes solving a differential equation. This is most commonly done using the FEM since these equations do not have a known analytical solution.” – [Wikipedia/topology optimization/problem statement](#)

1.2.2 Implementation methodologies

“There are various implementation methodologies that have been used to solve topology optimization problems.”

1.2.2.1 Discrete

“Solving topology optimization problems in a discrete sense is done by discretizing the design domain into finite elements. The material densities inside these elements are then treated as the problem variables. In this case material density of 1 indicates the presence of material, while 0 indicates an absence of material. Owing to the attainable topological complexity of the design being dependent on the number of elements, a large number is preferred. Large numbers of finite elements increases the attainable topological complexity, but come at a cost. 1stly, solving the FEM systems becomes more expensive. 2ndly, algorithms that can handle a large number (several thousands of elements is not uncommon) of discrete variables with multiple constraints are unavailable. Moreover, they are impractically sensitive to parameter variations. In literature problems with up to 30000 variables have been reported.” – [Wikipedia/topology optimization/implementation methodologies/discrete](#)

1.2.2.2 Solving the problem with continuous variables

“The earlier stated complexities with solving topology optimization problems using binary variables has caused the community to search for other options. One is the modeling of the densities with continuous variables. The material densities can now also attain values between 0 & 1. Gradient based algorithms that handle large amounts of continuous variables & multiple constraints are available. But the material properties have to be modeled in a continuous setting. This is done through interpolation. 1 of the most implemented interpolation methodologies is the *Solid Isotropic Material with Penalization* method (SIMP). This interpolation is essentially a power law $E = E_0 + \rho^p(E_1 - E_0)$. It interpolates the Young’s modulus of the material to the scalar selection field. The value of the penalization parameter p is generally taken between $[1, 3]$. This has been shown to confirm the micro-structure of the materials. In the SIMP method a lower bound on the Young’s modulus is added, E_0 , to make sure the derivatives of the objective function are nonzero when the density becomes 0. The higher the penalization factor, the more SIMP penalizes the algorithm in the use of non-binary densities. Unfortunately, the penalization parameter also introduces non-convexities.” – [Wikipedia/topology optimization/implementation methodologies/solving the problem with continuous variables](#)

1.2.2.3 Shape derivatives

“Topology optimization can be achieved by using shape derivatives.”

1.2.2.4 Topological derivatives

1.2.2.5 Level set

1.2.2.6 Phase field

1.2.2.7 Evolutionary structural optimization

1.2.2.8 Commercial software

“There are several commercial topology optimization software on the market. Most of them use topology optimization as a hint how the optimal design should look like, & manual geometry re-construction is required. There are a few solutions which produce optimal designs ready for Additive Manufacturing.” – [Wikipedia/topology optimization/implementation methodologies/commercial software](#)

1.2.3 Examples

1.2.3.1 Structural compliance

Fig. Checker Board Patterns are shown in this result. Fig. Topology optimization result when filtering is used. Fig. Topology optimization of a compliance problem.

“A stiff structure is one that has the least possible displacement when given certain set of boundary conditions. A global measure of the displacements is the **strain energy** (also called **compliance**) of the structure under the prescribed boundary conditions. The lower the strain energy the higher the stiffness of the structure. So, the objective function of the problem is to minimize the strain energy.

On a broad level, one can visualize that the more the material, the less the deflection¹ as there will be more material to resist the loads. So, the optimization requires an opposing constraint, the volume constraint. This is in reality a cost factor, as we would not want to spend a lot of money on the material. To obtain the total material utilized, an integration of the selection field over the volume can be done.

Finally the elasticity governing differential equations are plugged in so as to get the final problem statement.

$$\min_{\rho} \int_{\Omega} \frac{1}{2} \boldsymbol{\sigma} : \boldsymbol{\varepsilon} \, d\Omega \text{ subject to } \rho \in [0, 1], \int_{\Omega} \rho \, d\Omega \leq V^*, \nabla \cdot \boldsymbol{\sigma} + \mathbf{F} = \mathbf{0}, \boldsymbol{\sigma} = \mathbf{C} : \boldsymbol{\varepsilon}.$$

But, a straightforward implementation in the finite element framework of such a problem is still infeasible² owing to issues such as:

- **Mesh dependency** i.e., the design obtained on 1 mesh is not the one that will be obtained on another mesh. The features of the design become more intricate³ as the mesh gets refined.
- **Numerical instabilities.** The selection of region in the form of a chess board.

¹**deflection** [n] [uncountable, countable, usually singular] **deflection (of something)** a sudden change in the direction that something is moving in, usually after it has hit something; the act of causing something to change direction.

²**unfeasible** [a] not possible to do or achieve, OPPOSITE: **feasible**.

³**intricate** [a] having a lot of different parts & small details that fit together.

Some techniques such as [filtering](#) based on image processing are currently being used to alleviate⁴ some of these issues. Although it seemed like this was purely a heuristic approach for a long time, theoretical connections to nonlocal elasticity have been made to support the physical sense of these methods.” – [Wikipedia/topology optimization/examples/structural compliance](#)

1.2.3.2 Multiphysics problems

Fluid-structure-interaction. “[Fluid-structure-interaction](#) is a strongly coupled phenomenon & concerns the interaction between a stationary or moving fluid & an elastic structure. Many engineering applications & natural phenomenon are subject to fluid-structure interaction & to take such effects into consideration is therefore critical in the design of many engineering applications. Topology optimization for fluid structure interaction problems has been studied. Design solutions solved for different Reynolds numbers are shown below. The design solutions depend on the fluid flow with indicate that the coupling between the fluid & the structure is resolved in the design problems.” – [Wikipedia/topology optimization/examples/multiphysics problems/fluid-structure-interaction](#)

Fig. Design solutions for different Reynolds number for a wall inserted in a channel with a moving fluid. Fig. Sketch fo the well-known wall problem. The objective of the design problem is to minimize the structural compliance. Fig. Design evolution for a fluid-structure-interaction problem. The objective of the design problem is to minimize the structural compliance. The fluid-structure-interaction problem is modeled with Navier–Cauchy & NSEs.

Thermoelectric energy conversion. “[Thermoelectricity](#) is a multi-physic problem which concerns the interaction & coupling between electric & thermal energy in semi conducting materials. Thermoelectric energy conversion can be described by 2 separately identified effects: The Seebeck effect & the Peltier effect. The Seebeck effect concerns the conversion of thermal energy into electric energy & the Peltier effect concerns the conversion of electric energy into thermal energy. By spatially distributing 2 thermoelectric materials in a 2D design space with a topology optimization methodology, it is possible to exceed performance of the constitutive thermoelectric materials for [thermoelectric coolers](#) & [thermoelectric generators](#).” – [Wikipedia/topology optimization/examples/multiphysics problems/thermoelectric energy conversion](#)

1.2.3.3 3F3D Form Follows Force 3D Printing

“The current proliferation⁵ of 3D printer technology has allowed designers & engineers to use topology optimization techniques when designing new products. Topology optimization combined with 3D printing can result in less weight, improved structural performance & shortened design-to-manufacturing cycle. As the designs, while efficient, might not be realizable with more traditional manufacturing techniques.” – [Wikipedia/topology optimization/examples/3F3D Form Follows Force 3D Printing](#)

1.2.3.4 Design-dependent loads

“The direction, magnitude, & location of a design-dependent load alter with topology optimization iterations. Therefore, dealing with such loads in a TO setting is a challenging task. One can find novel methods to deal with such loads (e.g. pressure load, self-weight, etc.).” – [Wikipedia/topology optimization/examples/design-dependent loads](#)

Fig. A sketch of the design problem. The aim of the design problem is to spatially distribute 2 materials, Material A & Material B, to maximize a performance measure such as cooling power or electric power output. Fig. Design evolution for an off-diagonal thermoelectric generator. The design solution of an optimization problem solved for electric power output. The performance of the device has been optimized by distributing [Skutterudite](#) (yellow) & [bismuth telluride](#) (blue) with a density-based topology optimization methodology. The aim of the optimization problem is to maximize the electric power output of the thermoelectric generator. Fig. Design evolution for a thermoelectric cooler. The aim of the design problem is to maximize the cooling power of the thermoelectric cooler.

1.3 Wikipedia/Water Supply Network

“A *water supply network* or *water supply system* is a system of engineered [hydrologic](#) & [hydraulic](#) components that provide [water supply](#). A water supply system typically includes the following:

1. A [drainage basin](#) (see [water purification – sources of drinking water](#))
2. A [raw water](#) collection point (above or below ground) where the water accumulates, such as a [lake](#), a [river](#), or [ground-water](#) from an [underground aquifer](#). Raw water may be transferred using uncovered ground-level [aqueducts](#), covered [tunnels](#), or underground [water pipes](#) to water purification facilities.

⁴[alleviate](#) [v] **alleviate something** to make suffering or a problem less severe.

⁵[proliferation](#) [n] **1.** [uncountable, singular] **proliferation (of something)** a rapid increase in the number or amount of something; a large number of a particular thing; **2.** [uncountable] (*biology*) the rapid reproduction of a cell, part or organism.

3. **Water purification** facilities. Treated water is transferred using **water pipes** (usually underground).
4. Water storage facilities such as **reservoirs**, **water tanks**, or **water towers**. Smaller water systems may store the water in **cisterns** or **pressure vessels**. Tall buildings may also need to store water locally in pressure vessels in order for the water to reach the upper floors.
5. Additional water pressurizing components such as **pumping stations** may need to be situated at the outlet of underground or aboveground reservoirs or cisterns (if gravity flow is impractical).
6. A pipe network for distribution of water to consumers (which may be private houses or industrial, commercial, or institution establishments) & other usage points (such as **fire hydrants**)
7. Connections to the **sewers** (underground pipes, or aboveground **ditches** in some developing countries) are generally found downstream of the water consumers, but the sewer system is considered to be a separate system, rather than part of the water supply system.

Water supply networks are often run by **public utilities** of the **water industry**.” – **Wikipedia/water supply network**

1.3.1 Water abstraction & raw water transfer

“**Raw water** (untreated) is from a **surface water** source (such as an intake on a **lake** or a **river**) or from a **groundwater** source (such as a **water well** drawing from an underground **aquifer**) within the **watershed** that provides the **water resources**.

The raw water is transferred to the water purification facilities using uncovered aqueducts, covered tunnels or underground **water pipes**.” – **Wikipedia/water supply network/water abstraction & raw water transfer**

1.3.2 Water treatment

“Main article: **Wikipedia/water treatment**. Virtually all large systems must treat the water; a fact that is tightly regulated by global, state & federal agencies, such as the **World Health Organization** (WHO) or the **United States Environmental Protection Agency** (EPA). Water treatment must occur before the product reaches the consumer & afterwards (when it is discharged again). Water purification usually occurs close to the final delivery points to reduce pumping costs & the chances of the water becoming contaminated after treatment.

Traditional surface water treatment plants generally consists of 3 steps: clarification, filtration & disinfection. Clarification refers to the separation of particles (dirt, organic matter, etc.) from the water stream. Chemical addition (i.e., alum, ferric chloride) destabilizes the particle charges & prepares them for clarification either by setting or floating out of the water stream. Sand, anthracite or activated carbon filters refine the water stream, removing smaller particulate matter. While other methods of disinfection exist, the preferred method is via choline addition. Chlorine effectively kills bacteria & most viruses & maintains a residual to protect the water supply through the supply network.” – **Wikipedia/water supply network/water treatment**

1.3.3 Water distribution network

Fig. USA Not Combined City Water System. “Main article: **Wikipedia/water distribution system**. The product, delivered to the point of consumption, is called **potable water** if it meets the **water quality** standards required for human consumption.

The water in the supply network is maintained at positive **pressure** to ensure that water reaches all parts of the network, that a sufficient flow is available at every take-off point & to ensure that untreated water in the ground cannot enter the network. The water is typically pressurized by pumping the water into storage tanks constructed at the highest local point in the network. 1 network may have several such **service reservoirs**.

In small domestic systems, the water may be pressurized by a **pressure vessel** or even by an **underground cistern** (the latter however does need additional pressurizing). This eliminates the need of a water-tower or any other heightened water reserve to supply the water pressure.

These systems are usually owned & maintained by local governments, such as cities, or other public entities, but are occasionally operated by a commercial enterprise (see **water privatization**). Water supply networks are part of the master planning of communities, counties, & municipalities. Their planning & design requires the expertise of **city planners** & **civil engineers**, who must consider many factors, such as location, current demand, further growth, leakage, pressure, pipe size, pressure loss, fire fighting flows, etc. – using **pipe network analysis** & other tools.

As water passes through the distribution system, the water quality can degrade by chemical reactions & biological processes. **Corrosion** of metal pipe materials in the distribution system can cause the release of metals into the water with undesirable aesthetic & health effects. Release of **iron** from unlined iron pipes can result in customer reports of “red water” at the tap. Release of **copper** from **copper pipes** can result in customer reports of “blue water” &/or a metallic taste. Release

of **lead** can occur from the **solder** used to join copper pipe together or from **brass fixtures**. Copper & lead levels at the consumer's tap are regulated to protect consumer health.

Utilities will often adjust the chemistry of the water before distribution to minimize its corrosiveness. The simplest adjustment involves control of **pH** & **alkalinity** to produce a water that tends to passivate corrosion by depositing a layer of **calcium carbonate**. **Corrosion inhibitors** are often added to reduce release of metals into the water. Common corrosion inhibitors added to the water are **phosphates** & **silicates**.

Maintenance of a biologically safe drinking water is another goal in water distribution. Typically, a chlorine based **disinfectant**, such as **sodium hypochlorite** or **monochloramine** is added to the water as it leaves the treatment plant. Booster stations can be placed within the distribution system to ensure that all areas of the distribution system have adequate sustained levels of **disinfection**.” – [Wikipedia/water supply network/water distribution network](#)

Fig. The **Central Arizona Project Aqueduct** transfers untreated water. Fig. Most (treated) water distribution happens through underground pipes. Fig. Pressurizing the water is required between the small water reserve & the end-user.

1.3.3.1 Topologies

“Like electric power lines, roads, & microwave radio networks, water systems may have a **loop** or **branch** network topology, or a combination of both. The piping networks are circular or rectangular. If any 1 section of water distribution main fails or needs repair, that section can be isolated without disrupting all users on the network.

Most systems are divided into zones. Factors determining the extent or size of a zone can include hydraulics, **telemetry** systems, history, & population density. Sometimes systems are designed for a specific area then are modified to accommodate development. Terrain affects hydraulics & some forms of telemetry. While each zone may operate as a stand-alone system, there is usually some arrangements to interconnect zones in order to manage equipment failures or system failures.” – [Wikipedia/water supply network/water distribution network/topologies](#)

1.3.4 Water network maintenance

“Water supply networks usually represent the majority of assets of a water utility. Systematic documentation of maintenance works using a **computerized maintenance management system** (CMMS) is a key to a successful operation of a water utility.” – [Wikipedia/water supply network/water network maintenance](#)

1.3.5 Sustainable urban water supply

Fig. Clean drinking water is essential to human life.

“A sustainable urban water supply network covers all the activities related to provision of **potable water**. **Sustainable development** is of increasing importance for the water supply to urban areas. Incorporating innovative water technologies into **water supply** systems improves water supply from sustainable perspectives. The development of innovative water technologies provides flexibility to the water supply system, generating a fundamental & effective means of sustainability based on an integrated **real options** approach.

Water is an essential **natural resource** for human existence. It is needed in every industrial & natural process, e.g., it is used for **oil refining**, for **liquid-liquid extraction** in hydro-metallurgical processes, for cooling, for scrubbing in the iron & the steel industry, & for several operations in **food processing** facilities.

It is necessary to adopt a new approach to design urban water supply networks; **water shortages** are expected in the forthcoming decades & environmental regulations for water utilization & **waste-water** disposal are increasingly stringent.

To achieve a sustainable water supply network, new sources of water are needed to be developed, & to reduce environmental pollution.

The price of water is increasing, so less water must be wasted & actions must be taken to prevent pipeline leakage. Shutting down the supply service to fix leaks is less & less tolerated by consumers. A sustainable water supply network must monitor the freshwater consumption rate & the waste-water generation rate.

Many of the urban water supply networks in **developing countries** face problems related to **population increase**, **water scarcity**, & **environmental pollution**.” – [Wikipedia/water supply network/sustainable urban water supply](#)

1.3.5.1 Population growth

“In 1900 just 13% of the global population lived in cities. By 2005, 49% of the **global population** lived in urban areas. In 2030 it is predicted that this statistic will rise to 60%. Attempts to expand water supply by governments are costly & often not sufficient. The building of new illegal settlements makes it hard to map, & make connections to, the water supply, leads to inadequate water management. In 2002, there were 158 million people with inadequate **water supply**. An increasing number of people live in **slums**, in inadequate sanitary conditions, & are therefore at risk of **disease**.” – [Wikipedia/water supply network/sustainable urban water supply/population growth](#)

1.3.5.2 Water scarcity

“**Potable water** is not well distributed in the world. 1.8 million deaths are attributed to unsafe water supplies every year, according to the **WHO**. Many people do not have any access, or do not have access to quality & quantity of potable water, though water itself is abundant. Poor people in developing countries can be close to major rivers, or be in high rainfall areas, yet not have access to potable water at all. There are also people living where lack of water creates millions of deaths every year.

Where the water supply system cannot reach the slums, people manage to use **hand pumps**, to reach the pit wells, **rivers**, **canals**, **swamps** & any other source of water. In most cases the water quality is unfit for human consumption. The principal cause of water scarcity is the growth in demand. Water is taken from remote areas to satisfy the needs of urban areas. Another reason for water scarcity is **climate change**: **precipitation** patterns have changed; rivers have decreased their flow; **lakes** are drying up; & **aquifers** are being emptied.” – [Wikipedia/water supply network/sustainable urban water supply/water scarcity](#)

1.3.5.3 Governmental issues

“In developing countries many governments are **corrupt** & poor & they respond to these problems with frequently changing policies & non clear agreements. Water demand exceeds supply, & household & industrial water supplies are prioritized over other uses, which leads to **water stress**. Potable water has a price in the market; water often becomes a **business** for private companies, which earn a **profit** by putting a higher price on water, which imposes a barrier for lower-income people. The **Millennium Development Goals** propose the changes required.

Goal 6 of the United Nations’ **Sustainable Development Goals** is to “Ensure availability & sustainable management of water & sanitation for all”. This is in recognition of the human right to water & sanitation, which was formally acknowledged at the United Nations General Assembly in 2010, that “clean drinking water & sanitation are essential to the recognition of all human rights”. Sustainable water supply includes ensuring availability, accessibility, affordability & quality of water for all individuals.

In advanced economies, the problems are about optimizing existing supply networks. These economies have usually had continuing evolution, which allowed them to construct infrastructure to supply water to people. The **European Union** has developed a set of rules & policies to overcome expected future problems.

There are many international documents with interesting, but not very specific, ideas & therefore they are not put into practice. Recommendations have been made by the **United Nations**, such as the **Dublin Statement on Water & Sustainable Development**.” – [Wikipedia/water supply network/sustainable urban water supply/governmental issues](#)

1.3.6 Optimizing the water supply network

“The yield of a system can be measured by either its value or its net benefit. For a water supply system, the true value or the net benefit is a reliable water supply service having adequate quantity & good quality of the product. E.g., if the existing water supply of a city needs to be extended to supply a new **municipality**, the impact of the new branch of the system must be designed to supply the new needs, while maintaining supply to the old system.” – [Wikipedia/water supply network/optimizing the water supply network](#)

1.3.6.1 Single-objective optimization

“The design of a system is governed by multiple criteria, one being cost. If the benefit is *fixed*, the **least cost** design results in maximum benefit. However, the least cost approach normally results in a *minimum capacity* for a water supply network. A minimum cost model usually searches for the least cost solution (in pipe sizes), while satisfying the hydraulic constraints such as: required output pressures, maximum **pipe flow** rate & pipe flow velocities. The cost is a function of pipe diameters; therefore the **optimization** problem consists of finding a minimum cost solution by optimizing pipe sizes to provide the minimum acceptable capacity.” – [Wikipedia/water supply network/optimizing the water supply network/single-objective optimization](#)

1.3.6.2 Multi-objective optimization

“However, according to the authors of the paper entitled, “Method for optimizing design & rehabilitation of water distribution systems”, “the least capacity is not a desirable solution to a sustainable water supply network in a long term, due to the uncertainty of the future demand”. It is preferable to provide extra pipe capacity to cope with unexpected demand growth & with water outages. The problem changes from a single objective optimization problem (minimizing cost), to a multi-objective optimization problem (minimizing cost & maximizing flow capacity).” – [Wikipedia/water supply network/optimizing the water supply network/multi-objective optimization](#)

1.3.6.3 Weighted sum method

“To solve a multi-objective optimization problem, it is necessary to convert the problem into a single objective optimization problem, by using adjustments, such as a weighted sum of **objectives**, or an ε -constraint method. The weighted sum approach gives a certain weight to the different objectives, & then factors in all these weights to form a single objective function that can be solved by single factor optimization. This method is not entirely satisfactory, because the weights cannot be correctly chosen, so this approach cannot find the optimal solution for all the original objectives.” – [Wikipedia/water supply network/optimizing the water supply network/weighted sum method](#)

1.3.6.4 The constraint method

“The 2nd approach (the constraint method), chooses 1 of the objective functions as the single objective, & the other objective functions are treated as constraints with a limited value. However, the optimal solution depends on the pre-defined constraints limits.” – [Wikipedia/water supply network/optimizing the water supply network/the constraint method](#)

1.3.6.5 Sensitivity analysis

“The multiple objective optimization problems involve computing the **tradeoff** between the costs & benefits resulting in a set of solutions that can be used for sensitivity analysis & tested in different scenarios. But there is no single optimal solution that will satisfy the global optimality of both objectives. As both objectives are to some extent contradictory, it is not possible to improve 1 objective without sacrificing the other. It is necessary in some cases use a different approach. (e.g., **Pareto Analysis**), & choose the best combination.” – [Wikipedia/water supply network/optimizing the water supply network/sensitivity analysis](#)

1.3.6.6 Operational constraints

“Returning to the cost objective function, it cannot violate any of the operational constraints. Generally this cost is dominated by the energy cost for pumping. “The operational constraints include the standards of **customer service**, such as: the minimum delivered pressure, in addition to the physical constraints such as the maximum & the minimum water levels in storage tanks to prevent overtopping & emptying respectively.”

In order to optimize the operational performance of the water supply network, at the same time as minimizing the energy costs, it is necessary to predict the consequences of different pump & valve settings on the behavior of the network.

Apart from Linear & Nonlinear Programming, there are other methods & approaches to design, to manage & operate a water supply network to achieve sustainability – e.g., the adoption of **appropriate technology** coupled with effective strategies for operation & maintenance. These strategies must include effective management models, technical support to the householders & industries, sustainable financing mechanisms, & development of reliable **supply chains**. All these measures must ensure the following: system working lifespan; maintenance cycle; continuity of functioning; down time for repairs; water yield & water quality.” – [Wikipedia/water supply network/optimizing the water supply network/operational constraints](#)

1.3.7 Sustainable development

“In an unsustainable system there is insufficient maintenance of the water networks, especially in the major pipe lines in urban areas. The system deteriorates & then needs rehabilitation or renewal.” Fig. **Sustainable development in an urban water network**. “Householders & **sewage treatment** plants can both make the water supply networks more efficient & sustainable. Major improvements in **eco-efficiency** are gained through systematic separation of rainfall & wastewater. Membrane technology can be used for recycling wastewater.

The municipal government can develop a “Municipal Water Reuse System” which is a current approach to manage the rainwater. It applies a **water reuse** scheme for treated wastewater, on a municipal scale, to provide non-potable water for industry, household & municipal uses. This technology consists in separating the **urine** fraction of sanitary wastewater, & collecting it for recycling its **nutrients**. The **feces** & **graywater** fraction is collected, together with organic wastes from the households, using a **gravity sewer system**, continuously flushed with non-potable water. The water is treated **anaerobically** & the **biogas** is used for **energy production**.

The sustainable water supply system is an integrated system including water intake, water utilization, wastewater discharge & treatment & water **environmental protection**. It requires reducing **freshwater** & **groundwater** usage in all sectors of consumption. Developing sustainable water supply systems is a growing trend, because it serves people’s long-term interests. There are several ways to reuse & recycle the water, in order to achieve long-term sustainability, such as:

- Gray water re-use & treatment: **gray water** is a wastewater coming from **baths**, **showers**, sinks & **washbasins**. If this water is treated it can be used as a source of water for uses other than drinking. Depending on the type of gray water & its level of treatment, it can be re-used for **irrigation** & toilet flushing. According to an investigation about the impacts of domestic **grey water** reuse on public health, carried out by the New South Wales Health Center in Australia in the

year 2000, grey water contains less **nitrogen** & fecal pathogenic organisms than **sewage**, & the organic content of grey water decomposes more rapidly.

- Ecological treatment systems use little energy: there are many applications in gray water re-use, such as **reed beds**, soil treatment systems & plant filters. This process is ideal for gray water re-use, because of easier maintenance & higher removal rates of organic matter, **ammonia**, nitrogen & **phosphorus**.

Other possible approaches to scoping models for water supply, applicable to any urban area, include the following:

- **Sustainable drainage system**
- **Borehole** extraction
- Intercluster groundwater flow
- Canal & river extraction
- Aquifer storage
- A more user-friendly indoor water use

The **Dublin Statement on Water & Sustainable Development** is a good example of the new trend to overcome water supply problems. This statement, suggested by advanced economies, has come up with some principles that are of great significance to urban water supply. These are:

1. Fresh water is a finite & vulnerable resource, essential to sustain life, development & the environment.
2. Water development & management should be based on a participatory approach, involving users, planners & policy-makers at all levels.
3. Women play a central part in the provision, management & safeguarding of water. Institutional arrangements should reflect the role of women in water provision & protection.
4. Water has an **economic value** in all its competing uses & should be recognized as an economic good.

From these statements, developed in 1992, several policies have been created to give importance to water & to move urban **water system** management towards sustainable development. The **Water Framework Directive** by the **European Commission** is a good example of what has been created there out of former policies.” – [Wikipedia/water supply network/sustainable development](#)

1.3.8 Future approaches

“There is great need for a more sustainable water supply systems. To achieve sustainability several factors must be tackled at the same time: climate change, rising energy cost, & rising populations. All of these factors provoke change & put pressure on management of available water resources.

An obstacle to transforming conventional water supply systems, is the amount of time needed to achieve the transformation. More specifically, transformation must be implemented by municipal **legislation** bodies, which always need short-term solutions too. Another obstacle to achieving sustainability in water supply systems is the insufficient practical experience with the technologies required, & the missing know-how about the organization & the transition process.

Possible ways to improve this situation is simulating of the network, implementing **pilot projects**, learning from the costs involved & the benefits achieved.” – [Wikipedia/water supply network/future approaches](#)

Part I

Optimal Control

Chapter 2

Lions, 1971. Optimal Control of Systems Governed by PDEs

Introduction

1. “The development of a theory of optimal control (deterministic¹) requires the following initial data:

- (i) a *control* u belonging to some set \mathcal{U}_{ad} (the set of ‘admissible² controls’) which is at our disposition³,
- (ii) for a given control u , the state $y(u)$ of the system which is to be controlled is given by the solution of an equation (*)
 $\Lambda y(u) =$ given function of u , where Λ is an operator⁴ (assumed known) which specifies⁵ the system to be controlled (Λ is the ‘model⁶’ of the system⁷),
- (iii) the observation⁸ $z(u)$ which is a function of $y(u)$ (assumed to be known exactly; we consider only deterministic problems in this book),
- (iv) the “cost function” $J(u)$ (“economic⁹ function”) which is defined in terms of a numerical function $z \rightarrow \Phi(z) \geq 0$ on the “space of observations” by (**) $J(u) = \Phi(z(u))$. It is required to find (problem of the Calculus of Variations¹⁰)
 $\inf J(u), u \in \mathcal{U}_{\text{ad}}$.

The objectives of the theory are

¹**deterministic** [a] (*philosophy*) connected with the belief that people are not free to choose what they are like or how they behave, because these things are decided by their environment & other things over which they have no control.

²**admissible** [a] that can be allowed or accepted according to a set of rules, especially in a court of law, OPPOSITE: **inadmissible**.

³**disposition** [n] **1.** [countable, uncountable] the natural qualities of a person’s character; **2.** [countable] a tendency to behave in a particular way, or to have a particular opinion; **3.** [countable, uncountable] **disposition (of something)** (*specialist*) the way something is placed or arranged; the fact of something being placed somewhere; **4.** [countable, uncountable] (*law*) a formal act of giving property or money to somebody.

⁴**operator** [n] **1.** (often in compounds) a person or company that runs a particular business; **2.** (often in compounds) a person who operates equipment or a machine; **3.** (*mathematics*) a symbol or function which represents an operation in mathematics, e.g., \times , $+$.

⁵**specify** [v] to identify somebody/something clearly & definitely; to state a fact or something that is required clearly & exactly.

⁶**model** [n] **1.** a simple description, especially a mathematical one, of a group of complex systems or processes, used for understanding or explaining how something works; **2.** a way of doing something that others can copy or refer to; **3.** an object that is a copy of something, usually smaller than the original object; **4. model of something** a perfect example of something; **5.** a particular design or type of product; **6.** (in fashion & art) somebody who sits, stands or moves around in order to display clothes or so that somebody else can draw, paint or photograph them; [v] **1.** to create or use a description (especially a mathematical one), a computer program, a diagram or a copy of something, in order to explain or calculate something; **2. model something** to show somebody how to do something, especially how to behave well, SYNONYM: **simulate**; **model something on/after something** [phrasal verb] [usually passive] to make something so that it is like something else; to base something on something else.

⁷**system** [n] **1.** [countable] an organized way of doing something; an organized set of ideas or theories; **2.** [countable] a group of things that work together in a particular way or for a particular purpose; **3.** [countable] a human or animal body, or a part of it, when it is being thought of as the organs & processes that make it function; **4. (the system)** [singular] (*rather informal, usually disapproving*) the rules or people that control a country or an organization, especially when they seem to be unfair because you cannot change them.

⁸**observation** [n] **1.** [uncountable, countable] the act of watching somebody/something carefully for a period of time, especially to learn something; **2.** [uncountable] the ability to notice things, especially important details; **3.** [countable] **observation (about/on something)** a comment, especially based on something you have seen, heard or read, SYNONYM: **remark**.

⁹**economic** [a] **1.** [only before noun] connected with the trade, industry & development of wealth of a country, an area or a society; **2.** producing enough profit to continue; not costing much money, SYNONYM: **profitable**.

¹⁰**variation** [n] **1.** [countable, uncountable] a change or difference, especially in the amount or level of something, usually within particular limits; **2.** [uncountable] (*biology*) the fact of a living thing occurring in more than 1 different color or form; **3.** [countable] **variation (on something)** a thing that is different from other things in the same general group.

- (i) to obtain necessary (or possibly necessary & sufficient¹¹) conditions for u to be an extremum (or minimum),
- (ii) to study the structure & properties of the equations expressing these conditions (where the ‘model’ Λ naturally intervenes¹²),
- (iii) to obtain constructive¹³ algorithms¹⁴ amenable¹⁵ to numerical computations for the approximation of a (the) control $u \in \mathcal{U}_{\text{ad}}$ which determines the inf (such a control is termed an “optimal¹⁶ control”).

2. Clearly the development of such a theory depends on the model Λ in a fundamental manner¹⁷. The theory described in the works of Pontryagin–Boltyanskii–Gamkrelidze–Mischenko [1] & Hestenes [1] is concerned with the study of points (i) & (ii) of **1** in the case where Λ is a family of ordinary differential (or with delay¹⁸ or integro-differential) operators.

In a variety of applications, due to the complexity¹⁹ of the system to be controlled, it is often advantageous²⁰ to discard²¹ the above-mentioned mathematical model in favor of a model described by a family of partial differential operators (e.g., cf. Butkovskii [1], Wang [1] & the bibliography²² of these works). It is this situation that we propose²³ to investigate in this book^{24 25}. We thus consider systems whose state $y(u)$ is given by the solution of a *PDE* to which we must add appropriate *boundary conditions*²⁶ & in the case of evolution²⁷ equations *initial conditions*.

3. It is clear that unless we wish to restrict²⁸ ourselves to results which are purely²⁹ formal³⁰, the minimization³¹ of $(**)$ presupposes³² that the BVP $(*)$ is formulated³³ & solved in a precise³⁴ mathematical setting. The results that are needed in this direction are proved in this book for the case where Λ is an operator which is elliptic or parabolic or hyperbolic or

¹¹**sufficient** [a] enough for a particular purpose; as much as you need. In logic, a **sufficient condition** of a statement is a condition that, if true, makes the statement true. It is often combined with a **necessary condition**, which must be true in order for the statement to be true, OPPOSITE: **insufficient**.

¹²**intervene** [v] **1.** [intransitive] to become involved in a situation in order to improve it or stop it from getting worse; **2.** [intransitive] to happen in the time between events; **3.** [intransitive] to exist or be found in the space between things; **4.** [intransitive] to happen in a way that delays something or prevents it from happening.

¹³**constructive** [a] having a useful & helpful effect rather than being negative or with no purpose.

¹⁴**algorithm** [n] a process or set of rules to be followed when solving a particular problem, especially by a computer.

¹⁵**amenable** [a] **amenable to (doing) something** that you can treat in a particular way.

¹⁶**optimal** [a] [usually before noun] (also **optimum**) the best possible; producing the best possible results, SYNONYM: **ideal**.

¹⁷**manner** [n] **1.** [singular] the way that something is done or happens; **2.** [singular] the way that somebody behaves towards other people; **3.** (manners) [plural] behavior that is considered to be polite in a particular society or culture; **4.** (manners (of somebody/something)) [plural] the habits & customs of a particular group of people; **all manner of somebody/something** [idiom] many different types of people or things; **in the manner of somebody/something** [idiom] in a style that is typical of somebody/something.

¹⁸**delay** [n] **1.** [countable] a period or time by which something is slow or late; the period of time between 2 things happening; **2.** [uncountable] a situation in which something does not happen when it should; [v] **1.** [transitive] to not do something until a later time; **2.** [transitive, usually passive] if an event is delayed, it happens at a later time than is normal or expected.

¹⁹**complexity** [n] **1.** [uncountable] the state of being formed of many parts; the state of being difficult to understand; **2.** (complexities) [plural] complexity of something the features of a problem or situation that are difficult to understand.

²⁰**advantageous** [a] good or helpful to somebody in a particular situation, SYNONYM: **beneficial**, OPPOSITE: **disadvantageous**.

²¹**discard** [v] to get rid of something that you no longer want or need.

²²**bibliography** [n] (plural **bibliographies**) the list of books, etc. that have been used by somebody writing an article, essay, etc.; a list of books or articles about a particular subject or by a particular author.

²³**propose** [v] **1.** to suggest a plan or an idea for people to consider & decide on; **2.** to suggest an explanation of something for people to consider.

²⁴We hasten to add that only very partial results have been obtained in a number of directions.

²⁵**hasten** [v] **1.** [intransitive] **hasten to do something** to say or do something without delay; **2.** [transitive] **hasten something** (formal) to make something happen sooner or more quickly; **3.** [intransitive] + **adv./prep.** (literary) to go or move somewhere quickly, SYNONYM: **hurry**.

²⁶Control may be exercised through the boundary condition, which in fact is the situation generally encountered in practice.

²⁷**evolution** [n] [uncountable] **1.** (biology) the gradual development of living things over many years as they adapt to changes in their environment; **2.** the gradual development of something.

²⁸**restrict** [v] **1.** to limit or control the size, amount or range of something; **2.** to limit something to a particular time, place or group; **3.** **restrict somebody (from something/from doing something)** to not allow somebody to do something or to go somewhere; **4.** **restrict yourself/somebody/something (to something/to doing something)** to allow yourself or somebody to have, do or consider only a limited amount of something or to do only a particular kind of activity; **5.** **restrict something/somebody** to stop something/somebody from moving freely, SYNONYM: **impede**.

²⁹**purely** [adv] only; completely.

³⁰**formal** [a] **1.** following strict rules of how to do something; suitable for an official occasion, OPPOSITE: **informal**; **2.** (of speech or writing) suitable for official or serious situations, OPPOSITE: **informal**; **3.** (of education or training) received in a school, college or university rather than gained just through practical experience, OPPOSITE: **informal**; **4.** concerned with the form or structure of something rather than its content; **5.** concerned only with following rules.

³¹**minimization** [n] (British English also **minimisation**) [uncountable, countable, usually singular] the act of reducing something, especially something bad, to the lowest possible level.

³²**presuppose** [v] **1.** to accept that something is true & argue a case or take action on that basis, before it has been proved to be true, SYNONYM: **assume**, **presume**; **2.** **presuppose something** to require something or accept something as needing to exist.

³³**formulate** [v] **1.** **formulate something** to create or prepare something carefully, giving particular attention to the details; **2.** **formulate something** to express your ideas in carefully chosen words.

³⁴**precise** [a] **1.** clear & accurate, SYNONYM: **exact**; **2.** [only before noun] used to emphasize that something happens at a particular time or in a particular way; **to be (more) precise** [idiom] used to show that you are giving more detailed & accurate information about something you have just mentioned.

well-posed in the sense of Petrovsky. In order not to overburden³⁵ this work we have restricted ourselves to comparatively³⁶ simple examples. However the techniques we have used are quite general & hence more general problems can be solved using the same techniques. We refer the reader to Lions and Magenes, 1972, Chap. 6 where extensions to more general cases can be found.

Once we have in our possession³⁷ a good theory for the solution of (*), it remains to obtain & analyze the necessary (or necessary & sufficient in favorable cases) conditions for (**) to have a minimum. In this manner we are led to a number of BVPs which appear to be novel³⁸ in character. These problems however have a striking analogy³⁹ with “multiphase” & “unilateral⁴⁰” problems of mechanics⁴¹ – in particular in plasticity⁴².

4. We now give a brief analysis of the contents of the various chapters.

In Chap. 1 we study amongst other things, the minimization of positive definite⁴³ or semi-definite quadratic⁴⁴ forms defined on a closed, convex subset of a Hilbert space. Applications to unilateral problems are given. These unilateral problems are prototypes⁴⁵ of problems which we encounter in the sequel⁴⁶.

In Chap. 2 we study the optimal control of systems governed by elliptic equations. While in Chaps. 3–4 we examine the parabolic & hyperbolic or well-posed in the sense of Petrovsky cases. In each of these chapters we 1st consider the case of a linear system with a quadratic cost function & study the “unilateral problem” which this leads to. In Chaps. 3–4 we study in detail the “feedback problem” & the related integro-differential equation of Riccati type. We then study existence theorems for simple nonlinear systems (given the present state of the art of nonlinear PDEs, a general theory in this sense appears to be outside our reach for the moment) & 1st order necessary conditions.

Finally, in Chap. 5, we present various procedures of regularization, approximation & penalization. These procedures may be utilized in the numerical solution of optimal control problems which we have studied.

The chapter headings are the following:

- Chap. 1. Minimization of functions & unilateral BVPs.
- Chap. 2. Control of systems governed by elliptic PDEs.
- Chap. 3. Control of systems governed by parabolic PDEs.
- Chap. 4. Control of systems governed by hyperbolic PDEs or by equations well-posed in the sense of Petrovsky.
- Chap. 5. Regularization, approximation & penalization. Bibliography.

Each chapter begins with a detailed plan indicating the scope of the chapter & closes with bibliographic notes & indications on problems which are unsolved (of which there are many) or on aspects which have not been considered in this book. The whole subject is clearly in a process of evolution.

The contents of the book have developed from a course given at the Faculty of Sciences, University of Paris since the academic year 1965–1966. An abbreviated version of the book was presented in a summer course at the University of California, Los Angeles in Aug 1967 (the course was organized by A. V. Balakrishnan).” – Paris, Nov 1969, J. L. Lions, Lions, 1971, pp. 1–3

³⁵**overburden** [v] [usually passive] **overburden somebody/something (with something)** to give somebody/something more work, worry, etc. than they can deal with.

³⁶**comparatively** [adv] **1.** when measured or judged by how similar or different something is to something else, SYNONYM: **else**; **2.** connected with studying things to discover how they are similar or different.

³⁷**possession** [n] **1.** [uncountable] the state of having or owning something; **2.** [countable, usually plural] something that you own or have with you at a particular time; **3.** [countable] a country that is controlled or governed by another country; **4.** [uncountable] the state of having illegal drugs or weapons with you at a particular time; **5.** [uncountable] the situation when somebody’s mind is believed to be controlled by an evil spirit.

³⁸**novel** [n] a story long enough to fill a complete book, in which the characters & events are usually imaginary; [a] different from anything known before; new & interesting.

³⁹**analogy** [n] (plural **analogies**) [countable, uncountable] a comparison of 1 thing with another thing that has similar features, usually in order to explain it; a feature that is similar.

⁴⁰**unilateral** [a] **1.** [usually before noun] done by 1 member of a group or organization without the agreement of the other members; **2.** (*medical*) involving only 1 side of an organ or the body.

⁴¹**mechanics** [n] **1.** [uncountable] the science of movement & force; **2.** [plural] **mechanics of something** the way something works or is done.

⁴²**plasticity** [n] [uncountable] **1.** (*specialist*) the quality of being easily made into different shapes; **2.** (*biology*) the ability of a living thing to adapt to changes in its environment or differences between its various habitats.

⁴³**definite** [a] **1.** clearly stated or decided; sure or certain; **2.** clearly true or real; having a clear meaning; **3.** having an exact online or form that can be recognized easily.

⁴⁴**quadratic** [a] (*mathematics*) involving an unknown quantity that is multiplied by itself once only.

⁴⁵**prototype** [n] **1.** the 1st design of something from which other forms are copied or developed; **2.** **prototype (of something)** the 1st, original or typical form of something.

⁴⁶**sequel** [n] **1.** **sequel (to something)** a book, film, play, etc. that continues the story of an earlier one; **2.** [usually singular] **sequel (to something)** something that happens after an earlier event or as a result of an earlier event.

Principal Notations

“ $\mathbf{x} = (x_1, \dots, x_n)$ denotes the *space* variable; \mathbf{x} ranges in an open set $\Omega \subset \mathbb{R}^n$ with boundary Γ . t denotes time; in general $t \in (0, T)$, $T < \infty$. We set $Q := (0, T) \times \Omega$, $\Sigma = (0, T) \times \Gamma$. The *controls* (or *commands*⁴⁷) are, in general, denoted by u, v, w, \dots ; they are generally taken to be in a space \mathcal{U} (quite generally a Hilbert space on \mathbb{R}); \mathcal{U}_{ad} (= set of *admissible* controls) is a *closed, convex subset* of \mathcal{U} .

The *state* of the system is denoted by $y(v)$; in the *elliptic* case (Chap. 2) $y(v)$ is a function of $\mathbf{x} \in \Omega$, $y(\mathbf{x}, v)$; in the *evolution* case (Chaps. 3–4) $y(v)$ is a function of $\mathbf{x} \in \Omega$ & $t \in (0, T)$: $y(t, \mathbf{x}; v)$. The *observation* is denoted by $z(v) = Cy(v)$ (we do not study the case where there is noise present). The *cost function* (or *criterion*, or *economic function*) is denoted by $J(v)$. The $u \in \mathcal{U}_{\text{ad}}$ s.t. $J(u) \leq J(v)$, $\forall v \in \mathcal{U}_{\text{ad}}$ are the optimal controls. $p(v)$ denotes the *adjoint state*.

MAIN FUNCTION SPACES USED.

- $C^k(\Omega)$ = space of k -times continuously differentiable functions on $\overline{\Omega}$, k integer ≥ 0 .⁴⁸
- $\mathcal{D}(\Omega)$ = space of infinitely differentiable functions in Ω , with *compact support* in Ω , endowed with the inductive limit topology of L. Schwartz [1].
- $\mathcal{D}'(\Omega)$ = dual space of $\mathcal{D}(\Omega)$ = space of distributions on Ω .⁴⁹
If X is a Banach space, $\mathcal{D}'((0, T); X)$ denotes the space of distributions on $(0, T)$ with values in X (cf. Schwartz [3] & brief recapitulation⁵⁰ in Chap. 3, Sect. 1.1).
- $L^2(\Omega)$ = space (equivalence class) of functions square integrable on Ω .
- $H^m(\Omega)$ = (Sobolev [1] space of order m) = space of functions φ s.t. $\varphi \in L^2(\Omega)$, $\partial_{x_i} \varphi \in L^2(\Omega), \dots, D^\alpha \varphi \in L^2(\Omega)$, $\forall \alpha$, $|\alpha| \leq m$, $\alpha = (\alpha_1, \dots, \alpha_m)$, $|\alpha| = \sum_{i=1}^n \alpha_i$.
- $H_0^m(\Omega) = \{\varphi | \varphi \in H^m(\Omega), D^\alpha \varphi = 0 \text{ on } \Gamma, |\alpha| \leq m-1\}$.
- $H^s(\Omega)$ = *Fractional Sobolev space* of order s on Ω , = space of restrictions on Ω of functions of $H^s(\mathbb{R}^n)$ defined (by Fourier Transforms) in Chap. 1, (3.12).
- $L^2(S; E)$ = space (equivalence class) of functions defined on S (locally compact measure space with measure $\mu \geq 0$) with values in a Hilbert space E & s.t. $\int_S \|f(t)\|_E^2 d\mu(t) < \infty$. We use primarily $L^2(0, T; E)$, $d\mu(t) = dt$.
- $L^\infty(S; E)$ = space (equivalence class) of functions f defined on S with values in E & essentially bounded: $\|f(t)\|_E \leq \|f\|_{L^\infty(S, E)} < \infty$, a.e.” – Lions, 1971, pp. 4–5

2.1 Minimization of Functions & Unilateral BVPs

2.1.1 Minimization of Coercive Forms

2.1.1.1 Notation

“Let \mathcal{U} be a real Hilbert space. The elements of \mathcal{U} will be denoted by u, v, w, \dots . In the applications we have in mind in the following chapters. \mathcal{U} will be the space of controls.

In this chapter $\|\cdot\|$ will denote the norm on \mathcal{U} ; in general, if there is possible ambiguity, the norm in the space X will be denoted by $\|u\|_X$. For the moment, we shall assume that the following data is given:

- (i) a continuous bilinear form on \mathcal{U} , which is symmetric, $u, v \rightarrow \pi(u, v)$, $\pi(u, v) = \pi(v, u)$, $\forall u, v \in \mathcal{U}$,
- (ii) a continuous linear form on \mathcal{U} , $v \rightarrow L(v)$,

⁴⁷**command** [n] **1.** [uncountable] **command (of somebody/something)** control & authority over a situation or a group of people; **2.** [singular, uncountable] **command (of something)** your knowledge of something; your ability to do or use something, especially a language; **3.** [countable] an order given to a person or an animal; **4.** [countable] an instruction causing a computer to perform a function; **at your command** [idiom] if you have a skill or an amount of something at your command, you are able to use it well & completely; [v] **1.** [transitive] (of somebody in a position of authority) to tell somebody to do something, SYNONYM: **order**; **2.** [transitive, intransitive] **command (somebody/something)** to be in charge of a group of people in the army, navy, etc.; **3.** [transitive, no passive] (not used in the progressive tenses) **command something** to deserve & get something because of the special qualities you have; **4.** [transitive, no passive] (not used in the progressive tenses) **command something** to be in a strong enough position to have or get something; **5.** [transitive, no passive] (not used in the progressive tenses) **command something** to have something available for use; **6.** [transitive, no passive] (not used in the progressive tenses) **command something** to be in a position from where you can see or control something.

⁴⁸Clearly we have analogous notation for Q, Γ, Σ . All functions considered are real-valued.

⁴⁹In general, X' denotes the dual of X .

⁵⁰**recapitulation** [n] [countable, uncountable] (*formal*) (also **recap**) the act of repeating or giving a summary of what has already been said, decided, etc.

(iii) a closed, convex set \mathcal{U}_{ad} in \mathcal{U} .

In the applications considered in the sequel \mathcal{U}_{ad} will be the set of admissible controls. The quadratic functional (1.1) $J(v) = \pi(v, v) - 2L(v)$ is required to be minimized over the set \mathcal{U}_{ad} . – Lions, 1971, p. 6

2.1.1.2 The Case when π is Coercive

“ π is said to be *coercive* on \mathcal{U} if (1.2)

$$\pi(v, v) \geq c\|v\|^2, \quad \forall v \in \mathcal{U}, \quad c > 0.$$

We then have

Theorem 2.1. *Let $\pi(u, v)$ be a continuous symmetric bilinear form on \mathcal{U} satisfies (1.2). Then there exists a unique element $u \in \mathcal{U}_{\text{ad}}$ s.t. $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v)$.*

Proof. See Lions, 1971, pp. 7– □

Example 2.1. *Let us consider $\pi(u, v) = (u, v) = \text{scalar product in } \mathcal{U}$, $L(v) = (g, v)_{\mathcal{U}}$, where g is a given element in \mathcal{U} . Then, $J(v) = \|g - v\|_{\mathcal{U}}^2 - \|g\|_{\mathcal{U}}^2$ & the unique element $u \in \mathcal{U}_{\text{ad}}$ s.t. $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v)$ is characterized by $\|g - u\|_{\mathcal{U}} \leq \|g - v\|_{\mathcal{U}}$, $\forall v \in \mathcal{U}_{\text{ad}}$; u is thus the projection of g on \mathcal{U}_{ad} .*

Analysis of the Proof of Lions, 1971, Theorem 1.1. An analysis of the way the assumptions of Theorem 1.1 come into play in the proof of the theorem suggests the following remarks:

Remark 2.1. *Theorem 1.1 is true if we assume that the bilinear form $\pi(u, v)$ is defined on $\mathcal{U}_{\text{ad}} \times \mathcal{U}_{\text{ad}}$ & satisfies (1.2), $\forall v \in \mathcal{U}_{\text{ad}}$.*

The fact that the function $v \rightarrow J(v)$ is a quadratic form does not enter in any essential way in the proof of Theorem 1.1.

Remark 2.2. *Let $v \rightarrow J(v)$ be a convex function from $\mathcal{U}_{\text{ad}} \rightarrow \mathbb{R}$, s.t. (1.10)–(1.11)*

$$\begin{aligned} J(v) &\rightarrow +\infty \text{ as } \|v\| \rightarrow +\infty, \quad v \in \mathcal{U}_{\text{ad}}, \\ v &\rightarrow J(v) \text{ is strongly l.s.c.} \end{aligned}$$

Then there exists $u \in \mathcal{U}_{\text{ad}}$ s.t. (1.12) $J(u) = \inf_v J(v)$.

This remarks also applies to functions $v \rightarrow J(v)$ defined on $\mathcal{U}_{\text{ad}} \subset \mathcal{U}$, where \mathcal{U} is, e.g., a reflexive Banach space.

With hypotheses (1.10), (1.11) only, we do not necessarily have uniqueness; clearly we have uniqueness if we assume that the function $v \rightarrow J(v)$ is strictly convex.

Remark 2.3. *Assumption (1.10) is necessary to prove that every minimizing sequence is bounded (cf. 1.5). If \mathcal{U}_{ad} is also bounded then we may dispense⁵¹ with Assumption (1.10).” – Lions, 1971, pp. 6–8*

2.1.1.3 Characterization of the Minimizing Element. Variational Inequalities

Theorem 2.2. *Let the assumptions of Theorem 1.1 remain valid. The minimizing element u of \mathcal{U}_{ad} is characterized by (1.13) $\pi(u, v - u) \geq L(v - u)$, $\forall v \in \mathcal{U}_{\text{ad}}$.*

Proof. See Lions, 1971, p. 9. □

“Inequalities of the type given by (1.13) are termed “*variational inequalities*”.

Remark 2.4. *Let \mathcal{U} be a Hilbert space over \mathbb{C} (instead of \mathbb{R}) & let $\pi(u, v)$ be a sesquilinear hermitian form, i.e., $\pi(u, v) = \overline{\pi(v, u)}$, $\forall u, v \in \mathcal{U}$. Assuming that (1.2) is satisfied, Theorem 1.1 remains valid without any change. Replacing (1.13) by (1.18)*

$$\operatorname{Re} \pi(u, v - uu) \geq \operatorname{Re} L(v - u), \quad \forall v \in \mathcal{U}_{\text{ad}},$$

Theorem 1.2 remains valid.

⁵¹**dispense** [v] **1.** to provide something, usually something that is intended to help people; **2.** **dispense something** to prepare medicine & give it to people; **3.** **dispense something** (of a machine) to provide money, food, drink, etc.; **dispense with somebody/something** [phrasal verb] to not use or stop using somebody/something; to get rid of somebody/something.

Remark 2.5 (The Case $\mathcal{U}_{\text{ad}} = \mathcal{U}$ ⁵²). If $\mathcal{U}_{\text{ad}} = \mathcal{U}$, in (1.13) we may take $v = u \pm \varphi$, where φ is any element of \mathcal{U} & (1.13) becomes equivalent to (1.19)

$$\pi(u, \varphi) = L(\varphi), \quad \forall \varphi \in \mathcal{U}.$$

This is the Euler equation of the problem.

Remark 2.6 (The Case $\mathcal{U}_{\text{ad}} = \text{Cone}$). Let us suppose (1.20) $\mathcal{U}_{\text{ad}} = \text{closed convex cone with vertex at the origin}$. Then (1.13) is equivalent to (1.21)

$$\begin{cases} \pi(u, v) \geq L(v), \quad \forall v \in \mathcal{U}_{\text{ad}}, \\ \pi(u, u) = L(u). \end{cases}$$

In fact, in (1.13) we may replace v by $v + u$ which gives the 1st inequality in (1.21): putting $r = 0$ in (1.13), we obtain $\pi(u, v) \leq L(u)$ & hence the inequality $\pi(u, v) \leq L(u)$. Conversely, it is obvious that (1.21) implies (1.13).

Remark 2.7 (The case of a functional $v \rightarrow J(v)$ which is not necessarily quadratic). Suppose that the function (or functional) $v \rightarrow J(v)$ is differentiable⁵³ The proof of Theorem 1.2 is also applicable to

Theorem 2.3. Assume that the function $v \rightarrow J(v)$ is strictly convex, differentiable & satisfies (1.10) (the last hypothesis may be omitted if \mathcal{U}_{ad} is bounded). Then the unique element u in \mathcal{U}_{ad} satisfying $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v)$ is characterized by (1.22) $J'(u) \cdot (v - u) \geq 0, \forall v \in \mathcal{U}_{\text{ad}}$.

2.1.1.4 Alternative Form of Variational Inequalities

The following results are very useful in a technical sense:

Theorem 2.4. Let all the hypotheses of Theorem 1.3 be satisfied. Then the characterization (1.22) is equivalent to: (1.23) $J'(v) \cdot (v - u) \geq 0, \forall u \in \mathcal{U}_{\text{ad}}$.

Proof. See Lions, 1971, p. 11. □

1st assume that the following result, which is important in its own right, is true:

Theorem 2.5. Assume that the function $v \rightarrow J(v) : \mathcal{U} \rightarrow \mathbb{R}$ is convex & differentiable. Then the derivative $v \rightarrow J'(v) : \mathcal{U} \rightarrow \mathcal{U}'$ is monotone, i.e., (1.24) $(J'(v) - J'(w)) \cdot (v - w) \geq 0, \forall v, w \in \mathcal{U}$.

Proof. See Lions, 1971, p. 11. □

Remark 2.8. Summarizing: under the hypotheses of Theorem 1.3, we have 3 equivalent formulations of the problem:

- (i) the definition of the problem (when the minimum obtains): $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v), u \in \mathcal{U}_{\text{ad}},$
- (ii) $J'(u) \cdot (v - u) \geq 0, \forall v \in \mathcal{U}_{\text{ad}}, u \in \mathcal{U}_{\text{ad}},$
- (iii) $J'(v) \cdot (v - u) \geq 0, \forall u \in \mathcal{U}_{\text{ad}}, v \in \mathcal{U}_{\text{ad}}.$

2.1.1.5 Function f being the Sum of a Differentiable & Non-Differentiable Function

If we assume that the function $v \rightarrow J(v)$ is coercive, lower semi-continuous in the weak topology & strictly convex, then there exists a $u \in \mathcal{U}_{\text{ad}}$ s.t. $J(u) \leq J(v), \forall v \in \mathcal{U}_{\text{ad}}$. In this case it is clear that we cannot apply criteria (ii) & (iii) of Remark 1.8. However, these criteria are still applicable to the differentiable part of J . More precisely, we have the following result:

Theorem 2.6. Consider the function (1.26) $J(v) = J_1(v) + J_2(v)$ where we assume that the functions $J_i(v), i = 1, 2$, are continuous, convex, & lower semi-continuous in the weak topology. Further let $J(v) \rightarrow +\infty$ as $\|v\| \rightarrow +\infty, v \in \mathcal{U}_{\text{ad}}$. We assume that the function $v \rightarrow J_1(v)$ is differentiable, but J_2 is not necessarily differentiable. Finally assume that J is strictly convex. Then the unique element $u \in \mathcal{U}_{\text{ad}}$ s.t. $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v)$ is characterized by (1.27)

$$J'_1(u) \cdot (v - u) + J_2(v) - J_2(u) \geq 0, \quad \forall v \in \mathcal{U}_{\text{ad}}.$$

Proof. See Lions, 1971, pp. 12–13. □

Remark 2.9. Putting $J_2 = 0$, it is clear that Theorem 1.6 contains Theorems 1.3 & 1.4.

⁵²In control problems this corresponds to the case where there are no constraints.

⁵³Cf. J. Dieudonné [1], Chap. 8, Sect. 1.

Remark 2.10. Suppose that J is of the form (1.28) $J(v) = J_0(v) + J_1(v) + J_2(v)$, where the J_i 's satisfy the same hypotheses as in Theorem 1.6 & the functions J_0 & J_1 are differentiable. Then the unique element u s.t. $J(u) = \inf_{v \in \mathcal{U}_{\text{ad}}} J(v)$ is characterized by 1 of the following equivalent conditions: (1.29)–(1.30)

$$\begin{aligned} J'_0(u) \cdot (v - u) + J'_1(u) \cdot (v - u) + J_2(v) - J_2(u) &\geq 0, \quad \forall v \in \mathcal{U}_{\text{ad}}, \\ J'_0(u) \cdot (v - u) + J'_1(v) \cdot (v - u) + J_2(v) - J_2(u) &\geq 0, \quad \forall v \in \mathcal{U}_{\text{ad}}. \end{aligned}$$

Remark 2.11. In case we only have existence of the minimizing element u of \mathcal{U}_{ad} but not necessarily uniqueness, any 1 of the variational inequalities we have obtained characterizes the set of elements in \mathcal{U}_{ad} determining the minimum.

In applications to control problems an element $u \in \mathcal{U}_{\text{ad}}$ which determines the minimum is termed “optimal control”.

Remark 2.12. All the preceding results, without any change in their proofs are true when \mathcal{U} is a reflexive Banach space.

2.1.1.6 The Convexity Hypothesis on \mathcal{U}_{ad}

So far we have assumed that \mathcal{U}_{ad} is *convex*. The following development shows how we may obtain a (simple) necessary condition for extremality in case \mathcal{U}_{ad} is assumed to be only *closed* in \mathcal{U} .

Definition 2.1. Let \mathcal{U}_{ad} be closed & let $u \in \mathcal{U}_{\text{ad}}$. Define, (1.31)

$$\mathcal{C}(\mathcal{U}_{\text{ad}}; u) := \{w | w \in \mathcal{U}; \text{ there exists } u_n \in \mathcal{U}_{\text{ad}} \text{ \& } \lambda_n > 0, \text{ s.t. } u_n \rightarrow u \text{ \& } \lambda_n(u_n - u) \rightarrow w \text{ in } \mathcal{U}\}.$$

It may be easily verified that (1.32)

$$\mathcal{C}(\mathcal{U}_{\text{ad}}; u) \text{ is a closed cone with its vertex at } \{0\}.$$

We then have

Theorem 2.7. Let $v \rightarrow J(v)$ be a function which is differentiable & let u be an element (assumed to exist) of \mathcal{U}_{ad} s.t. $J(u) \leq J(v)$, $\forall v \in \mathcal{U}_{\text{ad}}$. Then (1.33) $J'(v) \cdot w \geq 0$, $\forall w \in \mathcal{C}(\mathcal{U}_{\text{ad}}; u)$.

” – Lions, 1971, pp. 9–

2.1.2 A Direct Solution of Certain Variational Inequalities

2.1.3 Examples

2.1.4 A Comparison Theorem

2.1.5 Non Coercive Forms

2.2 Control of Systems Governed by Elliptic PDEs

2.2.1 Control of Elliptic Variational Problems

2.2.2 1st Applications

2.2.3 A Family of Examples with $N = 0$ & \mathcal{U}_{ad} Arbitrary

2.2.4 Observation on the Boundary

2.2.5 Control & Observation on the Boundary. Case of the Dirichlet Problem

2.2.6 Constraints on the State

2.2.7 Existence Results for Optimal Controls

2.2.8 1st Order Necessary Conditions

2.3 Control of Systems Governed by Parabolic PDEs

2.4 Control of Systems Governed by Hyperbolic Equations or by Equations which are well Posed in the Petrowsky Sense

2.5 Regularization, Approximation & Penalization

2.5.1 Regularization

2.5.1.1 Parabolic Regularization

2.5.1.2 Approximation in Terms of Systems of Cauchy–Kowaleska Type

2.5.1.3 Penalization

Chapter 3

Tröltzsch, 2010. Optimal Control for PDEs

Preface

“The sections dealing with gradient methods were shortened in order to make space for *primal-dual active set strategies*; the exposition of the latter now leads to the systems of linear equations to be solved.” – Tröltzsch, 2010, Preface to the English edition, p. xi

“The mathematical optimization of process governed by PDEs has seen considerable progress in the past decade. Ever faster computational facilities & newly developed numerical techniques have opened the door to important practical applications in fields e.g. fluid flow, microelectronics¹, crystal² growth, vascular³ surgery⁴, & cardiac⁵ medicine, to name just a few. As a consequence, the communities of numerical analysts & optimizers have taken a growing interest in applying their methods to optimal control problems involving PDEs; at the same time, the demand from students for this expertise has increased, & there is a growing need for textbooks that provide an introduction to the fundamental concepts of the corresponding mathematical theory.” [...] “... the comprehensive text by J.-L. Lions Lions, 1971 covers much of the theory of linear equations & convex cost functionals.”

Tröltzsch, 2010 focuses “on basic concepts & notions e.g.:

- Existence theory for linear & semilinear PDEs
- Existence of optimal controls
- Necessary optimality conditions & adjoint equations
- 2nd-order sufficient optimality conditions
- Foundation of numerical methods

In this connection, we will always impose constraints on the control functions, & sometimes also on the state of the system under study. In order to keep the exposition to a reasonable length, we will not address further important subjects such as *controllability*, *Riccati equations*, *discretization*, *error estimates*, & *Hamilton–Jacobi–Bellman theory*.

The 1st part of the textbook deals with convex problems involving quadratic cost functionals & linear elliptic or parabolic equations. While these results are rather standard & have been treated comprehensively in Lions, 1971, they are well suited to facilitating the transition to problems involving semilinear equations. In order to make the theory more accessible to readers having only minor knowledge of these fields, some basic notions from functional analysis & the theory of linear elliptic & parabolic PDEs will also be provided.

The focus of the exposition is on nonconvex problems involving semilinear equations. Their treatment requires new techniques from analysis, optimization, & numerical analysis, which to a large extent can presently be found only in original papers. In particular, fundamental results due to E. Casas & J.-P. Raymond concerning the boundedness & continuity of solutions to semilinear equations will be needed.

¹**microelectronics** [n] [uncountable] the design, production & use of very small electronic circuits.

²**crystal** [n] **1.** [countable] a small piece of a substance with many even sides, that is formed naturally when the substance becomes solid; in chemistry, a **crystal** is any solid that has its atoms, ions or molecules arranged in an ordered, symmetrical way; **2.** [uncountable] a clear mineral, e.g. quartz, used in making decorative objects.

³**vascular** [a] [usually before noun] (*medical*) connected with or containing veins.

⁴**surgery** [n] **1.** [uncountable, countable] medical treatment of injuries or diseases that involves cutting open a person’s body, sewing up wounds, etc.; **2.** [countable] (*British English*) a place where a doctor sees patients; **3.** [countable] (*British English*) a time during which a doctor, an MP or another professional person is available to see people.

⁵**cardiac** [a] [only before noun] (*medical*) connected with the heart or heart disease; if somebody has a **cardiac arrest**, their heart suddenly stops temporarily or permanently.

This textbook is mainly devoted to the analysis of the problems, although numerical techniques will also be addressed. Numerical methods could easily fill another book. Our exposition is confined to brief introductions to the basis ideas, in order to give the reader an impression of how the theory can be realized numerically. Much attention will be paid to revealing hidden mathematical difficulties that, as experiences shows, are likely to be overlooked.” – Tröltzsch, 2010, Preface to the German edition, pp. xiii–xiv

3.1 Introduction & Examples

3.1.1 What is Optimal Control?

“The mathematical theory of optimal control has in the past few decades rapidly developed into an important & separate field of applied mathematics. 1 area of application of this theory lies in aviation⁶ & space technology: aspects of optimization come into play whenever the motion of an aircraft or a space vessel⁷ (which can be modeled by ODEs) has to follow a trajectory⁸ that is “optimal” in a sense to be specified.” – Tröltzsch, 2010, Sect. 1.1: *What is optimal control?*, p. 1

All the essential features of an *optimal control problem*:

- a *cost functional* to be minimized,
- an IVP for an ODE in order to determine the *state* y ,
- a *control function* u , &
- various constraints that have to be obeyed.

“The control u may be freely chosen within the given constraints, while the state is uniquely determined by the differential equation & the initial conditions. We have to choose u in such a way that the cost function is minimized. Such controls are called *optimal*.” [...] “The optimal control of ODEs is of interest not only for aviation & space technology. In fact, it is also important in fields e.g. robotics⁹, movement sequences in sports, & the control of chemical processes & power plants, to name just a few of the various applications. In many cases, however, the processes to be optimized can no longer be adequately modeled by ODEs; instead, PDEs have to be employed for their description. E.g., heat conduction¹⁰, diffusion¹¹, electromagnetic¹² waves, fluid flows, freezing processes, & many other physical phenomenon¹³ can be modeled by PDEs.

In these fields, there are numerous interesting problems in which a given cost functional has to be minimized subject to a differential equation & certain constraints being satisfied. The difference from the above problem “merely” consists of the fact that a PDE has to be dealt with in place of an ordinary one.” – Tröltzsch, 2010, pp. 2–3

Tröltzsch, 2010 discusses, “through examples in the form of mathematically simplified case studies, the optimal control of heating processes, 2-phase problems, & fluid flows”. Tröltzsch, 2010 focuses “on linear & semilinear elliptic & parabolic PDEs, since a satisfactory regularity theory is available for the solutions to such equations. This is not the case for hyperbolic equations. Also, the treatment of quasilinear PDEs is considerably more difficult, & the theory of their optimal control is still an open field in many respects.” [...] “... the Hilbert space setting suffices as a functional analytic framework in the case of linear-quadratic theory.” – Tröltzsch, 2010, p. 3

3.1.2 Examples of Convex Problems

3.1.2.1 Optimal boundary heating

See Tröltzsch, 2010, Subsect. 1.2.1, pp. 3–5.

Example 3.1 (Optimal boundary heating). *Consider a body heated or cooled which occupies the spatial domain $\Omega \subset \mathbb{R}^3$. Apply to its boundary Γ a heat source u (the control), which is constant in time but depends on the location \mathbf{x} on the boundary, i.e., $u = u(\mathbf{x})$. Aim: choose the control in such a way that the corresponding temperature distribution $y = y(\mathbf{x})$ in Ω (the state) is the best possible approximation to a desired stationary temperature distribution $y_\Omega = y_\Omega(\mathbf{x})$:*

$$\min J(y, u) := \frac{1}{2} \int_{\Omega} |y(\mathbf{x}) - y_\Omega(\mathbf{x})|^2 d\mathbf{x} + \frac{\lambda}{2} \int_{\Gamma} |u(\mathbf{x})|^2 ds(\mathbf{x}),$$

⁶**aviation** [n] [uncountable] the activity of designing, building & flying aircraft.

⁷**vessel** [n] **1.** a tube that carries blood through the body of a person or an animal, or liquid through the parts of a plant; **2.** (formal) a large ship or boat; **3.** (formal) a container used for holding liquids, e.g. a bowl or cup.

⁸**trajectory** [n] (plural **trajectories**) (specialist) **1.** the curved part of something that has been fired, hit or thrown into the air; **2.** the way in which a person, an event or a process develops over a period of time, often leading to a particular result.

⁹**robotics** [n] [uncountable] the science of designing & operating robots.

¹⁰**conduction** [n] [uncountable] (physics) the process by which heat or electricity passes along or through a material.

¹¹**diffusion** [n] [uncountable] **1.** the spreading of something more widely; **2.** the mixing of substances by the natural movement of their particles; **3.** the spreading of elements of culture from 1 region or group to another.

¹²**electromagnetic** [a] (physics) in which the electrical & magnetic properties of something are related.

¹³**phenomenon** [n] (plural **phenomena** a fact or an event in nature or society, especially one that is not fully understood.)

subject to the state equation:

$$\begin{cases} -\Delta y = 0, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = \alpha(u - y), & \text{on } \Gamma, \end{cases}$$

and the pointwise control constraints $u_a(\mathbf{x}) \leq u(\mathbf{x}) \leq u_b(\mathbf{x})$ on Γ . “Such pointwise bounds for the control are quite natural, since the available capacities for heating or cooling are usually restricted. The constant $\lambda \geq 0$ can be viewed as a measure of the energy costs needed to implement the control u . From the mathematical viewpoint, this term also serves as a regularization parameter; it has the effect that possible optimal controls show improved regularity properties.” [...] “The function α represents the heat transmission coefficient from Ω to the surrounding medium. The functional J to be minimized is called the cost functional. The factor $\frac{1}{2}$ appearing in it has no influence on the solution of the problem. It is introduced just for the sake of convenience: it will later cancel out a factor 2 arising from differentiation. We seek an optimal control $u = u(\mathbf{x})$ together with the associated state $y = y(\mathbf{x})$. The minus sign in front of the Laplacian Δ appears to be unmotivated at 1st glance. It is introduced because Δ is not a coercive operator, while $-\Delta$ is.” – Tröltzsch, 2010, p. 4

“Observe that in the above problem the cost functional is quadratic, the state is governed by a linear elliptic PDE, & the control acts on the boundary of the domain.”: thus have a *linear-quadratic elliptic boundary control problem*.

Remark 3.1 (Notations used in Tröltzsch, 2010). Denote the element of surface area by ds & the outward unit normal to Γ at $\mathbf{x} \in \Gamma$ by $\nu(\mathbf{x})$ ¹⁴.

Remark 3.2. “The problem is strongly simplified. Indeed, in a realistic model Laplace’s equation $\Delta y = 0$ has to be replaced by the stationary heat conduction equation $\nabla \cdot (a \nabla y) = 0$, where the coefficient a can depend on \mathbf{x} or even on y . If $a = a(y)$ or $a = a(\mathbf{x}, y)$, then the PDE is quasilinear. In addition, it will in many cases be more natural to describe the process by a time-dependent PDE.” – Tröltzsch, 2010, p. 4

Example 3.2 (Optimal heat source). Similarly, the control can act as a heat source in the domain Ω . Problems of this kind arise if the body Ω is heated by electromagnetic induction or by microwaves. Assuming at 1st that the boundary temperature vanishes, we obtain the following problem:

$$\min J(y, u) := \frac{1}{2} \int_{\Omega} |y(\mathbf{x}) - y_{\Omega}(\mathbf{x})|^2 d\mathbf{x} + \frac{\lambda}{2} \int_{\Omega} |u(\mathbf{x})|^2 d\mathbf{x},$$

subject to

$$\begin{cases} -\Delta y = \beta u, & \text{in } \Omega, \\ y = 0, & \text{on } \Gamma, \end{cases}$$

and $u_a(\mathbf{x}) \leq u(\mathbf{x}) \leq u_b(\mathbf{x})$ in Ω . Here, the coefficient $\beta = \beta(\mathbf{x})$ is prescribed. Observe that by the special choice $\beta = \chi_{\Omega_c}$ (where χ_E denotes the characteristic function of a set E), it can be achieved that u acts only in a subdomain $\Omega_c \subset \Omega$. This problem is a linear-quadratic elliptic control problem with distributed control. It can be more realistic to prescribe an exterior temperature y_a rather than assume that the boundary temperature vanishes. Then a better model is given by the state equation

$$\begin{cases} -\Delta y = \beta u, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = \alpha(y_a - y), & \text{on } \Gamma. \end{cases}$$

3.1.2.2 Optimal nonstationary boundary control

See Tröltzsch, 2010, pp. 5–6. “Let $\Omega \subset \mathbb{R}^3$ represent a potato that is to be roasted over a fire for some period of time $T > 0$.” Denote its temperature by $y = y(t, \mathbf{x})$, with $(t, \mathbf{x}) \in [0, T] \times \Omega$. “Initially, the potato has temperature $y_0 = y_0(\mathbf{x})$, & we want to serve it at a pleasant palatable¹⁵ temperature y_{Ω} at the final time T .” Write $Q := (0, T) \times \Omega$, $\Sigma := (0, T) \times \Gamma$. Then problem reads as follows:

$$\min J(y, u) := \frac{1}{2} \int_{\Omega} |y(T, \mathbf{x}) - y_{\Omega}(\mathbf{x})|^2 d\mathbf{x} + \frac{\lambda}{2} \int_0^T \int_{\Gamma} |u(t, \mathbf{x})|^2 d\Gamma dt,$$

subject to

$$\begin{cases} y_t - \Delta y = 0, & \text{in } Q, \\ \partial_{\mathbf{n}} y = \alpha(u - y), & \text{on } \Sigma, \\ y(0, \mathbf{x}) = y_0(\mathbf{x}), & \text{in } \Omega, \end{cases}$$

¹⁴NQBH: I prefer to use $\mathbf{n}(\mathbf{x})$, with “n” stands for “normal”, naturally & obviously.

¹⁵palatable [a] 1. (of food or drink) having a pleasant or acceptable taste; 2. palatable (to somebody) pleasant or acceptable to somebody, OPPOSITE: unpalatable.

& $u_a(t, \mathbf{x}) \leq u(t, \mathbf{x}) \leq u_b(t, \mathbf{x})$ on Σ . By continued turning of the spit¹⁶, we produce $u(t, \mathbf{x})$. The heating process has to be described by the *nonstationary heat equation*, which is a parabolic differential equation: thus have to deal with a *linear-quadratic parabolic boundary control problem*.

3.1.2.3 Optimal vibrations

“Suppose that a group of pedestrians crosses a bridge, trying to excite¹⁷ oscillations¹⁸ in it. This can be modeled (strongly abstracted) as follows: let $\Omega \subset \mathbb{R}^2$ denote the domain of the bridge, $y = y(t, \mathbf{x})$ its *transversal*¹⁹ *displacement*²⁰, $u = u(t, \mathbf{x})$ the *force density* acting in the vertical direction, & $y_d = y_d(t, \mathbf{x})$ a *desired evolution of the transversal vibrations*²¹. We then obtain the optimal control problem:

$$\min J(y, u) := \frac{1}{2} \int_0^T \int_{\Omega} |y(t, \mathbf{x}) - y_d(t, \mathbf{x})|^2 \, d\mathbf{x} \, dt + \frac{\lambda}{2} \int_0^T \int_{\Omega} |u(t, \mathbf{x})|^2 \, d\mathbf{x} \, dt,$$

subject to

$$\begin{cases} y_{tt} - \Delta y = u, & \text{in } Q, \\ y(0) = y_0, & \text{in } \Omega, \\ y_t(0) = y_1, & \text{in } \Omega, \\ y = 0, & \text{on } \Sigma, \end{cases}$$

and $u_a(t, \mathbf{x}) \leq u(t, \mathbf{x}) \leq u_b(t, \mathbf{x})$ in Q . This is a *linear-quadratic hyperbolic control problem with distributed control*.” [...] “Interesting control problems for oscillating elastic networks have been treated by Lagnese et al. [LLS94]. An elementary introduction to the controllability of oscillations can be found in [Kra95].

In the linear-quadratic case, the theory of hyperbolic problems has many similarities to the parabolic theory studied in Tröltzsch, 2010. However, the treatment of semilinear hyperbolic problems is much more difficult, since the smoothing properties of the associated solution operators are weaker. As a consequence, many of the techniques presented in Tröltzsch, 2010 fail in the hyperbolic case.” – Tröltzsch, 2010, pp. 6–7

3.1.3 Examples of Nonconvex Problems

“However, linear models do not suffice for many real-world phenomena. Instead, one often needs quasilinear or, much simpler, semilinear equations. Recall that a 2nd-order equation is called *semilinear* if the main parts (i.e., the expressions involving highest-order derivatives) of the differential operators considered in the domain & on the boundary are linear w.r.t. the desired solution. For such equations, the theory of optimal control is well developed.

Optimal control problems with semilinear state equations are, as a rule, nonconvex, even if the cost functional is convex.

“Associated optimal control problems can be obtained by prescribing a cost functional & suitable constraints.” – Tröltzsch, 2010, p. 7

3.1.3.1 Problems involving semilinear elliptic equations

Example 3.3 (Heating with radiation boundary condition). *If the heat radiation of the heated body is taken into account, then we obtain a problem with a nonlinear Stefan–Boltzmann boundary condition. If this case, the control u is given by the temperature of the surrounding medium:*

$$\begin{cases} -\Delta y = 0, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = \alpha(u^4 - y^4), & \text{on } \Gamma. \end{cases}$$

The nonlinearity y^4 occurs in the boundary condition, while the heat conduction equation itself is linear.

¹⁶**spit** [n] *in/from mouth* **1.** [uncountable] the liquid produced in your mouth, SYNONYM: **saliva**; **2.** [countable, usually singular] the act of spitting liquid or food out of your mouth; *piece of land* **3.** [countable] a long, thin piece of land that sticks out into the sea, a lake, etc.; *for cooking meat* **4.** [countable] a long, thin, straight piece of metal that you put through meat to hold & turn it while you cook it over a fire.

¹⁷**excite** [v] **1.** to make somebody feel a particular emotion or react in a particular way, SYNONYM: **arouse**; **2. excite somebody** to make somebody feel very pleased, interested or enthusiastic, especially about something that is going to happen; **3. excite somebody/something** to make somebody/something nervous, upset or active & unable to relax; **4. excite something** to produce a state of increased energy or activity in a physical or biological system, SYNONYM: **stimulate**; **5. excite something (physics)** to bring something to a state of higher energy.

¹⁸**oscillation** [n] **1.** [countable, uncountable] **oscillation (of something)** a regular movement between 1 position & another; **2.** [countable] **oscillation (between A & B)** a repeated change between different states, ideas, etc.; **3.** [countable] (*specialist*) regular variation in size, strength or position around a central point or value, especially of an electrical current or electric field.)

¹⁹**transversal** [n] a line that intersects a system of lines.

²⁰**displacement** [n] **1.** [uncountable] the act of displacing somebody/something; the process of being displaced; **2.** [uncountable, singular] **displacement (of something) (physics)** the distance between the final & initial (= 1st) positions of an object which has moved.

²¹**vibration** [n] [countable, uncountable] **1. vibration (of something)** a continuous shaking movement; **2. vibration (of something) (physics)** oscillation in a substance about its equilibrium state.

Example 3.4 (Simplified superconductivity).

Example 3.5 (Control of stationary flows).

3.1.3.2 Problems involving semilinear parabolic equations

3.1.4 Basic Concepts for the Finite-Dimensional Case

3.2 Linear-Quadratic Elliptic Control Problems

3.3 Linear-Quadratic Parabolic Control Problems

3.4 Optimal Control of Semilinear Elliptic Equations

3.5 Optimal Control of Semilinear Parabolic Equations

3.6 Optimization Problems in Banach Spaces

3.6.1 The Karush–Kuhn–Tucker Conditions

3.6.1.1 Convex problems

The Lagrange multiplier rule. “The formal Lagrange method, which was employed repeatedly in the previous chapters, has a rigorous mathematical foundation. In this section, we introduce the basics of this theory needed for understanding problems with state constraints. The corresponding proofs & further results can be found in texts dealing with optimization in general spaces. The theory of convex problems is described in Balakrishnan [Bal65], Barbu & Precupanu [BP78], & Ekeland & Temam [ET74]; nonconvex differentiable problems are treated in, e.g., Ioffe & Tihomirov [IT79], Jahn [Jah94], Luenberger [Lue69], & Tröltzsch [Trö84b].

There are numerous books dealing with the theory & numerical treatment of nonlinear differentiable finite-dimensional optimization problems. In this connection, we refer the interested reader to Alt [Alt02], Gill et al. [GMW81], Grossmann & Terno [GT97b], Kelley [Kel99], Luenberger [Lue84], Nocedal & Wright [NW99], Polak [Pol97], & Wright [Wri93], to name just a few.

In the following, we generally assume that U & Z are real Banach spaces, $G : U \rightarrow Z$ is in general a nonlinear mapping, & $C \subset U$ is a nonempty & convex set.

Definition 3.1 (Convex cone). *A convex set $K \subset Z$ is said to be a convex cone if $\lambda z \in K$ whenever $z \in K$ & $\lambda > 0$.*

Any convex cone induces a partial ordering \geq_K in the space Z :

Definition 3.2. *Let $K \subset Z$ be a convex cone. We write $z \geq_K 0$ iff $z \in K$. Analogously, we write $z \leq_K 0$ iff $-z \in K$.*

The elements in K are said to be *nonnegative*. Note, however, that nonnegativity in the sense of this definition does not imply the usual nonnegativity in the set of real numbers, as the following example shows.

Example 3.6. *Let $Z = \mathbb{R}^3$, & let $K = \{z \in \mathbb{R}^3; z_1 = 0, z_2 \leq 0, z_3 \geq 0\}$. Then K is evidently a convex cone, but $z \geq_K 0$ implies nonnegativity only for z_3 .*

The next definition enables us to introduce a notion of “nonnegativity” also in dual spaces. This notion will be needed for defining Lagrange multipliers, because they are elements of dual spaces.

Definition 3.3 (Dual cone). *Let $K \subset Z$ be a convex cone. Then the set $K^+ = \{z^* \in Z^*; \langle z^*, z \rangle_{Z^*, Z} \geq 0, \forall z \in K\}$ is called the dual cone of K .*

Example 3.7. (i) *Let $Z = L^2(\Omega)$ with a bounded domain $\Omega \subset \mathbb{R}^N$, & let $K = \{z \in L^2(\Omega); z(\mathbf{x}) \geq 0 \text{ for a.e. } \mathbf{x} \in \Omega\}$. Here, we have $Z = Z^*$ by the Riesz representation theorem & $K^+ = K$ according to Exercise 6.1.*

(ii) *Let Z be a Banach space & let $K = \{0\}$. Then $z \geq_K 0$ iff $z = 0$, & thus $K^+ = Z^*$; in fact, for any $z^* \in Z^*$ we have $\langle z^*, 0 \rangle_{Z^*, Z} = 0 \geq 0$.*

(iii) *If $K = Z$, then all elements of Z are nonnegative. Hence, $K^+ = \{0\}$ with the zero functional $0 \in Z^*$.*

Below we consider the following optimization problem in a Banach space: **(6.1)**

$$\min f(u), \quad G(u) \leq_K 0, \quad u \in C.$$

The constraints in (6.1) are viewed differently: as a “complicated” inequality $G(u) \leq_K 0$, which is to be eliminated by means of a Lagrange multiplier, & a “simple” constraint $u \in C$, which is accounted for explicitly. This motivates the following definition.

Definition 3.4 (Lagrange function). *The function $L : U \times Z^* \rightarrow \mathbb{R}$, **(6.2)** $L(u, z^*) = f(u) + \langle z^*, G(u) \rangle_{Z^*, Z}$, is called the Lagrange function. Any $(\bar{u}, z^*) \in U \times K^+$ satisfying the chain of inequalities **(6.3)** $L(\bar{u}, v^*) \leq L(\bar{u}, z^*) \leq L(u, z^*)$, $\forall u \in C$, $\forall v^* \in K^+$ is called a saddle point of L . If this is the case, z^* is said to be a Lagrange multiplier associated with \bar{u} .*

In the previous chapters, when dealing with the optimal control of PDEs we denoted the Lagrangian by \mathcal{L} . To facilitate the distinction, we use the letter L here. The existence of saddle points is most easily shown for *convex* optimization problems.

Definition 3.5. *Let U be a Banach space, & let the convex cone $K \subset Z$ induce the partial ordering \geq_K in the Banach space Z . An operator $G : U \rightarrow Z$ is said to be convex (w.r.t. \leq_K) if*

$$G(\lambda u + (1 - \lambda)v) \leq_K \lambda G(u) + (1 - \lambda)G(v), \quad \forall u, v \in U, \quad \forall \lambda \in (0, 1).$$

Evidently, every linear operator is convex. In the following, we write the strict inequality $z <_K 0$ iff $-z$ is an *interior point* of K , i.e., $z <_K 0 \Leftrightarrow -z \in \text{int } K$.

Theorem 3.1. *Suppose that a convex functional $f : U \rightarrow \mathbb{R}$, a convex operator $G : U \rightarrow Z$, & a solution \bar{u} to the problem (6.1) are given. Moreover, let there exist some $\tilde{u} \in C$ s.t. $G(\tilde{u}) <_K 0$, i.e., **(6.4)** $-G(\tilde{u}) \in \text{int } K$. Then there is some $z^* \in K^+$ s.t. (\bar{u}, z^*) is a saddle point of the Lagrangian L . In addition, we have the complementary slackness condition **(6.5)** $\langle z^*, G(\bar{u}) \rangle_{Z^*, Z} = 0$.*

The proof of the above theorem can be found in, e.g., Luenberger [Lue69]. In the literature, the condition (6.4) is usually referred to as the *Slater condition*. It can only be satisfied if the cone K has nonempty interior. This excludes, e.g., the case $K = \{0\}$, which corresponds to the equality constraint $G(u) = 0$. In this case, the above theorem fails to apply, but other existence results concerning Lagrange multipliers are available. The lack of interior points is a much more serious problem in the following situation.

Example 3.8. *Consider the natural nonnegative cone in $Z = L^2(0, 1)$, $K = \{z(\cdot) \in L^2(0, 1); z(x) \geq 0 \text{ for a.e. } x \in (0, 1)\}$. Quite unexpectedly, we have $\text{int } K = \emptyset$. How can this be possible? One is tempted to believe that, e.g., $z(x) \equiv 1$ is an interior point of K . Unfortunately, this is not true. In fact, the sequence $\{v_n\}_{n=1}^\infty \subset L^2(\Omega)$ with*

$$v_n(x) = \begin{cases} 1 & \text{in } \left[0, 1 - \frac{1}{n}\right], \\ -1 & \text{in } \left[1 - \frac{1}{n}, 1\right], \end{cases}$$

while obviously converging to z w.r.t. the L^2 norm, is not contained in K . Consequently, $z \notin \text{int } K$. This undesired behavior is simply a consequence of the fact that the L^2 norm, & likewise any other L^p norm with $1 \leq p < \infty$, measures an integral & not the maximal absolute value of a function. This fact constitutes a major obstacle in the treatment of optimization problems in function spaces.

Theorem 3.2. *Suppose that the mappings f & G in Theorem 6.1 are Gâteaux differentiable at \bar{u} . Then we have the variational inequality $D_u L(\bar{u}, z^*)(u - \bar{u}) \geq 0$, $\forall u \in C$.*

Here & in the following, D_u again denotes the partial Gâteaux or Fréchet derivative w.r.t. u . The assertion is an immediate consequence of the saddle point condition (6.3), which implies that \bar{u} solves the problem without the constraint $G(u) \leq_K 0$, namely $L(\bar{u}, z^*) = \min_{u \in C} L(u, z^*)$. The associated variational inequality reads, in explicit form,

$$f'(\bar{u})(u - \bar{u}) + \langle z^*, G'(\bar{u})(u - \bar{u}) \rangle_{Z^*, Z}, \quad \forall u \in C,$$

or, equivalently,

$$\langle f'(\bar{u}) + G'(\bar{u})^* z^*, u - \bar{u} \rangle_{U^*, U} \geq 0, \quad \forall u \in C.$$

In the unconstrained case where $C = U$, we get the equation $f'(\bar{u}) + G'(\bar{u})^* z^* = 0 \in U^*$.

Examples. We illustrate the application & limitations of the above theorems by means of simple examples that do not involve PDEs.

1-sided box constraints in $L^2(0,1)$. Let $u_d \in L^2(0,1)$ be given. We consider the minimization problem **(6.6)**

$$\min f(u) := \frac{1}{2} \int_0^1 |u(x) - u_d(x)|^2 dx \text{ subject to } u(x) \leq 0 \text{ for a.e. } x \in (0,1).$$

The above problem is a special case of problem (6.1), with the specifications $U = Z = L^2(0,1)$ & $G = I$ (the identity mapping). The associated convex cone K is the set of almost-everywhere nonnegative elements of $L^2(0,1)$, & we have $C = U$. The problem has a unique minimizer \bar{u} , namely, $\bar{u}(x) = \min\{u_d(x), 0\}$. We investigate whether there exists an associated Lagrange multiplier. The corresponding Lagrangian function reads

$$L(u, \mu) = f(u) + (\mu, G(u))_{L^2(0,1)} = \int_0^1 \left(\frac{1}{2} (u(x) - u_d(x))^2 + \mu(x)u(x) \right) dx.$$

Here, by the Riesz representation theorem, the functional $z^* \in Z^*$ has been identified with some function $\mu \in L^2(0,1)$. We search for a Lagrange multiplier $\mu \in L^2(0,1)$. Since $\text{int } K = \emptyset$, Theorem 6.1 does not apply. Instead, the Lagrange multiplier is constructed using a pointwise approach. To this end, recall that, owing to Lemma 2.21 on p. 63, we have the variational inequality

$$\int_0^1 (\bar{u}(x) - u_d(x))(u(x) - \bar{u}(x)) dx \geq 0, \quad \forall u(\cdot) \leq 0.$$

This can only be true if the implications

$$\begin{aligned} \bar{u}(x) < 0 &\Rightarrow \bar{u}(x) - u_d(x) = 0, \\ \bar{u}(x) = 0 &\Rightarrow \bar{u}(x) - u_d(x) \leq 0, \end{aligned}$$

are valid a.e. But then $\bar{u}(x) - u_d(x)$ must be nonpositive a.e. Since we have used arguments of this kind repeatedly in Chap. 2, we do not explain this in detail here. We now define $\mu(x) := -f'(\bar{u})(x) = -(\bar{u}(x) - u_d(x))$. Then, owing to the above implications, we have $\mu \geq 0$ as well as $\mu(x)\bar{u}(x) = 0$ for a.e. $x \in (0,1)$, which is the *pointwise form of the complementary slackness condition* (6.5). Finally, it follows from the definition of μ that $\mu = -f'(\bar{u})$, i.e., $f'(\bar{u}) + \mu = 0$, which, in turn, is equivalent to the equation $D_u L(\bar{u}, \mu) = 0$. Hence, the function μ defined above is a Lagrange multiplier.

2-sided box constraints in $L^2(0,1)$. We now consider the above minimization problem with the same functional f as in (6.6), but this time with constraints from both above & below: **(6.7)**

$$\min f(u) \text{ subject to } -1 \leq u(x) \leq 1 \text{ for a.e. } x \in (0,1).$$

Again, we put $C = U = L^2(0,1)$, & we cast the constraints in the form $u(x) - 1 \leq 0$, $-u(x) - 1 \leq 0$. We then have to choose $Z = L^2(0,1) \times L^2(0,1)$ & $K = L^2(0,1) \times L^2(0,1)_+$, where $L^2(0,1)_+$ denotes the set of a.e. nonnegative elements of $L^2(0,1)$. The convex operator $G : L^2(0,1) \rightarrow L^2(0,1) \times L^2(0,1)$ is defined by $G(u) := (u(\cdot) - 1, -u(\cdot) - 1)^\top$.

Although the function $\bar{u}(x) \equiv 0$ obeys both inequalities strictly, we again have $\text{int } K = \emptyset$ & thus cannot employ Theorem 6.1. However, the construction used in Sect. 1.4.7 works. The Lagrangian is now given by

$$L(u, \mu) = L(u, \mu_a, \mu_b) = \frac{1}{2} \|u - u_d\|_{L^2(0,1)}^2 + (-u - 1, \mu_a)_{L^2(0,1)} + (u - 1, \mu_b)_{L^2(0,1)}.$$

We make the pointwise definitions **(6.8)** $\mu_a(x) = (f'(x))_+ = (\bar{u}(x) - u_d(x))_+$, $\mu_b(x) = (f'(x))_- = (\bar{u}(x) - u_d(x))_-$, where, as usual, $z_+ = \frac{z+|z|}{2}$ & $z_- = \frac{|z|-z}{2}$. Obviously, μ_a & μ_b are nonnegative. The reader will be asked in Exercise 6.2 to check that the arguments from Sect. 1.4.7 carry over almost unchanged to give $D_u L(\bar{u}, \mu) = f'(\bar{u}) + \mu_b - \mu_a = 0$ & the slackness conditions $(-\bar{u} - 1, \mu_a)_{L^2(0,1)} = (\bar{u} - 1, \mu_b)_{L^2(0,1)} = 0$. Consequently, μ_a & μ_b are Lagrange multipliers for \bar{u} .

If we assume $u_d \in L^\infty(0,1)$, then both multipliers belong to $L^\infty(0,1)$. This nice byproduct of the pointwise construction follows from the fact that $\bar{u} - u_d \in L^\infty(0,1)$.

Remark 3.3. The problem with 2-sided constraints could also be considered in the space $L^\infty(0,1)$, since in this case every admissible control u is automatically bounded & measurable. Moreover, the cone K of nonnegative functions in $L^\infty(0,1)$ has interior points, & $\bar{u}(x) \equiv 0$ satisfies the Slater condition. Theorem 6.1 then yields the existence of Lagrange multipliers $\mu_a, \mu_b \in L^\infty(0,1)^*$. However, we do not gain much benefit from this result, since $L^\infty(0,1)^*$ is a space of continuous linear functionals that need not even be measures.

3.6.1.2 Differentiable problems

Lagrange multiplier rules & constraint qualifications. We now investigate the problem (6.1) without assuming f & G to be convex. We consider $\min f(u)$, $G(u) \leq_K 0$, $u \in C$, where C is still convex. Instead of convexity, we postulate the Fréchet differentiability of f & G . We use the same Lagrangian function $L = L(u, z^*)$ as in Sect. 6.1.1, but, in view of the nonconvexity, we can no longer expect a saddle point property to be valid. Therefore, Lagrange multipliers are defined in a slightly different way.

Definition 3.6 (Local solution). *Let $\bar{u} \in U$ be admissible. We call \bar{u} a local solution of the minimization problem (6.1) if there is some $\varepsilon > 0$ s.t. $f(\bar{u}) \leq f(u)$, $\forall u \in C$ with $G(u) \leq_K 0$ & $\|u - \bar{u}\|_U \leq \varepsilon$.*

Definition 3.7 (Lagrange multiplier). *Let \bar{u} be a local solution to the problem (6.1). Then any $z^* \in K^+$ satisfying the conditions (6.9)–(6.10)*

$$\begin{aligned} D_u L(\bar{u}, z^*)(u - \bar{u}) &\geq 0, \quad \forall u \in C, \\ \langle z^*, G(\bar{u}) \rangle_{Z^*, Z} &= 0 \end{aligned}$$

is called a Lagrange multiplier associated with \bar{u} .

In order that the existence of such a Lagrange multiplier be guaranteed, a so-called *constraint qualification* must be postulated. Since such a condition involves the locally optimal control itself, it usually cannot be verified without knowledge of this function. There are various constraint qualifications. A rather general one, which suffices for our purposes, is the *Zowe–Kurcyusz condition* (see Zowe & Kurcyusz [ZK79]).

Definition 3.8. *Suppose that $\bar{u} \in C$ with $G(\bar{u}) \leq_K 0$ is given. We call the sets*

$$C(\bar{u}) := \{\alpha(u - \bar{u}); \alpha \geq 0, u \in C\}, \quad K(\bar{z}) := \{\beta(z - \bar{z}); \beta \geq 0, z \in K\}$$

the conical hulls to C & K at \bar{u} & \bar{z} , respectively. The condition (6.11)

$$G'(\bar{u})C(\bar{u}) + K(-G(\bar{u})) = Z$$

is called the Zowe–Kurcyusz constraint qualification.

The above relation is obviously equivalent to saying that for any $z \in Z$ the equation (6.12) $\alpha G'(\bar{u})(u - \bar{u}) + \beta(v + G(\bar{u})) = z$ is solvable with suitable $u \in C$, $v \geq_K 0$, $\alpha \geq 0$, & $\beta \geq 0$. Recall that $v \geq_K 0$ iff $v \in K$.

Theorem 3.3. *Let \bar{u} be a local solution to problem (6.1), & let f & G be continuously Fréchet differentiable in an open neighborhood of \bar{u} . If the constraint qualification (6.11) holds, then there exists a Lagrange multiplier $z^* \in Z^*$ associated with \bar{u} . Moreover, the set of Lagrange multipliers associated with \bar{u} is bounded.*

The proof of this multiplier rule is due to Zowe & Kurcyusz [ZK79]. From (6.9) it follows that (6.13)

$$\langle f'(\bar{u}) + G'(\bar{u})^* z^*, u - \bar{u} \rangle_{U^*, U} \geq 0, \quad \forall u \in C.$$

Remark 3.4. *Sometimes it is difficult or even meaningless to establish $G'(\bar{u})^*$ in explicit form. Then (6.13) is replaced by the equivalent inequality*

$$f'(\bar{u})(u - \bar{u}) + \langle z^*, G'(\bar{u})(u - \bar{u}) \rangle_{Z^*, Z} \geq 0, \quad \forall u \in C.$$

Example 3.9. *Of particular interest is the minimization problem with both equality & set constraints: (6.14) $\min f(u)$, $G(u)$, $u \in C$, where f, G , & C are defined as before. In this special case, the constraint qualification (6.11) reads (6.15) $G'(\bar{u})C(\bar{u}) = Z$. If it is satisfied, then a Lagrange multiplier $z^* \in Z^*$ exists s.t. the variational inequality (6.13) is valid. The complementary slackness condition (6.10) is meaningless for equality constraints.*

In Sect. 6.1.3, we will apply this result to the special case of $G(u) = Ay - Bv = 0$, where $A : Y \rightarrow Y^$ is a continuously invertible operator representing an elliptic differential operator, y denotes the state, & $v \in V_{\text{ad}} \subset V$ is the control.*

In this case, we have $Z := Y^$, $U := Y \times V$, & $C := Y \times V_{\text{ad}}$. The constraint qualification is always satisfied, since the equation $G'(\bar{u})(u - \bar{u}) = A(y - \bar{y}) + B(v - \bar{v}) = z$ is solvable for any $z \in Z = Y^*$ with $v = \bar{v}$ & $y = A^{-1}z + \bar{y}$. The element $u - \bar{u} = (y - \bar{y}, v - \bar{v})$ belongs to the cone $C(\bar{u})$.*

Discussion of the Zowe–Kurcyusz constraint qualification. In the following, we illustrate the application of condition (6.11) for various types of constraints, 1st in the general situation, & then for pointwise constraints in function spaces.

Pure equality constraints $G(u) = 0$. With $C = U$ & $K = \{0\}$, (6.11) becomes **(6.16)** $G'(\bar{u})U = Z$. In other words, the operator $G'(\bar{u})$ must be surjective. This surjectivity requirement comes from the classical Lagrange multiplier rule for equality constraints. The relation (6.13) attains the form **(6.17)** $f'(\bar{u}) + G'(\bar{u})^* z^* = 0$.

Inequality constraints. Let the constraints be given as in (6.1). If the minimizer \bar{u} satisfies $G(\bar{u}) <_K 0$, i.e., if $-G(\bar{u}) \in \text{int } K$, then the constraint qualification (6.11) is fulfilled (Exercise 6.3). Since the constraint is not active, this case is not interesting.

The following *linearized Slater condition* is sufficient for the Zowe–Kurcyusz constraint qualification (6.11) to hold: **(6.18)**

$$\boxed{\exists \tilde{u} \in C : G(\bar{u}) + G'(\bar{u})(\tilde{u} - \bar{u}) <_K 0.}$$

This is easily seen: the Zowe–Kurcyusz condition postulates for any $z \in Z$ the existence of constants $\alpha \geq 0$, $\beta \geq 0$ & elements $k \in K$, $u \in C$ s.t. the equation $\alpha G'(\bar{u})(u - \bar{u}) + \beta(k + G(\bar{u})) = z$ is valid. To show this, put $\alpha = \beta$, $u = \tilde{u}$, & $\bar{z} = G(\bar{u}) + G'(\bar{u})(\tilde{u} - \bar{u})$. Then the above equation reduces to $\alpha(\bar{z} + k) = z$ &, since K is a cone, to $\alpha\bar{z} + q = z$, with $q \in K$. Now we choose α so large that $z - \alpha\bar{z} \geq_K 0$. This is possible, because by (6.18) \bar{z} lies in the interior of $-K$. With this, we satisfy the above condition with the choice $q = z - \alpha\bar{z} \geq_K 0$.

If both K & C have interior points, then (6.18) is equivalent to the following condition (cf. Penot [Pen82]): **(6.19)** $\exists h \in \text{int } C(\bar{u}) : G(\bar{u}) + G'(\bar{u})h <_K 0$.

Equality & inequality constraints. Suppose the constraints have the form $G_1(u) = 0$, $G_2(u) \leq_K 0$, $u \in C$. Then the following condition is sufficient for (6.11) to hold (cf. [HPUU09], Lemma 1.14): $G'_1(\bar{u})$ is surjective, & **(6.20)** $\exists h \in C(\bar{u}) : G'_1(\bar{u})h = 0$, $G_2(\bar{u}) + G'_2(\bar{u})h <_K 0$.

As the following examples will show, the applicability of the Zowe–Kurcyusz constraint qualification to inequality constraints in function spaces is essentially restricted to cones of nonnegative functions with nonempty interior.

1-sided box constraints for u . We begin our analysis with a problem involving 1-sided constraints: **(6.21)**

$$\min f(u) := \int_{\Omega} \psi(x, u(x)) \, dx, \quad u(x) \leq u_b(x) \text{ for a.e. } x \in \Omega.$$

Here, $u_b \in L^\infty(\Omega)$ is given, & the function ψ is sufficiently smooth & satisfies a suitable growth condition in order to guarantee that the given integral functional f be continuous differentiable in $U = L^2(\Omega)$. The above minimization problem is of the form $\min f(u)$, $G(u) <_K 0$, with $G(u)(x) := u(x) - u_b(x)$. As an affine continuous operator, G is differentiable from U into $Z = U$. The cone K is given by the set of a.e. nonnegative elements of $L^2(\Omega)$.

In this case, the Zowe–Kurcyusz constraint qualification (6.11) is satisfied: in view of $C = L^2(\Omega)$, we have $C(\bar{u}) = L^2(\Omega)$. & since $G'(\bar{u})$ is the identity mapping, for any $z \in L^2(\Omega)$ there is some $u \in L^2(\Omega) = C(\bar{u})$ s.t. $G'(\bar{u})u = z$: one simply chooses $u = z$.

Consequently, Theorem 6.3 may be applied in $L^2(\Omega)$, where, in view of the Riesz representation theorem, every $z^* \in L^2(\Omega)^*$ can be identified with some $\mu \in L^2(\Omega)$. Hence, for any local solution \bar{u} there exists some a.e. nonnegative multiplier $\mu \in L^2(\Omega)$ s.t. $f'(\bar{u}) + G'(\bar{u})^* \mu = 0$. We may identify $f'(\bar{u})$ with the function $\psi_u(\cdot, \bar{u}(\cdot)) \in L^2(\Omega)$, & $G'(\bar{u})^*$ is the identity operator in $L^2(\Omega)$. We therefore find that $\psi_u(x, \bar{u}(x)) + \mu(x) = 0$, $\mu(x) \geq 0$, for a.e. $x \in \Omega$.

In this example, the Zowe–Kurcyusz condition was applicable even though the cone of nonnegative functions in $L^2(\Omega)$ had empty interior. This is in a certain sense an exceptional case. Alternatively, we could have constructed the multiplier directly as in (6.8) by setting $\mu(x) = (\psi_u(x, \bar{u}(x)))_-$.

2-sided box constraints for u . We consider the same problem as above, but this time with the 2-sided control constraints $u_a(x) \leq u(x) \leq u_b(x)$ for a.e. $x \in \Omega$, with bounded & measurable functions $u_a \leq u_b$. We fit these constraints into the abstract framework of (6.1) by choosing the operator G to be of the form $G(u) = (u_a - u, u - u_b)^\top$.

Evidently, G is a continuously differentiable mapping from $L^2(\Omega)$ into $L^2(\Omega) \times L^2(\Omega)$. However, the Zowe–Kurcyusz constraint qualification cannot be directly satisfied in the form (6.11), as can be shown with a little effort. Again, we have the problem that the cone of nonnegative functions in $L^2(\Omega)$ has empty interior. It would also not be helpful to work in $L^\infty(\Omega)$ instead, since then we would at best obtain measures as multipliers for the control constraints. As in (6.8), a possible way out is to define Lagrange multipliers by $\mu_a(x) := \psi_u(x, \bar{u}(x))_+$, $\mu_b(x) := \psi_u(x, \bar{u}(x))_-$, with which the Karush–Kuhn–Tucker conditions are fulfilled.

2nd-order optimality conditions. The scope of the Karush–Kuhn–Tucker theory in Banach spaces also encompasses 2nd-order necessary & sufficient optimality conditions. For illustration, we only discuss the problem (6.14): $\min f(u)$, $G(u) = 0$,

$u \in C$, additionally assuming that f & G are twice continuously Fréchet differentiable. Suppose that \bar{u} satisfies, together with $z^* \in Z^*$, the 1st-order necessary condition (6.22)

$$f'(\bar{u})(u - \bar{u}) + \langle z^*, G'(\bar{u})(u - \bar{u}) \rangle_{Z^*, Z} \geq 0, \quad \forall u \in C.$$

Moreover, let there exist some $\delta > 0$ s.t. (6.23)

$$L''(\bar{u}, z^*)[u, u] := f''(\bar{u})[u, u] + \langle z^*, G''(\bar{u})[u, u] \rangle_{Z^*, Z} \geq \delta \|u\|_U^2$$

for all $u \in C(\bar{u})$ s.t. (6.24) $G'(\bar{u})u = 0$.

Lemma 3.1. *If \bar{u} is admissible for problem (6.14) & the conditions (6.22)–(6.24) are fulfilled, then \bar{u} is locally optimal for (6.14).*

These 2nd-order sufficient optimality conditions follow from general results due to Maurer & Zowe [MZ79, Mau81]. The lemma applies only to problems in which the 2-norm discrepancy²² does not play a role. Since we have not provided any further information concerning the structure of the set C , we are not in a position to define & make use of strongly active constraints. The conditions above are thus too restrictive. In the case of inequality constraints of the form $G(u) \leq_K 0$, 1st-order sufficient optimality conditions can also be employed; see [MZ79]. For PDEs with state constraints, refer to [CDIRT08]. In the case of pointwise constraints in function spaces, usually strongly active sets in the sense of Dontchev et al. [DHPY95] are used for this purpose.

3.6.1.3 A semilinear elliptic problem

Let $\Omega \subset \mathbb{R}^N$, $N \leq 3$, be a bounded Lipschitz domain. For given $v \in L^2(\Omega)$, we consider the elliptic BVP

$$\begin{cases} -\Delta y + y + y^3 = v, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = 0, & \text{on } \Gamma. \end{cases}$$

As shown on pp.181–183, this problem is easy to handle in the state space $Y = H^1(\Omega)$. Introducing the mapping $A : Y \rightarrow Y^*$ generated by the elliptic operator $-\Delta + I$, the *Nemytskii operator* $\Phi : Y \rightarrow V = L^2(\Omega)$, $y(\cdot) \mapsto y(\cdot)^3$, & the embedding operator $B : L^2(\Omega) \rightarrow Y^*$, we can transform the above BVP into the equation $Ay + B\Phi(y) = Bv$ in Y^* .

In the following, we are going to demonstrate how Theorem 6.3 & Lemma 6.4 can be applied to a corresponding optimal control problem. To this end, we study the minimization of

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

subject to the above elliptic state problem & to the control constraint $-1 \leq v(x) \leq 1$ for almost every $x \in \Omega$. With the embedding operator $E_Y : H^1(\Omega) \rightarrow L^2(\Omega)$ & the admissible set

$$V_{\text{ad}} = \{v \in L^2(\Omega); -1 \leq v(x) \leq 1 \text{ for a.e. } x \in \Omega\},$$

we obtain the problem (6.25)

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

subject to (6.26)

$$Ay + B(\Phi(y) - v) = 0, \quad v \in V_{\text{ad}}.$$

Obviously, this is a special case of the problem (6.14), $\min J(u)$, $G(u) = 0$, $u \in C$, with the specifications $U := Y \times V$, $u := (y, v)$, $G : Y \rightarrow Y^*$, $G(u) := Ay + B(\Phi(y) - v)$, & $C := Y \times V_{\text{ad}}$.

1st-order necessary conditions. skipped Tröltzsch, 2010, pp. 336–337 ...

Remark 3.5. *In applying the general result Theorem 6.3 in function spaces, usually a compromise has to be made between 2 conflicting restraints: in order that the constraint qualification be valid that, at the same time, the nonlinearities be differentiable, the range space Z should not be too large; on the other hand, Z also should not be too small, since otherwise the dual space Z^* becomes too large in the sense that it contains functions of low regularity that can no longer be interpreted as (weakly differentiable) solutions to adjoint problems.*

²²**discrepancy** [n] (plural **discrepancies**) a difference between 2 or more things that should be the same.

2nd-order sufficient conditions. Since the functional J & the Neymyskii operator Φ are twice continuously Fréchet differentiable in $H^1(\Omega) \times L^2(\Omega)$ & $H^1(\Omega)$, respectively, so is the Lagrangian. Thus, in view of Lemma 6.4, the following condition is sufficient for local optimality: the pair (\bar{u}, \bar{v}) satisfies both the 1st-order necessary conditions & the definiteness condition

$$L''(\bar{y}, \bar{v}, p)[(y, v), (y, v)] \geq \delta \left(\|y\|_{H^1(\Omega)}^2 + \|v\|_{L^2(\Omega)}^2 \right)$$

for all pairs (y, v) satisfying the BVP

$$\begin{cases} -\Delta y + y + 3\bar{y}^2 y = v, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = 0, & \text{on } \Gamma. \end{cases}$$

Then \bar{v} is locally optimal in the sense of the norm of $L^2(\Omega)$. The above definiteness condition is already valid if we merely have, with a modified δ , $L''(\bar{y}, \bar{v}, p)[(y, v), (y, v)] \geq \delta \|v\|_{L^2(\Omega)}^2$. The explicit expression for the 2nd derivative L'' is

$$L''(\bar{y}, \bar{v}, p)[(y, v), (y, v)] = \|y\|_{L^2(\Omega)}^2 + \lambda \|v\|_{L^2(\Omega)}^2 - 6 \int_{\Omega} p \bar{y} y^2 \, dx.$$

3.6.2 Control Problems with State Constraints

“State constraints naturally arise in many applications. A typical example is that of heating problems in which the temperature is forbidden to exceed or fall short of certain prescribed threshold values. Such problems raise interesting, & in parts still unsolved, mathematical questions. Here, we shall only briefly address some basic ideas in order to enable the reader to consult the relevant literature for a more in-depth study. For a comprehensive treatment of the elliptic case, we refer the reader to Neittaanmäki et al. [NST06]. For simplicity, we confine ourselves to elliptic problems; the theory for parabolic problems is quite similar.

The necessary optimality conditions to be proved below may also be derived from Pontryagin’s maximum principle for state-constrained elliptic problems; for this purpose, the maximum condition is transformed into a variational inequality. In the case of boundary controls, the corresponding maximum principle was proved by Alibert & Raymond [AR97] & by Casas [Cas93]. The same applies to state-constrained parabolic problems, which were treated in Casas [Cas97] & in Raymond & Zidani [RZ99]. However, the proof of Pontryagin’s maximum principle is very technical, while the optimality conditions to be presented here can be obtained much more simply by means of the Lagrange method in Banach spaces. This technique was employed also in the works of Casas [Cas86] & Tröltzsch [Trö84b]. In the following, we apply it to derive 1st-order necessary conditions. We do not pursue 2nd-order sufficient conditions, referring the reader to the papers [CTU00] and [RT00]. Thus far, 2nd-order sufficient conditions for problems with pointwise state constraints in the whole domain could only be shown for low-dimensional domains; see Casas et al. [CDIRT08].” – Tröltzsch, 2010, Sect. 6.2, pp. 338–339

3.6.2.1 Convex problems

An elliptic problem with pointwise state constraints. Let $\Omega \subset \mathbb{R}^N$ denote a bounded Lipschitz domain. We consider the optimal control problem (6.28)

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

subject to (6.29)

$$\begin{cases} -\Delta y + y = u, & \text{in } \Omega, \\ \partial_{\mathbf{n}} y = 0, & \text{on } \Gamma, \end{cases}$$

& the constraints (6.30)

$$\begin{aligned} u_a(x) &\leq u(x) \leq u_b(x) \text{ for a.e. } x \in \Omega, \\ y(x) &\leq 0, \quad \forall x \in \bar{\Omega}. \end{aligned}$$

p. 339

3.7 Supplementary Results on PDEs

Chapter 4

Shape Optimization

4.1 Introduction

“Shape Optimization was introduced around 1970 by Jean Céa Céa, Gioan, and Michel, 1973, who understood, after several engineering studies [127, 12, 35, 110, 102, 83, 84, 7],¹ the future issues in the context of optimization problems. At that time, he proposed a list of open problems at the French National Colloquium in Numerical Analysis. These new problems were formulated in terms of minimization of functionals (referred as *open loop control* or *passive control*) governed by partial differential BVPs where the control variable was the geometry of a given boundary part [103, 76]. From the beginning, the terminology *shape optimization* was not connected to the structural mechanical sciences in which elasticity & optimization of the compliance played a central role. Furthermore, these research studies were mainly addressed in the context of the numerical analysis of the FEMs.

At the same time, there was some independent close results concerning fluid mechanics by young researchers e.g. O. Pironneau [123, 124, 78], Ph. Morice [107] & also several approaches related to perturbation theory by P.R. Garabedian [74, 75] & D.D Joseph [91, 92].

Very soon, it appeared that the shape control of BVPs was at the crossroads of several disciplines such as PDE analysis, non-autonomous semi-group theory, numerical approximation (including FEMs), control & optimization theory, geometry & even physics. Indeed several classical modeling in both structural & fluid mechanics (among other fields) needed to be extended. An illustrative example concerns a very *popular* problem in the 80’s concerning the thickness optimization of a plate modeled by the classical Kirchoff biharmonic equation. This kind of solid model is based on the assumption that the thickness undergoes only small variations. Therefore, many pioneering works were violating the validity of this assumption, leading to strange results, e.g., the work presented in the Iowa NATO Study [85] stating the existence of optimal beams having *zero cross section* values.

In the *branch* which followed the passive control approach, we shall mention the work of G. Chavent [32, 34] based on the theory of distributed system control introduced by J.-L. Lions Lions, 1971. Those results did not address optimization problems related to the domain but instead related to the coefficients inside the PDE. At that time, it was hoped that the solution of elliptic problems would be continuous w.r.t. the weak convergence of the coefficients. It appeared that this property was not achieved by this class of problem² At that point a main *bifurcation* arose with the homogenization approach [10] which up to some point was considered as a part of the *Optimal Design* theory.

The mathematical analysis of shape optimization problems began with the correct definition of derivatives of functionals & functions w.r.t. the domain, together with the choice of tangential space to the family of shapes. Following the very powerful theory developed by J. Nečas [117], the role of bilipschitzian mapping was emphasized for Sobolev spaces defined in moving domains based on the Identity perturbation method [115, 106, 134]. Concerning the large domain deformation viewpoint the previous approach led to the incremental domain evolution methods [143].

After 1975, the 2nd author introduced [145] an asymptotic analysis for domain evolution using classical geometrical flows which are intrinsic tools for manifolds evolutions & gave existence results for the so-called *shape differential equation* (see also [79]). At that period, applications focused more on sensitivity analysis problems than on asymptotic analysis of domains evolution. In 1972, A.M. Micheletti introduced in parallel [105, 104] a metric based on the Identity perturbation method thanks to the use of differentiable mappings, in order to study eigenvalues perturbation problems. The associated topology was extended by M. Delfour et. al [52] & turns out to be the same as the one induced by the continuity along flow field deformations [147].

¹(see refs. in Moubachir and Jean-Paul Zolésio, 2006)

²Indeed, in his thesis [33], G. Chavent referred to such a result to appear in a work by F. Murat [113]. That paper [111] appeared but as a counterexample to the expected continuity property. He showed on a 1D simple example that with weak *oscillating* convergence of the coefficients, the associated solution was converging to another problem in which the new coefficients were related to the limit of the *inverse* coefficients associated to the original problem [112, 114].

The systematic use of flow mapping & intrinsic geometry through the fundamental role of the oriented distance function [47, 50] led to the revised analysis of the elastic shell theory [48, 49, 25, 26, 27, 28], of the boundary layer theory [3] or of the manifold derivation tools [53].

The use of both Bounded Variation (BV) analysis & the notion of Cacciopoli sets led to the 1st compactness method for domain sequences & several extensions to more regular boundaries were done through the use of different concepts such as *fractal boundaries*, *density parameter* [23, 20, 21, 19] or *Sobolev domains* [50].

At that point, an other important *bifurcation point* in that theory occurred with the relaxation theory & the Special Bounded Variation (SBV) analysis which was particularly well adapted for image segmentation problem [6]. At the opposite, the capacity constraint for Dirichlet boundary conditions led to a fine analysis initiated in [18] & is still going on for cracks analysis.

The method of large evolution based on the flow mapping (known from 1980 as the *speed method* [150]) turns to be the natural setting for weak evolution of geometry allowing topological changes through the convection of either characteristic functions or oriented distance functions.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, pp. 1–3

4.1.1 Classical & moving shape analysis

“The classical shape analysis investigates the effects of perturbations of the geometry in terms of continuity, differentiability & optimization of quantities related to the state of a system defined in that geometry. In this case, the geometry is usually perturbed thanks to a map involving a scalar parameter usually referred to as a fictitious time. On the contrary, the moving shape analysis deals with systems that are intrinsically defined on a moving geometry. Hence, we shall deal with sensitivity analysis w.r.t. a continuous famil of shapes over a given time period. In this context, if we consider the geometry in a space-time configuration, the moving shape analysis may also be referred to as a *non-cylindrical shape analysis*³.

A 1st issue in this analysis is to model the evolution of the geometry. This is a common topic with the classical shape analysis. There exists many ways to build families of geometries. E.g., a domain can be made variable by considering its image by a family of diffeomorphisms parametrized by the time parameter as it happens frequently in mechanics for the evolution of continuous media. This way of defining the motion of domains avoids a priori the modification of the underlying topology. This change of topology can be allowed by using the characteristic function of families of sets or the level set of a space-time scalar function.” Refer to Delfour and J.-P. Zolésio, 2001; Delfour and J.-P. Zolésio, 2011 for a complete review on this topic. “In Chap. 2, we shall deal with the particular problem of defining in a weak manner the convection of a characteristic function in the context of the *speed method* developed in Zolésio’s PhD thesis [147].

In numbers of applications, we shall consider a state variable associated to a system which is a solution of a PDE defined inside the moving domain over a given time period. Hence, we need to analyze the solvability of this non-cylindrical PDE system before going further. Here, again this topic has been already studied since it enters the classical shape analysis problem while introducing a perturbed state defined in the moving domain parametrized by the fictitious time parameter. Furthermore, this solvability analysis has been performed in numbers of mathematical problems involving moving domains.” Refer to [135, 51] for some particular results in the context of the classical shape analysis. Also refer to the extensive literature concerning the analysis of PDE systems defined in moving domains, e.g., [96, 126, 132, 62, 130, 88, 128, 70, 100, 11].

“Contrary to the last topic, very few references exist for the sensitivity analysis w.r.t. the perturbation of the evolution of the moving geometry. Early studies have been conducted in [90, 151, 158, 141, 120, 43, 142, 2] for specific hyperbolic & parabolic linear problems. An important step was performed in [154, 155] where Zolésio established the derivative of integrals over a moving domain w.r.t. its associated Eulerian velocity. These results were applied in order to study variational principles for an elastic solid under large displacements & the incompressible Euler equation. This work was generalized in [58, 59].” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Sect. 1.1, pp. 3–4

4.1.2 Fluid-Structure Interaction Problems

“A general fluid–solid model consists of an elastic solid either surrounded by a fluid (aircrafts, automobiles, bridge decks, ...) or surrounding a fluid flow (pipelines, arteries, reservoir tanks, ...). Here the motion of the interface between the fluid & the solid is part of the unknown of the coupled system. It is a free boundary problem that can be solved by imposing continuity properties through the moving interface (e.g., the kinematic continuity of the velocities & the kinetic continuity of the normal stresses). This model has been intensively studied in the last 2 decades on the level of its mathematical solvability [87, 54, 82, 39, 80, 131, 15, 9, 138, 41], its numerical approximation [89, 55, 119, 118, 122, 95, 67], its stability [73, 38, 64, 65] & more recently on its controllability [66, 109]. In this lecture note, we will restrict ourselves to viscous Newtonian incompressible fluid flows described by the NSEs in space dimension 2 or 3. The case of a compressible Newtonian fluid can be incorporated in the present framework with the price of a heavier mathematical analysis (solvability, non-differentiability around shocks ...).”

³The notion of tube (non-cylindrical evolution domains) was also independently introduced by J.P. Aubin via the concept of *abstract mutations* [8].

Goal. “To solve inverse or control problems based on the previous general fluid-solid model. As an example, we think to decrease the drag of a car inside the atmospheric air flow by producing specific vibrations on its body using smart materials such as piezoelectrical layers. In this example, the control variable can be chosen as the electrical energy input evolution inside the piezoelectrical device & the objective is to decrease the drag which is a function of the coupled fluid-structure state (the air & the body of the car) & this state depends on the control variable. In order to build a control law for the electrical input, we need to characterize the relationship between the drag function & the control variable on the level of its computation & its variations.

As an other example, we can think of the problem of aeroelastic stability of structures. Both authors have been dealing with such a problem in the context of the stability analysis against wind loads of bridge decks. In [108], it has been suggested that such a problem can be set as the inverse problem consisting in recovering the smallest upstream wind speed that leads to the worst bridge deck vibrations. In this example, the decision variable can be chosen as the upstream wind speed & the objective is to increase a functional based on the vibration amplitude history of the bridge deck during a given characteristic time period which is a function of the coupled fluid-structure state (the wind flow & the bridge deck) which is also a function of the decision variable. Again, in order to recover the wind speed history, we need to characterize the relationship between the objective functional & the decision variable on the level of its computation & its variations.

In order to characterize the sensitivity of the objective functional w.r.t. the control variable, it is obvious that we need to characterize the sensitivity of the coupled fluid-structure state w.r.t. the control variable. Here we recall that the coupled fluid-structure state is the solution of a system of PDEs that are coupled through continuity relations defined on the moving interface (the fluid-structure interface). The key point towards this sensitivity analysis is to investigate the sensitivity of the fluid state, which is an Eulerian quantity, w.r.t. the motion of the solid, which is a Lagrangian quantity. This task falls inside the moving shape analysis framework described earlier. Indeed the fluid state is the solution of system of nonlinear PDEs defined in a moving domain. The boundary of this moving domain is the solid wall. Then using the tools developed in [59], it has been possible to perform in [58] the moving shape sensitivity analysis in the case of a Newtonian incompressible fluid inside a moving domain driven by the non-cylindrical NSEs.

All the previous results use a parametrization of the moving domain based on the Lagrangian flow of a given velocity field. Hence, the design variable is the Eulerian velocity of the moving domain, allowing topology changes while using the associated level set formulation. In [13, 14], the author used a non-cylindrical identity perturbation technique. It consists in perturbing the space-time identity operator by a family of diffeomorphism. Then, this family is chosen as the design parameter. It is a Lagrangian description of the moving geometry, which a priori does not allow topology changes but which leads to simpler sensitivity analysis results which are still comparable with the one obtained by the non-cylindrical *speed method*. In [57], the authors came back to the dynamical shape control of the Navier-Stokes & recovered the results obtained in [58] using the Min-Max principle allowing to avoid the state differentiation step w.r.t. the velocity of the domain.

Now, we come back to the original problem consisting in the sensitivity analysis of the coupled fluid-structure state w.r.t. the control variable. Using the chain rule, the derivative of the coupled state w.r.t. the control variable involves the partial derivative of the fluid state w.r.t. the motion of the fluid-structure interface already characterized in [58, 57]. Hence, again using a Lagrangian penalization technique, already used & justified in [45, 46], it has been possible to perform in [109] the sensitivity analysis of a simple fluid-structure interaction problem involving a rigid solid within an incompressible flow of a Newtonian fluid w.r.t. the upstream velocity field. As already mentioned, this simple model is particularly suited for bridge deck aeroelastic stability analysis [121].” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Sect. 1.2, pp. 4–6

- “Moubachir and Jean-Paul Zolésio, 2006, Chap. 2 furnishes a simple illustration to some of the moving shape analysis results reported in the core of the lecture note. We deal with a simple inverse problem arising in phase change problems consisting in recovering the moving interface at the isothermal interface between a solid & liquid phase from measurements of the temperature on a insulated fixed part of the solid boundary. We use a least-square approach & we show how to compute the gradient of the least-square functional w.r.t. the velocity of the moving interface. It involves an adjoint state problem together with an adjoint transverse state, which is the novelty of the moving shape analysis compared to the classical one.
- In Moubachir and Jean-Paul Zolésio, 2006, Chap. 3, we consider the weak Eulerian evolution of domains through the convection, generated by a non-smooth vector field \mathbf{V} , of measurable sets. The introduction of transverse variations enables the derivation of functionals associated to evolution tubes. We also introduce Eulerian variational formulations for the minimal curve problem. These formulations involve a geometrical adjoint state λ which is backward in time & is obtained thanks to the use of the so-called *transverse field* \mathbf{Z} .
- In Moubachir and Jean-Paul Zolésio, 2006, Chap. 4, we recall the concept of shape differential equation developed in [145, 147]. Here, we present a simplified version & some applications in 2D which enable us to reach the time asymptotic result. Furthermore, we introduce the associated level set formulation whose speed vector version was already contained in [149].
- In Moubachir and Jean-Paul Zolésio, 2006, Chap. 5, we deal with a challenging problem in fluid mechanics which consists in the control of a Newtonian fluid flow thanks to the velocity evolution law of a moving wall. Here, the

optimal control problem has to be understood as the open loop version, i.e., it consists in minimizing a given objective functional w.r.t. the velocity of the moving wall. This study is performed within the non-cylindrical Eulerian moving shape analysis described in Chaps. 2–3. We focus on the use of a Lagrangian penalization formulation in order to avoid the fluid state differentiation step.

- In Moubachir and Jean-Paul Zolésio, 2006, Chap. 6, we introduce the Lagrangian moving shape analysis framework. It differs from the Eulerian one from the fact that the design variable is the diffeomorphism that parametrizes the moving geometry. The sensitivity analysis is simpler since it does not involve the transverse velocity field. We apply these tools in order to deal with the control of a Newtonian fluid flow thanks to the displacement evolution law of a moving wall.
- Moubachir and Jean-Paul Zolésio, 2006, Chap. 7 moves to inverse problems related to fluid-structure interaction systems. Here, we consider a 2D elastic solid with rigid displacements inside the incompressible flow of a viscous Newtonian fluid. We try to recover informations about the inflow velocity field from the partial measurements of the coupled fluid-structure state. We use a least-square approach together with a Lagrangian penalization technique. We derive the structure of the gradient w.r.t. the inflow velocity field of a given cost function. Using the Min-Max principle, the cost function gradient reduces to the derivative of the Lagrangian w.r.t. the inflow velocity at the saddle point. This saddle point is solution of 1st order optimality conditions. We use non-cylindrical Eulerian derivatives to compute the partial derivative of the Lagrangian functional w.r.t. the solid state variables, involved in the optimality system.
- In Moubachir and Jean-Paul Zolésio, 2006, Chap. 8 we extend the results of Chap. 7, to the case of an elastic solid under large displacements inside an incompressible fluid flow. The main difference with the previous case is the use of a non-cylindrical Lagrangian shape analysis for establishing the KKT system. It forms the adjoint counterpart of the sensitivity analysis conducted in [66].” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Sect. 1.3, pp. 6–7

“... we shall describe the different steps encountered while designing a complex fluid-structure interaction system. Indeed, let us consider a mechanical system that consists of a solid & a fluid interacting with each other. We would like to increase the performances of this system. These performances have to be quantitatively translated inside a cost function that we have to optimize w.r.t. some parameters that we will call the control variables. In the sequel, we will describe different control situations:

1. *Control of a fluid flow around a fixed body*: it consists in trying to modify the fluid flow pattern around a fixed body using a boundary control which can act e.g. by blowing or suctioning the fluid at some part of the solid boundary. The control law will be designed in order to match some efficiency goals using the minimization of a cost functional.
2. *Shape design of a fixed solid inside a fluid flow*: in this case, the control is the shape of the body. We would like to find the best shape satisfying some geometrical constraints that will optimize some cost functionals. This problem is somewhat classical in the aeronautical field, but it requires some subtle mathematical tools that we will quickly recall.
3. *Dynamical shape design of a solid inside a fluid flow*: the novelty compared to the last item is that the shape is moving & we are looking for the best evolution of this shape that both satisfies some geometrical constraints & optimizes some cost functionals. This is a rather natural technique in order to control a fluid flow pattern, but still its design requires some new mathematical tools that will be sketched in this introduction & more detailed in the core of this lecture note.
4. *Control of an elastic solid inside a fluid flow*: this is the most complex & most realistic situation where both the fluid & the solid have their own dynamics which are coupled through the fluid-solid interface. Then, we would like to control or optimize the behavior of this coupled system thanks to boundary conditions. The mathematical analysis of this situation uses the whole framework introduced previously. This is a challenging problem, both on the mathematical point of view & on the technological side. The goal of this book is to partially answer to some issues related to this problem.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Sect. 1.4, pp. 7–8

4.1.2.1 Control of a fluid flow around a fixed body

The objective functional. “A common topic in the optimization & control field of PDE systems is the choice of appropriate cost functionals, i.e., meeting both our objectives & the mathematical requirements that guarantee the convergence to at least 1 optimum parameter. This functional can depend both on the state variables (\mathbf{u}, p) & on the control parameter \mathbf{g} .” [...] “More generally, we can consider any cost functionals that are twice-differentiable w.r.t. their arguments.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Subsect. 1.4.1, pp. 9–10

The control problem. “Our goal is now furnish the 1st-order optimality conditions associated to the optimization problem. These conditions are very useful since they are the basis in order to build both a rigorous mathematical analysis & gradient-based optimization algorithms.

There exists 2 main methods in order to derive these conditions: the 1st one is based on the differentiability of the state variables w.r.t. the control parameter & the 2nd one relies on the existence of Lagrangian multipliers.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Subsect. 1.4.1, p. 10

Sensitivity. “Let us consider a *control point* $\mathbf{g} \in \mathcal{U}$, then the cost functional $j(\mathbf{g})$ is Fréchet differentiable w.r.t. \mathbf{g} Abergel and Temam, 1990, [71] & its directional derivative is given by (1.9)

$$\langle j'(\mathbf{g}), \mathbf{h} \rangle = \langle \partial_{(\mathbf{u}, p)} J[(\mathbf{u}, p)(\mathbf{g})], (\mathbf{u}', p')(\mathbf{g}; \mathbf{h}) \rangle, \text{ where } (\mathbf{u}', p')(\mathbf{g}; \mathbf{h}) := \frac{d}{d\mathbf{g}}(\mathbf{u}, p)(\mathbf{g}) \cdot \mathbf{h}$$

stands for the directional derivative of $(\mathbf{u}, p)(\mathbf{g})$ w.r.t. \mathbf{g} .” [...]

“Then the 1st-order optimality condition writes (1.11) $\langle j'(\mathbf{g}), \mathbf{h} \rangle = 0, \forall \mathbf{h} \in \mathcal{U}$. I.e., the set of optimal controls is contained in the set of critical points for the cost function $j(\mathbf{g})$. However, we would like to obtain an expression of this condition avoiding the direction $\mathbf{h} \in \mathcal{U}$. To this end, we introduce the *adjoint variable* (\mathbf{v}, π) solution of the *adjoint linearized Navier–Stokes system* (1.12). Consequently, we are able to identify the gradient of the cost function as the trace on Γ^c of the *adjoint normal stress tensor*, i.e., (1.13)

$$\nabla j(\mathbf{g}) = {}^* \gamma_{(0, \tau) \times \Gamma^c} [\sigma(\mathbf{v}, \pi) \cdot \mathbf{n}].$$

This formal proof provides the basic steps needed in order to build a gradient-based optimization method associated to the control problem $\min_{\mathbf{g} \in \mathcal{U}} j(\mathbf{g})$.

An alternative approach consists in avoiding the derivation of the fluid state (\mathbf{u}, p) w.r.t. the control \mathbf{g} thanks to the introduction of a Lagrangian functional that includes not only the cost functional but also the state equation, (1.14)

$$\mathcal{L}(\psi, r, \phi, q; \mathbf{g}) = J(\psi, r) + \langle e(\psi, r; \mathbf{g}), (\phi, q) \rangle,$$

where $\langle e(\mathbf{u}, p; \mathbf{g}), (\phi, q) \rangle$ stands for the weak form of the state equation NSEs (1.8), e.g.,

$$\begin{aligned} \langle e(\psi, r; \mathbf{g}), (\phi, q) \rangle &= \int_{(0, \tau) \times \Omega_f} [-\psi \cdot \partial_t \phi + (D\psi \cdot \psi) \cdot \phi - \nu \psi \cdot \Delta \phi + \psi \cdot \nabla q - r \nabla \cdot \phi] + \int_{(0, \tau) \times \Gamma^c} \mathbf{g} \cdot (\sigma(\phi, q) \cdot \mathbf{n}) d\Gamma dt \\ &+ \int_{(0, \tau) \times \partial D} \mathbf{u}_\infty \cdot (\sigma(\phi, q) \cdot \mathbf{n}) + \int_{\Omega_f} \psi(\tau) \cdot \phi(\tau) - \int_{\Omega_f} \mathbf{u}_0 \cdot \phi(0). \end{aligned}$$

Hence the control problem $\min_{\mathbf{g} \in \mathcal{U}} j(\mathbf{g})$ is equivalent to the min-max problem, (1.15)

$$\min_{\mathbf{g} \in \mathcal{U}} \min_{(\psi, r)} \max_{(\phi, q)} \mathcal{L}(\psi, r, \phi, q; \mathbf{g}).$$

For every control $\mathbf{g} \in \mathcal{U}$, it can be proven that the min-max problem,

$$\min_{(\psi, r)} \max_{(\phi, q)} \mathcal{L}(\psi, r, \phi, q; \mathbf{g})$$

admits a unique saddle-point $(\mathbf{u}, p; \mathbf{v}, \pi)$ which are solutions of the systems (1.8)–(1.12). Finally the 1st-order optimality for the problem (1.15) writes (1.16)

$$\partial_{\mathbf{g}} \mathcal{L}(\mathbf{u}, p, \mathbf{v}, \pi; \mathbf{g}) = 0$$

which turns out to be equivalent to (1.13). Then, we can think to solve the optimality condition (1.11), using a continuous iterative method. Indeed let us introduce a scalar parameter $s \geq 0$, & a *control variable* $\mathbf{g}(s)$ that is differentiable w.r.t. s . Hence using the differentiability of $J(\mathbf{g})$, we get

$$J(\mathbf{g}(r)) - J(\mathbf{g}(0)) = \int_0^r \langle \nabla J(\mathbf{g}(s)), \mathbf{g}'(s) \rangle_{\mathcal{U}^*, \mathcal{U}} ds.$$

Let us choose the control s.t. (1.17)

$$\mathbf{g}'(s) + \mathbb{A}^{-1}(s) \nabla J(\mathbf{g}(s)) = 0, \quad s \in (0, r),$$

where \mathbb{A} stands for an appropriate duality operator, then the functional writes

$$J(\mathbf{g}(r)) - J(\mathbf{g}(0)) = - \int_0^r |\nabla J(\mathbf{g}(s))|^2 ds.$$

I.e., the control law (1.17) leads to a functional’s decrease & is referred to as a continuous gradient based optimization method. Using a discretization of the parameter s leads to a standard gradient-based method such as the conjugate-gradient or the quasi-Newton method depending on the choice of $\mathbb{A}(s)$.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Subsect. 1.4.1, pp. 11–12

4.1.2.2 Shape design of a fixed solid inside a fluid flow

“We again consider the situation where a fixed solid is surrounded by a fluid flow. The shape control consists in finding the optimal shape of the solid that reduces some objective functional (e.g., the drag) under some perimeter, volume or curvature constraints. This optimization is an open-loop control since the shape of the obstacle is time-independent.”

The speed method. “Here the space of shapes is no more a linear space & the associated differential calculus becomes more tricky. Our goal is to build *gradient-based methods* in order to find the optimal shape, i.e., we would like to solve the following problem, (1.18) $\min_{\Omega \in \mathcal{A}} J(\Omega)$. In order to carry out the sensitivity analysis of functionals depending on the shape of the solid Ω , we introduce a family of pertubated domains $\Omega_s \subset D$ parametrized by a scholar parameter $0 \leq s \leq \varepsilon$. These domains are the images of the original domain Ω through a given family of smooth maps⁴ $\mathbf{T}_s : \overline{D} \rightarrow \overline{D}$, i.e. $\Omega_s = \mathbf{T}_s(\Omega)$, $\Gamma_s = \mathbf{T}_s(\Gamma)$.

2 major classes of such mappings are given by:

- the *identity perturbation method* ([116, 125]), $\mathbf{T}_s = \mathbf{I} + s\boldsymbol{\theta}$, where $\boldsymbol{\theta} : \overline{D} \rightarrow \overline{D}$.
- the *speed method* [145] Pironneau, 1984, where the transformation is the flow associated to a given velocity field $\mathbf{V}(s, \mathbf{x})$,

$$\begin{cases} \partial_s \mathbf{T}_s(\mathbf{x}) = \mathbf{V}(s, \mathbf{T}_s(\mathbf{x})), & (s, \mathbf{x}) \in (0, \varepsilon) \times D, \\ \mathbf{T}_{s=0}(\mathbf{x}) = \mathbf{x}, & \mathbf{x} \in D. \end{cases}$$

In order for \overline{D} to be globally invariant under $\mathbf{T}_s(\mathbf{V})$ we need to impose the following *viability conditions*, $\mathbf{V}(s, \mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) = 0$, $\mathbf{x} \in \partial D$.

Let us consider the family of functionals $J(\Omega_s)$ that depends on the shapes Ω_s , e.g., the work to overcome the drag exerted by the fluid on the solid boundary, (1.19)

$$J_{\text{drag}}(\Omega) = \int_{(0, \tau) \times \Gamma} (\mathbf{u} - \mathbf{u}_\infty) \cdot \sigma(\mathbf{u}, p) \cdot \mathbf{n} \, d\Gamma \, dt.$$

This functional depends on Ω not only because it is an integral over the boundary Γ , but also because it involves the solution (\mathbf{u}, p) of the Navier–Stokes system, (1.20), that depends on Ω .

To perform our sensitivity analysis, we choose to work in the framework of the *speed method*⁵ We define the Eulerian derivative of the shape functional $J(\Omega)$ at point Ω in the direction of the vector field $\mathbf{V} \in \mathcal{V}$ as the limit,

$$dJ(\Omega; \mathbf{V}) = \lim_{s \downarrow 0} \frac{J(\Omega_s(\mathbf{V})) - J(\Omega)}{s},$$

where \mathcal{V} is a linear space⁶ If this limit exists & is finite $\forall \mathbf{V} \in \mathcal{V}$ & the mapping $\mathcal{V} \rightarrow \mathbb{R}$, $\mathbf{V} \mapsto dJ(\Omega; \mathbf{V})$ is linear & continuous, then the functional $J(\Omega)$ is said to be *shape differentiable*.

Actually if $J(\Omega)$ is shape differentiable, then its Eulerian derivative only depends on $\mathbf{V}(0)$ & there exists a distribution $\mathbf{G}(\Omega) \in \mathcal{D}(D; \mathbb{R}^3)'$ that we call the *shape gradient* s.t.

$$dJ(\Omega; \mathbf{V}) = \langle \mathbf{G}(\Omega), \mathbf{V}(0) \rangle, \quad \forall \mathbf{V} \in \mathcal{V}.$$

In the sequel, we shall use the notation $\nabla J(\Omega) := \mathbf{G}(\Omega)$. In the case of smooth domain, the gradient is only supported on the boundary Γ & depends linearly on the normal vector field \mathbf{n} . This result, called the *structure theorem*⁷, is recalled as follows,

Theorem 4.1 (Shape derivative structure theorem). *Let $J(\cdot)$ be a differentiable shape functional at every shape Ω of class \mathcal{C}^{k+1} for $k \geq 0$ with shape gradient $\mathbf{G}(\Omega) \in \mathcal{D}(D; \mathbb{R}^3)'$. In this case, the shape gradient has the following representation,*

$$\mathbf{G}(\Omega) = {}^* \gamma_\Gamma(g\mathbf{n}),$$

where $g(\Gamma) \in \mathcal{D}^{-k}(\Gamma)$ stands for a scalar distribution & ${}^* \gamma_\Gamma$ stands for the adjoint trace operator⁸.

⁴Typically we have the following Lipschitz regularity assumptions:

$$\begin{aligned} \mathbf{T}(\cdot, \mathbf{x}) &\in \mathcal{C}^1([0, \varepsilon]; \mathbb{R}^3), \quad \forall \mathbf{x} \in D, \quad \|\mathbf{T}(\cdot, \mathbf{x}) - \mathbf{T}(\cdot, \mathbf{y})\|_{\mathcal{C}^0([0, \varepsilon]; \mathbb{R}^3)} \leq C \|\mathbf{x} - \mathbf{y}\|_{\mathbb{R}^3}, \\ \mathbf{T}^{-1}(\cdot, \mathbf{x}) &\in \mathcal{C}^0([0, \varepsilon]; \mathbb{R}^3), \quad \forall \mathbf{x} \in D, \quad \|\mathbf{T}^{-1}(\cdot, \mathbf{x}) - \mathbf{T}^{-1}(\cdot, \mathbf{y})\|_{\mathcal{C}^0([0, \varepsilon]; \mathbb{R}^3)} \leq C \|\mathbf{x} - \mathbf{y}\|_{\mathbb{R}^3}, \end{aligned}$$

where $C > 0$.

⁵which leads, at least for the 1st order terms, to the same results as the identity perturbation framework [51].

⁶e.g., $\mathcal{V} := \{\mathbf{V} \in \mathcal{C}^0(0, \varepsilon; \mathcal{C}^1(D; \mathbb{R}^3)), \nabla \cdot \mathbf{V} = 0 \text{ in } D, \langle \mathbf{V}, \mathbf{n} \rangle = 0 \text{ on } \partial D\}$.

⁷Refer to Delfour and J.-P. Zolésio, 2001, Theorem 3.5 for the case of non-smooth domains.

⁸i.e.,

$$\langle {}^* \gamma_\Gamma(g\mathbf{n}), \mathbf{V} \rangle_{\mathcal{D}(D; \mathbb{R}^3)', \mathcal{D}(D; \mathbb{R}^3)} = \langle g, \mathbf{V} \cdot \mathbf{n} \rangle_{\mathcal{D}'(\Gamma), \mathcal{D}(\Gamma)}.$$

This result is easier to understand when $g(\Gamma)$ is integrable over Γ , that is to say, $g \in L^1(\Gamma)$. Indeed in this case, it means that the directional shape derivative can always be written as follows,

$$dJ(\Omega; \mathbf{V}) = \int_{\Gamma} g \mathbf{V} \cdot \mathbf{n} d\Gamma.$$

Basic shape derivative calculus. The notion of Eulerian derivative for shape functionals can be extended to functions defined on Banach or Hilbert spaces built on smooth domains Ω . Hence, a function $y \in H(\Omega)$ admits a *material derivative* at Ω in the direction $\mathbf{V} \in \mathcal{V}$ if the following limit

$$\dot{y}(\Omega; \mathbf{V}) := \lim_{s \rightarrow 0} \frac{1}{s} [y(\Omega_s(\mathbf{V})) \circ \mathbf{T}_s(\mathbf{V}) - y(\Omega)]$$

admits a limit in the Hilbert space⁹ $H(\Omega)$.

Endowed with the following definition, it is possible to derive the Eulerian shape derivative of the following functionals,

$$J(\Omega) = \int_{\Omega} y(\Omega) d\Omega.$$

If $y(\Omega)$ is weakly shape differentiable in $L^1(\Omega)$, then the functional $J(\Omega)$ is shape differentiable & its directional derivative writes,

$$dJ(\Omega; \mathbf{V}) = \int_{\Omega} [\dot{y}(\Omega; \mathbf{V}) + y(\Omega) \nabla \cdot \mathbf{V}(0)] d\Omega.$$

In order to apply the structure theorem, it is useful to define the notion of shape derivative for functions. Hence, if $y \in H(\Omega)$ admits a material derivative $\dot{y}(\Omega; \mathbf{V}) \in H(\Omega)$ & $\nabla y \cdot \mathbf{V}(0) \in H(\Omega)$ for all $\mathbf{V} \in \mathcal{V}$, we define the *shape derivative* as

$$y'(\Omega; \mathbf{V}) = \dot{y}(\Omega; \mathbf{V}) - \nabla y(\Omega) \cdot \mathbf{V}(0).$$

In this case, the Eulerian shape derivative of $J(\Omega)$ takes the following form,

$$dJ(\Omega; \mathbf{V}) = \int_{\Omega} [y'(\Omega; \mathbf{V}) + \nabla \cdot (y(\Omega) \mathbf{V}(0))] d\Omega.$$

If Ω is class \mathcal{C}^k with $k \geq 1$, then using the Stokes formula, we get

$$dJ(\Omega; \mathbf{V}) = \int_{\Omega} y'(\Omega; \mathbf{V}) + \int_{\Gamma} y(\Omega) \mathbf{V} \cdot \mathbf{n} d\Gamma.$$

Remark 4.1. In the case where $y(\Omega) = Y|_{\Omega}$, where $Y \in H(D)$ with $\Omega \subset D$, its shape derivative is zero since $\dot{y}(\Omega; \mathbf{V}) = \nabla Y \cdot \mathbf{V}$. Hence,

$$dJ(\Omega; \mathbf{V}) = \int_{\Gamma} y(\Omega) \mathbf{V} \cdot \mathbf{n} d\Gamma.$$

This is a simple illustration of the structure theorem.

In the case of functionals involving integration over the boundary Γ , we need to introduce the notion of *material derivative* on Γ .

Let $z \in W(\Gamma)$ where $W(\Gamma)$ is an Hilbert space of functions (e.g., $W^{m,p}(\Gamma)$) defined over Γ . It is said that it admits a material derivative in the direction $\mathbf{V} \in \mathcal{V}$, if the following limit,

$$\dot{z}(\Gamma; \mathbf{V}) := \lim_{s \rightarrow 0} \frac{1}{s} [z(\Gamma_s(\mathbf{V})) \circ \mathbf{T}_s(\mathbf{V}) - z(\Gamma)]$$

admits a limit in the Hilbert space $W(\Gamma)$.

As a consequence of this definition, it is possible to derive the Eulerian shape derivative of the following functional,

$$J(\Gamma) = \int_{\Gamma} z(\Gamma) d\Gamma.$$

⁹e.g., $W^{m,p}(\Omega)$ or $L^2(0, \tau; W^{m,p}(\Omega))$.

If $z(\Gamma)$ is weakly shape differentiable in $L^1(\Omega)$, then the functional $J(\Gamma)$ is shape differentiable & its directional derivative writes

$$dJ(\Gamma; \mathbf{V}) = \int_{\Gamma} [\dot{z}(\Gamma; \mathbf{V}) + z(\Gamma) \operatorname{div}_{\Gamma} \mathbf{V}(0)] d\Gamma,$$

where

$$\operatorname{div}_{\Gamma} \mathbf{V} := \gamma_{\Gamma}[\nabla \cdot \mathbf{V} - (\mathbf{D}\mathbf{V} \cdot \mathbf{n}) \cdot \mathbf{n}]$$

stands for the *tangential divergence*.

As in the previous case, it is also possible to introduce the notion of shape derivative for $z(\Gamma)$. Let Ω be of class \mathcal{C}^k with $k \geq 2$. If $z \in W(\gamma)$ admits a material derivative $\dot{z}(\Gamma; \mathbf{V}) \in W(\Gamma)$ & $\nabla_{\Gamma} y \cdot \mathbf{V}(0) \in W(\Gamma)$ for all $\mathbf{V} \in \mathcal{V}$, we define the *shape derivative* as

$$z'(\Gamma; \mathbf{V}) = \dot{z}(\Gamma; \mathbf{V}) - \nabla_{\Gamma} z(\Gamma) \cdot \mathbf{V}(0) \text{ where } \nabla_{\Gamma} z = \nabla Z|_{\Gamma} - (\nabla Z \cdot \mathbf{n})\mathbf{n}$$

stands for the *tangential gradient* & Z is any smooth extension of z inside Ω .

Using the above definition, it is possible to transform the expression of the differential as follows,

$$dJ(\Gamma; \mathbf{V}) = \int_{\Gamma} [z'(\Gamma; \mathbf{V}) + H z(\Gamma) \mathbf{V}(0) \cdot \mathbf{n}] d\Gamma,$$

where H stands for the *mean curvature* of Γ .

Remark 4.2. In the case where $z(\Gamma) = y(\Omega)|_{\Gamma}$, the Eulerian derivative takes the following form,

$$dJ(\Gamma; \mathbf{V}) = \int_{\Gamma} [y'(\Omega; \mathbf{V}) + (\nabla y(\Omega) \cdot \mathbf{n} + H z(\Gamma) \mathbf{V}(0) \cdot \mathbf{n})] d\Gamma.$$

Application to shape design. Thanks to the framework introduced previously, it is possible to build a complete sensitivity analysis of shape functionals. Coming back to our optimal shape problem, we can state the following: the shape gradient for the tracking functional

$$J(\Omega) = \int_{(0,\tau) \times \Omega} (\mathbf{u} - \mathbf{u}_d)^2 d\mathbf{x} dt$$

is given by (1.21)

$$\nabla J(\Omega) = {}^* \gamma_{\Gamma}[\sigma(\mathbf{v}, \pi) \cdot \mathbf{n}]$$

where (\mathbf{u}, p) is a solution of system (1.20) associated to the shape Ω & (\mathbf{v}, π) is solution of the adjoint system (1.22).

The associated shape differential equation. Now as in the previous section, we can choose to solve the 1st-order optimality equation (1.21) using a continuous gradient-based method. I.e., we write

$$J(\Omega_r(\mathbf{V})) - J(\Omega_0) = \int_0^r \langle \nabla J(\Omega_s(\mathbf{V})), \mathbf{V}(s) \rangle ds.$$

Then solving the equation (1.23)

$$\nabla J(\mathbf{V}(s)) + \mathbb{A}^{-1}(s) \cdot \mathbf{V}(s) = 0, \quad s \in (0, +\infty)$$

leads to a decrease of the functional $J(\Omega_s(\mathbf{V}))$. The equation (1.23) is referred to as the *shape differential equation* & some of its properties are studied in Moubachir and Jean-Paul Zolésio, 2006, Chap. 4. Notably, we study its solvability in the case of smooth shape functionals. We also prove some results concerning the asymptotic behavior of the solution of this equation, which hold essentially when the shape gradient has some continuity properties for an *ad-hoc* shape topology¹⁰.

¹⁰The Hausdorff-complementary topology.

The level-set framework. In Chap. 4, we also relate the shape differential equation to the Hamilton–Jacobi equation involved in the level-set setting. The level-set setting consists in parametrizing the perturbed domain Ω_s as the positiveness set of a scalar function $\Phi : (0, \varepsilon) \times \overline{D} \rightarrow \mathbb{R}$,

$$\Omega_s = \Omega_s(\Phi) := \{\mathbf{x} \in D, \Phi(s, \mathbf{x}) > 0\},$$

& its boundary is the zero-level set,

$$\Gamma_s = \Gamma_s(\Phi) := \{\mathbf{x} \in D, \Phi(s, \mathbf{x}) = 0\}.$$

This parametrization & the one introduced in the *speed method* can be linked thanks to the following identity,

$$\mathbf{V}(s) = -\partial_s \Phi(s) \frac{\nabla \Phi(s)}{\|\nabla \Phi(s)\|^2}.$$

Both frameworks are equivalent if $\Phi(s)$ belongs to set of functions without steps, i.e., $\|\nabla \Phi(s)\|$ is different from zero a.e. in D . We show how to build without step functions & we study the shape differential equation in this setting.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Subsect. 1.4.2, pp. 13–19

4.1.2.3 Dynamical shape design of a solid inside a fluid flow

“... we consider that the shape of the solid is moving.” **Goals.** to control this motion in order to optimize some objective functionals; to build gradient-based methods in order to find the optimal shape dynamic, i.e., we would like to solve the following problem, (1.24) $\min_{Q \in \mathcal{E}} J(Q)$ where $Q \in \mathcal{E}$ is a smooth evolution set, which means

$$Q := \bigcup_{t \in (0, \tau)} \{t\} \times \Omega_t,$$

where Ω_t is a smooth domain of \mathbb{R}^3 with boundary Γ_t . The set,

$$\Sigma := \bigcup_{t \in (0, \tau)} \{t\} \times \Gamma_t$$

stands for the *non-cylindrical lateral boundary*. We call the set Q a *tube*.

The \mathbb{R}^{d+1} -approach. “The optimal control of moving domain is a problem which is relevant of classical (nonlinear) control theory as well as of classical shape theory. In fact on the pure theoretical level, the dynamical shape control theory can be viewed as an application of the shape optimization theory for space-time manifolds.” Fig. Non-cylindrical space-time domain. “Indeed the dynamical shape control consists in finding the optimal evolution of a spatial domain Ω_t in \mathbb{R}^d . Let us consider the mapping $\mathcal{S} : t \in \mathbb{R} \rightarrow \Omega_t \in \mathcal{P}(\mathbb{R}^d)$ where $\mathcal{P}(\mathbb{R}^d)$ stands for the set of parts inside \mathbb{R}^d . Usually we would like to minimize some cost functional,

$$j(\mathcal{S}) = \int_0^\tau J(t, \mathcal{S}(t)) dt.$$

Obviously, it is equivalent to the problem of finding the optimal tube

$$Q = \bigcup_{0 < t < \tau} \{t\} \times \Omega_t \in \mathbb{R}^{d+1}.$$

In fact the tube Q is the graph in $\mathbb{R} \times \mathcal{P}(\mathbb{R}^d) \subset \mathbb{R}^{d+1}$ of the shape mapping \mathcal{S} . As for usual mappings defined from \mathbb{R} in some space E , the graph $G \subset \mathbb{R} \times E$ & of course any subset G is not a graph. Now under simple conditions on that set G , it becomes a graph. In the same way, any subset $Q \in \mathbb{R} \times \mathcal{P}(\mathbb{R}^d)$ will not be a tube. Intuitively we would say that we require some *causality* in the evolution of the set Ω_t .

When the boundary of the set Ω_t is smooth enough¹¹, the idea is to avoid the normal field ν to the lateral boundary Σ of the tube to be strictly vertical. To handle non-smooth situations, we adopt an Eulerian viewpoint that associates to each tube Q the non-empty closed convex set of speed vector fields \mathbf{V} which transport (in a weak sense) the characteristic function of the moving domain. When we consider the tube Q as a subset of \mathbb{R}^{d+1} , the control problems becomes a usual shape optimization problem (as far as no real time consideration enters). The sensitivity analysis is then classically performed by considering *horizontal* vector fields $\tilde{\mathbf{Z}}(s, t, \mathbf{x}) = (0, \mathbf{Z}(s, t, \mathbf{x})) \in \mathbb{R}^{d+1}$ where s is the *perturbation parameter* of the tube.

¹¹say there exists a tangent space.

Then the $(d+1)$ -dimensional shape optimization analysis fully applies & the so-called *Shape Differential Equation* furnishes descent direction, i.e., it furnishes the existence of a vector field \mathbf{Z}^* s.t., for some $\alpha > 0$,

$$J(Q_s) \leq J(Q) - \alpha \int_0^s \|\mathbf{Z}^*(\sigma)\|^2 d\sigma, \quad \forall s > 0.$$

We show that the existence of that field \mathbf{Z}^* induces the existence of a usual vector field $\mathbf{V}(t, \mathbf{x}) \in \mathbb{R}^d$ which builds that tube, i.e., $\Omega_t = \mathbf{T}(\mathbf{V})(\Omega_0)$.

The \mathbb{R}^d -approach. In order to carry out the sensitivity analysis of functionals depending on the tube Q , we assume that the domains are the images of the domain $\Omega_0 := \Omega_{t=0}$ through a given family of smooth maps $\mathbf{T}_t : \bar{D} \rightarrow \bar{D}$, i.e., $\Omega_t = \mathbf{T}_t(\Omega_0)$, $\Gamma_t = \mathbf{T}_t(\Gamma_0)$. 2 major class of such mappings are given by:

- the *Lagrangian parametrization* $\mathbf{T}_t = \boldsymbol{\theta}(t, \cdot)$ where $\boldsymbol{\theta} : (0, \tau) \times \bar{D} \rightarrow \bar{D}$. In this case, the minimization problem (1.24) can be transformed as **(1.25)** $\min_{\boldsymbol{\theta} \in \Theta} J(Q(\boldsymbol{\theta}))$.
- the *Eulerian parametrization*, where the transformation is the flow associated to a given velocity field $\mathbf{V}(t, \mathbf{x})$,

$$\begin{cases} \partial_t \mathbf{T}_t(\mathbf{x}) = \mathbf{V}(t, \mathbf{T}_t(\mathbf{x})), & (t, \mathbf{x}) \in (0, \tau) \times D, \\ \mathbf{T}_{t=0}(\mathbf{x}) = \mathbf{x}, & \mathbf{x} \in D. \end{cases}$$

In this case, the minimization problem (1.24) can be transformed as **(1.26)** $\min_{\mathbf{V} \in \mathcal{V}} J(Q(\mathbf{V}))$.

Existence of tubes. In the smooth case, the existence of tubes follows the Cauchy–Lipschitz theory on differential equations [147], Delfour and J.-P. Zolésio, 2001. In the non-smooth case, the Lipschitz regularity of the velocities \mathbf{V} can be weakened using the equations satisfied by the characteristic functions $\xi(t, \mathbf{x})$ associated to the domain $\Omega_t(\mathbf{V})$, **(1.27)**

$$\begin{cases} \partial_t \xi + \nabla \xi \cdot \mathbf{V} = 0, & (0, \tau) \times D, \\ \xi_{t=0} = \chi_\Omega, & D. \end{cases}$$

We shall consider velocity fields s.t. $\mathbf{V} \in L^1(0, \tau; L^2(D; \mathbb{R}^d))$ & the divergence positive part $(\nabla \cdot \mathbf{V})^+ \in L^1(0, \tau; L^\infty(D))$. In this case, using a Galerkin approximation & some energy estimates, we are able to derive an existence result of solutions with initial data given in $H^{-1/2}(D)$. For the time being, no uniqueness result has been obtained for this smoothness level.

Actually, when the field \mathbf{V} & its divergence are simply L^1 functions, the notion of weak solutions associated to the convection problems (1.27) does not make sense. In this case, the correct modeling tool for shape evolution is to introduce the product space of elements $(\xi = \xi^2, \mathbf{V})$ equipped with a parabolic BV like topology for which the constraint (1.27) defines a closed subset \mathcal{T}_Ω which contains the weak closure of smooth elements

$$\mathcal{T}_\Omega := \{(\chi_\Omega \circ \mathbf{T}_t^{-1}(\mathbf{V}), \mathbf{V}); \mathbf{V} \in \mathcal{U}_{\text{ad}}\}.$$

This approach consists in handling characteristic functions $\xi = \xi^2$ which belongs to $L^1(0, \tau; \text{BV}(D))$ together with vector fields $\mathbf{V} \in L^2(0, \tau; L^2(D; \mathbb{R}^d))$ solution of problem (1.27). For a given element $(\xi, \mathbf{V}) \in \mathcal{T}_\Omega$, we consider the set of fields \mathbf{W} s.t. $(\xi, \mathbf{W}) \in \mathcal{T}_\Omega$, we consider the set of fields \mathbf{W} s.t. $(\xi, \mathbf{W}) \in \mathcal{T}_\Omega$. It forms a closed convex set, noted \mathcal{V}_ξ . Hence, we can define the unique minimal norm energy element \mathbf{V}_ξ in the convex set \mathcal{V}_ξ . For a given tube ξ , the element \mathbf{V}_ξ is the unique (with minimal norm) vector field associated to ξ via the convection equation (1.27).

We choose to adopt a different point of view inspired by the optimization problems framework. Indeed, our final goal is to apply the weak set evolution setting to the control problem arising in various fields such as free boundary problems or image processing. The usual situation can be described as follows. Let us consider a given smooth enough functional $J(\xi, \mathbf{V})$. We would like to solve the following optimization problem **(1.28)** $\inf_{(\xi, \mathbf{V}) \in \mathcal{T}_\Omega} J(\xi, \mathbf{V})$. The space \mathcal{U}_{ad} is a space of smooth velocities.

In most situations, such a problem does not admit solutions & we need to add some regularization terms to ensure its solvability. Consequently, we shall introduce different penalization terms which furnish compactness properties of the minimizing sequences inside an ad-hoc weak topology involving bounded variation constraints. Then, the new problem writes **(1.29)**

$$\inf_{(\xi, \mathbf{V}) \in \mathcal{T}_\Omega} J(\xi, \mathbf{V}) + F(\xi, \mathbf{V}).$$

The *penalization term* $F(\xi, \mathbf{V})$ can be chosen using several approaches:

- We can 1st consider the time-space perimeter of the lateral boundary Σ of the tube, developed in [155]. This approach easily draws part of the variational properties associated to the bounded variation functions space framework. In particular, it uses the compactness properties of tube family with bounded perimeters in \mathbb{R}^{d+1} . Nevertheless, this method leads to heavy variational analysis developments.
- We can rather consider the time integral of the spatial perimeter of the moving domain which builds the tube, as introduced in [157]. We shall extend these results to the case of vector fields living in $L^2((0, \tau) \times D; \mathbb{R}^d)$. In this case, only existence results for solutions of the convection equation can be handled & the uniqueness property is lost.

Tube derivative. In this paragraph, we are interested in differentiability properties of integrals defined over moving domains,

$$J(Q(\mathbf{V})) = \int_{Q(\mathbf{V})} f(\mathbf{V}) \, d\mathbf{x} \, dt.$$

The transverse map \mathcal{T}_ρ^t associated to 2 vector fields $(\mathbf{V}, \mathbf{W}) \in \mathcal{U}$ is defined as follows,

$$\begin{aligned} \mathcal{T}_\rho^t : \overline{\Omega_t} &\rightarrow \overline{\Omega_t^\rho} := \overline{\Omega_t(\mathbf{V} + \rho \mathbf{W})} \\ \mathbf{x} &\mapsto T_t(\mathbf{V} + \rho \mathbf{W}) \circ T_t(\mathbf{V})^{-1}. \end{aligned}$$

Remark 4.3. *The transverse map allows us to perform sensitivity analysis on functions defined on the unperturbed domain $\Omega_t(\mathbf{V})$.*

The following result states that the transverse map \mathcal{T}_ρ^t can be considered as a dynamical flow w.r.t. the perturbation variable ρ – Moubachir and Jean-Paul Zolésio, 2006, Chap. 1, Subsect. 1.4.3, pp. 19–23 [skipped pp. 23–31]

4.2 Inverse Stefan Problem

“... we consider the identification of a moving boundary that represents the isothermal interface between a solid phase & a liquid phase, from measurements on a fixed part of the solid boundary. This problem is referred in the literature as the *inverse Stefan problem* [61, 144]. We make use of the *transverse derivative concepts* introduced in [154, 155].” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 2, p. 33

4.2.1 The inverse problem setting

“For a given evolution of the *melting interface*, we consider the solution of the heat equation (2.3) & we consider its trace on the fixed boundary Σ^s .” Fig. 2.2. **Non-cylindrical space-time domain.** “On the mathematical viewpoint, we introduce the *observation space* $\mathcal{O} := L^2((0, \tau); L^2(\Gamma^s))$ & the *observation operator* (2.5) $\mathfrak{D} : \mathcal{U}_{\text{ad}} \rightarrow \mathcal{O}$, $\mathbf{V} \mapsto \mathfrak{D}(\mathbf{V}) := \gamma_{\Gamma^s}(y(\mathbf{V}))$ where $y(\mathbf{V})$ stands for the solution of (2.3) & γ_{Γ^s} is the zero order trace operator on Σ^s .

The inverse Stefan problem consists in recovering the evolution of the melting front $\Gamma^f(t)$ from the knowledge of the temperature on the fixed solid boundary Γ^s . I.e., for a given temperature $y_d \in L^2((0, \tau); L^2(\Gamma^s))$, we look for $\mathbf{V} \in \mathcal{U}_{\text{ad}}$ s.t. (2.6) $\mathfrak{D}(\mathbf{V}) = y_d$, in \mathcal{O} . It is a nonlinear ill-posed inverse problem that can be solved using a least-square minimization problem regularized thanks to a Tikhonov zero order term. Hence we look for the solution \mathbf{V} of the following optimization problem (2.7)

$$\min_{\mathbf{V} \in \mathcal{U}_{\text{ad}}} \frac{1}{2} \|\mathfrak{D}(\mathbf{V}) - y_d\|_{\mathcal{O}}^2 + \frac{\alpha}{2} \|\mathbf{V}\|_{\mathcal{U}_{\text{ad}}}^2,$$

with $\alpha > 0$.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 2, Sect. 2.2, p. 35

4.2.2 The Eulerian derivative & the transverse field

“A possible choice in order to solve the above minimization problem is to use a gradient based method such as the *conjugate gradient method*. Hence, we need to evaluate the gradient w.r.t. \mathbf{V} of the functional (2.13) $j(\mathbf{V}) := \frac{1}{2} \|\mathfrak{D}(\mathbf{V}) - y_d\|_{\mathcal{O}}^2$ where, for the sake of simpleness, we have dropped the regularizing term $\frac{\alpha}{2} \|\mathbf{V}\|_{\mathcal{U}_{\text{ad}}}^2$. Let us choose a perturbation direction $\mathbf{W} \in \mathcal{U}_{\text{ad}}$. We would like to compute the *directional derivative* of j , (2.14)

$$[\mathbf{D}_{\mathbf{V}}[j](\mathbf{V})] \cdot \mathbf{W} := \lim_{\rho \rightarrow 0} \frac{1}{\rho} (j(\mathbf{V} + \rho \mathbf{W}) - j(\mathbf{V})).$$

Then the goal is to evaluate the directional derivative of the element $y(\mathbf{V})$ which is a solution of the moving heat equation (2.3). In order to do so, we write the associated variational formulation satisfied $y(\mathbf{V})$, (2.15)

$$\int_0^\tau \int_{\Omega_t(\mathbf{V})} [\partial_t y(\mathbf{V}) \phi(\mathbf{V}) + \nabla y(\mathbf{V}) \cdot \nabla \phi(\mathbf{V})] d\mathbf{x} dt = 0, \quad \forall \phi(\mathbf{V}) \in L^2((0, \tau); H_{0, \Gamma_t(\mathbf{V})}^1(\Omega_t(\mathbf{V}))),$$

where we have set w.l.o.g., $(f, y_d, y_0) = (0, 0, 0)$ together with $\Omega_t(\mathbf{V}) := \Omega^s(t)$ & $\Gamma_t(\mathbf{V}) := \Gamma^f(t)$. Looking at (2.15), it is clear that we need to establish how to differentiate the generic term $J(\mathbf{V}) = \int_0^\tau \int_{\Omega_t(\mathbf{V})} f(\mathbf{V}) d\mathbf{x} dt$ w.r.t. \mathbf{V} . To this end, we introduce the *perturbed moving domain* $\Omega_t(\mathbf{V} + \rho \mathbf{W}) := T_t(\mathbf{V} + \rho \mathbf{W})(\Omega_0)$. This family generates a *perturbed tube*

$$Q(\mathbf{V} + \rho \mathbf{W}) := \bigcup_{0 \leq t \leq \tau} (\{t\} \times \Omega_t(\mathbf{V} + \rho \mathbf{W}))$$

as described in Fig. 2.3. Perturbed tube. Since the function $f(\mathbf{V})$ is defined on the *non-cylindrical reference tube* $Q(\mathbf{V})$, it is natural to introduce the transformation between $Q(\mathbf{V})$ & $Q(\mathbf{V} + \rho \mathbf{W})$. A canonical choice is furnished by

$$\begin{aligned} \mathcal{T}^t(\rho; \mathbf{x}) &: \Omega_t(\mathbf{V}) \rightarrow \Omega_t(\mathbf{V} + \rho \mathbf{W}) \\ \mathbf{x} &\mapsto \mathcal{T}^t(\rho; \mathbf{x}) := [\mathbf{T}_t(\mathbf{V} + \rho \mathbf{W}) \circ \mathbf{T}_t(\mathbf{V})^{-1}](\mathbf{x}). \end{aligned}$$

Hence, the perturbed functional can be written as follows,

$$J(\mathbf{V} + \rho \mathbf{W}) = \int_0^\tau \int_{\Omega_t(\mathbf{V} + \rho \mathbf{W})} f(\mathbf{V} + \rho \mathbf{W}) d\mathbf{x} dt = \int_0^\tau \int_{\mathcal{T}_\rho^t(\Omega_t(\mathbf{V}))} f(\mathbf{V} + \rho \mathbf{W}) d\mathbf{x} dt = \int_0^\tau \int_{\Omega_t(\mathbf{V})} (\det D\mathcal{T}_\rho^t) f(\mathbf{V} + \rho \mathbf{W}) \circ \mathcal{T}_\rho^t d\mathbf{x} dt,$$

where we have performed a transport into the *moving reference domain* $\Omega_t(\mathbf{V})$. Now we shall need to differentiate the terms inside the integral w.r.t. ρ at point $\rho = 0$. The easiest way to do so is to connect this problem to the classical shape derivative calculus handled inside the speed method framework [147, 135]. I.e., we need to identify a transverse velocity field that may generate the transverse map $\mathcal{T}^t(\rho; \mathbf{x})$ as the solution of a dynamical system w.r.t. the parameter $\rho \in [0, \rho_0]$. Actually, it can be proven that $\mathbf{T}(\mathbf{V} + \rho \mathbf{W})$ is continuously differentiable¹² w.r.t. ρ & that the transverse map \mathcal{T}_ρ^t can be considered as the flow w.r.t. ρ of the transverse vector field

$$\mathcal{Z}(\rho; (t, \mathbf{x})) := [\partial_\rho \mathcal{T}^t(\rho)] \circ \mathcal{T}^t(\rho)^{-1}(\mathbf{x}) = [\partial_\rho \mathbf{T}(\mathbf{V} + \rho \mathbf{W})] \circ \mathbf{T}(\mathbf{V} + \rho \mathbf{W})^{-1}(\mathbf{x}).$$

p. 39

Fig. 2.4. Transverse map.

” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 2, Sect. 2.3, pp. 37–

4.3 Dynamical Shape Control of NSEs

Moubachir and Jean-Paul Zolésio, 2006, Chap. 5 “deals with the analysis of an *inverse dynamical shape problem* involving a fluid inside a moving domain. This type of inverse problem happens frequently in the design & the control of many industrial devices such as aircraft wings, cable-stayed bridges, automobile shapes, satellite reservoir tanks & more generally of systems involving fluid-solid interactions.

The control variable is the shape of the moving domain, & the objective is to minimize a given cost functional that may be chosen by the designer.

On the theoretical level, early works concerning optimal control problems for general parabolic equations written in non-cylindrical domains have been considered in [43, 29, 30, 142, 2]. In [140, 151, 152], the stabilization of structures using the variation of the domain has been addressed. The basic principle is to define a map sending the non-cylindrical domain into a cylindrical one. This process leads to the mathematical analysis of non-autonomous PDE’s systems.

Recently, a new methodology to obtain *Eulerian derivatives* for non-cylindrical functionals has been introduced in [157, 156, 58]. This methodology was applied in [59] to perform dynamical shape control of the non-cylindrical NSEs where the evolution of the domain is the control variable. Hence the classical optimal shape optimization theory has been extended to deal with non-cylindrical domains.”

Aim. “review several results on the dynamical shape control of the Navier–Stokes system & suggest an alternative treatment using the Min-Max principle [45, 46]. Despite its lack of rigorous mathematical justification in the case where the Lagrangian functional is not convex, we shall show how this principle allows, at least formally, to bypass the tedious computation of the state differentiability w.r.t. the shape of the moving domain.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 5, p. 109

¹²in $\mathcal{Z}_{\text{ad}} := \mathcal{C}^0([0, \tau]; (\mathcal{C}^{k-1}(\overline{D}))^d)$.

4.3.1 Problem Statement

“Let us consider a moving domain $\Omega_t \in \mathbb{R}^d$. We introduce a diffeomorphic map sending a fixed reference domain Ω_0 into the physical configuration Ω_t at time $t \geq 0$. W.l.o.g., we choose the reference configuration to be the physical configuration at initial time $\Omega_{t=0}$. Hence we define a map $T_t \in \mathcal{C}^1(\overline{\Omega_0})$ s.t. $\overline{\Omega_t} = T_t(\overline{\Omega_0})$, $\overline{\Gamma_t} = T_t(\overline{\Gamma_0})$. We set $\Sigma := \bigcup_{0 < t < T} (\{t\} \times \Gamma_t)$, $Q := \bigcup_{0 < t < T} (\{t\} \times \Omega_t)$. The map T_t can be actually defined as the flow of a particular vector field, as described in the following lemma:

Theorem 4.2 (ref. 147). $\overline{\Omega_t} = T_t(V)(\overline{\Omega_0})$, $\overline{\Gamma_t} = T_t(V)(\overline{\Gamma_0})$ where $T_t(V)$ is the solution of the following dynamical system:

$$\begin{aligned} T_t(V) : \Omega_0 &\rightarrow \Omega \\ x_0 &\mapsto x(t, x_0) := T_t(V)(x_0) \end{aligned}$$

with (5.1)

$$\begin{cases} \frac{dx}{d\tau} = V(\tau, x(\tau)), & \tau \in [0, T], \\ x(\tau = 0) = x_0, & \text{in } \Omega_0. \end{cases}$$

The fluid filling Ω_t is assumed to be a viscous incompressible Newtonian fluid. Its evolution is described by its velocity \mathbf{u} & its pressure p . The couple (\mathbf{u}, p) satisfies the classical NSEs written in non-conservative form (5.2)

$$\begin{cases} \partial_t \mathbf{u} + \mathbf{Du} \cdot \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = 0, & Q(V), \\ \nabla \cdot \mathbf{u} = 0, & Q(V), \\ \mathbf{u} = V, & \Sigma(V), \\ \mathbf{u}(t = 0) = \mathbf{u}_0, & \Omega_0, \end{cases}$$

where ν stands for the kinematic viscosity. The quantity $\sigma(\mathbf{u}, p) = -p\mathbf{I} + \nu(\mathbf{Du} + {}^*\mathbf{Du})$ stands for the *fluid stress tensor* inside Ω_t , with $(\mathbf{Du})_{i,j} = \partial_j u_i$. We are interested in solving the following minimization problem: (5.3) $\min_{V \in \mathcal{U}} j(V)$ where $j(V) = J_V(\mathbf{u}(V), p(V))$ with $(\mathbf{u}(V), p(V))$ is a weak solution of problem (5.2) & $J_V(\mathbf{u}, p)$ is a real functional of the following form: (5.4)

$$J_V(\mathbf{u}, p) = \frac{\alpha}{2} \|\mathcal{B}\mathbf{u}\|_{Q(V)}^2 + \frac{\gamma}{2} \|\mathcal{K}V\|_{\Sigma(V)}^2,$$

where $\mathcal{B} \in \mathcal{L}(\mathcal{H}, \mathcal{H}^*)$ is a general linear differential operator satisfying the following identity, (5.5)

$$\langle \mathcal{B}\mathbf{u}, \mathbf{v} \rangle + \langle \mathbf{u}, \mathcal{B}^*\mathbf{v} \rangle = \langle \mathcal{B}_\Sigma \mathbf{u}, \mathbf{v} \rangle_{L^2(\Sigma)},$$

where $\mathcal{H} = \{\mathbf{v} \in L^2(0, T; (H_0^1(\text{div}, \Omega_t(V)))^d)\}$ & $\mathcal{K} \in \mathcal{L}(\mathcal{U}, L^2(\Sigma(V)))$ is a general linear differential operator satisfying the following identity, (5.6)

$$\langle \mathcal{K}\mathbf{u}, \mathbf{v} \rangle_{L^2(\Sigma)} + \langle \mathbf{u}, \mathcal{K}^*\mathbf{v} \rangle_{L^2(\Sigma)} = \langle \mathcal{K}_\Sigma \mathbf{u}, \mathbf{v} \rangle_{L^2(\Sigma)}.$$

The main difficulty in dealing with such a minimization problem is related to the fact that integrals over the domain $\Omega_t(V)$ depend on the control variable V . This point will be solved by using the Arbitrary Lagrange–Euler (ALE) map $T_t(V)$ introduced previously. The purpose of this chapter is to prove using several methods the following result,

Theorem 4.3 (Main result). *For $V \in \mathcal{U}$ & Ω_0 of class \mathcal{C}^2 , the functional $j(V)$ possesses a gradient $\nabla j(V)$ which is supported on the moving boundary $\Gamma_t(V)$ & can be represented by the following expression, (5.7)*

$$\nabla j(V) = -\lambda \mathbf{n} - \sigma(\varphi, \pi) \cdot \mathbf{n} + \alpha \mathcal{B}_\Sigma \mathbf{u} + \gamma [-\mathcal{K}^* \mathcal{K}V + \mathcal{K}_\Sigma \mathcal{K}V],$$

where (φ, π) stands for the adjoint fluid state solution of the following system, (5.8)

$$\begin{cases} -\partial_t \varphi - \mathbf{D}\varphi \cdot \mathbf{u} + {}^*\mathbf{Du} \cdot \varphi - \nu \Delta \varphi + \nabla \pi = -\alpha \mathcal{B}^* \mathcal{B}\mathbf{u}, & Q(V), \\ \nabla \cdot \varphi = 0, & Q(V), \\ \varphi = 0, & \Sigma(V), \\ \varphi(T) = 0, & \Omega_T, \end{cases}$$

& λ is the adjoint transverse boundary field, solution of the tangential dynamical system, (5.9)

$$\begin{cases} -\partial_t \lambda - \nabla_\Gamma \lambda \cdot V - (\nabla \cdot V)\lambda = f, & (0, T), \\ \lambda(T) = 0, & \Gamma_T(V), \end{cases}$$

with (5.10)

$$f = [-(\sigma(\varphi, \pi) \cdot \mathbf{n}) + \alpha \mathcal{B}_\Sigma \mathcal{B}\mathbf{u}] \cdot (\mathbf{DV} \cdot \mathbf{n} - \mathbf{Du} \cdot \mathbf{n}) + \frac{1}{2} [\alpha |\mathcal{B}\mathbf{u}|^2 + \gamma H |\mathcal{K}V|^2].$$

Example 4.1. We set $(\mathcal{B}, \mathcal{B}^*, \mathcal{B}_\Sigma) = (\mathbf{I}, -\mathbf{I}, 0)$, $(\mathcal{K}, \mathcal{K}^*, \mathcal{K}_\Sigma) = (\mathbf{I}, -\mathbf{I}, 0)$. I.e., we consider the cost functional, (5.11)

$$J_V(\mathbf{u}, p) = \frac{\alpha}{2} \|\mathbf{u}\|_{L^2(Q(V))}^2 + \frac{\gamma}{2} \|V\|_{L^2(\Sigma(V))}^2.$$

Then its gradient is given by (5.12)

$$\nabla j(V) = -\lambda \mathbf{n} - \sigma(\boldsymbol{\varphi}, \pi) \cdot \mathbf{n} + \gamma V,$$

where $(\boldsymbol{\varphi}, \pi)$ stands for the adjoint fluid state solution of the following system (5.13)

$$\begin{cases} -\partial_t \boldsymbol{\varphi} - \mathbf{D}\boldsymbol{\varphi} \cdot \mathbf{u} + {}^*\mathbf{D}\mathbf{u} \cdot \boldsymbol{\varphi} - \nu \Delta \boldsymbol{\varphi} + \nabla \pi = \alpha \mathbf{u}, & Q(V), \\ \nabla \cdot \boldsymbol{\varphi} = 0, & Q(V), \\ \boldsymbol{\varphi} = \mathbf{0}, & \Sigma(V), \\ \boldsymbol{\varphi}(T) = \mathbf{0}, & \Omega_T, \end{cases}$$

$\mathcal{E} \lambda$ is the adjoint transverse boundary field, solution of the tangential dynamical system, (5.14)

$$\begin{cases} -\partial_t \lambda - \nabla_\Gamma \lambda \cdot V - (\nabla \cdot V) \lambda = f, & (0, T), \\ \lambda(T) = 0, & \Gamma_T(V), \end{cases}$$

with (5.15)

$$f = -\nu(\mathbf{D}\boldsymbol{\varphi} \cdot \mathbf{n}) \cdot (\mathbf{D}V \cdot \mathbf{n} - \mathbf{D}\mathbf{u} \cdot \mathbf{n}) + \frac{1}{2}(\alpha + \gamma H)|V|^2.$$

Example 4.2. We set $(\mathcal{B}, \mathcal{B}^*, \mathcal{B}_\Sigma) = (\text{curl}, \text{curl}, \wedge \mathbf{n})$, $(\mathcal{K}, \mathcal{K}^*, \mathcal{K}_\Sigma) = (\mathbf{I}, -\mathbf{I}, 0)$, (5.16)

$$J_V(\mathbf{u}, p) = \frac{\alpha}{2} \|\text{curl } \mathbf{u}\|_{L^2(Q(V))}^2 + \frac{\gamma}{2} \|V\|_{L^2(\Sigma(V))}^2.$$

Then its gradient is given by (5.17)

$$\nabla j(V) = -\lambda \mathbf{n} - \sigma(\boldsymbol{\varphi}, \pi) \cdot \mathbf{n} + \alpha(\text{curl } \mathbf{u}) \wedge \mathbf{n} + \gamma V,$$

where $(\boldsymbol{\varphi}, \pi)$ stands for the adjoint fluid state solution of the following system, (5.18)

$$\begin{cases} -\partial_t \boldsymbol{\varphi} - \mathbf{D}\boldsymbol{\varphi} \cdot \mathbf{u} + {}^*\mathbf{D}\mathbf{u} \cdot \boldsymbol{\varphi} - \nu \Delta \boldsymbol{\varphi} + \nabla \pi = -\alpha \Delta \mathbf{u}, & Q(V), \\ \nabla \cdot \boldsymbol{\varphi} = 0, & Q(V), \\ \boldsymbol{\varphi} = \mathbf{0}, & \Sigma(V), \\ \boldsymbol{\varphi}(T) = \mathbf{0}, & \Omega_T, \end{cases}$$

$\mathcal{E} \lambda$ is the adjoint transverse boundary field, solution of the tangential dynamical system, (5.19)

$$\begin{cases} -\partial_t \lambda - \nabla_\Gamma \lambda \cdot V - (\nabla \cdot V) \lambda = f, & (0, T), \\ \lambda(T) = 0, & \Gamma_T(V), \end{cases}$$

with

$$f = [-\nu \mathbf{D}\boldsymbol{\varphi} \cdot \mathbf{n} + \alpha(\text{curl } \mathbf{u}) \wedge \mathbf{n}] \cdot (\mathbf{D}V \cdot \mathbf{n} - \mathbf{D}\mathbf{u} \cdot \mathbf{n}) + \frac{1}{2} [\alpha |\text{curl } \mathbf{u}|^2 + \gamma H|V|^2].$$

In the next section, we introduce several concepts closely related to shape optimization tools for moving domain problems. We also recall *elements of tangential calculus* that will be used through this chapter. Then we treat successively the following points,

1. In Sect. 5.5, we choose to prove the differentiability of the fluid state (\mathbf{u}, p) w.r.t. the design variable V . The directional shape derivative $(\mathbf{u}', p')(V) \cdot W$ is then used to compute the directional derivative $j'(V) \cdot W$ of the cost functional $j(V)$. Using the adjoint state $(\boldsymbol{\varphi}, \pi)(V)$ associated to $(\mathbf{u}', p')(V)$ & the adjoint field Λ associated to the *transverse field* Z_t introduced in sect. 5.3, we are able to furnish an expression of the gradient $\nabla j(V)$ which is a distribution supported by the moving boundary $\Gamma_t(V)$.

2. In Sect. 5.6, we choose to bypass the computation of the state shape derivative $(\mathbf{u}', p')(V) \cdot W$, by using a Min-Max formulation of problem (5.3) & a transport technique. The state & multiplier spaces are chosen in order to be independent on the scalar perturbation parameter used in the computation of the derivative of the Lagrangian functional w.r.t. V . This method directly furnishes the fluid state & transverse field adjoint systems & the resulting gradient $\nabla j(V)$.
3. In Sect. 5.7, we again use a Min-Max strategy coupled with a state & multiplier functional space embedding. I.e., the state & multiplier variables live in the hold-all domain D . Hence the derivative of the Lagrangian functional w.r.t. V only involves terms coming from the flux variation through the moving boundary $\Gamma_t(V)$. This again leads to the direct computation of the fluid state & transverse field adjoints & consequently to the gradient $\nabla j(V)$. – Moubachir and Jean-Paul Zolésio, 2006, Chap. 5, Sect. 5.2, pp. 110–114

4.3.2 Elements of Non-cylindrical Shape Calculus

“This section introduces several concepts that will be intensively used through this chapter. It concerns the differential calculus of integrals defined on moving domains or boundaries w.r.t. their support.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 5, Sect. 5.3, p. 114

4.3.2.1 Non-cylindrical speed method

“In this paragraph, we are interested in differentiability properties of integrals defined over moving domains,

$$J_1(\Omega_t) = \int_{\Omega_t} f(\Omega_t) d\Omega, \quad J_2(\Gamma_t) = \int_{\Gamma_t} g(\Gamma_t) d\Gamma.$$

The behavior of J_1 & J_2 while perturbing their moving support highly depends on the regularity in space & time of the domains. In this work, we choose to work with domains Ω_t that are images of a fixed domain Ω_0 through an ALE map $T_t(V)$ as introduced in the 1st section. Hence, the design parameter is no more the support Ω_t but rather the velocity field $V \in \mathcal{U} := \mathcal{C}([0, T]; (W^{k, \infty}(D))^d)$ that builds the support. This technique has the advantage to transform shape calculus into classical differential calculus on vector spaces [157, 59]. For another choice based on the *non-cylindrical identity perturbation*, the reader is referred to the Moubachir and Jean-Paul Zolésio, 2006, Chap. 6.”

Transverse applications.

Definition 4.1 (Transverse map). *The transverse map \mathcal{T}_ρ^t associated to 2 vector fields $(V, W) \in \mathcal{U}$ is defined as follows,*

$$\begin{aligned} \mathcal{T}_\rho^t : \overline{\Omega_t} &\rightarrow \overline{\Omega_t^\rho} := \overline{\Omega_t(V + \rho W)} \\ \mathbf{x} &\mapsto T_t(V + \rho W) \circ T_t(V)^{-1}. \end{aligned}$$

Remark 4.4. *The transverse map allows us to perform sensitivity analysis on functions defined on the unperturbed domain $\Omega_t(V)$.*

The following result states that the transverse map \mathcal{T}_ρ^t can be considered as a dynamical flow w.r.t. the perturbation variable ρ ,

Theorem 4.4 (ref. 156). *The Transverse map \mathcal{T}_ρ^t is the flow of a transverse field \mathcal{Z}_ρ^t defined as follows (5.20)*

$$\mathcal{Z}_\rho^t := \mathcal{Z}^t(\rho, \cdot) = \left(\frac{\partial \mathcal{T}_\rho^t}{\partial \rho} \right) \circ (\mathcal{T}_\rho^t)^{-1},$$

i.e., is the solution of the following dynamical system:

$$\begin{aligned} T_t^\rho(\mathcal{Z}_\rho^t) : \overline{\Omega_t} &\rightarrow \overline{\Omega_t^\rho} \\ \mathbf{x} &\mapsto \mathbf{x}(\rho, \mathbf{x}) := T_t^\rho(\mathcal{Z}_\rho^t)(\mathbf{x}) \end{aligned}$$

with (5.21)

$$\begin{cases} \frac{d\mathbf{x}(\rho)}{d\rho} = \mathcal{Z}^t(\rho, \mathbf{x}(\rho)), & \rho \geq 0, \\ \mathbf{x}(\rho = 0) = \mathbf{x}, & \text{in } \Omega_t(V). \end{cases}$$

Since, we will mainly consider derivatives of perturbed functions at point $\rho = 0$, we set $Z_t := \mathcal{Z}_{\rho=0}^t$. A fundamental result lies in the fact that Z_t can be obtained as the solution of a linear time dynamical system depending on the vector fields $(V, W) \in \mathcal{U}$,

Theorem 4.5. *The vector field Z_t is the unique solution of the following Cauchy problem, (5.22)*

$$\begin{cases} \partial_t Z_t + [Z_t, V] = W, & (0, T) \times D, \\ Z_{t=0} = 0, & D, \end{cases}$$

where $[Z_t, V] := DZ_t \cdot V - DV \cdot Z_t$ stands for the Lie bracket of the pair (Z_t, V) .

Shape derivative of non-cylindrical functionals. The main theorem of this section uses the notion of a non-cylindrical material derivative that we recall here,

Definition 4.2. *The derivative w.r.t. ρ at point $\rho = 0$ of the following composed function,*

$$\begin{aligned} f^\rho : [0, \rho_0] &\rightarrow H(\Omega_t(V)) \\ \rho &\mapsto f(V + \rho W) \circ \mathcal{T}_\rho^t \end{aligned}$$

$\dot{f}(V; W)$ is called the non-cylindrical material derivative of $f(V)$ at point $V \in U$ in the direction $W \in \mathcal{U}$. We shall use the notation,

$$\dot{f}(V) \cdot W = \dot{f}(V; W) := \left. \frac{d}{d\rho} f^\rho \right|_{\rho=0}.$$

With the above definition, we can state the differentiability properties of non-cylindrical integrals w.r.t. their moving support,

Theorem 4.6 (ref. 59). *For a bounded measurable domain Ω_0 with boundary Γ_0 , let us assume that for any direction $W \in U$ the following hypothesis holds,*

- (i) *$f(V)$ admits a non-cylindrical material derivative $\dot{f}(V) \cdot W$ then $J_1(\cdot)$ is Gâteaux differentiable at point $V \in \mathcal{U}$ & its derivative is given by the following expression, (5.23)*

$$J'_1(V) \cdot W = \int_{\Omega_t(V)} [\dot{f}(V) \cdot W + f(V) \nabla \cdot Z_t] d\Omega.$$

Furthermore, if

- (ii) *$f(V)$ admits a non-cylindrical shape derivative given by the following expression, (5.24)*

$$f'(V) \cdot W = \dot{f}(V) \cdot W - \nabla f(V) \cdot Z_t,$$

then (5.25)

$$J'_1(V) \cdot W = \int_{\Omega_t(V)} [f'(V) \cdot W + \nabla \cdot (f(V) Z_t)] d\Omega.$$

Furthermore, if Ω_0 is an open domain with a Lipschitzian boundary Γ_0 , then (5.26)

$$J'_1(V) \cdot W = \int_{\Omega_t(V)} f'(V) \cdot W d\Omega + \int_{\Gamma_t(V)} f(V) \langle Z_t, \mathbf{n} \rangle d\Gamma.$$

Remark 4.5. *The last identity will be of great interest while trying to prove a gradient structure result for general non-cylindrical functionals.*

It is also possible to establish a similar result for integrals over moving boundaries. For that purpose, we need to define the non-cylindrical tangential material derivative,

Definition 4.3. *The derivative w.r.t. ρ at point $\rho = 0$ of the following composed function,*

$$\begin{aligned} g^\rho : [0, \rho_0] &\rightarrow H(\Gamma_t(V)) \\ \rho &\mapsto g(V + \rho W) \circ \mathcal{T}_\rho^t \end{aligned}$$

is called the non-cylindrical material derivative of the function $g(V) \in H(\Gamma_t(V))$ in the direction $W \in \mathcal{U}$. We shall use the notation

$$\dot{g}(V) \cdot W = \dot{g}(V; W) := \left. \frac{d}{d\rho} g^\rho \right|_{\rho=0}.$$

This concept is involved in the differentiability property of boundary integrals,

Theorem 4.7. *For a bounded measurable domain Ω_0 with boundary Γ_0 , let us assume that for any direction $W \in U$ the following hypothesis holds,*

- (i) *$g(V)$ admits a non-cylindrical material derivative $\dot{g}(V) \cdot W$ then $J_2(\cdot)$ is Gâteaux differentiable at point $V \in \mathcal{U}$ & its derivative is given by the following expression, (5.27)*

$$J'_2(V) \cdot W = \int_{\Gamma_t(V)} [\dot{g}(V) \cdot W + g(V) \operatorname{div}_\Gamma Z_t] d\Gamma.$$

Furthermore, if

- (ii) *$g(V)$ admits a non-cylindrical shape derivative given by the following expression, (5.28)*

$$g'(V) \cdot W = \dot{g}(V) \cdot W - \nabla_\Gamma g(V) \cdot Z_t,$$

then (5.29)

$$J'_2(V) \cdot W = \int_{\Gamma_t(V)} [\tilde{g}'(V) \cdot W + Hg(V) \langle Z_t, \mathbf{n} \rangle] d\Gamma,$$

where H stands for the additive curvature. Furthermore, if $g(V) = \tilde{g}(V)|_{\Gamma_t(V)}$ with $\tilde{g} \in H(\Omega_t(V))$, then (5.30)

$$J'_2(V) \cdot W = \int_{\Gamma_t(V)} [g'(V) \cdot W + (\nabla \tilde{g}(V) \cdot \mathbf{n} + Hg(V)) \langle Z_t, \mathbf{n} \rangle] d\Gamma.$$

Adjoint transverse field. It is possible to define the solution of the adjoint transverse system,

Theorem 4.8. *For $F \in L^2(0, T; (H^1(D))^d)$, there exists a unique field $\Lambda \in \mathcal{C}^0([0, T]; (L^2(D))^d)$ solution of the backward dynamical system, (5.31)*

$$\begin{cases} -\partial_t \Lambda - \operatorname{DA} \cdot V - {}^* \operatorname{DV} \cdot \Lambda - (\nabla \cdot V) \Lambda = F, & (0, T), \\ \Lambda(T) = 0, \end{cases}$$

Remark 4.6. *The field Λ is the dual variable associated to the transverse field Z_t & is the solution of the adjoint problem associated to the transverse dynamical system.*

In this chapter, we shall deal with a specific RHS F of the form $F(t) = {}^* \gamma_{\Gamma_t(V)}(f(t)\mathbf{n})$. In this case, the adjoint field Λ is supported on the moving boundary $\Gamma_t(V)$ & has the following structure,

Theorem 4.9 (ref. 59). *For $F(t) = {}^* \gamma_{\Gamma_t(V)}(f(t)\mathbf{n})$, with $f \in L^2(0, T; L^2(\Gamma_t(V)))$, the unique solution Λ of the problem is given by the following identity, (5.32)*

$$\Lambda = (\lambda \circ p) \nabla_{\chi_{\Omega_t(V)}} \in \mathcal{C}^0([0, T]; (H^1(\Gamma_t))^d),$$

where $\lambda \in \mathcal{C}^0([0, T]; H^1(\Gamma_t))$ is the unique solution of the following boundary dynamical system, (5.33)

$$\begin{cases} -\partial_t \lambda - \nabla_\Gamma \lambda \cdot V - (\nabla \cdot V) \lambda = f, & (0, T) \\ \lambda(T) = 0, & \Gamma_t(V), \end{cases}$$

p is the canonical projection on $\Gamma_t(V)$ & $\chi_{\Omega_t(V)}$ is the characteristic function of $\Omega_t(V)$ inside D .

Gradient of non-cylindrical functionals. In the next sections, we will often deal with boundary integrals of the following forms,

$$K = \int_0^T \int_{\Gamma_t(V)} E \langle Z_t, \mathbf{n} \rangle$$

with $E \in L^2(0, T; \Gamma_t(V))$ & Z_t is the solution of the transverse equation (5.22). The following result allows us to eliminate the auxiliary variable Z_t inside the functional K ,

Theorem 4.10 (ref. 59). *For any $E \in L^2(0, T; \Gamma_t(V))$ & $(V, W) \in \mathcal{U}$, the following identity holds, (5.34)*

$$\int_0^T \int_{\Gamma_t(V)} E \langle Z_t, \mathbf{n} \rangle = - \int_0^T \int_{\Gamma_t(V)} \lambda \langle W, \mathbf{n} \rangle,$$

where $\lambda \in \mathcal{C}^0([0, T]; H^1(\Gamma_t))$ is the unique solution of problem (5.33) with $f = E$.

4.3.3 Elements of tangential calculus

“In this section, we review basic *elements of differential calculus* on a \mathcal{C}^k -submanifold with $k \geq 2$ of codimension 1 in \mathbb{R}^d . The following approach avoids the use of local bases & coordinates by using the intrinsic tangential derivative.” – Moubachir and Jean-Paul Zolésio, 2006, Chap. 5, Sect. 5.4, pp. 119

4.3.3.1 Oriented distance function

Chapter 5

Topology Optimization

Bibliography

- Abergel, F. and R. Temam (1990). “On Some Control Problems in Fluid Mechanics”. In: *Theoret. Comput. Fluid Dynamics* 1.6, pp. 303–325. DOI: [10.1007/BF00271794](https://doi.org/10.1007/BF00271794). URL: <https://doi.org/10.1007/BF00271794>.
- Céa, Jean, Alain Gioan, and Jean Michel (1973). “Quelques résultats sur l’identification de domaines”. In: *Calcolo* 10, pp. 207–232. ISSN: 0008-0624. DOI: [10.1007/BF02575843](https://doi.org/10.1007/BF02575843). URL: <https://doi.org/10.1007/BF02575843>.
- Delfour, M. C. and J.-P. Zolésio (2001). *Shapes and geometries*. Vol. 4. Advances in Design and Control. Analysis, differential calculus, and optimization. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, pp. xviii+482. ISBN: 0-89871-489-3.
- (2011). *Shapes and geometries*. Second. Vol. 22. Advances in Design and Control. Metrics, analysis, differential calculus, and optimization. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, pp. xxiv+622. ISBN: 978-0-898719-36-9. DOI: [10.1137/1.9780898719826](https://doi.org/10.1137/1.9780898719826). URL: <https://doi.org/10.1137/1.9780898719826>.
- Lions, J.-L. (1971). *Optimal control of systems governed by partial differential equations*. Die Grundlehren der mathematischen Wissenschaften, Band 170. Translated from the French by S. K. Mitter. Springer-Verlag, New York-Berlin, pp. xi+396.
- Lions, J.-L. and E. Magenes (1972). *Non-homogeneous boundary value problems and applications. Vol. I*. Die Grundlehren der mathematischen Wissenschaften, Band 181. Translated from the French by P. Kenneth. Springer-Verlag, New York-Heidelberg, pp. xvi+357.
- Moubachir, Marwan and Jean-Paul Zolésio (2006). *Moving shape analysis and control*. Vol. 277. Pure and Applied Mathematics (Boca Raton). Applications to fluid structure interactions. Chapman & Hall/CRC, Boca Raton, FL, pp. xx+291. ISBN: 978-1-58488-611-2; 1-58488-611-0. DOI: [10.1201/9781420003246](https://doi.org/10.1201/9781420003246). URL: <https://doi.org/10.1201/9781420003246>.
- Pironneau, Olivier (1984). *Optimal shape design for elliptic systems*. Springer Series in Computational Physics. Springer-Verlag, New York, pp. xii+168. ISBN: 0-387-12069-6. DOI: [10.1007/978-3-642-87722-3](https://doi.org/10.1007/978-3-642-87722-3). URL: <https://doi.org/10.1007/978-3-642-87722-3>.
- Tröltzsch, Fredi (2010). *Optimal control of partial differential equations*. Vol. 112. Graduate Studies in Mathematics. Theory, methods and applications, Translated from the 2005 German original by Jürgen Sprekels. American Mathematical Society, Providence, RI, pp. xvi+399. ISBN: 978-0-8218-4904-0. DOI: [10.1090/gsm/112](https://doi.org/10.1090/gsm/112). URL: <https://doi.org/10.1090/gsm/112>.